

Name:

Student ID:

Midterm Exam
Computer Architecture (263-2210-00L)
ETH Zürich, Fall 2017

Prof. Onur Mutlu

Problem 1 (30 Points):	
Problem 2 (80 Points):	
Problem 3 (90 Points):	
Problem 4 (40 Points):	
Problem 5 (70 Points):	
Problem 6 (90 Points):	
Problem 7 (70 Points):	
Problem 8 (BONUS: 80 Points):	
<hr/>	
Total (550 (470 + 80 bonus) Points):	

Examination Rules:

1. Written exam, 180 minutes in total.
2. No books, no calculators, no computers or communication devices. 6 pages of handwritten notes are allowed.
3. Write all your answers on this document, space is reserved for your answers after each question. Blank pages are available at the end of the exam.
4. Clearly indicate your final answer for each problem. Answers will only be evaluated if they are readable.
5. Put your Student ID card visible on the desk during the exam.
6. If you feel disturbed, immediately call an assistant.
7. Write with a black or blue pen (no pencil, no green or red color).
8. Show all your work. For some questions, you may get partial credit even if the end result is wrong due to a calculation mistake.
9. Please write your initials at the top of every page.

Tips:

- **Be cognizant of time.** Do not spend too much time on one question.
- **Be concise.** You may be penalized for verbosity.
- **Show work when needed.** You will receive partial credit at the instructors' discretion.
- **Write legibly.** Show your final answer.

Initials: _____

Computer Architecture

December 7th, 2017

This page intentionally left blank

1 Emerging Memory Technologies [30 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) [15 points] 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) [15 points] 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.

2 Cache Performance Analysis [80 points]

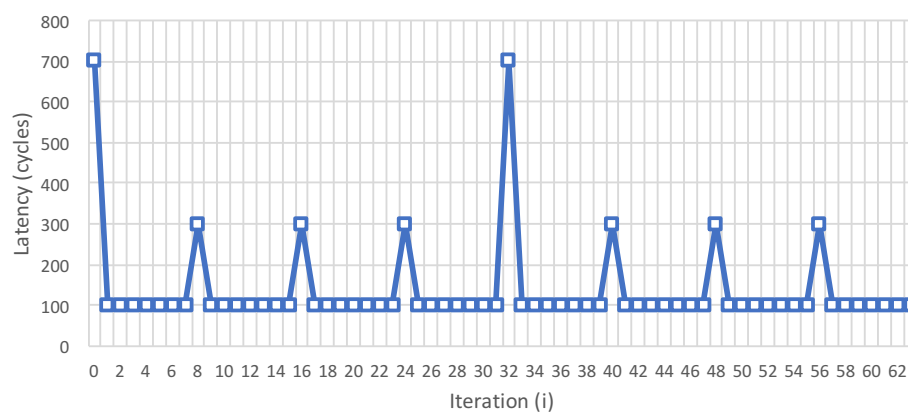
We are going to microbenchmark the cache hierarchy of a computer with the following two codes. The array `data` contains 32-bit unsigned integer values. For simplicity, we consider that accesses to the array latency bypass all caches (i.e., latency is *not* cached). `timer()` returns a timestamp in cycles.

```
(1) j = 0;
    for (i=0; i<size; i+=stride){
        start = timer();
        d = data[i];
        stop = timer();
        latency[j++] = stop - start;
    }

(2) for (i=0; i<size1; i+=stride1){
        d = data[i];
    }
    j = 0;
    for (i=0; i<size2; i+=stride2){
        start = timer();
        d = data[i];
        stop = timer();
        latency[j++] = stop - start;
    }
```

The cache hierarchy has two levels. L1 is a 4kB set associative cache.

