# MUIC: Computer System and Architecture

Date: Thursday, March 31, 2022 Instructor: Rachata Ausavarungnirun

Name:	Monthon	Kraitheerawut	-	
		Problem 1 (35 Points):		
		Problem 2 (15 Points):		
		Problem 3 (30 Points):		
		Problem 4 (20+5 Points):		
		Total (100+5 Points):		

#### Instructions:

- 1. This exam lasts 24 hours. This is an open-book, open-everything exam. Every single question can be answered through materials we have covered.
- 2. Clearly indicate your final answer for each problem.
- 3. Please show your work when needed.
- 4. Please write your initials at the top of every page.
- 5. Please make sure that your answers to all questions (and all supporting work that is required) are contained in the space required.
- 6. **DO NOT CHEAT.** If I catch you cheating in any shape or form, you will be penalized based on our plagiarism policy (N\* 10% of your total grade, where N is the number of times you plagiarized previously).

## Tips:

- Read everything. Read all the questions on all pages first and formulate a plan.
- Be cognizant of time. Do not spend too much time on one question.
- Be concise. You will be penalized for verbosity and unnecessarily long answers.
- Show work when needed. You will receive partial credit at the instructors' discretion.
- Write legibly. Show your final answer.

## 1. Potpourri [35 points]

#### (a) ISA vs. uArch [7 points]

For each of the following, circle if the concept is a part of the ISA versus the microarchitecture (circle only one):

Number of cycles needed for the adder

Circle one: ISA Microarchitecture

Program counter

Circle one: (ISA) Microarchitecture

Pipeline registers

Circle one: ISA Microarchitecture

Number of registers

Circle one: ISA Microarchitecture

Cache replacement policy

Circle one: (ISA) Microarchitecture

ADD instruction

Circle one: (ISA) Microarchitecture

The branch predictor

Circle one: ISA Microarchitecture

#### (b) Performance Measurement [10 points]

You and your friends both buy a CPU with 3.5 GHz clock with the exact same motherboard, DRAM and SSD. However, when running the exact same assembly code (Program A), your friend's CPU is somehow 20% faster than your machine. What happen here? Explain your reasoning.

The reason might be Program A strelf as it might need something like 5/0 operation to a parate which needs Buses for sending /recessor, fits for 1th CPU, and another reason might be that my CPU have lesson core than my friend.

After noticing your friends computer runs Program A faster then your computer, you decide to buy the CPU with your same spec except this time there are 16 cores and a high-end GPU instead of 4 cores. Yet again, your friends' computer is still 20% faster. What can happen in this case? Please explain your reasoning and assumptions in detail.

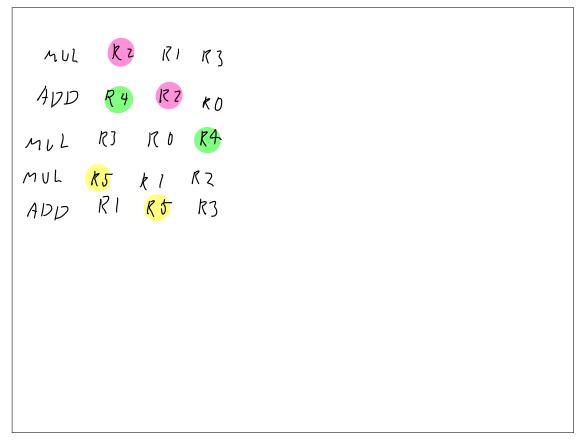
Program A might not work well against CPU with many cores as it might not optimize well with CPU with lb bores.

#### (c) Design Tradeoffs [6 points]

What is the key benefit of branch prediction over stalling?

Without branch prediction, the processor would have to wait until the condition surprincition has passed the execute stage define the next instruction can enter the feter stage in the pipeline. It attemps to avoid writing time by messing whenther the condition jump is not likely to be taken or not. Faster performance

Please provide a simple code example in MIPS ISA showing data dependency assuming a 5-stage pipeline? You must highlight the dependency.



#### (d) DRAM Microarchitecture [4 points]

List one difference between a DRAM bank and a DRAM channel.

DRAM Chunel is a connection path between the Memory Controller and prim module

DRAM Bank is a set of independent array inside a DRAM chip Each Bank of menory is an independent array that can be in different phases of a date access/refush yele.

What is the benefit of having more DRAM channels?

