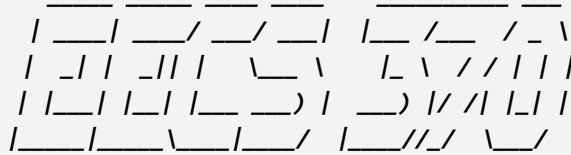


Final Exam-Key



EECS 370 Winter 2023: Introduction 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: Key

Name: _____

Uniqname: _____

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

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

Exam Directions:

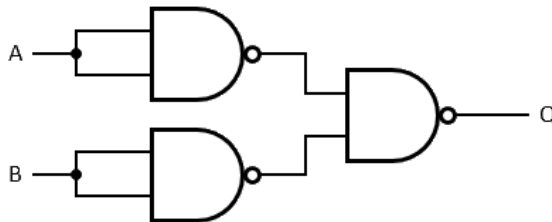
- You have **120 minutes** to complete the exam. There are **9** questions in the exam on **6** pages (so 12 “sides” total). **Please flip through your exam to ensure you have all 6 pages.**
- You must show your work to be eligible for partial credit.
- Write legibly and dark enough for the scanners to read your answers.
- **Write your unickname on the line provided at the top of each page. Do this at the beginning of the exam; you will NOT be given time to do it at the end.**

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 electronic devices with an internet connection are strictly forbidden.

1.	<div data-bbox="276 226 570 270">Multiple choice</div> <div data-bbox="1295 226 1435 270">[15 pts]</div> <div data-bbox="276 300 844 331">Completely fill in the circle of the <i>best</i> answer.</div>
-----------	--

- Under what conditions would you expect a write-through cache to have a lower number of bytes transferred between the cache and memory than a write-back cache? **[2]**
 - ☐ The program has low spatial locality but high temporal locality
 - ☐ The program has high spatial locality but low temporal locality
 - ☐ The program has high spatial locality and high temporal locality
 - ☒ The program has low spatial locality and low temporal locality
 - ☐ Never
- Functions are surprisingly difficult for the branch predictors we've discussed to deal with. What is it about functions that typically cause problems for those predictors? **[2]**
 - ☐ It is often hard to predict the "direction" of a function call.
 - ☐ It is often hard to predict the "direction" of the return from a function.
 - ☒ It is often hard to predict the target of the return from a function.
 - ☐ It is often hard to predict the target of a function call.
 - ☐ The branches associated with function calls and returns have very little spatial locality.
- Which of the following formulas is equivalent to the circuit below? **[2]**



- ☐ $Q = A \text{ nor } B$
- ☐ $Q = A \text{ and } B$
- ☒ $Q = A \text{ or } B$
- ☐ $Q = \text{not}(A) \text{ or } \text{not}(B)$
- ☐ None of the above

4. When comparing direct-mapped caches to fully-associative caches that otherwise have identical parameters, which of the following would be expected to be true? **[3]**
- a) Direct-mapped caches will have a lower hit latency, fully-associative caches will have a higher hit rate
 - b) Direct-mapped caches will require more index bits, fully-associative caches will have more tag bits.
 - c) Direct-mapped caches will require fewer block offset bits, fully-associative caches will have more LRU bits.
- ☐ Only a
- ☐ Only b
- ☐ Only c
- ☒ Only a and b
- ☐ Only b and c
- ☐ Only a and c
- ☐ All of a, b, and c.

For the next two questions, assume you have a byte-addressable, 256-byte virtually addressed cache with 16-byte blocks. Assume all entries in the cache start as “invalid” and addresses are 16-bits. All loads and stores are to 4-byte values.