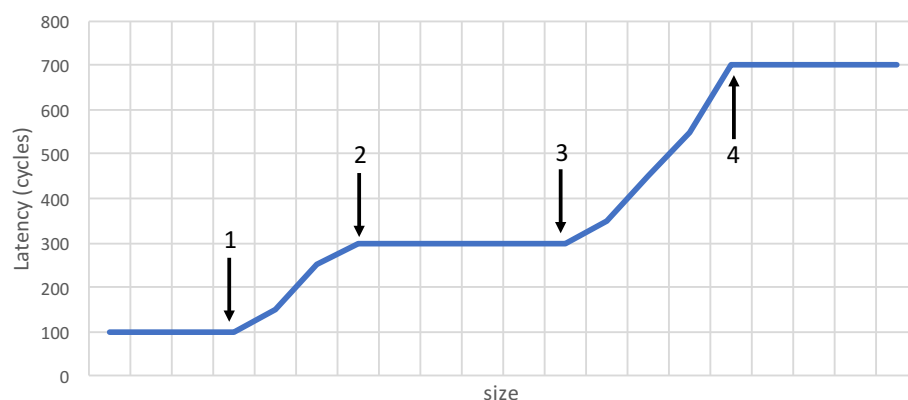
- (a) [15 points] When we run code (1), we obtain the latency values in the following chart for the first 64 reads to the array `data` (in the first 64 iterations of the loop) with `stride` equal to 1. What are the cache block sizes in L1 and L2?



- (b) [20 points] Using code (2) with `stride1 = stride2 = 32`, `size1 = 1056`, and `size2 = 1024`, we observe `latency[0] = 300` cycles. However, if `size1 = 1024`, `latency[0] = 100` cycles. What is the maximum number of ways in L1? (Note: The replacement policy can be either FIFO or LRU).

- (c) [20 points] We want to find out the exact replacement policy, assuming that the associativity is the maximum obtained in part (b). We first run code (2) with `stride1 = 32`, `size1 = 1024`, `stride2 = 64`, and `size2 = 1056`. Then (after resetting `j`), we run code (1) with `stride = 32` and `size = 1024`. We observe `latency[1] = 100` cycles. What is the replacement policy? Explain. (Hint: The replacement policy can be either FIFO or LRU. You need to find the correct one and explain).

- (d) [25 points] Now we carry out two consecutive runs of code (1) for different values of `size`. In the first run, `stride` is equal to 1. In the second run, `stride` is equal to 16. We ignore the latency results of the first run, and average the latency results of the second run. We obtain the following graph. What do the four parts shown with the arrows represent?



Before arrow 1:

Between arrow 1 and arrow 2:

Between arrow 2 and arrow 3:

Between arrow 3 and arrow 4:

After arrow 4:

Explain as needed (if you need more):

3 GPUs and SIMD [90 points]

We define the *SIMD utilization* of a program run on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program. As we saw in lecture and practice exercises, the SIMD utilization of a program is computed across the *complete run* of the program.

The following code segment is run on a GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A, B, and C are already in vector registers, so there are no loads and stores in this program. (Hint: Notice that there are 6 instructions in each thread.) A warp in the GPU consists of 32 threads, and there are 32 SIMD lanes in the GPU. Please assume that all values in arrays B and C have magnitudes less than 10 (i.e., $|B[i]| < 10$ and $|C[i]| < 10$, for all i).

```
for (i = 0; i < 1008; i++) {  
    A[i] = B[i] * C[i];  
    if (A[i] < 0) {  
        C[i] = A[i] * B[i];  
        if (C[i] < 0) {  
            A[i] = A[i] + 1;  
        }  
        A[i] = A[i] - 2;  
    }  
}
```

Please answer the following six questions.

- (a) [10 points] How many warps does it take to execute this program?

- (b) [10 points] What is the *maximum* possible SIMD utilization of this program?

- (c) [20 points] Please describe what needs to be true about arrays B and C to reach the *maximum* possible SIMD utilization asked in part (b). (Please cover all possible cases in your answer)

- (d) [10 points] What is the *minimum* possible SIMD utilization of this program?

