

Final Examination
Design of Digital Circuits (252-0028-00L)
ETH Zürich, Spring 2017

Professors Onur Mutlu and Srdjan Capkun

Problem 1 (70 Points):	
Problem 2 (50 Points):	
Problem 3 (40 Points):	
Problem 4 (40 Points):	
Problem 5 (60 Points):	
Problem 6 (60 Points):	
<hr/>	
Total (320 Points):	

Examination Rules:

1. Written exam, 90 minutes in total.
2. No books, no calculators, no computers or communication devices. Six pages of hand-written 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.

This page intentionally left blank

1 Potpourri

1.1 Processor Design [20 points]

Circle the lines including terms that are compatible with each other and it makes sense for a processor design to include both.

- superscalar execution — in-order execution **2 points**
- superscalar execution — out-of-order execution **2 points**
- single-cycle machine — branch prediction **2 points**
- reservation station — microprogramming **2 points**
- fine-grained multithreading — single-core processor **2 points**
- Tomasulo's algorithm — in-order execution **2 points**
- precise exceptions — out-of-order instruction retirement **2 points**
- branch prediction — fine-grained multithreading **2 points**
- direct-mapped cache — LRU replacement policy **2 points**
- fine-grained multithreading — pipelining **2 points**

1.2 Pipelining [6 points]

What are the three major causes of pipeline stalls?

Data/Control Flow dependences (other possible answer: Data flow dependence)
2 points

Multi-cycle operations (other possible answer: Control flow dependence) **2 points**

Resource contention **2 points**

1.3 Caches I [5 points]

Please reason about the following statements about a possible processor cache one can design.

Can a cache be 5-way set associative?

YES

NO

Explain your reasoning. Be concise. Show your work.

Answer: we just need 5 tag comparators.

Explanation: Nothing wrong with a non-power-of-two associativity.

1.4 Caches II [10 points]

Assume a processor where instructions operate on 8-byte operands. An instruction is also encoded using 8 bytes. Assume that the designed processor implements a 16 kilo-byte, 4-way set associative cache that contains 1024 sets.

How effective is this cache? Explain your reasoning. Be concise. Show your work.

Answer:

1) The cache requires two accesses to be effective. (5 points)

2) The cache cannot exploit spatial locality. (5 points)

Explanation: The cache has $4 * 1024 = 4096$ cache lines in total. That means, each cache line is $16KB/4096 = 4$ bytes. With 4-byte cache lines, each operand and each instruction needs to be stored in two cache lines, which will require 2 accesses to the cache for each load/store operation and instruction fetches. The cache cannot exploit spatial locality, but only can provide benefit by exploiting temporal locality (albeit requiring two accesses).

1.5 Performance Analysis [15 points]

A multi-cycle processor executes *arithmetic instructions* in **5 cycles**, *branch instructions* in **4 cycles** and *memory instructions* in **10 cycles**. You have a program where 30% of all instructions are arithmetic instructions, 35% of *all instructions* are memory instructions, and the rest are branch instructions. You figured out that the processor cannot execute the program fast enough to meet your performance goals. Your goal is to reduce the execution time of this program by at least 10%. Hence, you decide to change the processor design to improve the performance of **arithmetic instructions**.

In the new processor design, **at most** how many cycles should the execution of **a single arithmetic instruction** take to reduce the execution time of the *entire program* by **at least** 10%? Show your work.

Answer: 2 cycles. (10 points)

Explanation: Let the total number of instructions be X .

