

CMU 18-447 INTRODUCTION TO COMPUTER ARCHITECTURE, SPRING 2015

MIDTERM EXAM 1

DATE: FRI., 3/20

INSTRUCTOR: ONUR MUTLU

TAS: RACHATA AUSAVARUNGNIRUN, KEVIN CHANG, ALBERT CHO, JEREMIE KIM, CLEMENT LOH

Name:

Problem 1 (80 Points):	<input type="text"/>
Problem 2 (40 Points):	<input type="text"/>
Problem 3 (40 Points):	<input type="text"/>
Problem 4 (60 Points):	<input type="text"/>
Problem 5 (30 Points):	<input type="text"/>
Problem 6 (30 Points):	<input type="text"/>
Problem 7 (60 Points):	<input type="text"/>
Bonus (60 Points):	<input type="text"/>
Total (340 + 60 Points):	<input type="text"/>

Instructions:

1. This is a closed book exam. You are allowed to have one letter-sized cheat sheet.
2. No electronic devices may be used.
3. This exam lasts 1 hour and 50 minutes.
4. Clearly indicate your final answer for each problem.
5. Please show your work when needed.
6. Please write your initials at the top of every page.
7. Please make sure that your answers to all questions (and all supporting work that is required) are contained in the space required.

Tips:

- **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.

Initials: _____

1. Potpourri [80 points]

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

Number of threads in fine-grained multithreading

Circle one: ISA Microarchitecture

Number of DRAM banks

Circle one: ISA Microarchitecture

Vectored interrupts

Circle one: ISA Microarchitecture

Number of entries in reservation stations

Circle one: ISA Microarchitecture

Number of entries in the reorder buffer

Circle one: ISA Microarchitecture

Number of entries in the architectural register file

Circle one: ISA Microarchitecture

Number of entries in the physical register file

Circle one: ISA Microarchitecture

Number of sets in the L3 cache

Circle one: ISA Microarchitecture

The page table base register of the executing process

Circle one: ISA Microarchitecture

- (b) Why does ISA change more slowly than microarchitecture?

- (c) A program is written in C. We execute this program on two different computers:

Computer A: has a processor that implements the x86 ISA and has 3 GHz clock frequency

Computer B: has a processor that implements the x86 ISA and has 3 GHz clock frequency

When we execute this program and measure its cycles per instruction (CPI) in x86 instructions, we find the following result:

On Computer A: CPI is equal to 10

On Computer B: CPI is equal to 8

What can you say about on which computer (A or B) this program runs faster?

Initials: _____

Explain and show all your work below:

- (d) You are designing an ISA that uses delayed branch instructions. You are trying to decide how many instructions to place into the branch delay slot. How many branch delay slots would you need for the following different implementations? Explain your reasoning briefly.

An in-order processor where conditional branches resolve during the 4th stage:

An out-of-order processor with 64 unified reservation station entries where conditional branches resolve during the 2nd cycle of branch execution. The processor has 15 pipeline stages until the start of the execution stages.

- (e) What three key pieces of information does the compiler not know when performing instruction scheduling?

Initials: _____

- (f) In class, we discussed the concept of traces and trace scheduling.

What is a trace?

Now suppose we make each trace atomic, as we also discussed in class. What is the benefit of making each trace atomic? Explain.

What is the disadvantage?

- (g) Assume we have an ISA with virtual memory. It is byte-addressable and its address space is 64 bits. The physical page size is 8KB. The size of the physical memory we use in a computer that implements the ISA is 1 TB (2^{40} bytes).

Assume the demand paging system we would like to design for this uses the perfect LRU algorithm to decide what to evict on a page fault. What is the minimum number of bits that the operating system needs to keep to implement this perfect LRU algorithm? Show your work.

Initials: _____

2. Register Renaming [40 points]

In this problem, we will give you the state of the Register Alias Table (RAT), Reservation Stations (RS), and Physical Register File (PRF) for a Tomasulo-like out-of-order execution engine.

The out-of-order machine in this problem has the following characteristics:

- The processor is fully pipelined with four stages: Fetch, decode, execute, and writeback.
- For all instructions, fetch takes 1 cycle, decode takes 1 cycle, and writeback takes 1 cycle.
- The processor implements ADD and MUL instructions only. Both the adder and multiplier are fully pipelined. ADD instructions take 3 cycles and MUL instructions take 4 cycles in the execute stage. Note that the adder and multiplier have separate common data buses (CDBs), which allow both the adder and multiplier to broadcast results in the same cycle.
- An instruction always allocates the first reservation station that is available (in top-to-bottom order) at the required functional unit.