- (e) [20 points] Please describe what needs to be true about arrays B and C to reach the *minimum* possible SIMD utilization asked in part (d). (Please cover all possible cases in your answer)

- (f) [20 points] Now consider a GPU that employs *Dynamic Warp Formation (DWF)* to improve the SIMD utilization. As we discussed in the class, DWF dynamically merges threads executing the same instruction (after branch divergence). What is the maximum achievable SIMD utilization using DWF? Explain your answer (Hint: The *maximum* SIMD utilization can happen under the conditions you found in part (e)).

Initials: _____

Computer Architecture

December 7th, 2017

A large, empty rectangular box with a thin black border, occupying the upper half of the page. It is intended for a student to draw a diagram or write notes related to the exam question.

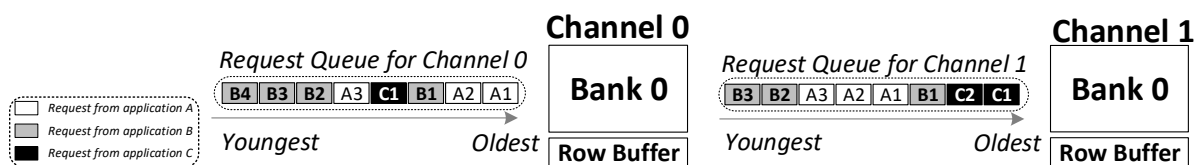
4 Memory Scheduling [40 points]

4.1 Basics and Assumptions

To serve a memory request, the memory controller issues one or multiple DRAM commands to access data from a bank. There are four different DRAM commands as discussed in class.

- **ACTIVATE:** Loads the row (that needs to be accessed) into the bank's row-buffer. This is called opening a row. (**Latency: 15ns**)
- **PRECHARGE:** Prepares the bank for the next access by closing the row in the bank (and making the row buffer empty). (**Latency: 15ns**)
- **READ/WRITE:** Accesses data from the row-buffer. (**Latency: 15ns**)

The diagrams below show the snapshots of memory controller's *request queues* at time 0, i.e., t_0 , when applications A, B, and C are executed together on a multi-core processor. Each application runs on a separate core but shares the memory subsystem with other applications. Each request is color-coded to denote the application to which it belongs. Additionally, each request is annotated with a number that shows the order of the request among the set of enqueued requests of the application to which it belongs. For example, A3 means that this is the third request from application A enqueued in the request queue. Assume all memory requests are reads and a read request is considered to be served when the READ command is complete (i.e., 15 ns after the request's READ command is issued).



Assume also the following:

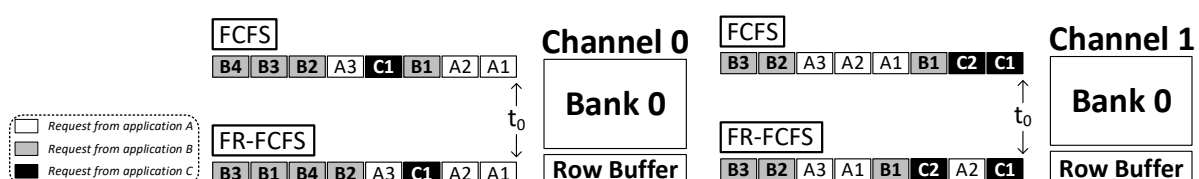
- The memory system has two DRAM channels, one DRAM bank per channel, and four rows per bank.
- All the row-buffers are closed (i.e., empty) at time 0.
- All applications start to stall at time 0 because of memory.
- No additional requests from any of the applications arrive at the memory controller.
- An application (A, B, or C) is considered to be stalled until *all* of its memory requests (across all the request buffers) have been served.

4.2 Problem Specification

The below table shows the stall time of applications A, B, and C with the FCFS (First-Come, First-Served) and FR-FCFS (First-Ready, First-Come, First-Served) scheduling policies.