The processor will execute the program in:

$$5 * 0.3 * X + 4 * 0.35 * X + 10 * 0.35 * X = 6.4 * X \text{ cycles.}$$

To improve the execution time by 10%, the program should complete in:

$$6.4 * X * 0.1 = 0.64 * X \text{ less cycles.}$$

The cost of executing the arithmetic instructions was $5 * 0.3 * X = 1.5 * X$ cycles. To improve the program's performance by 10%, the arithmetic instructions should complete execution at least in $1.5 * X - 0.64 * X = 0.86 * X$ cycles. Hence, $A * 0.3 * X \leq 0.86 * X$, $A \leq 2.87$, where A is the new number of cycles that the processor should execute an arithmetic instruction in. So, to improve the overall performance by 10%, an arithmetic instruction needs to execute 3 cycles faster. Hence, it should take at most $5 - 3 = 2$ cycles. (a correct explanation that proves the student's understanding may receive 13/14 points.)

1.6 Microprogrammed Design [4 points]

In lecture, we discussed a design principle for microprogrammed processors. We said that it is a good design principle to generate the control signals for cycle $N + 1$ in cycle N .

Why is this a good design principle? Be concise in your answer.

Answer: Likely keeps the critical path short (it follows the critical path design principle).

Explanation: By generating the control signals in advance, we can make the critical path of the circuit likely shorter. Shorter critical path can increase the frequency of the processor.

1.7 Processor Performance [10 points]

Assume that we test the performance of two processors, A and B, on a benchmark program. We find the following about each:

- Processor A has a CPI of 2 and executes 4 Billion Instructions per Second.
- Processor B has a CPI of 1 and executes 8 Billion Instructions per Second.

Which processor has higher performance on this program? Circle one.

Recall that CPI stands for Cycles Per Instruction.

- A. Processor A
- B. Processor B
- C. They have equal performance
- D. Not enough information to tell

Explain concisely your answer in the box provided below. Show your work.

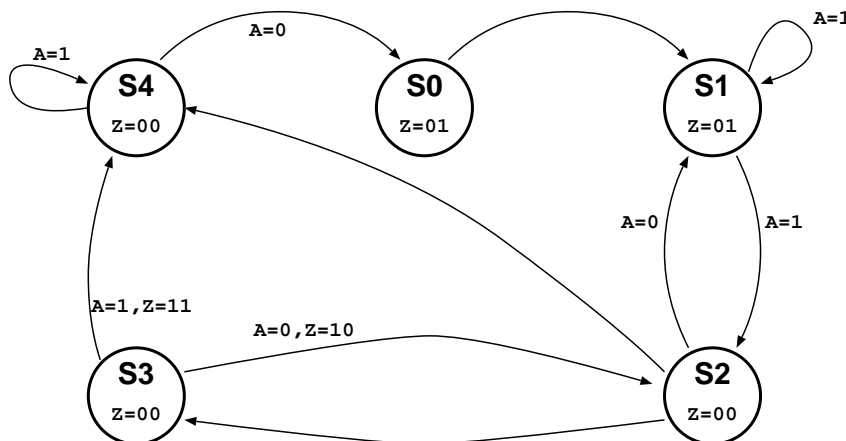
Answer: Neither of these metrics nor their combination provide execution time.

Explanation: Although information about the CPI and the instructions/second is provided, it is not enough to reason about the processors' performance. The processors may support different Instruction Set Architectures, in which case the benchmark program will be compiled into a different assembly code. The fact that one of the processors execute more instructions per second does not necessarily mean that the processor makes more progress on the benchmark.

2 Finite State Machines

This question has three parts.

- (a) [20 points] An engineer has designed a deterministic finite state machine with a one-bit input (A) and a two-bit output (Z). He started the design by drawing the following state transition diagram:



Although the exact functionality of the FSM is not known to you, there are **at least three mistakes** in this diagram. Please list **all** the mistakes.

