

CS152 Computer Architecture and
Engineering

Solution

Assigned 01/18/2022 ISAs, Microprogramming and Pipelining
Problem Set #1, Version (1.2) Due February 3

<http://inst.eecs.berkeley.edu/~cs152/sp22>

The problem sets are intended to help you learn the material, and we encourage you to collaborate with other students and to ask questions in discussion sections and office hours to understand the problems. However, each student must turn in their own solution to the problems.

The problem sets also provide essential background material for the exam and the midterms. The problem sets will be graded primarily on an effort basis, but if you do not work through the problem sets yourself you are unlikely to succeed on the exam or midterms! We will distribute solutions to the problem set on the day after the deadline to give you feedback.

Assignments must be submitted through **Gradescope** by **11:59pm PST** on the specified due date. Refer to Piazza for the entry code to join the CS152 Gradescope. Late submissions will not be accepted, except for extreme circumstances and with prior arrangement.

Name:

SID:

Problem 1: CISC, RISC, Accumulator, and Stack: Comparing ISAs

In this problem, your task is to compare four different ISAs: x86 (a CISC architecture with variable-length instructions), RISC-V (a load-store, RISC architecture with 32-bit instructions in its base form), a stack-based ISA, and an accumulator-based ISA.

Problem 1.A CISC

Let us begin by considering the following C code, which (inefficiently) rotates the bits in a 32-bit value by n times.

```
unsigned int rotate(unsigned int x, unsigned int n) {
    unsigned int msb;

    while (n != 0) {
        msb = x >> 31;
        x = (x << 1) | msb;
        n--;
    }
    return x;
}
```

Using `gcc` and `objdump` on an x86 machine, we see that the above loop compiles to the following x86 instruction sequence. On entry to this code, register `%eax` contains `x` and register `%ecx` contains `n`. Throughout parts (a–d), we will ignore what happens in the `done` label and return statement.

```
loop:    test   %ecx,%ecx (2 bytes)
        jz     done      (2 bytes)
        mov   %eax,%ebx (2 bytes)
        shr  $31,%ebx   (2 bytes)
        shl  $1,%eax   (2 bytes)
        or   %ebx,%eax (2 bytes)
        dec  %ecx      (1 byte)
        jmp  loop      (2 bytes)

done:    ...
```

The meanings and instruction lengths of the instructions used above are given in the following table. Registers are denoted with $R_{\text{SUBSCRIPT}}$, register contents with $\langle R_{\text{SUBSCRIPT}} \rangle$.

Instruction	Operation	Length
<code>mov R_{SRC}, R_{DEST}</code>	$\langle R_{\text{DEST}} \rangle = \langle R_{\text{SRC}} \rangle$	2 bytes
<code>test R_{SRC1}, R_{SRC2}</code>	<code>temp = <R_{SRC1}> & <R_{SRC2}></code> Set flags based on value of temp	2 bytes
<code>dec R_{DEST}</code>	$\langle R_{\text{DEST}} \rangle = \langle R_{\text{DEST}} \rangle - 1$	1 byte
<code>shl \$imm8, R_{DEST}</code>	$\langle R_{\text{DEST}} \rangle = \langle R_{\text{DEST}} \rangle \ll \text{imm32}$	2 bytes

<code>shr \$imm8, R_DEST</code>	$\langle R_DEST \rangle = \langle R_DEST \rangle \gg imm32$	2 bytes
<code>or R_SRC, R_DEST</code>	$\langle R_DEST \rangle = \langle R_DEST \rangle \langle R_SRC \rangle$	2 bytes
<code>jmp label</code>	jump to the address specified by <code>label</code>	2 bytes
<code>jz label</code>	if ($ZF == 1$), jump to the address specified by <code>label</code>	2 bytes

Notice that the jump instruction `js` (jump if negative) depends on `ZS`, which is a status flag. Status flags are set by the instruction preceding the jump, based on the result of the computation. Some instructions, like the `test` instruction, perform a computation and set status flags, but do not return any result. The meanings of the status flags are given in the following table:

Name	Purpose	Condition Reported
ZF	Zero	Result is zero

How many bytes is the program? For the above x86 assembly code, how many bytes of instructions need to be fetched if $x = 0x01020304$ and $n = 5$? Assuming 32-bit data values, how many bytes of data memory need to be loaded? Stored?