Suppose the pipeline is initially empty and the machine fetches exactly 5 instructions. The diagram below shows the snapshot of the machine at a particular point in time.

Register Alias Table

ID	V	Tag
R0	1	P0
R1	1	P8
R2	0	P15
R3	1	P3
R4	0	P10
R5	1	P5
R6	1	P12
R7	1	P7

Physical Register File

ID	V	Data	ID	V	Data
P0	1	1	P8	1	87
P1	0	10	P9	1	90
P2	0	2	P10	1	11
P3	1	30	P11	0	110
P4	0	3	P12	1	37
P5	1	50	P13	0	130
P6	0	5	P14	1	17
P7	1	70	P15	1	159

ADD Reservation Station

ID	V	Tag	V	Tag	Dest. Tag
A	0	-	0	-	-
B	1	P15	1	P5	P14

ADD CDB

MUL Reservation Station

ID	V	Tag	V	Tag	Dest. Tag
X	1	P8	1	P12	P15
Y	1	P8	1	P15	P9
Z	0	P14	0	P14	P10

MUL CDB

Initials: _____

- (a) Your first task is to use only the supplied information to draw the data flow graph for the five instructions which have been fetched. Label nodes with the operation (+ or *) being performed and edges with the architectural register alias numbers (e.g., R0).

- (b) Now, use the data flow graph to fill in the table below with the five instructions being executed on the processor in program order. The source registers for the first instruction can be specified in either order. Give instructions in the following format: “opcode, destination, source1, source2.”

OP	Dest	Src 1	Src 2
			R1

- (c) Now show the full pipeline timing diagram below for the sequence of five instructions that you determined above, from the fetch of the first instruction to the writeback of the last instruction. Assume that the machine stops fetching instructions after the fifth instruction.

As we saw in class, use F for fetch, D for decode, En to signify the nth cycle of execution for an instruction, and W to signify writeback. You may or may not need all columns shown. Finally, identify the cycle after which the snapshot of the microarchitecture was taken. Shade the corresponding cycle in the last row of the table.

Cycle:	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Instruction 1																			
Instruction 2																			
Instruction 3																			
Instruction 4																			
Instruction 5																			
Snapshot cycle																			

3. Branch Prediction [40 points]

Assume the following piece of code that iterates through two large arrays, j and k , each populated with completely (i.e., truly) random positive integers. The code has two branches (labeled B1 and B2). When we say that a branch is *taken*, we mean that the code *inside* the curly brackets is executed. Assume the code is run to completion without any errors (there are no exceptions). For the following questions, assume that this is the only block of code that will ever be run, and the loop-condition branch (B1) is resolved first in the iteration before the if-condition branch (B2). N and X are unspecified non-zero integers.

```
for (int i = 0; i < N; i++) {           /* B1 */
    if (i % X == 0) {                   /* B2 */
        j[i] = k[i] - i;               /* TAKEN PATH for B2 */
    }
}
```

You are running the above code on a machine with a two-bit global history register (GHR) shared by all branches, which starts at *Strongly Not Taken* ($2'b00$). Each pattern history table entry (PHTE) contains a 2-bit saturating counter, which is initialized to *Strongly Taken* ($2'b11$)

The saturating counter values are as follows:

- $2'b00$ - Strongly Not Taken
- $2'b01$ - Weakly Not Taken
- $2'b10$ - Weakly Taken
- $2'b11$ - Strongly Taken

- (a) Assuming that N is larger than 10 (ten), after running the loop for 10 iterations, you observe that the branch predictor mispredicts 0% of the time. What is the value of X ? Explain your reasoning.

Initials: _____

- (b) What is the prediction accuracy of the branch predictor if $N = 20$ and $X = 2$? Explain your answer. You can leave your answer as a fraction.

- (c) Assuming that N is larger than 10 (ten), after running this code until it exits the loop, you observe that the state of the branch predictor is as follows:

GHR: 2'b00

PHTE:

NN:2'b11

NT:2'b00

TN:2'b10

TT:2'b11

What can you tell about the value of X and N ? Explain your reasoning.

Initials: _____

4. GPUs and SIMD [60 points]

We define the *SIMD utilization* of a program running on a GPU as the fraction of SIMD lanes that are kept busy with *active threads*.