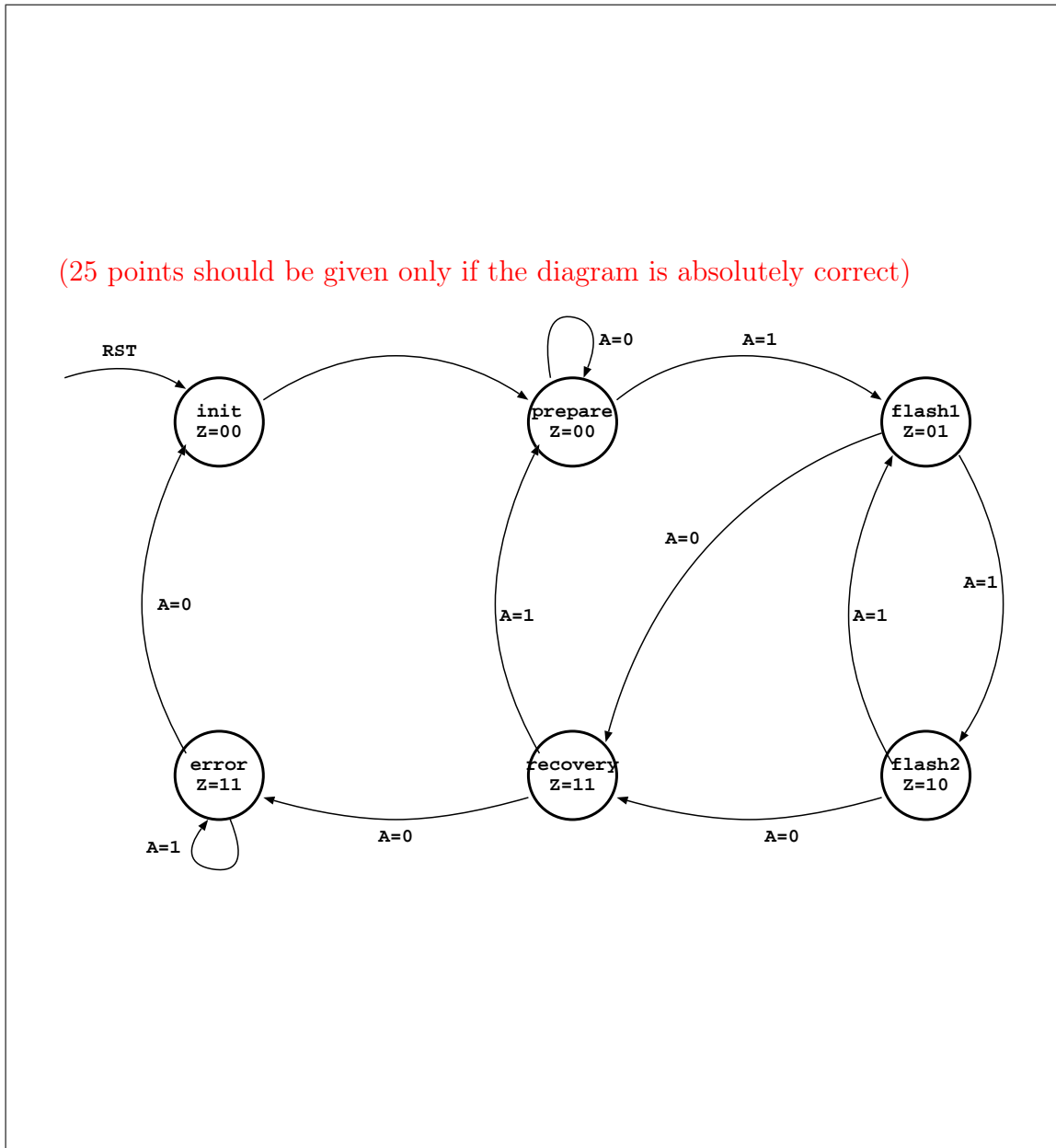
There are four problems with this diagram

- Most states have a Moore labelling (output state in the bubble), one has a Mealy type labelling (output given with input transitions) (5 points)
- There are two different transitions both with $A = 1$ from state $S1$. What will happen with $A = 0$ is missing (5 points)
- There are two different transitions from state $S2$, without labeling which input triggers them (5 points)
- There is no reset state (5 points)