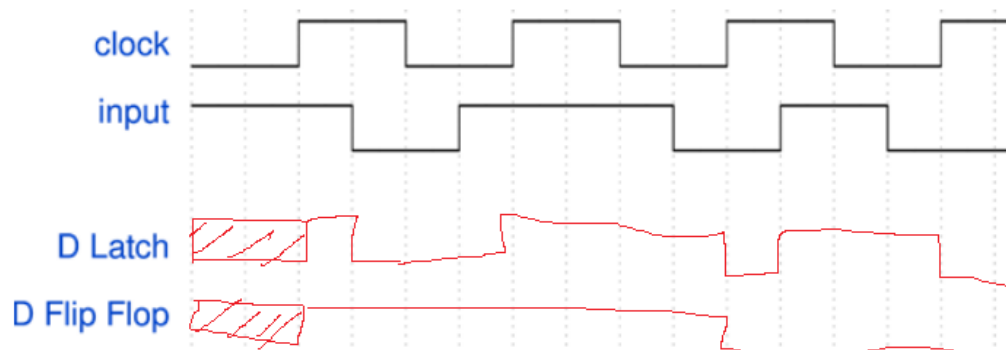
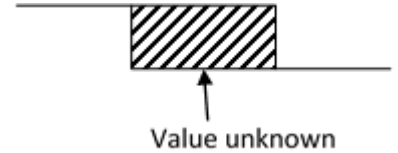
5. For which of the following access patterns will a direct-mapped cache will get a better hit-rate than a two-way associative cache using LRU replacement? **[3]**
- ☐ 0x0000, 0x0010, 0x0020, 0x0000
 - ☐ 0x0000, 0x0080, 0x0080, 0x0001
 - ☒ 0x0000, 0x0180, 0x0080, 0x0000
 - ☐ 0x0000, 0x0060, 0x0080, 0x0001
 - ☐ None of the above
6. For which of the following access patterns will a fully-associative cache using LRU replacement will get a better hit-rate than a two-way associative cache? **[3]**
- ☒ 0x0000, 0x0100, 0x0200, 0x0000
 - ☐ 0x0000, 0x0001, 0x0200, 0x0002
 - ☐ 0x0000, 0x0010, 0x0020, 0x0000
 - ☐ 0x0000, 0x0410, 0x0020, 0x0000
 - ☐ None of the above

2.	True or False	[13 pts]
	Complete the following true or false questions.	

- (1) A clock with a 2ns period has a frequency of 200MHz. ☐ True
☒ False
- (2) The number of LRU bits required for a set associative cache depends on cache associativity. ☒ True
☐ False
- (3) An XOR gate can be created using only AND gates. ☐ True
☒ False
- (4) A multi-level page table can take up more memory space than a single level page table. ☒ True
☐ False
- (5) In the 3C's cache model, a "compulsory" cache miss can sometimes be avoided by changing the cache's total size while holding block size constant. ☐ True
☒ False
- (6) A virtually-addressed cache doesn't need to access the TLB to see if the cache will get a hit or a miss. ☒ True
☐ False
- (7) Virtual address space is generally limited to the amount of DRAM on a computer. ☐ True
☒ False
- (8) Dennard scaling is the claim that each transistor will generally use the same amount of power no matter how small the transistor is. ☐ True
☒ False
- (9) The size of a virtual page can be larger than the physical page. ☐ True
☒ False
- (10) A 2-bit branch predictor can get about a 0% hit rate on a branch that alternates between taken and not-taken forever (so T, N, T, N, T, N...) ☒ True
☐ False
- (11) In the 3 C's model, you would expect to be able to reduce the number of conflict misses by increasing the associativity of the cache. ☒ True
☐ False
- (12) Tags in TLBs are derived from virtual page numbers. ☒ True
☐ False
- (13) If you increase the page size while holding DRAM's size constant, you would expect the number of physical pages to increase. ☐ True
☒ False

3.	Older Stuff [9 pts]
	Things from days of yore.

1. Complete the timing diagram below for both a D latch and a rising-edge triggered D flip-flop. If a value is unknown, indicate that clearly using the notation shown. Assume there is no meaningful delay. **[4]**



2. Using 2's complement notation, write the *8-bit binary* representation of -22. **[2]**

0b11101010

0001_0110=22, so -22 is 1110_1001 +1 = 1110_1010

3. Consider an 8-bit floating point format based on the IEEE standard where the most significant bit is used for the sign (as the IEEE format), the next 3 bits are used for the exponent and the last 4 bits are used for the mantissa. The scheme uses "biased 3" to represent the exponent (rather than biased 127 used for a 32-bit IEEE floating point number) and has an implicit one just like the IEEE format. This scheme is called "VSF" (very short float).