Bytes in program = $7 * 2 \text{ bytes} + 1 \text{ byte} = 15 \text{ bytes}$

Instruction bytes fetched = $5 * 15 \text{ bytes (loop iterations)} + 4 \text{ bytes (final test and jz)} = 79 \text{ bytes}$

Date loaded/stored = 0 bytes

Problem 1.B RISC

Translate each of the x86 instructions in the following table into one or more RISC-V instructions. Place the `loop` label where appropriate. You should use the minimum number of instructions needed to translate each x86 instruction. (You are allowed to replace *multiple* x86 instructions with a single RISC-V instructions). Assume that `x1` contains `x` upon entry, and `x2` should receive `n`. If needed, use `x4` as a condition register, and `x6`, `x7`, etc., for temporaries. You should not need to use any floating-point registers or instructions in your code. A description of the RISC-V instruction set architecture can be found in the class website, resources page.

Note: It is possible to replace the loop in 1.A with an $O(1)$ non-loop-based solution. For this problem, we want you to use the more inefficient $O(N)$ loop-based solution.

x86 instruction	Label	RISC-V instruction sequence
test %ecx,%ecx	loop:	
jz done		beq x2, x0 done
mov \$eax,%ebx		
shr \$31,%ebx		srli x6, x1, 31
shl \$1,%eax		slli x2, x2, 1
or %ebx,%eax		or x1, x6, x1
dec %ecx		addi x2, x2, -1
jmp loop		jal x0, loop
...	done:	...

How many bytes is the RISC-V program using your direct translation? How many bytes of RISC-V instructions need to be fetched for $x = 0x01020304$ and $n = 5$ with your direct translation? Assuming 32-bit data values, how many bytes of data memory need to be loaded? Stored?

The answer for this question may be slightly different, depending on how exactly students translated their programs.

Bytes in program = $6 * 4$ bytes = 24 bytes

Instruction bytes fetched = $5 * 24$ bytes (loop iterations) + 4 bytes (final beq) = 124 bytes

Bytes loaded/stored = 0 bytes

Problem 1.C**Stack**

In a stack architecture, all operations occur on top of the stack. Only push and pop access memory, and all other instructions remove their operands from the stack and replace them with the result. The hardware implementation we will assume for this problem set uses stack registers for the topmost two entries; accesses that involve deeper stack positions (e.g., pushing or popping something when the stack has more than two entries) use an extra memory reference. Assume each instruction occupies three bytes if it takes an address or label; other instructions occupy one byte.

Instruction	Definition
PUSH <i>addr</i>	load value at <i>addr</i> ; push value onto stack
POP <i>addr</i>	pop stack; store value to <i>addr</i>
OR	pop two values from the stack; OR them; push result onto stack
SHL	pop value from top of stack; shift left by 1; push result onto stack
SIGN	pop value from top of stack; shift right by 31; push result onto stack
DEC	pop value from top of stack; decrement value by 1; push result onto stack
BEQZ <i>label</i>	pop value from stack; if it's zero, branch to <i>label</i> ; else, continue with next instruction
BNEZ <i>label</i>	pop value from stack; if it's not zero, branch to <i>label</i> ; else, continue with next instruction
JUMP <i>label</i>	continue execution at location <i>label</i>

Translate the `rotate` loop to the stack ISA. You are permitted to change the sequence of instructions from (a) and (b). Assume that when we reach the loop, `n` is at the top of the stack and `x` is underneath it. At the end of the loop, the stack should contain only `x` at the top. Assume that memory starting at address `0x8000` (to fit within a 2-byte address specifier) is available to use as temporary storage.

How many bytes is your program? How many bytes of instructions need to be fetched for `x = 0x01020304` and `n = 5` with your translation? Assuming 32-bit data values, how many bytes of data memory need to be loaded? Stored? Would the number of bytes loaded and stored change if the stack could fit 8 entries in registers?

Answers may be slightly different depending on how students translated their programs.

```
pop 0x8000 # [n,x]
pop 0x8004 # [x]; n is at 0x8000; x is at 0x8004
```

```
push 0x8000 # []
```

```
loop: beqz done # [n]
```

```

push 0x8004 # []
shl # [x]
push 0x8004 # [x<<1]
sign # [x, x<<1]

or # [msb(x), x<<1]

pop 0x8004 # [x']
push 0x8000 # []
dec # [n]
jump loop # [n-1]
done: push 0x8004 # []

```

Bytes in program = 10 * 3 bytes (instructions with addresses/labels) + 4 * 1 byte (instructions *without* addresses/labels) = 34 bytes

Instruction bytes fetched = 9 bytes (prologue) + 5 * 22 bytes (loop iterations) + 3 bytes (final `beqz`) + 3 bytes (epilogue) = 125 bytes

Bytes loaded = 1 push * 4 bytes (prologue) + 5 * 3 pushes * 4 bytes (loop iterations) + 1 push * 4 bytes (epilogue) = 68 bytes

Bytes stored = 2 pops * 4 bytes (prologue) + 5 * 1 pop * 4 bytes (loop iterations) = 28 bytes

The solution above doesn't use more than two stack entries, so there would be no change if more stack entries were added. But if your solution does use more than two stack entries, then adding more registers may eliminate implicit loads.

