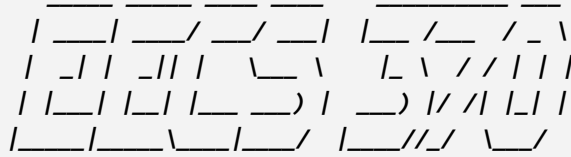


# Final Exam **ANSWER KEY**



## EECS 370 Winter 2022: Intro to Computer Organization

You are to abide by the University of Michigan College of Engineering Honor Code. Please sign below to signify that you have kept the honor code pledge:

***I have neither given nor received aid on this exam,  
nor have I concealed any violations of the Honor Code.***

Signature: \_\_\_\_\_

Name: \_\_\_\_\_

Uniquename: \_\_\_\_\_

First/Last name of person sitting to your **Right**  
(Write ⊥ if you are at the end of the row) \_\_\_\_\_

First/Last name of person sitting to your **Left**  
(Write ⊥ if you are at the end of the row) \_\_\_\_\_

### Exam Directions:

- You have **120 minutes** to complete the exam. There are **8** questions in the exam on **15** pages (double-sided). **Please flip through your exam to ensure you have all 15 pages.**
- You must show your work to be eligible for partial credit!
- Write legibly and dark enough for the scanners to read your answers.
- **Write your uniquename on the line provided at the top of each page.**

### Exam Materials:

- You are allotted **one 8.5 x 11 double-sided** note sheet to bring into the exam room.
- You are allowed to use calculators that do not have an internet connection. All other electronic devices, such as cell phones or anything or calculators with an internet connection, are strictly forbidden and usage will result in an Honor Code violation.

1. Short Questions	_____ / 12 pts
2. Branch Prediction	_____ / 11 pts
3. LC2K Pipeline Datapath Performance	_____ / 10 pts
4. The 3 C's of Caches	_____ / 9 pts
5. New LC2K Pipeline Datapath	_____ / 15 pts
6. Cache Locality	_____ / 15 pts
7. VM Simulation	_____ / 16 pts
8. VM Performance / Cache Performance	_____ / 12 pts
TOTAL _____ / 100 pts	

1.	<b>Short Questions</b> <span style="float: right;"><b>[12 pts]</b></span>
	Complete the following true/false and short answer questions

**True/False Questions [5 pts]****Circle One:**

- (a) Using the speculate-and-squash method to resolve control hazards can result in the same CPI as the detect-and-stall method. True / False
- (b) Increasing the associativity of a cache can reduce capacity misses. True / False
- (c) If a machine has infinite physical memory, virtual memory is no longer needed. True / False
- (d) A multi-level page table can consume more space than a single-level page table. True / False
- (e) Virtual address to physical address translation can happen simultaneously with the cache access. True / False

**Short Answer Questions**

- (f) **[2 pts]** Consider a cache with a given size. You are asked to change the associativity and block size to reduce the tag area overhead:

i) How would you change the associativity? Circle around your choice.

Increase Associativity      ,      **Decrease Associativity**

ii) How would you adjust the block size? Circle around your choice.

**Increase Block Size**      ,      Decrease Block Size

- (g) [3 pts] Consider the following assembly code being simulated on the **5-stage** LC2K pipeline datapath that uses **detect and forward** for data hazards, **speculate not-taken and squash** for control hazards, and **internal forwarding** for register file.

i) Circle around the registers of the executed instructions that are going to use forwarded-data from pipeline registers (including the internal forwarding of register file):

**Example:**      `nor`      5      6      7

1	<code>lw</code>	0	1	10
2	<code>lw</code>	0	2	10
3	<code>beq</code>	<u>1</u>	<u>2</u>	0
4	<code>add</code>	1	2	3
5	<code>nor</code>	1	<u>3</u>	4
6	<code>sw</code>	<u>3</u>	<u>4</u>	10
7	<code>halt</code>			

ii) How many cycles would it take to complete the program?