Write the *binary* encoding of -2.5 as a VSF number. **[3]**

0b11000100

1 100 0100=-1.01*2¹=-2.5

4.	Pipeline stalls and forwarding [8 pts]
	Determine data hazards and avoidance methods in a pipeline

Consider a **5-stage** LC2K pipeline datapath that uses **detect-and-forward** to resolve data hazards, **detect-and-stall** to resolve control hazards (no branch prediction), and has **internal forwarding** for its register file.

Determine the number of pipeline stalls for each of the following benchmarks. Also, specify all (could be more than one) the instructions that receive forwarded data by shading the circles. **Ignore and do NOT specify instructions that received data through register internal forwarding.**

Benchmark 1:

<input type="radio"/>	beq	5	6	end //Not Taken
<input type="radio"/>	lw	0	1	data1
<input type="radio"/>	lw	0	1	data2
<input type="radio"/>	nor	3	4	7
<input checked="" type="radio"/>	add	2	1	3

of Stalls : 3

Benchmark 2:

<input type="radio"/>	add	2	2	3
<input type="radio"/>	lw	0	2	str
<input checked="" type="radio"/>	lw	2	7	data3
<input checked="" type="radio"/>	beq	1	7	if //Not Taken
<input type="radio"/>	add	7	5	1

of Stalls : 5

Benchmark 3:

<input type="radio"/>	add	3	3	2
<input type="radio"/>	nor	4	5	6
<input type="radio"/>	lw	0	1	data4
<input checked="" type="radio"/>	sw	0	1	data5
<input type="radio"/>	add	1	1	2

of Stalls : 1

Benchmark 4:

<input type="radio"/>	nor	4	5	6
<input type="radio"/>	add	1	2	3
<input type="radio"/>	lw	1	2	data6
<input checked="" type="radio"/>	add	2	2	4
<input type="radio"/>	nor	2	2	3

of Stalls : 1

5.	Pipeline Performance	[15 pts]
	Perform performance calculations on a given pipeline with a cache	

Consider the following 5-stage LC2K pipeline:

- **Detect-and-forward** is used to handle data hazards.
- **Speculate-and-squash** is used to handle control hazards.
- When a lw or sw instruction accesses the memory system in the MEM stage, either
 - There is a **cache hit** (95% of time time) and the pipeline **does not stall**
 - or, there is a **cache miss**, which causes the pipeline to **stall for 20 cycles** while the main memory is accessed.
- 1% of instruction fetches from an instruction cache are cache misses and result in a stall of 20 cycles.
- Throughout this problem you are to assume that no sources of stalls will overlap.

Say we run a program with the following characteristics:

lw	30%
sw	10%
add/nor	40%
beq	20%

- 25% of each type of instruction that writes to a register (**lw, add, nor**) is immediately followed by an instruction that depends on it.
- 5% of instructions that write to a register (**lw, add, nor**) are immediately followed by an independent instruction, and then immediately followed by a dependent instruction.
- 35% of branches are mispredicted.