Problem 1.D**Accumulator**

In an accumulator ISA, one operand is implicitly a specific register (the same for all instructions), called the accumulator. To make programming easier, we will consider a modified architecture that has a secondary accumulator to hold an additional value. Assume each instruction occupies three bytes if it takes an address or label; other instructions occupy one byte.

Instruction	Definition
LOAD <i>addr</i>	load value at <i>addr</i> into the primary accumulator
STORE <i>addr</i>	store the primary accumulator's value to <i>addr</i>
OR <i>addr</i>	OR the value at <i>addr</i> with the value in the primary accumulator
SHL	left-shift the value in the primary accumulator by one bit
SIGN	logical right-shift the value in the primary accumulator by 31 bits
INC	increment the primary accumulator by 1
DEC	decrement the primary accumulator by 1
SWAP	swap the values in the primary and secondary accumulators
ZERO	zero the value in the primary accumulator
BEQZ <i>label</i>	branch to <i>label</i> if the primary accumulator holds a zero value
BNEZ <i>label</i>	branch to <i>label</i> if the primary accumulator holds a non-zero value
JUMP <i>label</i>	continue execution at location <i>label</i>

Notice that all instructions operate on the primary accumulator. Also note that there are no register specifiers in this architecture; *addr* and *label* represent memory addresses. Translate the `rotate` loop to use this ISA. Assume that `x` initially held at address `0x8000`, and `n` is initially held at address `0x8004`. You are permitted to write temporary variables to any addresses above `0x8000`. You should return `x` in the **primary** accumulator.

How many bytes is your program? How many bytes of instructions need to be fetched for `x = 0x01020304` and `n = 5` with your translation? Assuming 32-bit data values, how many bytes of data memory need to be loaded? Stored?

```

LOAD 0x8000 # (p: ?, s: ?) 1
SWAP # (p: x, s: ?)
LOAD 0x8004 # (p: ?, s: x) 2

loop: BEQZ done # (p: n, s: x) 3
      DEC # (p: n, s: x)
      SWAP # (p: n', s: x)

      SHL # (p: x, s: n-1)
      STORE 0x8008 # (p: x<<1, s: n') 4

```

LOAD 0x8000 # (p: x<<1, s: n') 5
SIGN # (p: x, s: n-1)

OR 0x8008 # (p: msb(x), s: n') 6
STORE 0x8000 # (p: x', s: n') 7

SWAP # (p: x', s: n)
JUMP loop # (p: n', s: x') 8

done: SWAP # (p: n, s: x)

Bytes in program = 8 * 3 bytes (instructions with addresses/labels) + 7 * 1 byte (instructions
without addresses/labels) = 31 bytes

Instruction bytes fetched = 7 bytes (prologue) + 5 * 23 bytes (loop iterations) + 3 bytes (final
BEQZ) + 1 byte (epilogue) = 126 bytes

Data bytes loaded = 2 * 4 bytes (prologue) + 5 * 2 * 4 bytes (loop iterations) = 48 bytes

Data bytes stored = 5 * 2 * 4 (loop iterations) = 40 bytes

Problem 1.E

Conclusions

In just a few sentences, compare the four ISAs you have studied with respect to code size, number of instructions fetched, and data memory traffic. Which one would you choose if you were to build a specialized processor to execute the code in this program, and why?

- Static code size: CISC < RISC < (Stack \approx Accumulator)
- Dynamic code size: CISC < (RISC \approx Stack \approx Accumulator)
- Data memory traffic: (CISC \approx RISC) < Accumulator < Stack
 - If your code is not well-matched for a stack machine, even accumulator machines can be more efficient
- We would choose CISC if we wanted to minimize bandwidth and memory storage requirements
 - Another ISA choice is also acceptable as long as the student provides a reasonable explanation

Problem 1.F

Optimization

To get more practice with RISC-V, optimize the code from part B so that fewer dynamic instructions are executed on average and the frequency of taken branches is minimized. There are solutions more efficient than simply translating each individual x86 instruction as you did in part (b). Your solution should contain commented assembly code, a brief explanation of your optimizations, and a short analysis of the savings you obtained.