Scheduling	Application A	Application B	Application C
FCFS	195 ns	285 ns	135 ns
FR-FCFS	135 ns	225 ns	90 ns

The diagrams below show the scheduling order of requests for Channel 0 and Channel 1 with the FCFS and FR-FCFS scheduling policies.



What are the numbers of row hits and row misses for each DRAM bank with either of the scheduling policies? Show your work.

Channel 0, hits:

Channel 0, misses:

Channel 1, hits:

Channel 1, misses:

Extra space for explanation (use only if needed):

5 Branch Prediction [70 points]

A processor implements an *in-order* pipeline with *12 stages*. Each stage completes in a single cycle. The pipeline stalls on a conditional branch instruction until the condition of the branch is evaluated. However, you *do not* know at which stage the branch condition is evaluated. Please answer the following questions.

- (a) [15 points] A program with 1000 dynamic instructions completes in 2211 cycles. If 200 of those instructions are conditional branches, at the end of which pipeline stage the branch instructions are resolved? (Assume that the pipeline does not stall for any other reason than the conditional branches (e.g., data dependencies) during the execution of that program.)

- (b) In a new, higher-performance version of the processor, the architects implement a *mysterious* branch prediction mechanism to improve the performance of the processor. They keep the rest of the design exactly the same as before. The new design with the mysterious branch predictor completes the execution of the following code in 115 cycles.

```
MOV R1, #0 // R1 = 0

LOOP_1:
    BEQ R1, #5, LAST // Branch to LAST if R1 == 5
    ADD R1, R1, #1    // R1 = R1 + 1
    MOV R2, #0        // R2 = 0
LOOP_2:
    BEQ R2, #3, LOOP_1 // Branch to LOOP_1 if R2==3.
    ADD R2, R2, #1     // R2 = R2 + 1
    B LOOP_2           // Unconditional branch to LOOP_2

LAST:
    MOV R1, #1        // R1 = 0
```

Assume that the pipeline never stalls due to a data dependency. Based on the given information, determine which of the following branch prediction mechanisms could be the *mysterious* branch predictor implemented in the new version of the processor. For each branch prediction mechanism below, you should circle the configuration parameters that makes it match the performance of the mysterious branch predictor.

i) [15 points] **Static Branch Predictor**

Could this be the mysterious branch predictor?

YES

NO

If YES, for which configuration below is the answer *YES*? Pick an option for each configuration parameter.

i. Static Prediction Direction

Always taken

Always not taken

Explain:

ii) [15 points] **Last Time Branch Predictor**

Could this be the mysterious branch predictor?

YES

NO

If YES, for which configuration is the answer *YES*? Pick an option for each configuration parameter.

i. Initial Prediction Direction

Taken

Not taken

ii. Local for each branch instruction (PC-based) or global (shared among all branches) history?

Local

Global

Explain:

iii) [10 points] **Backward taken, Forward not taken (BTFN)**

Could this be the mysterious branch predictor?

YES

NO

Explain:

iv) [15 points] **Two-bit Counter Based Prediction** (using saturating arithmetic)

Could this be the mysterious branch predictor?

YES

NO

If YES, for which configuration is the answer *YES*? Pick an option for each configuration parameter.

i. Initial Prediction Direction

00 (Strongly not taken)

01 (Weakly not taken)

10 (Weakly taken)

11 (Strongly taken)

ii. Local for each branch instruction (i.e., PC-based, without any interference between different branches) or global (i.e., a single counter shared among all branches) history?

Local

Global

Explain:

6 SIMD [90 points]

We have two SIMD engines: 1) a traditional vector processor and 2) a traditional array processor. Both processors can support a vector length up to 16.

All instructions can be fully pipelined, the processor can issue one vector instruction per cycle, and the pipeline does not forward data (no chaining). For the sake of simplicity, we ignore the latency of the pipeline stages other than the execution stages (e.g, decode stage latency: 0 cycles, write back latency: 0 cycles, etc).