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 and B, are already in the vector registers so there are no loads and stores in this program. Hint: Notice that there are 2 instructions in each thread. A warp in this GPU consists of 32 threads, and there are 32 SIMD lanes in the GPU. Assume that each instruction takes the same amount of time to execute.

```
for (i = 0; i < N; i++) {  
    if (A[i] % 3 == 0) {        // Instruction 1  
        A[i] = A[i] * B[i];    // Instruction 2  
    }  
}
```

- (a) What's the minimum number of bits required to encode the warp ID in order to execute this program? Please leave the answer in terms of N .

- (b) Assume integer array A has a repetitive pattern which consists of 24 ones followed by 8 zeros, and integer array B has a different repetitive pattern which consists of 48 zeros followed by 64 ones. What is the SIMD utilization of this program?

- (c) Is it possible for this program to yield a SIMD utilization of 100% (circle one)?

YES NO

If YES, what should be true about arrays A for the SIMD utilization to be 100%?

If NO, explain why not.

What should be true about array B?

Initials: _____

(d) Is it possible for this program to yield a SIMD utilization of 56.25% (circle one)? *Hint: $56.25\% = 36/64$.*

YES NO

If YES, what should be true about arrays A for the SIMD utilization to be 56.25%?

What should be true about arrays B?

If NO, explain why not.

(e) Is it possible for this program to yield a SIMD utilization of 50% (circle one)?

YES NO

If YES, what should be true about arrays A for the SIMD utilization to be 50%?

What should be true about arrays B?

If NO, explain why not.

Initials:

In lecture, we learned a technique called dynamic warp formation, which tries to improve SIMD utilization, by merging threads executing the same instruction together. The key idea of dynamic warp formation is to move an "active" thread from one warp to another warp such that SIMD utilization is improved and all threads can access their registers.

Consider the following three warps X, Y, and Z that are executing the same code segment specified in the beginning of this question. Assume that the vectors we provide below specify the "active mask", i.e., whether or not the instruction should be executed by each thread in the warp: 1 means the instruction should be executed, 0 means it should not be executed. Assume each warp is at the same Program Counter.

```
Warp X = {10000000000000000000000000000010}  
Warp Y = {10000000000000000000000000000001}  
Warp Z = {01000000000000000000000000000000}
```

- (f) Suppose that you perform dynamic warp formation on these three warps. What are the resulting active masks for each of the newly formed warps X', Y', and Z'? Explain your reasoning as necessary.

- (g) Given the original specification for arrays A and B (integer array A has a repetitive pattern which consists of 24 ones followed by 8 zeros, and integer array B has a different repetitive pattern which consists of 48 zeros followed by 64 ones). Is it possible for this program to yield a better SIMD utilization than the baseline when dynamic warp formation is used?

Initials: _____

5. Caches [30 points]

A byte-addressable system with 16-bit addresses ships with a three-way set associative, write-back cache. The cache implements a true LRU replacement policy using the minimum number of replacement policy bits necessary to implement it. The tag store requires a total of 264 bits of storage. What is the block size of the cache? (Hint: $264 = 2^8 + 2^3$)

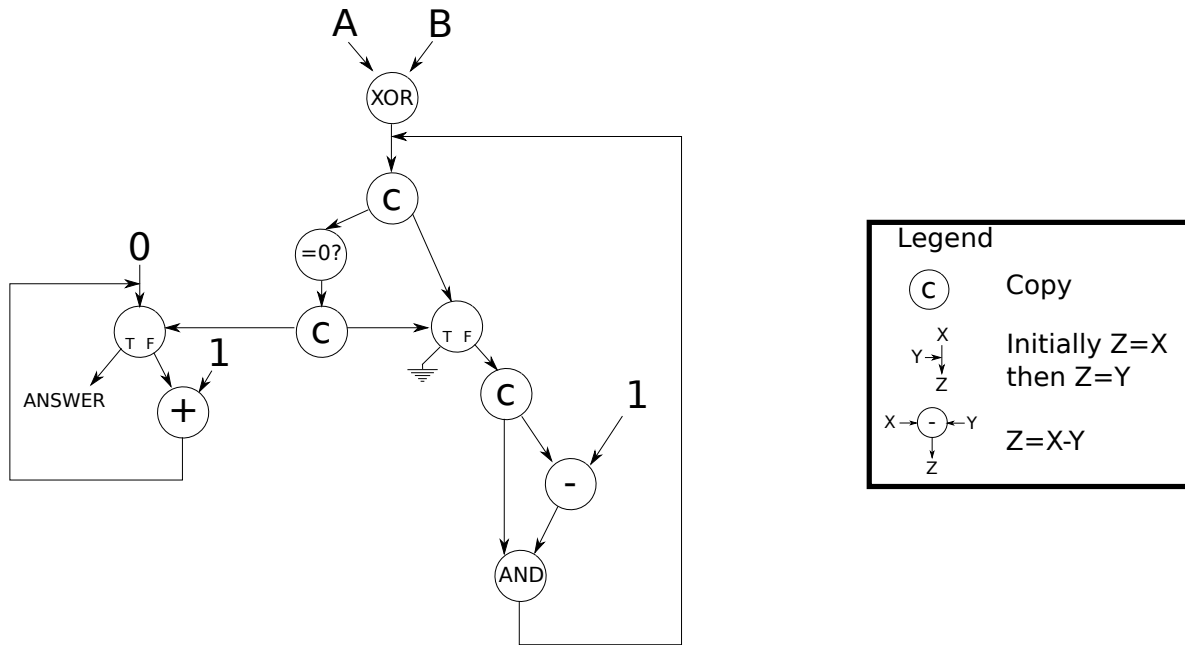
Answer:

Show all your work.

Initials: _____

6. Dataflow [30 points]

Here is a dataflow graph representing a dataflow program:



Note that the inputs, A and B, are non-negative integers.

What does the dataflow program do? Specify clearly in less than 15 words.

Initials: _____

7. VLIW and Instruction Scheduling [60 points]

Explain the motivation for VLIW in one sentence.

You are the human compiler for a VLIW machine whose specifications are as follows:

- There are **3 fully** pipelined functional units (ALU, MU and FPU).
 - Integer Arithmetic Logic Unit (ALU) has a 1-cycle latency.
 - Memory Unit (MU) has a 2-cycle latency.
 - Floating Point Unit (FPU) has a 3-cycle latency, and can perform either FADD or FMUL (floating point add / floating point multiply) on **floating point registers**.
 - This machine has **only** 4 integer registers (r1 .. r4) and 4 floating point registers (f1 .. f4)
 - The machine does not implement hardware interlocking or data forwarding.
- (a) For the given assembly code on the next page, fill **Table 1** (on the next page) with the appropriate VLIW instructions for only one iteration of the loop (The C code is also provided for your reference). Provide the VLIW instructions that lead to the **best** performance. Use the minimum number of VLIW instructions. Table 1 should **only** contain instructions provided in the assembly example. For all the instruction tables, show the NOP instructions you may need to insert. Note that BNE is executed in the **ALU**.

The base addresses for A, B, C are stored in r1, r2, r3 respectively. The address of the last element in the array C[N-1] is stored in r4, where N is an integer multiplier of 10! (read: 10 factorial).

Three extra tables are available for you to work with in the scratchpad, for the entirety of this question.

Initials: _____

C code	Assembly Code
float A[N];	loop: LD f1, 0 (r1)
float C[N];	LD f2, 0 (r2)
int B[N];	FMUL f1, f1, f1
... <i>//code to initialize A and B</i>	FADD f1, f1, f2
for (int i=0; i<N; i++)	ADDI r3, r3, 4
C[i] = A[i] * A[i] +B[i];	ST f1, -4 (r3)
	ADDI r1, r1, 4
	ADDI r2, r2, 4
	BNE r3, r4, loop

VLIW Instruction	ALU	MU	FPU
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			

Table 1

What is the performance in Ops/VLIW instruction (Operations/VLIW instruction) for this design? An operation here refers to an instruction (in the Assembly Code), excluding NOPs.

Initials: _____

- (b) Assume now we decide to unroll the loop once. Fill **Table 2** with the new VLIW instructions. You should optimize for latency first, then instruction count. **You can choose to use different offsets, immediates and registers, but you may not use any new instructions.**

VLIW Instruction	ALU	MU	FPU
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			

Table 2

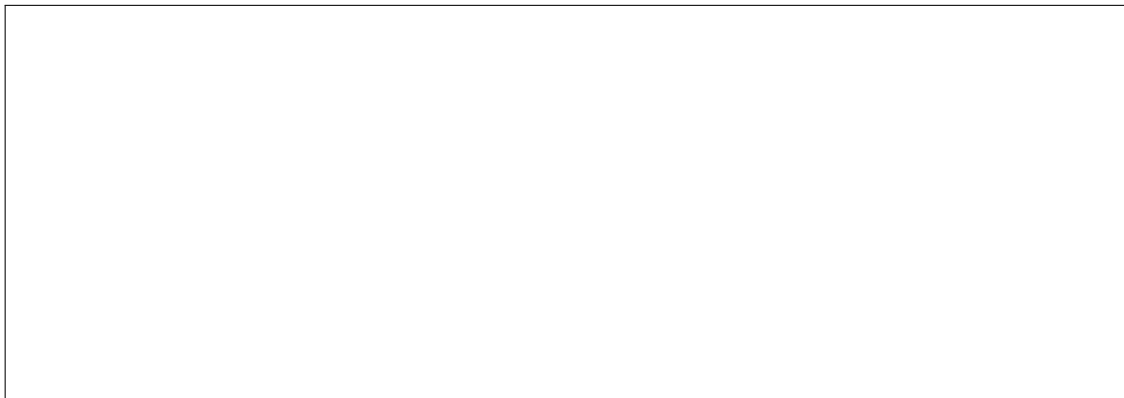
What is the performance in Ops/VLIW instruction for this design?