Note: It is possible to replace the loop in 1.A with an $O(1)$ non-loop-based solution. For this problem, we want you to use the more inefficient $O(N)$ loop-based solution.

Common optimizations may include:

- Loop unrolling
 - Reduces the loop overhead
- Loop inversion: translating the while loop to a do-while loop
 - Eliminates the unconditional jump

Problem 2: Microprogramming and Bus-based Architectures

In this problem, we explore microprogramming by writing microcode for the bus-based implementation of the RISC-V machine described in Handout #1 (Bus-Based RISC-V Implementation). Read the instruction fetch microcode in Table H1-3 of Handout #1. Make sure that you understand how different types of data and control transfers are achieved by setting the appropriate control signals before attempting this problem.

The final solution should be as elegant and efficient as possible with respect to the number of microinstructions used.

Problem 2.A

Implementing SUBLEQ

For this problem, you are to implement a new kind of arithmetic instruction, **MODULOM**. The new instruction has the following format:

MODULOM rd, rs1, rs2

MODULOM performs the following operation: The memory word at the address in *rs1* is divided by the memory word at the address in *rs2*, and the *remainder* is stored in address in the memory word at the address in *rd*.

$$M[rd] \leftarrow M[rs1] \% M[rs2]$$

Your CPU's ALU does **not** have support for a remainder or a division operation. Fortunately, the modulo operation can also be implemented as a loop, as illustrated below:

```
unsigned int modulo(unsigned int x, unsigned int y) {
    while (x >= y) x -= y;
    return x;
}
```

This loop *is* realizable with the microcode of your CPU.

Fill in Worksheet 2.A with the microcode for MODULOM. Use *don't cares* (*) for fields where it is safe to use don't cares. Study the hardware description well, and make sure all your microinstructions are legal.

Please comment your code clearly. If the pseudo-code for a line does not fit in the space provided, or if you have additional comments, you may write in the margins so long as you do it neatly. Your code should exhibit “clean” behavior and not modify *rd*, *rs1*, *rs2*, or other general-purpose architectural registers while executing the instruction.

Finally, make sure that the instruction fetches the next instruction (i.e., by doing a microbranch to FETCH0 as discussed above) once the result has been saved to $M[rd]$.

You may want to consult the microcode found in the micro-coded processor provided in Lab1, which can be viewed at `lab1/generators/riscv-sodor/src/main/scala/rv32_ucose/microcode.scala` for guidance. Warning: While that microcode passes all provided assembly tests and benchmarks, no guarantees to the optimality of that code are assured, and there may still be bugs in the provided implementation.

We will accept any reasonable solution, even if it is different from the one on the next page.

State	PseudoCode	ldIR	Reg Sel	Reg Wr	en Reg	ldA	ldB	ALUOp	en ALU	ld MA	Mem Wr	en Mem	Imm Sel	en Imm	μBr	Next State
FETCH0:	MA <- PC; A <- PC	*	PC	0	1	1	*	*	0	1	0	0	*	0	N	*
	IR <- Mem	1	*	0	0	0	*	*	0	0	0	1	*	0	S	*
	PC <- A+4	0	PC	1	0	0	*	INC_A_4	1	*	0	0	*	0	D	*
...																
NOP0:	microbranch back to FETCH0	*	*	0	0	*	*	*	0	*	0	0	*	0	J	FETCH0
MODULOM 0:	MA <- R[rs1]	0	rs1	0	1	*	*	*	0	1	0	0	*	0	N	*
	A <- Mem	0	*	0	0	1	*	*	0	0	0	1	*	0	S	
	MA <- R[rs2]	0	rs2	0	1	0	*	*	0	1	0	0	*	0	N	*
	B <- Mem	0	*	0	0	0	1	*	0	0	0	1	*	0	S	
LOOP	if (A < B) goto DONE	0	rd	0	1	0	0	SLTU	0	1	0	0	*	0	NZ	DONE
	MA <- R[rd]															
	A <- A – B goto LOOP	*	*	0	0	1	0	SUB	1	0	0	0	*	0	J	LOOP
DONE	Mem <- A	*	*	0	0	0	*	COPY_A	1	0	1	0	*	0	S	
	goto FETCH0	*	*	0	0	*	*	*	0	*	0	0	*	0	J	FETCH0

In this question we ask you to implement a useful vector instruction to find the *largest number in a vector of unsigned integers*. This instruction has the same format as other arithmetic (R-type) instructions in RISC-V:

MAXV rd, rs1, rs2

The MAXV instruction takes a pointer to the beginning of a vector in memory (*rs1*) and a pointer to the end of a vector in memory (*rs2*), and it returns in register *rd* the largest number in that register. Your code is permitted to modify register *rs1* during the execution of this instruction.