We implement the following instructions in both designs, with their corresponding execution latencies:

Operation	Description	Name	Latency of a single operation (VLEN=1)
VADD	$VDST \leftarrow VSRC1 + VSRC2$	vector add	5 cycles
VMUL	$VDST \leftarrow VSRC1 * VSRC2$	vector mult.	15 cycles
VSHR	$VDST \leftarrow VSRC \gg 1$	vector shift	1 cycles
VLD	$VDST \leftarrow \text{mem}[SRC]$	vector load	20 cycles
VST	$VSRC \rightarrow \text{mem}[DST]$	vector store	20 cycles

- All the vector instructions operate with a vector length specified by VLEN. The VLD instruction loads VLEN consecutive elements from the DST address specified by the value in the VDST register. The VST instruction stores VLEN elements from the VSRC register in consecutive addresses in memory, starting from the address specified in DST.
- Both processors have eight vector registers (VR0 to VR7) which can contain up to 16 elements, and eight scalar registers (R0 to R7). The entire vector register needs to be ready (i.e., populated with all VLEN elements) before any element of it can be used as part of another operation.
- The memory can sustain a throughput of one element per cycle. The memory consists of 16 banks that can be accessed independently. A single memory access can be initiated in each cycle. The memory can sustain 16 parallel accesses if they all go to different banks.

(a) [10 points] Which processor (array or vector processor) is more costly in terms of chip area? Explain.

(b) [25 points] The following code takes 52 cycles to execute on the vector processor:

```
VADD VR2 ← VR1, VR0
VADD VR3 ← VR2, VR5
VMUL VR6 ← VR2, VR3
```

What is the VLEN of the instructions? Explain your answer.

How long would the same code execute on an array processor with the same vector length?

- (c) [25 points] The following code takes 94 cycles to execute on the vector processor:

```
VLD  VR0 ← mem[R0]
VLD  VR1 ← mem[R1]
VADD VR2 ← VR1, VR0
VSHR VR2 ← VR2
VST  VR2 → mem[R2]
```

Assume that the elements loaded in VR0 are all placed in different banks, and that the elements loaded into VR1 are placed in the same banks as the elements in VR0. Similarly, the elements of VR2 are stored in different banks in memory. What is the VLEN of the instructions? Explain your answer.

- (d) [30 points] We replace the memory with a new module whose characteristics are unknown. The following code (the same as that in (c)) takes 163 cycles to execute on the vector processor:

```
VLD  VR0 ← mem[R0]
VLD  VR1 ← mem[R1]
VADD VR2 ← VR1, VR0
VSHR VR2 ← VR2
VST  VR2 → mem[R2]
```

The VLEN of the instructions is 16. The elements loaded in VR0 are placed in consecutive banks, the elements loaded in VR1 are placed in consecutive banks, and the elements of VR2 are also stored in consecutive banks. What is the number of banks of the new memory module? Explain.

7 In-DRAM Bitmap Indices [70 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) [15 points] What is the total number of DRAM rows that are occupied by u users and w weeks?

(b) [20 points] What is the throughput in users/second of the *Ambit*-based implementation?

- (c) [20 points] What is the throughput in users/second of the CPU implementation?

- (d) [15 points] 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 BONUS: 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) [15 points] What is the execution time of the above piece of code in cycles?

- (b) [25 points] 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) [15 points] 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) [15 points] 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) [10 points] What other changes could be made to the cache design to improve the performance of the single-core processor on this program?

Initials: _____

Computer Architecture

December 7th, 2017

- SCRATCHPAD -

Initials: _____

Computer Architecture

December 7th, 2017

- SCRATCHPAD -

Initials: _____

Computer Architecture

December 7th, 2017

- SCRATCHPAD -

Initials: _____

Computer Architecture

December 7th, 2017

- SCRATCHPAD -