Initials: _____

- (c) Assume now we have **unlimited registers** and the loop is fully optimized (unrolled to the **best** performance possible). What is the performance in Ops/cycle for this design? Show your work and explain **clearly** how you arrived at your answer. You are not required to draw any tables, but you may choose to do so to aid your explanation. You will receive **zero** credit for a correct answer without any explanation. (Hint: trace the dependent instructions)



What is the performance bottleneck for this code and why? Explain.



Initials: _____

8. [Bonus] Mystery Instruction [60 points]

A pesky engineer implemented a mystery instruction on the LC-3b. It is your job to determine what the instruction does. The mystery instruction is encoded as:

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
1010				DestR				SR1			0	0	0	SR2		

The instruction is only defined if the value of SR2 is greater than the value of SR1.

The modifications we make to the LC-3b datapath and the microsequencer are highlighted in the attached figures (see the four pages at the end of this question). We also provide the original LC-3b state diagram, in case you need it. (As a reminder, the selection logic for SR2MUX is determined internally based on the instruction.)

The additional control signals are

LD_TEMP1/1: NO, YES

LD_TEMP2/1: NO, YES

GateTEMP3/1: NO, LOAD

Reg_IN_MUX/1: BUS, Mystery2 – (Assume BUS is asserted if this signal is not specified)

Mystery_MUX/2: SR2MUX, PASS_1 (outputs value 1), PASS_0 (outputs value 0)

Additional Signals for ALUK: PASS_B (outputs the value from input B), SUB (A-B)

Also note that both of DRMUX and SR1MUX can now choose DR, SR1, and SR2

COND/4:

COND₀₀₀₀ ;Unconditional

COND₀₀₀₁ ;Memory Ready

COND₀₀₁₀ ;Branch

COND₀₀₁₁ ;Addressing mode

COND₀₁₀₀ ;Mystery 1

COND₁₀₀₀ ;Mystery 2

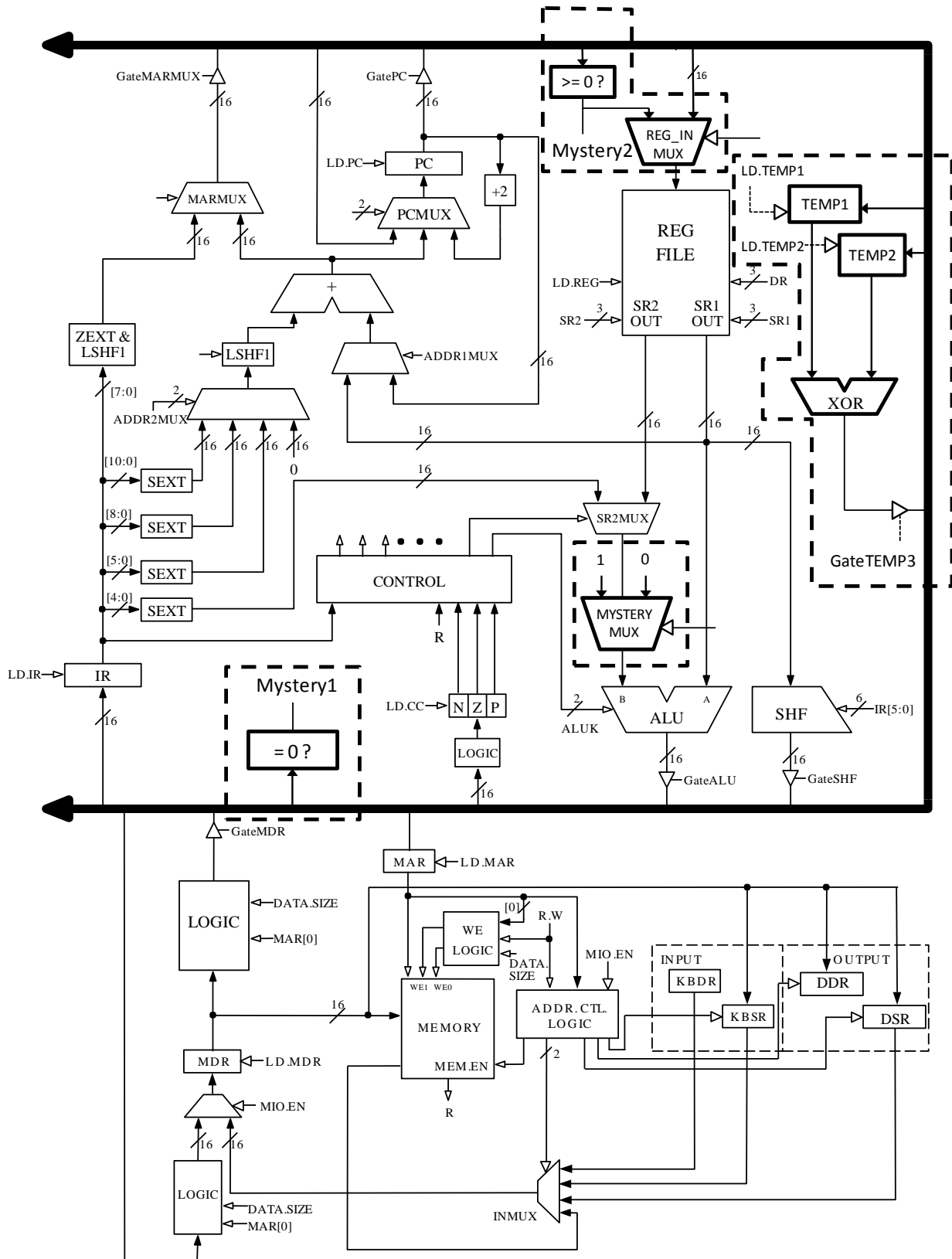
Initials: _____

The microcode for the instruction is given in the table below.