- (b) [25 points] After learning from his mistakes, your colleague has proceeded to write the following Verilog code for a much better (and **different**) FSM. The code has been verified for syntax errors and found to be OK.

```
1  module fsm (input CLK, RST, A, output [1:0] Z);
2
3  reg [2:0] nextState, presentState;
4
5  parameter start = 3'b000;
6  parameter flash1 = 3'b010;
7  parameter flash2 = 3'b011;
8  parameter prepare = 3'b100;
9  parameter recovery = 3'b110;
10 parameter error = 3'b111;
11
12 always @ (posedge CLK, posedge RST)
13     if (RST) presentState <= start;
14     else presentState <= nextState;
15
16 assign Z = (presentState == recovery) ? 2'b11 :
17             (presentState == error) ? 2'b11 :
18             (presentState == flash1) ? 2'b01 :
19             (presentState == flash2) ? 2'b10 : 2'b00;
20
21 always @ (presentState, A)
22     case (presentState)
23         start : nextState <= prepare;
24         prepare : if (A) nextState <= flash1;
25         flash1 : if (A) nextState <= flash2;
26                 else nextState <= recovery;
27         flash2 : if (A) nextState <= flash1;
28                 else nextState <= recovery;
29         recovery : if (A) nextState <= prepare;
30                 else nextState <= error;
31         error : if (~A) nextState <= start;
32         default : nextState <= presentState;
33     endcase
34
35 endmodule
```

Draw a proper state transition diagram that corresponds to the FSM described in this Verilog code.



- (c) [5 points] Is the FSM described by the previous Verilog code a Moore or a Mealy FSM? Why?

Moore, the output Z only depends on the state (*presentState*) and not on the input (A).

3 Verilog

Please answer the following four questions about Verilog.

- (a) [10 points] Does the following code result in a sequential circuit or a combinational circuit? Please explain why.

```
1 module one (input clk, input a, input b, output reg [1:0] q);
2   always @ (*)
3     if (b)
4       q <= 2'b01;
5     else if (a)
6       q <= 2'b10;
7 endmodule
```

Answer and concise explanation:

This code results in a sequential circuit because a latch is required to store the old value of `q` if both conditions are **not** satisfied.

- (b) [10 points] What is the value of the output `z` if the input `c` is 10101111 and `d` is 01010101?

```
1 module two (input [7:0] c, input [7:0] d, output reg [7:0] z);
2   always @ (c,d)
3     begin
4       z = 8'b00000001;
5       z[7:5] = c[5:3];
6       z[4] = d[7];
7       z[3] = d[7];
8     end
9 endmodule
```

Please answer below. Show your work.

10100001. Last assignment of a bit overrides all previous assignments.

- (c) [10 points] Is the following code correct? If not, please explain the mistake and how to fix it.

```
1 module mux2 ( input [1:0] i, input sel, output z );
2   assign z= (sel) ? i[1]:i[0];
3 endmodule
4
5 module three ( input [3:0] data, input sel1, input sel2, output z );
6
7   wire m;
8
9   mux2 i0 (.i(data[1:0]), .sel(sel1), .z(m[0]) );
10  mux2 i1 (.i(data[3:2]), .sel(sel1), .z(m[1]) );
11  mux2 i2 (.i(m), .sel(sel2), .z(z) );
12
13 endmodule
```

Answer and concise explanation:

No. The wire m is declared to be only 1-bit wide but it needs to be 2-bit wide.

- (d) [10 points] Does the following code correctly implement a multiplexer?

```
1 module four (input sel, input [1:0] data, output reg z);
2   always@(sel)
3   begin
4     if(sel == 1'b0)
5       z = data[0];
6     else
7       z = data[1];
8   end
9 endmodule
```

Answer and concise explanation:

No. The input data is not in the sensitivity list and therefore changes to the input would not be reflected in the output z.