- 1) Complete the equation for CPI below using data given above. It is fine to leave your answer as an equation that can be plugged into a calculator. **[3]**

$$\begin{aligned}
 \text{CPI} = & \quad 1 \\
 & + \quad \underline{0.21} \quad (\text{increase due to control hazards}) \quad .2 * .35 * 3 \\
 & + \quad \underline{0.075} \quad (\text{increase due to data hazards}) \quad .3 * .25 * 1 \\
 & + \quad \underline{0.6} \quad (\text{increase due to cache misses}) \quad (.01 + .05 * (.3 + .1)) * 20
 \end{aligned}$$

- 2) Say with a process shrink (i.e. the transistors are made smaller) we increase the clock frequency of the processor by a factor of 2 (including the cache but not the memory). In order to make this work, we had to split the execution stage into two stages (EX1 and EX2, where all ALU operations finish in EX2). Branches still resolve in the MEM stage. If we run the same program from part (a) on our new pipeline, what is the new CPI? It is fine to leave your answer as an equation that can be plugged into a calculator. [9]

$$\begin{aligned}
 \text{CPI} = & \quad 1 \\
 & + \quad \underline{0.28} \quad (\text{increase due to control hazards}) \quad .2 \cdot .35 \cdot 4 \\
 & + \quad \underline{0.265} \quad (\text{increase due to data hazards}) \quad .3 \cdot .25 \cdot 2 + .3 \cdot .05 \cdot 1 + .4 \cdot .25 \cdot 1 \\
 & + \quad \underline{1.2} \quad (\text{increase due to cache misses}) \quad (.01 + .05 \cdot (.3 + .1)) \cdot 40
 \end{aligned}$$

- 3) Say for part 1) you had found a CPI of 2.0 and for part 2) you had found a CPI of 3.0. What would be the speedup¹ after our process shrink expressed as a percentage? It is fine to leave your answer as an equation that can be plugged into a calculator. [3]

$$\underline{133} \% \quad \text{Old} = 2.0 \cdot 1; \text{ New} = 3.0 \cdot .5. \text{ Speedup} = 2 / 1.5$$

¹Recall that in general, if we say Processor A is 50% of the speed of processor B, we mean A is half as fast. Which is the same as saying that A takes twice as long to do the task. If we say Processor A is 300% the speed of processor B, that means it is 3 times as fast.

6.	Cache Basics [8 pts]
	Just the facts

Indicate the number of bits used for the index and block offset of each of the following caches. Assume the address size (physical and virtual) is 32-bits and that memory is byte addressable.

1) A 128KB, 4-way associative cache with 32-byte blocks. Index: 10 Offset: 5

2) A 1MB, direct-mapped cache with 16-byte blocks. Index: 16 Offset: 4

3) A 48KB, 3-way associative cache with 8-byte blocks. Index: 11 Offset: 3

4) A 1MB fully-associative cache with 128-byte blocks. Index: 0 Offset: 7

7.	<div style="display: flex; justify-content: space-between;"> Examining the Memory System Bit by Bit [10 pts] </div> <div style="border: 1px solid black; padding: 5px; margin-top: 5px;">Working with the data</div>
-----------	--

Consider a byte-addressable architecture with 12-bit virtual addresses. The system has a maximum of 1KB of physical memory with 16-byte page sizes. The system has a 16-byte 2-way set associative physically-addressed cache with a 2-byte block size and a fully-associative TLB with 4 entries. The TLB, cache contents, and a *snippet* of the single level page table are provided below:

Page table

VPN	PPN	Valid
0x00	0x01	1
0x01	0x03	0
0x02	0x0F	1
0x03	0x00	1
0x04	0x22	1
0x05	0x1A	0
0x06	0x31	0
0x07	0x13	1

VPN	PPN	Valid
0x08	0x25	1
0x09	0x26	1
0x0A	0x3A	1
0x0B	0x3A	0
0x0C	0x1B	1
0x0D	0x1C	1
0x0E	0x27	1
0x0F	0x1F	0

TLB

Tag	PPN	Valid
0x02	0x0F	1
0x2B	0x1A	1
0x06	0x0A	0
0x0D	0x1C	1

Cache

Set Index	Tag	Valid	Byte0	Byte 1
0	0x37	1	0xDE	0xAD
	0x1A	0	0x12	0xB0
1	0x65	1	0x0A	0xC1
	0x4F	1	0x99	0x1F
2	0x1C	1	0x84	0x92
	0x00	1	0xBE	0xCF
3	0x7B	1	0xCC	0xA0
	0x0A	1	0x45	0x67

Say that the processor reads one byte from virtual address **0x0EB**. Recall that the cache is physically addressed. Answer the following questions. Provide all numeric answers in hex. If the answer is unknown, write “unknown”. Note that early errors on this problem could cause later answers to also be wrong.