Answer: 7 + 1 + 3 + 4 = 15

- (h) [2 pts] Consider a byte-addressable system with a 128 B fully-associative cache that has a block size of 32 B. The following lines of code are executed:

```

1  int data[20];
2  for(unsigned int i = 0; i < 20; i++) {
3      data[i] = data[i] + i;
4  }
```

How many bytes would be written into main memory for data array accesses, in each case of write-through and write-back policy for the cache. You can assume data starts from address 0 and all dirty blocks are written *in their entirety* back to cache at the end of the program.

Write-through: 80 B

Write-back: 96 B

<b>2.</b>	<b>Branch Prediction</b> <span style="float: right;"><b>[11 pts]</b></span>
	Simulate an LC2K assembly program and answer questions about its branches

Consider the following LC2K assembly program running on our 5-stage pipelined datapath from lecture. This program **counts the number of zeros** in the binary array **Arr** and stores that count in **r5**. Assume all registers are initialized to zero. Answer the following question about the program's branch decisions.

1		lw	0	1	One	//r1 = 1
2		lw	0	2	Len	//r2 = Len
3		lw	0	3	Zero	//r3 = loop counter i
4	loop	beq	3	2	end	//loop terminates if i==Len
5		lw	3	4	Arr	//load Arr[i]
6		beq	4	1	cont	//if Arr[i]==1, skip Line 7
7		add	5	1	5	//increment r5
8	cont	add	3	1	3	//increment i
9		beq	0	0	loop	
10	end	halt				
11	Zero	.fill	0			
12	One	.fill	1			
13	Len	.fill	3			
14	Arr	.fill	1			
15		.fill	1			
16		.fill	0			

- (a) [7 pts] Write the sequence of branch decisions for each **beq** instruction, using **T** to represent “taken” and **N** to represent “not taken.” (You might not need all the boxes.)

Line 4 **beq** (Loop termination condition)

N	N	N	T								
---	---	---	---	--	--	--	--	--	--	--	--

Line 6 **beq** (If-Condition)

T	T	N									
---	---	---	--	--	--	--	--	--	--	--	--

Line 9 **beq** (Loop condition)

T	T	T									
---	---	---	--	--	--	--	--	--	--	--	--

Combined sequence, as seen globally

N	T	T	N	T	T	N	N	T	T		
---	---	---	---	---	---	---	---	---	---	--	--

- (b) [4 pts] What percentage of the total branch decisions are correctly predicted when using the following prediction schemes? (You can leave your answers as fractions.)

Partial credits are given based on incorrect prediction patterns of part a.

Predict always not taken (N)

\_\_\_\_\_ 4/10 \_\_\_\_\_

Predict backwards taken (T), forwards not taken (N)

\_\_\_\_\_ 7/10 \_\_\_\_\_

Predict using a local (i.e., per **beq** instruction) 2-bit branch predictor initialized to the starting state “strongly taken”

\_\_\_\_\_ 6/10 \_\_\_\_\_

3.	<b>LC2K Pipeline Datapath Performance</b> <span style="float: right;"><b>[10 pts]</b></span>
	Calculate the CPI and execution time for a pipelined datapath

Consider an LC2K assembly program running on our 5-stage pipelined datapath from lecture, which has internal forwarding for the register file. The program has the following instruction breakdown:

- R-type instructions (**add** and **nor**)                      50%
- Load instructions (**lw**)                                      25%
- Store instructions (**sw**)                                        15%
- Branch instructions (**beq**)                                    15%

In the program, 60% of branches are not taken (N). Additionally, 20% of the R-type instructions are followed immediately by a dependent R-type instruction (e.g., **add 1 1 2**, **add 2 2 3**) and 80% of load instructions are followed immediately by a dependent R-type instruction (e.g., **lw 0 3 0**, **add 3 3 4**). There are no other data dependencies.