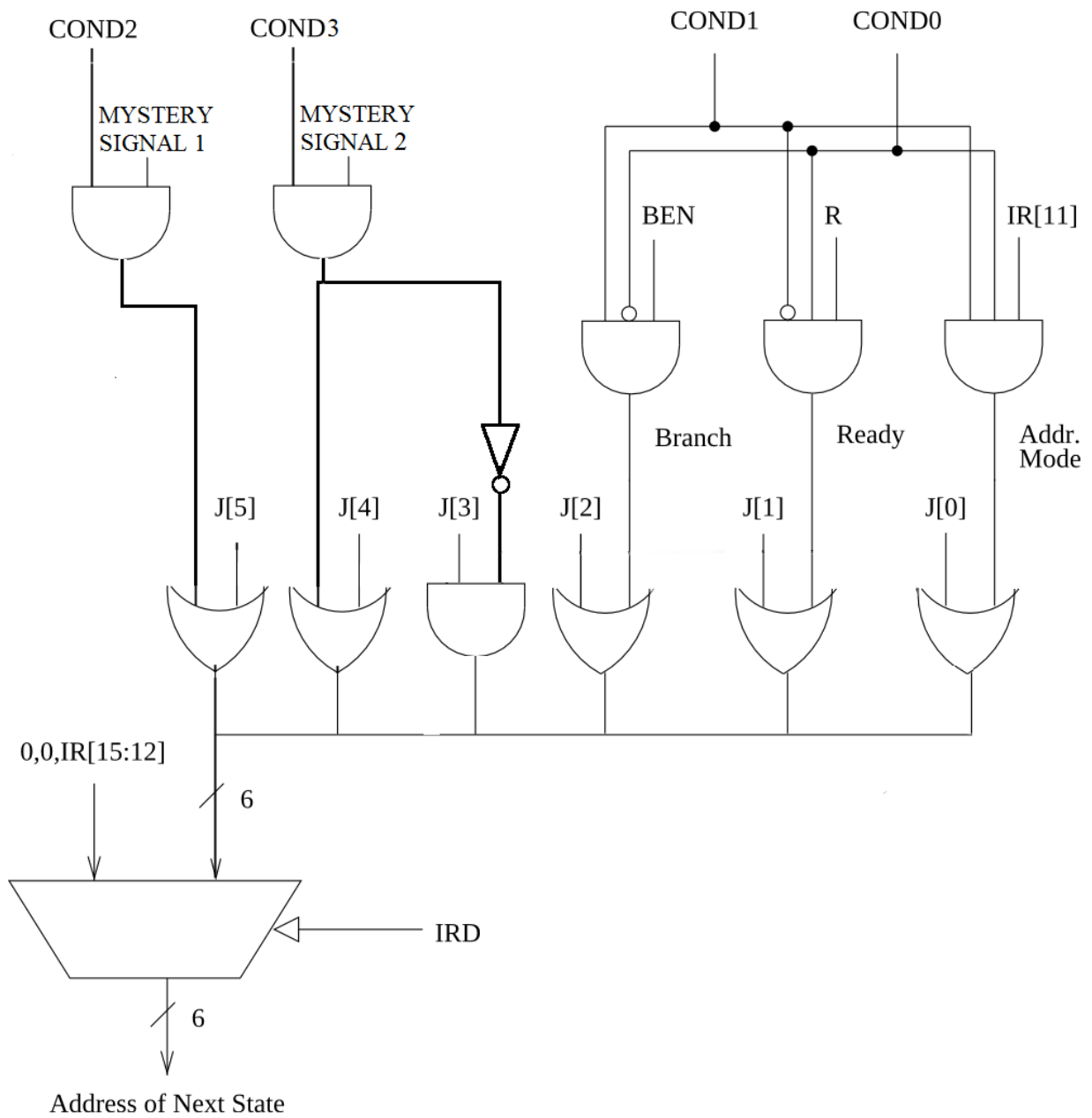
State	Cond	J	Asserted Signals
001010 (10)	COND ₀₀₀₀	001011	LD.REG, DRMUX = DR(IR [11:9]), GateALU, ALUK = PASS_B, MYSTERY_MUX = PASS_0
001011 (11)	COND ₀₀₀₀	110001	LD.MAR, SR1MUX = SR1(IR[8:6]), ADDR1MUX = SR1OUT, ADDR2MUX = 0, MARMUX = ADDER, GateMARMUX
110001 (49)	COND ₀₀₀₁	110001	LD.MDR, MIO.EN, DATA.SIZE=BYTE, R.W = R
110011 (51)	COND ₀₀₀₀	100100	GateMDR, LD.TEMP1, DATA.SIZE=BYTE
100100 (36)	COND ₀₀₀₀	100101	LD.MAR, SR1MUX = SR2(IR[2:0]), ADDR1MUX = SR1OUT, ADDR2MUX = 0, MARMUX = ADDER, GateMARMUX
100101 (37)	COND ₀₀₀₁	100101	LD.MDR, MIO.EN, DATA.SIZE=BYTE, R.W = R
100111 (39)	COND ₀₀₀₀	101000	GateMDR, LD.TEMP2, DATA.SIZE=BYTE
101000 (40)	COND ₀₁₀₀	010010	GateTEMP3
110010 (50)	COND ₀₀₀₀	101001	LD.REG, DRMUX = SR1(IR[8:6]), GateALU, ALUK = ADD, SR1MUX = SR1(IR[8:6]), MYSTERY_MUX = PASS_1
101001 (41)	COND ₀₀₀₀	101010	LD.REG, DRMUX = SR2(IR[2:0]), GateALU, ALUK = SUB, SR1MUX = SR2 (IR[2:0]), MYSTERY_MUX = PASS_1
101010 (42)	COND ₁₀₀₀	001011	LD.REG, DRMUX = DR (IR[11:9]), Reg.IN_MUX = MYSTERY2, GateALU, ALUK = SUB, SR1MUX = SR1(IR[8:6]), MYSTERY_MUX = SR2MUX

Describe what this instruction does. Show your work for partial credit.

