

Solver: Dang Xuan Vu

Email Address: xuanvu001@e.ntu.edu.sg

1.

- a) Assumed that the program comprises of n instructions with execution time of 1 unit

Original execution time: $n \times 1 = n$ (units)

Enhanced execution time: $\frac{n}{2} \times \frac{1}{10} + \frac{n}{2} \times \frac{1}{5} = \frac{3n}{20}$ (units)

$$\rightarrow \text{Speedup} = \frac{\text{Original exe. time}}{\text{Enhanced exe. time}} = \frac{20}{3} \sim 6.67$$

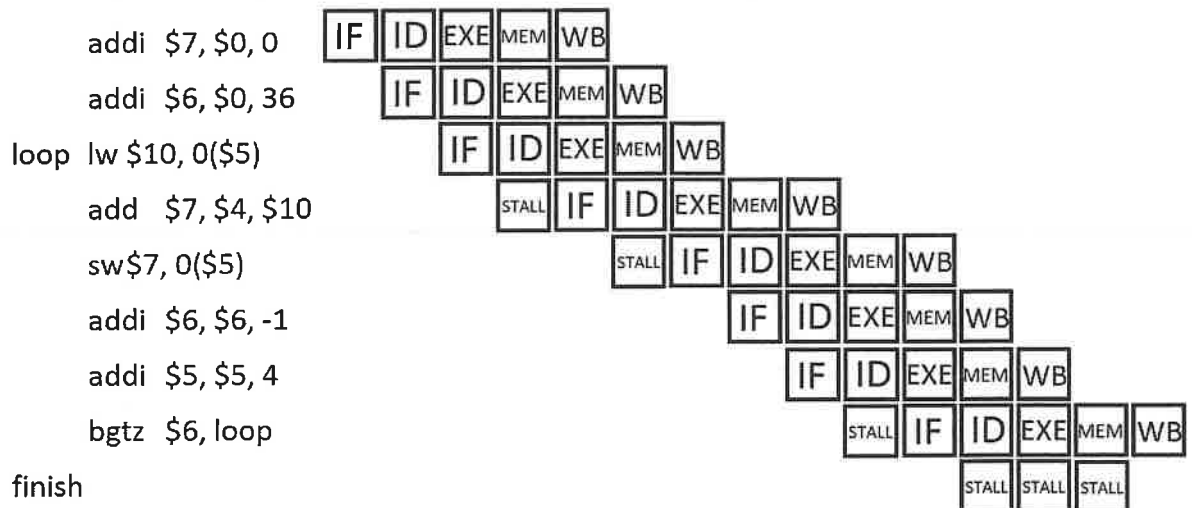
b)

- (i) Types of hazards:

- Data hazard
- Control hazard

- (ii) Assume that data is needed at ID stage and PC update for branch is after the MEM stage.

According to the diagram, there are 6 stall cycles for each loop



- (iii) Code segment with loop unrolling:

```

addi $7, $0, 0
addi $8, $0, 0
addi $6, $0, 18
loop lw $10, 0($5)
    lw $11, 1($5)
    add $7, $4, $10
    add $8, $4, $11
    addi $6, $6, -1
    sw $7, 0($5)
    bgtz $6, loop
    addi $5, $5, 4
    sw $8, 1($5)
finish
    
```

\rightarrow Total reduction of stall cycles: 5 cycles/loop

a) Structure of a JUMP instruction is:

➔ Maximum forwarding in byte is $0x1FFFFFF = 2^{25} - 1$ (bytes)

b)

$$0x7FFFFFFF7 + 0xFFFFFFFF9 = 0x7FFFFFF0$$
$$0x7FFFFFFF7 + 0x00000009 = 0x80000000$$

The diagram illustrates a processor architecture with the following components and connections:

- PC (Program Counter):** Provides an input to the **IMEM** and the **ADDER**.
- ADDER:** A 1-bit adder that takes the PC input and a constant '1' to calculate the next PC value.
- IMEM (Instruction Memory):** Receives the PC input and outputs the **inst** (instruction) to the **RF**.
- CU (Control Unit):** Receives the **opcode** from the **inst** and controls the **RF** and **ALU**.
- RF (Register File):** Receives **raddr1**, **raddr2**, and **waddr** from the **CU**. It outputs **rdata1** and **rdata2** to the **ALU**. It also receives **wdata** from the **ALU**.
- ALU (Arithmetic Logic Unit):** Takes **rdata1** and **rdata2** from the **RF** and performs operations based on the **aluop** (arithmetic logic operation) from the **CU**. It outputs the result to the **wdata** bus.