4 Boolean Logic and Truth Tables

You will be asked to derive the Boolean Equations for two 4-input logic functions, X and Y. Please use the Truth Table below for the following three questions.

Inputs				Outputs	
A_3	A_2	A_1	A_0	X	Y
0	0	0	0	1	0
0	0	0	1	1	0
0	0	1	0	1	0
0	0	1	1	1	0
0	1	0	0	1	0
0	1	0	1	1	1
0	1	1	0	1	0
0	1	1	1	0	0
1	0	0	0	1	0
1	0	0	1	1	0
1	0	1	0	1	1
1	0	1	1	1	0
1	1	0	0	1	0
1	1	0	1	1	0
1	1	1	0	0	0
1	1	1	1	0	0

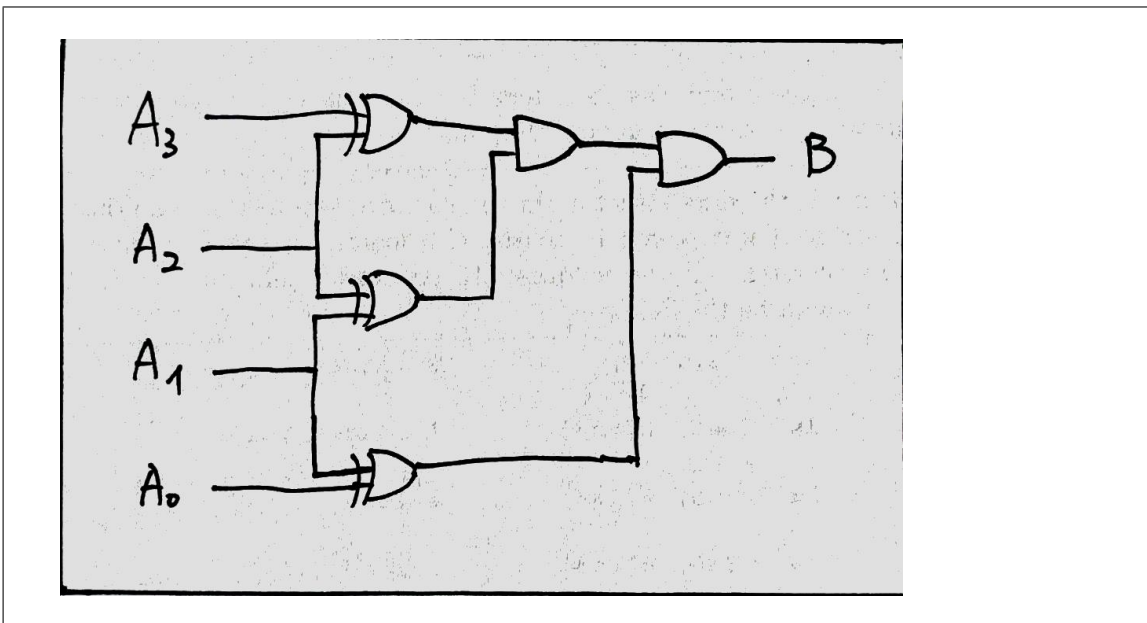
- (a) [15 points] The output X is *one* when the input does **not** contain 3 consecutive 1's in the word A_3, A_2, A_1, A_0 . The output X is *zero*, otherwise. **Fill in the truth table on the previous page and write the Boolean equation in the box below for X using the *Sum of Products* form. (No simplification needed.)**

$$X = (A_3 + \overline{A_2} + \overline{A_1} + \overline{A_0}) \cdot (\overline{A_3} + \overline{A_2} + \overline{A_1} + A_0) \cdot (\overline{A_3} + \overline{A_2} + \overline{A_1} + \overline{A_0})$$

- (b) [15 points] The output Y is *one* when no two adjacent bits in the word A_3, A_2, A_1, A_0 are the same (e.g., if A_2 is 0 then A_3 and A_1 cannot be 0). The output Y is *zero*, otherwise (for example 0000). **Fill in the truth table on the previous page and write the Boolean equation in the box below for Y using the *Sum of Products* form. (No simplification needed.)**

$$Y = \overline{A_3}A_2\overline{A_1}A_0 + A_3\overline{A_2}A_1\overline{A_0}$$

- (c) [10 points] Please represent the circuit of Y using only 2-input XOR and AND gates.