- (a) **[4 pts]** What is the CPI of the program when the processor uses **detect-and-stall** for data hazards and **speculate-and-squash** for control hazards with a branch predictor that always predicts not taken (N). Please show your work for partial credit.

$$\begin{aligned}
 \text{CPI} &= 1 + \text{stalls for data hazards} && + \text{stalls for control hazards} \\
 &= 1 + (0.5 \cdot 0.2 + 0.25 \cdot 0.8)(2 \text{ cycles}) && + (0.15 \cdot 0.4)(3 \text{ cycles}) \\
 &= 1 + (0.5 \cdot 0.2 \cdot 2 \text{ cycles}) + (0.25 \cdot 0.8 \cdot 2 \text{ cycles}) && + (0.15 \cdot 0.4 \cdot 3 \text{ cycles}) \\
 &= 1 + 0.6 + 0.18 = \mathbf{1.78}
 \end{aligned}$$

- (b) **[4 pts]** What is the CPI of the program when the processor uses **detect-and-forward** for data hazards and **speculate-and-squash** for control hazards with a branch predictor that always predicts not taken (N). Please show your work for partial credit.

$$\begin{aligned}
 \text{CPI} &= 1 + \text{stalls for data hazards} && + \text{stalls for control hazards} \\
 &= 1 + (0.25 \cdot 0.8)(1 \text{ cycles}) && + (0.15 \cdot 0.4)(3 \text{ cycles}) \\
 &= 1 + 0.2 + 0.18 = \mathbf{1.38}
 \end{aligned}$$

\*\*\* Term for control hazards can be transferred from (a) with no additional loss of points

- (c) **[2 pts]** Given that the processor frequency is 1MHz and the program executes one million instructions, how many seconds faster is the execution time of the program when using detect-and-forward? Show your work by calculating the execution time in seconds for parts (a) and (b), then finding the difference between the two.

$$\text{Ex. time for (a): } (1.78 \text{ cycles/insn}) \cdot (1,000,000 \text{ insn}) / (10^6 \text{ cycles/second}) = 1.78 \text{ s}$$

$$\text{Ex. time for (b): } (1.38 \text{ cycles/insn}) \cdot (1,000,000 \text{ insn}) / (10^6 \text{ cycles/second}) = 1.38 \text{ s}$$

$$\text{Difference: } 1.78 - 1.38 = \mathbf{0.4 \text{ s}}$$

4.	<b>The 3 C's of Caches</b> <span style="float: right;"><b>[9 pts]</b></span>
	Determine if the memory references result in compulsory, capacity, or conflict misses

Consider an 8-bit processor with byte-addressable memory that has a cache with the following configuration:

- **Cache size** 32 B
- **Block size** 16 B
- **Associativity** Direct mapped (1-way)
- **Replacement Policy** Least-recently used (LRU)

Assume the cache is initially empty. Fill out the table given the provided sequence of memory references. If the access was a miss, fill in the bubble (○) to indicate the type of miss. If a block was evicted because of the access, write the tag of the block that was evicted.

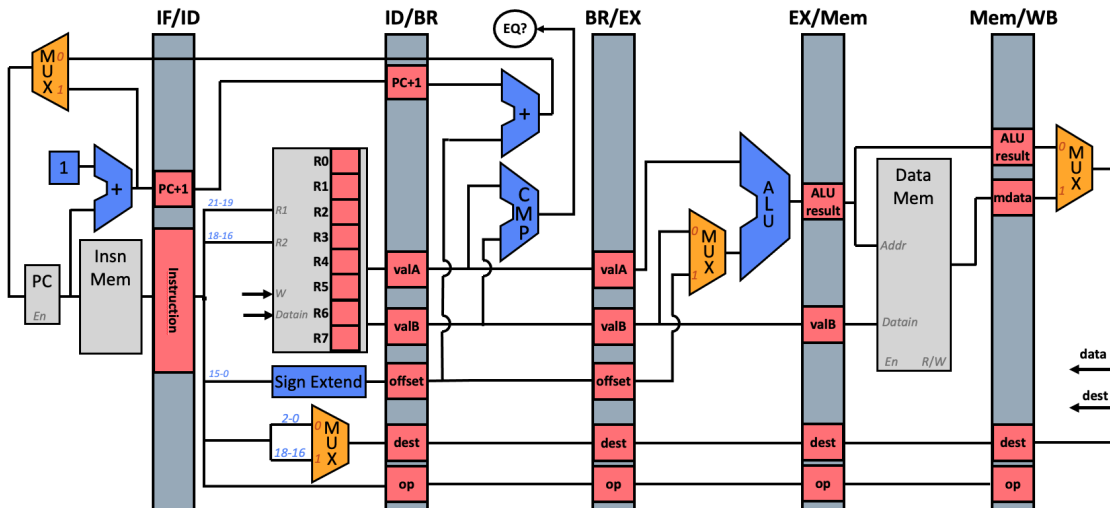
Reference Address in Binary (0b)	Tag	Set Index	Hit or Miss?	Type of Miss (if applicable)	Tag of evicted block (if applicable)
0b1001 0110	0b 100	0b 1	Miss	<input checked="" type="radio"/> <u>Compulsory</u> <input type="radio"/> Capacity <input type="radio"/> Conflict	N/A
0b1110 1111	0b 111	0b 0	Miss	<input checked="" type="radio"/> <u>Compulsory</u> <input type="radio"/> Capacity <input type="radio"/> Conflict	N/A
0b1001 0010	0b 100	0b 1	Hit	<input type="radio"/> Compulsory <input type="radio"/> Capacity <input type="radio"/> Conflict	N/A
0b0110 1011	0b 011	0b 0	Miss	<input checked="" type="radio"/> <u>Compulsory</u> <input type="radio"/> Capacity <input type="radio"/> Conflict	0b 111
0b1110 0000	0b 111	0b 0	Miss	<input type="radio"/> Compulsory <input checked="" type="radio"/> <u>Capacity</u> <input type="radio"/> Conflict	0b 011
0b0110 1000	0b 011	0b 0	Miss	<input type="radio"/> Compulsory <input type="radio"/> Capacity <input checked="" type="radio"/> <u>Conflict</u>	0b 111