PC: Program Counter

IMEM: Instruction Memory

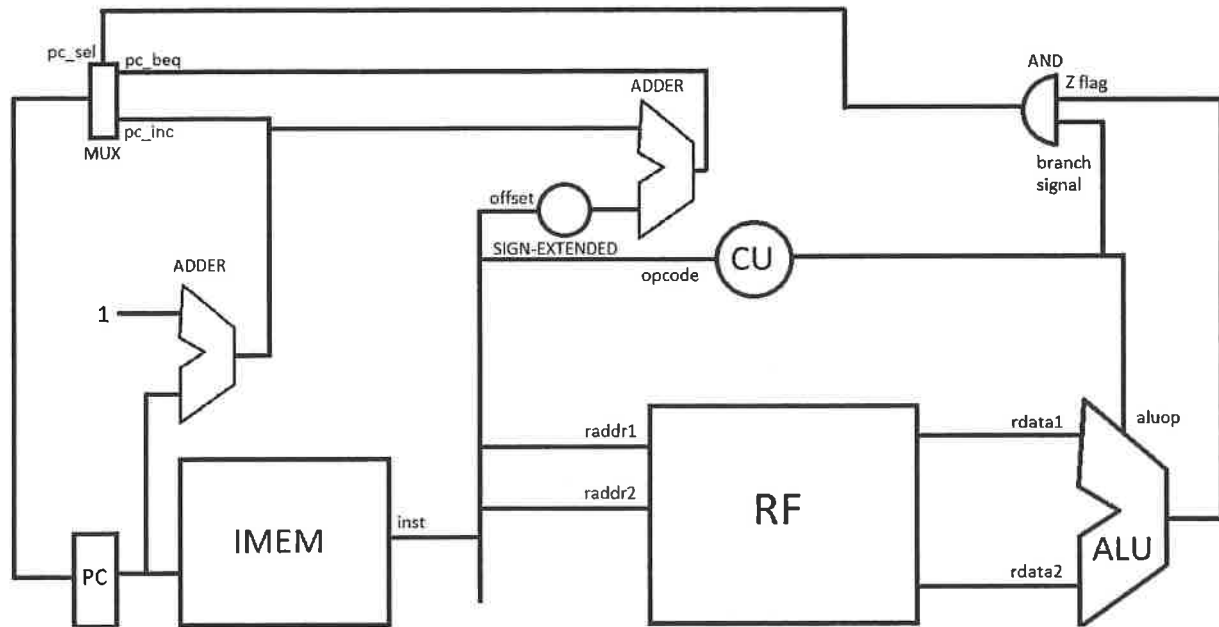
ALU: Arithmetic – Logic Unit

RF: Register File

CU: Control Unit

- 1 MUX to choose the destination register for the Register File (\$rd/\$rt)
- 1 MUX to choose the 2nd operand for the ALU (\$rt/imm)
- 1 MUX to choose the value source for Register File writeback (ALU/DM)

d) I-format BEQ instructions datapath



LEGEND:

PC: Program Counter

IMEM: Instruction Memory

ALU: Arithmetic – Logic Unit

RF: Register File

CU: Control Unit

Execution of BRANCH instructions:

- Instruction is fetched & decoded in the first two stages
- ALU compares the content of \$rs & \$rt
- The offset is left-shifted by 2 to convert byte-addressable to word-addressable
- If the result is equal (\$rs == \$rt), PC = PC (incremented) + offset (sign-extended)
- Otherwise if not equal, PC = PC (incremented)

3.

a)

(i) Prediction Accuracy = $4/12 = 33.33\%$

Branch Prediction	Branch Outcome	Hit/Miss
N	T	M
T	T	H
T	T	H
T	N	M
N	T	M
T	N	M
N	T	M
T	T	H
T	T	H
T	N	M
N	T	M
T	N	M

(ii) Methods to eliminate control hazards:

- Traditional
 - Stall the CPU for as many cycles as the branch instruction requires for PC updating, with minimum cost of 2 stall cycles
- Branch Prediction
 - Predict the outcome of the branch then compare with actual result for precision, can either be static (taken/not taken) or dynamic (multiple-bit predictor)
- Delayed Branching
 - Apply out-of-order execution to insert independent instructions subsequent to control instructions to minimize waste cycles

b) Parallelism Features:

- Instruction-level (Superscalar/VLIW processors)
 - Multiple independent instructions are identified and grouped to be executed concurrently in different functional units in a single processor
- Data-level (Vector/array processors)
 - Same operation is performed on multiple data values concurrently in multiple processing units
- Thread-level (Multi-core/processor systems)
 - More than one independent threads/tasks are executed simultaneously

c) Motivation of using GPU for high performance GPP computing:

- Data-level parallelism by massive data parallel stream processing
- Thread-level parallelism by using thousands of threads running in hundreds of cores
- High cost-effective due to GPU high performance
- Utilizes the already available hardware

Behavior of GPU compared with GPP:

- High latency & throughput vs. low latency & throughput

d) Relation between operating voltage & clock frequency:

$$f \propto V$$

when the threshold voltage (V_{th}) is negligible compared to the operating voltage (the full formula can be found in p7, Module 1.4 Power dissipation in processors).

Dynamic power consumption formula:

$$P_{dynamic} = ACV^2f \propto ACf^3$$

where A: switching factor (workload equivalent)
 C: total load capacitance
 V: operating voltage
 f: clock frequency

Compare to quad-core running at frequency f , a quad-core processor running at frequency $f/4$ has equal switching factor ($A_f = A_{f/4}$) and same chip area, or total capacitance ($C_f = C_{f/4}$). Dynamic power consumption in comparison:

$$P_{\frac{f}{4}} = A_{\frac{f}{4}} C_{\frac{f}{4}} \left(\frac{1}{4}f\right)^3 = \frac{1}{64} A_f C_f f^3 = \frac{1}{64} P_f = \frac{1}{16} P_{single}$$

Hence, for similar performance, using quad-core processors with sufficient independent threads available can greatly reduce the dynamic power consumption by 16 times compared to normal single-core processors.

4.

a)

(i) Address Format for Direct Mapped Cache:

$$[tag] :: [index] :: [offset] - [7:5] :: [4:2] :: [1:0]$$

Address Reference	Tag	Index	Offset
0x12	000	100	10
0x01	000	000	01
0x09	000	010	01
0x05	000	001	01
0x22	001	000	10
0x28	001	010	00
0x05	000	001	01
0x01	000	000	01
0x0A	000	010	10
0x21	001	000	01

Address Format for 4-way Set Associative Cache:

$$[tag] :: [index] :: [offset] - [7:3] :: [2] :: [1:0]$$

Address Reference	Tag	Index	Offset
0x12	00010	0	10
0x01	00000	0	01
0x09	00001	0	01
0x05	00000	1	01
0x22	00100	0	10
0x28	00101	0	00
0x05	00000	1	01
0x01	00000	0	01
0x0A	00001	0	10
0x21	00100	0	01

(ii) Cache Content for Direct Mapped Cache:

Index	Content
000	000
	001
	000
	001
001	000
010	000
	001
	000
011	
100	000
101	

110	
111	

→ Hit-rate = $1/10 = 10\%$

Cache Content for 4-way Set Associative Cache:

Index	Way-1	Way-2	Way-3	Way-4
0	00010 00101	00000	00001	00100
1	00000			

→ Hit-rate = $4/10 = 40\%$

(iii) Cache Content for Fully Associative Cache:

Way-1	Way-2	Way-3	Way-4	Way-5	Way-6	Way-7	Way-8
000100	000000	000010	000001	001000	001010		

→ Hit-rate = $4/10 = 40\%$

→ The Fully Associative scheme does not show any advantage compared with the 4-way Set Associative, as the memory content for latter reference stay in the cache without being cached out

b) Number of page table entries:

$$\frac{\text{Virtual memory space}}{\text{Page size}} = \frac{2^{16}}{64} = 2^{10} \text{ (entries)}$$

Single-level page table size:

$$\text{Number of entries} \times \text{Entry size} = 2^{10} \text{ (entries)} \times 4 \left(\frac{B}{\text{entry}} \right) = 2^{12} B = 4MB$$

For reporting of errors and errata, please visit pypdiscuss.appspot.com

Thank you and all the best for your exams! ☺