5 Tomasulo's Algorithm

In this problem, we consider an in-order fetch, out-of-order dispatch, and in-order retirement execution engine that employs Tomasulo's algorithm. This engine behaves as follows:

- The engine has four main pipeline stages: Fetch (F), Decode (D), Execute (E), and Write-back (W).
- The engine can fetch one instruction per cycle, decode one instruction per cycle, and write back the result of one instruction per cycle.
- The engine has two execution units: 1) an adder for executing ADD instructions and 2) a multiplier for executing MUL instructions.
- The execution units are fully pipelined. The adder has two stages (E1-E2) and the multiplier has four stages (E1-E2-E3-E4). Execution of each stage takes one cycle.
- The adder has a two-entry reservation station and the multiplier has a four-entry reservation station.
- An instruction always allocates the first available entry of the reservation station (in top-to-bottom order) of the corresponding execution unit.
- Full data forwarding is available, i.e., during the last cycle of the E stage, the tags and data are broadcast to the reservation station and the Register Alias Table (RAT). For example, an ADD instruction updates the reservation station entries of the dependent instructions in E2 stage. So, the updated value can be read from the reservation station entry in the next cycle. Therefore, a dependent instruction can potentially begin its execution in the next cycle (after E2).
- The multiplier and adder have separate output data buses, which allow both the adder and the multiplier to update the reservation station and the RAT in the same cycle.
- An instruction continues to occupy a reservation station slot until it finishes the Write-back (W) stage. The reservation station entry is deallocated after the Write-back (W) stage.

5.1 Problem Definition

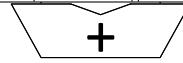
The processor is about to fetch and execute *six* instructions. Assume the *reservation stations (RS)* are all initially empty and the initial state of the *register alias table (RAT)* is given below in Figure (a). Instructions are fetched, decoded and executed as discussed in class. At some point during the execution of the six instructions, a snapshot of the state of the RS and the RAT is taken. Figures (b) and (c) show the state of the RS and the RAT at the snapshot time. A dash (–) indicates that a value has been cleared. A question mark (?) indicates that a value is unknown.

Reg	Valid	Tag	Value
R0	1	-	1900
R1	1	-	82
R2	1	-	1
R3	1	-	3
R4	1	-	10
R5	1	-	5
R6	1	-	23
R7	1	-	35
R8	1	-	61
R9	1	-	4

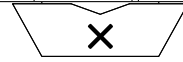
Reg	Valid	Tag	Value
R0	1	?	1900
R1	0	Z	?
R2	1	?	12
R3	1	?	3
R4	1	?	10
R5	0	B	?
R6	1	?	23
R7	0	H	?
R8	1	?	350
R9	0	A	?

(a) Initial state of the RAT (b) State of the RAT at the snapshot time

ID	V	Tag	Value
A	1	?	350
B	0	A	?