Partial credits are given based on incorrect calculation of tag and set index bits.

**5. New LC2K Pipeline Datapath****[15 pts]**

Identify the hazards and pipeline stalls in the new LC2K pipelined-datapath

Consider the following revised LC2K pipelined-datapath, in which a new pipeline stage called **Branch Resolve (BR)** is added between **ID** and **EX** stage, so we can resolve control hazards earlier in the pipeline. The new **BR** stage includes an adder to calculate the branch target as well as a comparator to check the equality of regA and regB values.



- (a) [6 pts] Assume that this new pipeline design has **only the original data forwarding paths** of the 5-stage LC2K pipeline (i.e.,  $EX/MEM \rightarrow EX$ ,  $MEM/WB \rightarrow EX$ , and internal forwarding for the register file). Compute how many cycles we will need to stall the pipeline for each of the following scenarios of instructions containing data dependencies.

In the table, "Distance" denotes the distance between dependent instructions. When Distance = 0, the instructions are adjacent (e.g., add 1 1 2, nor 1 2 3). When the Distance = 1, there is one instruction separating the dependent instructions (eg., add 1 1 2, <unrelated instruction>, nor 1 2 3), and so on.

Source of Data Dependency	Dependent Instruction	# of stalled cycles		
		Distance = 0	Distance = 1	Distance = 2
R-type (add 1 1 2)	R/I-type (nor 1 2 3)	0 cycles	0 cycles	0 cycles
R-type (add 1 1 2)	Branch (beq 1 2 end)	3	2	1
Load (lw 0 2 0)	R/I-type (nor 1 2 3)	1	0	0
Load (lw 0 2 0)	Branch (beq 1 2 end)	3	2	1



- (b) [3 pts] If we allow new data-forwarding paths, what is the **minimum set** of forwarding paths that are necessary to **minimize** the number of stalling cycles for part (a)? (You might not need all the paths.)

	<u>Source Pipeline Register</u>	→	<u>Destination Stage</u>
<b>Example:</b>	<u>EX/MEM</u>	→	<u>EX</u>
<b>New Forwarding Path 1:</b>	<u>EX/MEM</u>	→	<u>BR</u>
<b>New Forwarding Path 2:</b>	<u>MEM/WB</u>	→	<u>BR</u>
<b>New Forwarding Path 3:</b>	<u>                    </u>	→	<u>                    </u>

- (c) [5 pts] With the original data forwarding paths of the 5-stage LC2K pipeline and the **new data forwarding paths** from part (b), compute how many cycles we will need to stall the pipeline for the following scenarios of instructions containing data dependencies.