As each PRAN chands is physicisty independent from other channels and has its own nearly controller that send series of DRAM cornands to read/vite data to DRAM cell so everything vill be faster

### (e) Simple Branch Prediction [8 points]

What is the accuracy of a branch predictor is we always take the branch?

$$B_1 = 100 \%$$

$$B_2 = 50\%$$
Within  $0 - 100$  there are 50 even overless

## 2. 5-stage Pipeline [15 points]

In this question, we will use the code below as a reference code we are running. For the instruction below, assume the instruction format is [operation name] [dest], [src1], [src2]. Assume that for every operation, fetch, decode, execute, memory and writeback stages take one cycle.

```
ADD R3, R1, R6 | MUL R2, R0, R3 | MUL R4, R2, R3 | ADD R6, R2, R4 | MUL R6, R4, R4 | MUL R6, R4 | MUL
```

Assume **no pipelining**, and add instructions takes 5 cycles total while multiply takes 6 cycles total. How many cycles would it take to finish running the code above?

Now, let's assume a 5-stage pipeline with no data forwarding and the adder takes 1 cycle to execute, the multiplier takes two cycles to execute, all other stages take one cycle to execute, fill in the table below for how each instruction runs and what stage would each instruction be at each of these clock cycle. Put in **F** for fetch, **D** for decode, **E** for execute, **M** for memory and **W** for writeback. Leave the block blank if the instruction is stalling the pipeline.

Cycle:	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
MUL R3, R1, R6	F	D	$\epsilon$	6	M	W														
ADD R2, R0, R3		F	D				6	M	W											
MUL R4, R2, R3			F				P			6	6	M	W							
ADD R6, R1, R0							F			D	6		М	W						
MUL R6, R4, R4										F	D			6	6	M	W			

Finally, let's assume a 5-stage pipeline with data forwarding, fill in the table below for how each instruction runs and what stage would each instruction be at each of these clock cycle. Put in **F** for fetch, **D** for decode, **E** for execute, **M** for memory and **W** for writeback.

Cycle:	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
MUL R3, R1, R6	F	D	6	Μ	W															
ADD R2, R0, R3		F	D		$\epsilon$	M	٧													
MUL R4, R2, R3			F		D	E	6	M	V											
ADD R6, R1, R0					F	D	6		M	W										
MUL R6, R4, R4						6	b	6	6	M	W									

### 3. Caching in Multicore [30 points]

Assume that we have a processor with 8 cores and 1-level 2kB, 2-way cache associative *shared* cache with a block size of 32 bytes, and assume that an integer is 8 bytes and an address is 32 bits. The array j and k are also 8 bytes and we are running the code below.

What is the cache hit rate of the code above assuming that this is the only thread running on the system? Show your work and make sure you specify the tag, index bits.

Index = 
$$69_{2}(32) = 5$$
 bits

bit  $75t_{2} = \frac{30,000}{40,000}$ 

this is  $\frac{3}{4} = 75\%$ 

Total access

Therefore  $600 = \frac{20,000}{4} = \frac{20,000}{40,000} = \frac{10,000}{40,000}$ 

hit =  $\frac{30,000}{4}$ 

Your friend suggest that you can improve the performance of your program by using pthread. So, you ended up spawning four threads, each of which are responsible for a portion of the loop. You also ensure that you load balance the work in each threads (i.e., each thread is responsible for exactly 5000 loop iterations). While this is likely make the program much faster, you observe that the cache hit rate drops.

What can contribute to the drop in the cache hit rage? (Hint: Our cache is a shared cache)

Generally cache hit rate generally drop for the type of access, size of cache and frequency of consistency checks and in this question 4 threads are trying to rues the cache which might be the reason in the drop of hit rate

After running this multithreaded version of the code multiple times, you observe a variation of the cache miss rate. Is this possible or are you dreaming? If this is possible, please explain why and provide the minimum cache hit rate that your new code can actually see. If you are dreaming, please provide a reason to convince yourself so that you wake up from this nightmare.

possible because multithreaded rule each thread runs
parallel to each other and create in additional stress on
nevery hierarchy caused by the interference army threads

#### 4. DRAM [20+5 points]

In class, we learn of a few things in DRAM. One of which is the DRAM banks and the DRAM row.

Why do we need a row buffer and what is the effect of having the row buffer to the amount of time it takes to access DRAM?

Let's go back to our computer in Question 3 Assume that we have a 1-level 4kB, 4-way cache associative cache with a block size of 32 bytes. Let's also assume a row buffer is 4kB.

Let's assume I have a series of 50 data accesses that results in cache misses from each of the two CPU cores (Core A and Core B), making it 100 data accesses, all of which are cache misses and all go to the same DRAM channel and DRAM bank. Let's also assume that a row buffer hit (i.e., access to the address that is already on the row buffer) takes 15 ns while a row buffer miss (an access to an address not on the row buffer) takes 50 ns. Let's also assume we do not have to worry about any other DRAM timing parameters such as DRAM refreshes.

What is the minimum and the maximum time it takes to process these 100 data accesses? Why?

nitials:								
xtra credit:	5 points]	If the size	of my rov	w buffer ch	ange from	4kB to 2kl	3.	
hat is the r	ninimum a	and the ma	aximum ti	me it take	s to proces	s these 100	data acces	sses? Why