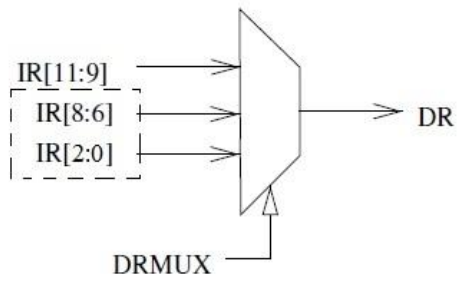
Initials: _____



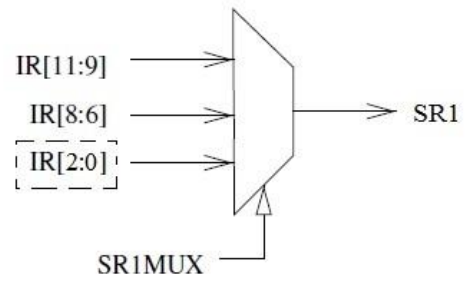
Initials: _____



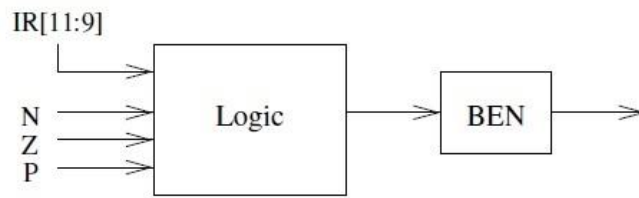
Initials: _____



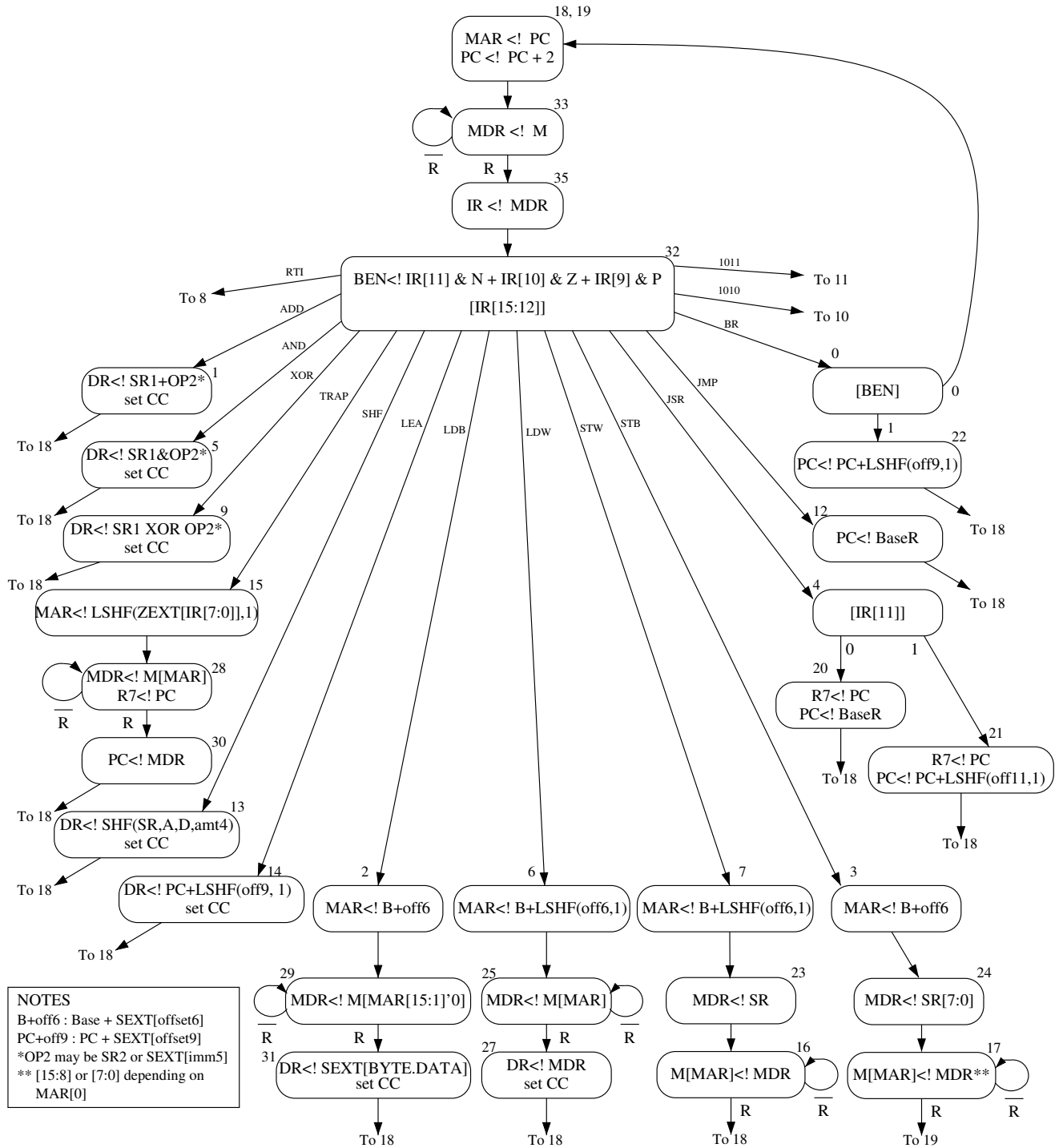
(a)



(b)



(c)



Initials: _____

Stratchpad

VLIW Instruction	ALU	MU	FPU
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			

Initials: _____

VLIW Instruction	ALU	MU	FPU
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			

Initials: _____

VLIW Instruction	ALU	MU	FPU
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			