Source of Data Dependency	Dependent Instruction	# of stalled cycles		
		Distance = 0	Distance = 1	Distance = 2
R-type (add 1 1 2)	R/I-type (nor 1 2 3)	0 cycles	0 cycles	0 cycles
R-type (add 1 1 2)	Branch (beq 1 2 end)	1	0	0
Load (lw 0 2 0)	R/I-type (nor 1 2 3)	1	0	0
Load (lw 0 2 0)	Branch (beq 1 2 end)	2	1	0

- (d) [1 pts] With the new pipeline design, how many cycles of stalls are we going to have in the pipeline in case of a branch misprediction?

Stalled Cycles: 2 cycles

<b>6.</b>	<b>Cache Locality</b>	<b>[15 pts]</b>
	Trace memory accesses and determine the ideal cache layout for a C++ program	

Consider the following convolutional kernel executing on a system with a byte-addressable memory. Load/store instructions operate with 4 B granularity (i.e., each load/store instruction returns a 4 B value from the cache to the processor, resulting in a *single* cache hit or miss).

1	#define N 3
2	int Out[N], Input[N], Mask[3];
3	
4	//initialization code
5	
6	for(int i=0; i<N; i++)
7	{
8	dot=0;
9	for(int j=0; j<3; j++)
10	{
11	if((i+j-1)>=0 && (i+j-1)<N)
12	dot += Input[i+j-1] * Mask[j];
13	}
14	Out[i] = dot;
15	}

For this program:

- All variables except the arrays **Out**, **Input**, and **Mask** are mapped to registers. Hence, only these arrays require loads/stores to the cache.
- **Input** starts at address 0x100, **Mask** starts at address 0x200, and **Output** starts at address 0x400.
- On line 14, the **Input** array is accessed before the **Mask** array.

(a) [2.5 pts] Write down the memory addresses for the array accesses in the program. For example, the iteration  $i=0$  accesses  $\text{Input}[0]$ ,  $\text{Mask}[1]$ ,  $\text{Input}[1]$ ,  $\text{Mask}[2]$ ,  $\text{Out}[0]$  and the memory addresses for this iteration are completed in the table below.

for $i=0$	0x100, 0x204, 0x104, 0x208, 0x400
for $i=1$	0x100, 0x200, 0x104, 0x204, 0x108, 0x208, 0x404
for $i=2$	0x104, 0x200, 0x108, 0x204, 0x408

Total number of memory accesses: \_\_\_\_\_ 17 \_\_\_\_\_

- (b) [3 pts] Say our processor has a fully associative cache size of 16 B and block size of 8 B with LRU replacement policy. Determine hits/misses for the above program. Assume the memory address has 12 bits. Show your work.

\* # hits and misses are graded based on the address pattern of part a.

i=0: M M H M M

i=1: M M H H M M M

i=2: M M M H M

# of Hits: \_\_\_\_\_ 4 \_\_\_\_\_  
 # of Misses: \_\_\_\_\_ 13 \_\_\_\_\_

- (c) [4.5 pts] What is the minimum size of a fully associative cache with the same block size as above (8 B) that can eliminate all misses except compulsory misses. Determine hits/misses for this new cache. **Note that cache size does not need to be a power of 2.** Show your work.

There are 6 compulsory misses, so the minimum cache size would have 11 hits and 6 misses. This could be achieved by 5 blocks, note that we don't need to have all 6 blocks at the same time in the cache, due to lack of reuse in Out array.

i=0: M M H M M

i=1: H H H H M H H

i=2: H H H H M

Cache Size (B): \_\_\_\_\_ 40 \_\_\_\_\_  
 # of Hits: \_\_\_\_\_ 11 \_\_\_\_\_  
 # of Misses: \_\_\_\_\_ 6 \_\_\_\_\_

- (d) [5 pts] Say you can increase the cache size to 64 B. What cache configuration for this increased size will result in the lowest number of cache misses? Determine hits/misses for your new cache configuration. Show your work.

We need to minimize the compulsory misses, so the largest block size needs to be used. A block size of 16 would make sure that we would only have 1 inevitable compulsory miss per array. The associativity of the cache should be fully-associative, otherwise we would face conflict misses due to addresses being mapped to the same set.