ID	V	Tag	Value
-	-	-	-
T	1	?	10
H	1	?	35
Z	1	?	82

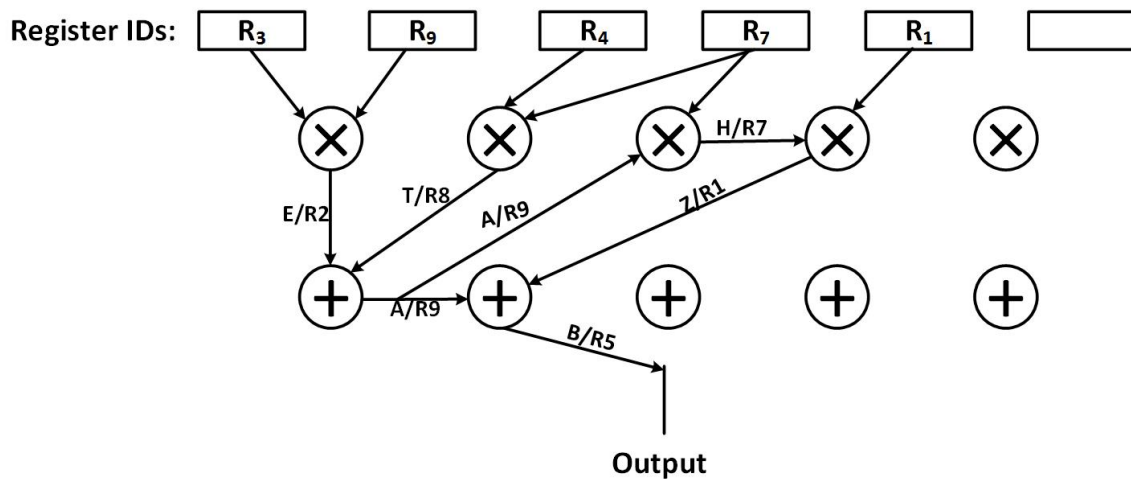


(c) State of the RS at the snapshot time

5.2 Questions

5.2.1 Data Flow Graph [40 points]

Based on the information provided above, identify the instructions and complete the dataflow graph below for the six instructions that have been fetched. Please appropriately connect the nodes using edges and specify the direction of each edge. Label each edge with the destination architectural register and the corresponding Tag. *Note that you may not need to use all registers and/or nodes provided below. (40 points if everything is correct. Deduct 2 points per mistake.)*



5.2.2 Program Instructions [20 points]

Fill in the blanks below with the six-instruction sequence in program order. When referring to registers, please use their architectural names (R0 through R9). Place the register with the smaller architectural name on the left source register box.

For example, ADD R8 \leftarrow R1, R5. (20 points if everything is correct.)

MUL	R2	\leftarrow	R3	,	R9
MUL	R8	\leftarrow	R4	,	R7
ADD	R9	\leftarrow	R2	,	R8
MUL	R7	\leftarrow	R7	,	R9
MUL	R1	\leftarrow	R1	,	R7
ADD	R5	\leftarrow	R1	,	R9

6 GPUs and SIMD

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 64 threads, and there are 64 SIMD lanes in the GPU. Please assume that all values in array B have magnitudes less than 10 (i.e., $|B[i]| < 10$, for all i).

```
for (i = 0; i < 1024; i++) {
    A[i] = B[i] * B[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 five questions.

- (a) [5 points] How many warps does it take to execute this program?

Warps = (Number of threads) / (Number of threads per warp)
Number of threads = 2^{10} (i.e., one thread per loop iteration).
Number of threads per warp = $64 = 2^6$ (given).
Warps = $2^{10}/2^6 = 2^4$

- (b) [5 points] What is the maximum possible SIMD utilization of this program?

100%

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

B:

For every 64 consecutive elements: every value is 0, every value is positive, or every value is negative. Must give all three of these.

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

Answer: 132/384

Explanation: The first two lines must be executed by every thread in a warp (64/64 utilization for each line). The minimum utilization results when a single thread from each warp passes both conditions on lines 2 and 4, and every other thread fails to meet the condition on line 2. The thread per warp that meets both conditions, executes lines 3-6 resulting in a SIMD utilization of 1/64 for each line. The minimum SIMD utilization sums to $(64 * 2 + 1 * 4) / (64 * 6) = 132/384$

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

B:

Exactly 1 of every 64 consecutive elements must be negative. The rest must be zero. This is the only case that this holds true.