- 1) In hex, what is the virtual page number associated with that address? [2] 0x**E**_____
- 2) Is this a TLB hit? [1] ☐ Yes ☒ No
- 3) In hex, what is the physical page number associated with this access? [2] 0x**27**_____
- 4) What set index could hold the data? [2] 0x**1**_____
- 5) What is the value of that byte of memory being accessed? [3] 0x**1F**_____

8.	Hierarchical Page tables [10 pts]
	Clearly write answers in the blanks provided.

Consider a 42-bit byte-addressable system that supports virtual memory with the following specifications:

- A page is 2 KB.
- The page table is hierarchical with 3 levels.
- The first-level page table occupies exactly 1 page of memory.
- Each second-level page table occupies exactly 2 pages of memory.
- All page table entries are 4 bytes .
- A maximum of 32 GB of physical memory can be installed.

Provide a number (rather than an equation) in each blank for parts 1-6. Something like 2^{20} is fine. $2^{22}/4$ is not.

1. How many bits are used for the page offset? [1] 11 bits
2. How many virtual pages exist in this system? [1] 2^{31} pages
3. How many physical pages can exist in the system? [1] 2^{24} pages
4. How many bits in a virtual address are used to index the first-level page table? [1] 9 bits
5. How many bits in a virtual address are used to index a second-level page table? [1] 10 bits
6. How many bits in a virtual address are used to index a third-level page table? [2] 12 bits
7. In the worst case, how many pages would this page table occupy? [3] $1+2^{10}+2^{22}$ pages (You can write an "equation", rather than number, for this one)

9.	Cache Analysis [12 pts]
	Deeper thoughts about caches

- 1) Consider a 1024-byte **direct-mapped cache** with a block size of 32 bytes. The cache starts empty and P, Q, and R are addresses. You are given the following stream of address references. Cache misses are marked as “M”, while hits are marked as “H”.

Address	P	Q	R	P	Q
Cache Access	M	M	H	M	M

- a. From the above you can be **sure** of which of the following? Select *all* that are true. [3]
- ☐ P and Q are in the same cache block
 - ☐ P and R are in the same cache block
 - ☒ R and Q are in the same cache block
 - ☐ You cannot be sure of any of the three above.
- b. From the above you can be **sure** of which of the following? Select *all* that are true. [3]
- ☒ P and Q have the same line index, but are in different cache blocks
 - ☒ P and R have the same line index, but are in different cache blocks
 - ☐ R and Q have the same line index, but are in different cache blocks
 - ☐ You cannot be sure of any of the three above.
- 2) Consider a 1024-byte **2-way associative cache** with a block size of 32 bytes. The cache starts empty and A, B, C, D, and E are addresses. You are given the following stream of address references. Cache misses are marked as “M”, while hits are marked as “H”..

Address	A	B	C	D	E	D	C	B	A
Cache Access	M	M	H	M	M	H	M	H	H

- a. From the above which of the following **could** be true? Select *all* that could be true. [3]
- ☐ A and B are in the same cache block
 - ☒ A and C are in the same cache block
 - ☒ B and C are in the same cache block
 - ☐ D and E are in the same cache block
- b. From the above which of the following **could** be true? Select *all* that could be true. [3]
- ☐ A and B map to the same set, but different cache blocks
 - ☐ B and C map to the same set, but different cache blocks
 - ☒ B and D map to the same set, but different cache blocks
 - ☒ D and E map to the same set, but different cache blocks