Associativity: \_\_\_\_\_ Fully-associative \_\_\_\_\_

^^^ Just writing "1" without explicitly saying "1 set" or "fully associative" is not counted

Block Size (B): \_\_\_\_\_ 16 \_\_\_\_\_  
 # of Hits: \_\_\_\_\_ 14 \_\_\_\_\_  
 # of Misses: \_\_\_\_\_ 3 \_\_\_\_\_

<b>7.</b>	<b>Virtual Memory and Page Table</b> <span style="float: right;"><b>[16 pts]</b></span>
	Simulate virtual to physical page translation and calculate page table overheads

Consider a 32-bit processor that has a byte-addressable memory with the following configuration:

- **Page Size** 256 B
- **Virtual Memory Size** 16 KB
- **Physical Memory Size** 2 KB (8 pages)
- **Page Replacement Policy** Least-recently used (LRU)
- **Page Table Layout** Single-level page table
- **Page Table Size** 256 B
- **Page Table Entry Size** 4 B

Physical page 0x0 is reserved for the operating system (OS) and cannot be replaced. Similarly, physical pages 0x1 and 0x2 are reserved for the current processes' page tables and cannot be replaced. On a page fault, the page table is updated before allocating a physical page. If more than one free page is available, the smallest physical page number is chosen.

Assume that two processes with Process ID (PID) 370 and 281, respectively, have been running. The initial state of physical memory is shown below:

Physical Page # (PPN)	Memory Contents
0x0	Reserved for OS
0x1	PID 370 Page Table
0x2	PID 281 Page Table
0x3	PID 370: VPN 0x3
0x4	
0x5	PID 281: VPN 0x2
0x6	
0x7	PID 370: VPN 0xF

- (a) [10 pts] Complete the following table for the given sequence of virtual address requests. Please express your solution in hexadecimal.

Time	PID	Virtual Address (VA)	Virtual Page # (VPN)	Physical Page # (PPN)	Page Fault? (Y/N)	Physical Address (PA)
0	370	0x31A	0x3	0x3	N	0x31A
1	370	0xA02	0xA	0x4	Y	0x402
2	281	0x31A	0x3	0x6	Y	0x61A
3	370	0xF00	0xF	0x7	N	0x700
4	281	0x618	0x6	0x5 (evicted)	Y	0x518

- (b) Suppose the page size is reduced to **16 B** and the single-level page table is replaced by a **3-level page table**, in which the size of each 2nd level page table is **2 pages** and the size of each 3rd level page table is **4 pages**. Answer the following questions about this new configuration. (Note that the page table entry size is still 4 B)

- i) [2 pts] How many index bits are required for the 2nd level and the 3rd level page table, respectively?

2nd level:  $\log_2 (2 \text{ pages} * 16 \text{ bytes per page} / 4 \text{ bytes per entry}) = 3 \text{ bits}$

3rd level:  $\log_2 (4 \text{ pages} * 16 \text{ bytes per page} / 4 \text{ bytes per entry}) = 4 \text{ bits}$

2nd level:   3   index bits

3rd level:   4   index bits

- ii) [2 pts] What is the **minimum** storage space this 3-level page table would occupy in bytes?

Best case: only the 1st level page table is allocated.

Index bits for the 1st level page table:  $\log_2 (16KB) - 3 \text{ bits (2nd level)} - 4 \text{ bits (3rd level)} - 4 \text{ bits (offset)} = 3 \text{ bits}$

Total size =  $2^3 * 4 = 32 \text{ bytes}$

- iii) [2 pts] What is the **maximum** storage space this 3-level page table would occupy in bytes?