For this problem, each vector element will be a 32-bit unsigned number. You can assume that the address in *rs2* is larger than the address in *rs1*.

Your task is to fill out Worksheet 2.B for the MAXV instruction. You should try to optimize your implementation for the minimal number of cycles necessary and for which signals can be set to don't-cares.

We will accept any reasonable solution, even if it is different from the one on the next page.

We will accept solutions for which *rs2* points to the last value in the vector, as well as solutions for which *rs2* points *after* the last value in the vector

State	PseudoCode	ldIR	Reg Sel	Reg Wr	en Reg	ldA	ldB	ALUOp	en ALU	ld MA	Mem Wr	en Mem	Imm Sel	en Imm	μBr	Next State
FETCH0:	MA <- PC; A <- PC	*	PC	0	1	1	*	*	0	1	0	0	*	0	N	*
	IR <- Mem	1	*	0	0	0	*	*	0	0	0	1	*	0	S	*
	PC <- A+4	0	PC	1	0	0	*	INC_A_4	1	*	0	0	*	0	D	*
...																
NOP0:	microbranch back to FETCH0	*	*	0	0	*	*	*	0	*	0	0	*	0	J	FETCH0
MAXV0:	A <- B	0	*	0	0	1	0	COPY_B	1	*	0	0	*	0	N	
	R[rd] <- A – B (== 0)	0	rd	1	0	*	*	SUB	1	*	0	0	*	0	N	
LOOP	A, MA <- R[rs1]	0	rs1	0	1	1	*	*	0	1	0	0	*	0	N	
	B <- R[rs2]	0	rs2	0	1	0	1	*	0	0	0	0	*	0	N	
	If NOT (A < B) goto FETCH0	0	rd	0	1	0	1	SLTU	0	0	0	0	*	0	EZ	FETCH0
	B <- R[rd]															
	R[rs1] <- A + 4	0	rs1	1	0	*	0	INC_A_4	1	0	0	0	*	0	N	
	A <- MEM	0	*	0	0	1	0	*	0	0	0	1	*	0	S	
	if (A < B) goto LOOP	0	*	0	0	0	0	SLTU	0	*	0	0	*	0	NZ	LOOP
	R[rd] <- B goto LOOP	0	rd	1	0	*	*	*	0	*	0	0	*	0	J	LOOP

Problem 2.C

Instruction Execution Times

How many cycles does it take to execute the following instructions on the microcoded RISC-V implementation? Use the states and control signals from Handout #1 (or Lab 1, in `lab1/generators/riscv-sodor/src/main/scala/rv32_ucose/microcode.scala`) and assume that memory does not assert its busy signal.

Instruction	Cycles
<code>SUB x3, x2, x1</code>	$3 + 3 = 6$
<code>ANDI x2, x1, #4</code>	$3 + 3 = 6$
<code>LW x1, 0(x2)</code>	$3 + 5 = 8$
<code>BNE x1, x2, label # (x1 == x2)</code>	$3 + 4 = 7$
<code>BNE x1, x2, label # (x1 != x2)</code>	$3 + 3 + 4 = 10$
<code>BEQ x1, x2, label # (x1 != x2)</code>	$3 + 4 = 7$
<code>BEQ x1, x2, label # (x1 == x2)</code>	$3 + 3 + 4 = 10$
<code>J label</code>	$3 + 6 = 9$
<code>JAL label</code>	$3 + 6 = 9$
<code>JALR x1</code>	$3 + 6 = 9$
<code>AUIPC x1, #128</code>	$3 + 4 = 7$

The answers below are derived from the microcoded processor from the handout.