Worst case: all pages at different levels are occupied

Total Size =  $2^3 * 4 \text{ (first level)} + 2^3 * 2^3 * 4 \text{ (second level)} + 2^3 * 2^3 * 2^4 * 4 \text{ (third level)}$

$= 2^5 + 2^8 + 2^{12} \text{ bytes}$

$= 4384 \text{ bytes}$

8.	<b>VM + Cache Performance</b> <span style="float: right;"><b>[12 pts]</b></span>
	Optimize the virtual memory system based on latency-cost tradeoff

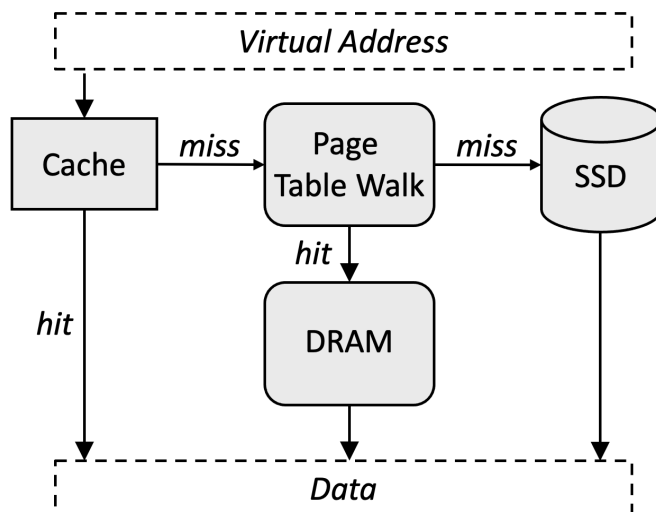
You have been tasked to configure the memory system for a custom processor designed by a startup company. The company has placed the following limitations on the design:

- The cache should be **virtually addressed** (virtual to physical address translation occurs in case of the cache miss)
- Virtual memory should use a single-level page table
- There is no TLB in the memory system.

The company is considering the following technologies the memory system components:

Component	Hit Rate	Technology	Latency	Cost
Cache	80%	SRAM-A	2 ns	\$20
		SRAM-B	1 ns	\$40
DRAM (Main Memory)	90%	DDR3	20 ns	\$15
		DDR4	10 ns	\$45
SSD (Disk)	100%	SATA 3.0	800 ns	\$20

- (a) [6 pts] Considering this memory system, write down the coefficients for the average memory access time (AMAT) equation. The variables  $LC$ ,  $LD$ , and  $LS$  represent the latencies for the cache, DRAM, and SSD components, respectively. You can assume that any necessary updates to cache, page table, and DRAM due to data miss would happen in parallel to data retrieval (No impact on AMAT). Please show your work.



$$\begin{aligned}
 \text{AMAT} &= LC + && \text{[cache access]} \\
 &0.2 * ( && \text{[cache Miss]} \\
 &0.9 * (LD + LD) + && \text{[PT Hit]} \\
 &0.1 * (LD + LS) && \text{[PT Miss]} \\
 &) \\
 &= LC + 0.38 * LD + 0.02 * LS
 \end{aligned}$$

$$\text{AMAT} = \alpha \times LC + \beta \times LD + \gamma \times LS$$

$$\alpha = \underline{1} \quad , \quad \beta = \underline{0.38} \quad , \quad \gamma = \underline{0.02}$$

(b) [4 pts] The company has set the following restrictions for the memory system:

- Average memory access time (AMAT) must be **less than 25 ns**.
- The memory system components must **cost less than \$100 in total**.

Circle the ideal technologies for the cache and DRAM to **minimize the latency while meeting the cost budget**. Please show your work.

Cache: SRAM-A , SRAM-B

DRAM: DDR3 , DDR4

SSD: SATA 3.0

Here are the breakdown of all configurations:

SSD	DRAM	Cache	Latency	Cost
SATA 3.0	DDR3	SRAM-A	25.6	55
<b>(1) SATA 3.0</b>	<b>DDR3</b>	<b>SRAM-B</b>	<b>24.6</b>	<b>75</b>
<b>(2) SATA 3.0</b>	<b>DDR4</b>	<b>SRAM-A</b>	<b>21.8</b>	<b>85</b>
SATA 3.0	DDR4	SRAM-B	20.8	105

Out of the both options, option 2 minimizes the latency

\*\*\* Graded based on the equation in part (a) \*\*\*

(c) [2 pts] Using the equation in part (a), what would the AMAT and total cost be for the technologies selected in part (b)?

\*\*\* Graded based on the selections to part (b).\*\*\*

AMAT = 21.8 , Cost = 85