Instruction	Cycles	Summary (not including fetch and dispatch)
<code>SUB x3, x2, x1</code>	$3 + 3 = 6$	1) $A \leftarrow R[x2]$; 2) $B \leftarrow R[x1]$; 3) $R[x3] \leftarrow A - B$
<code>ANDI x2, x1, #4</code>	$3 + 3 = 6$	1) $A \leftarrow R[x1]$; 2) $B \leftarrow Imm$; 3) $R[x2] \leftarrow A \&^1 B$
<code>LW x1, 0(x2)</code>	$3 + 4 = 7$ or $3 + 5 = 8$	1) $A \leftarrow R[x2]$; 2) $B \leftarrow Imm$; 3) $MA \leftarrow A + B$; 4) $R[x1] \leftarrow Mem$; 5) $\mu Br J FETCH0^2$
<code>BNE x1, x2, label # (x1 == x2)</code>	$3 + 3 = 6$	1) $A \leftarrow R[x1]$; 2) $B \leftarrow R[x2]$; 3) $A - B$; $\mu Br EZ FETCH0$; $B \leftarrow Imm^4$
<code>BNE x1, x2, label # (x1 != x2)</code>	$3 + 3 + 4 = 9$	4) $A \leftarrow PC$; 5) $A \leftarrow A - 4$; 6) $PC \leftarrow A + B$
<code>BEQ x1, x2, label # (x1 != x2)</code>	$3 + 3 = 6$	1) $A \leftarrow R[x1]$; 2) $B \leftarrow R[x2]$; 3) $A - B$; $\mu Br NZ FETCH0$; $B \leftarrow Imm^4$
<code>BEQ x1, x2, label # (x1 == x2)</code>	$3 + 3 + 4 = 9$	4) $A \leftarrow PC$; 5) $A \leftarrow A - 4$; 6) $PC \leftarrow A + B$
<code>J label</code>	$3 + 3 = 6$	1) $R[rd]^5 \leftarrow PC$; 2) $B \leftarrow Imm$; 3) $PC \leftarrow A^3 + B$
<code>JAL label</code>	$3 + 3 = 6$	Same as above
<code>JALR x1</code>	$3 + 4 = 7$	1) $R[rd] \leftarrow PC$; 2) $A \leftarrow R[x1]$; 3) $B \leftarrow Imm$; 4) $PC \leftarrow A + B$
<code>AUIPC x1, #128</code>	$3 + 2 = 5$	1) $B \leftarrow Imm$; 2) $R[x1] \leftarrow A^3 +^1 B$

⁰ Terminal microinstructions are assumed to have a $\mu Br J$ back to `FETCH0` unless stated otherwise.

¹ These operations were not provided in the handout, but it is reasonable to assume they can occur in a single ALU op.

² The wording of the question is ambiguous as to whether “memory does not assert its busy signal” implies that the machine has guaranteed single-cycle memory versus you can exclude stall cycles for purposes of the cycle accounting. The former would permit the μBr to occur in parallel with the load. The safer, intended answer assumes the latter case and separates the μBr and the memory op.

³ The A register contains PC after fetch, whereas PC is speculatively set to PC+4. Thus, we reuse A to speed up AUIPC and JAL instructions, eliding the need to load the PC into A and decrement it by 4. (Conversely, this cannot be avoided in taken conditional branches.)

⁴ Speculatively load the branch offset into B to shave a cycle off a taken conditional branch. This will just be discarded in fetch if the branch is not taken.

⁵ J is encoded as JAL with $rd=x0$.

Which instruction takes the most cycles to execute? Which instruction takes the fewest cycles to execute?

Most cycles: Taken branch (BEQ, BNE)

Fewest cycles: Arithmetic operations (SUB, ANDI, etc.)

Problem 3: 6-Stage Pipeline

In this problem, we consider a modification to the fully bypassed 5-stage RISC-V processor pipeline presented in Lecture 3. Our new processor has a data cache with a two-cycle latency. To accommodate this cache, the memory stage is pipelined into two stages, M1 and M2, as shown in Figure 1-A. Additional bypasses are added to keep the pipeline fully bypassed.

Suppose we are implementing this 6-stage pipeline in a technology in which register file ports are inexpensive but bypasses are costly. We wish to reduce cost by removing some of the bypass paths, but without increasing CPI. The proposal is for all integer arithmetic instructions to write their results to the register file at the end of the Execute stage, rather than waiting until the Writeback stage. A second register file write port is added for this purpose. Remember that register file writes occur on each rising clock edge, and values can be read in the next clock cycle. The proposed change is shown in Figure 1-B.

In this problem, assume that the only exceptions that can occur in this pipeline are illegal opcodes (detected in the Decode stage) and invalid memory address (detected at the start of the M2 stage). Additionally, assume that the control logic is optimized to stall only when necessary. **You may ignore branch and jump instructions in this problem.**

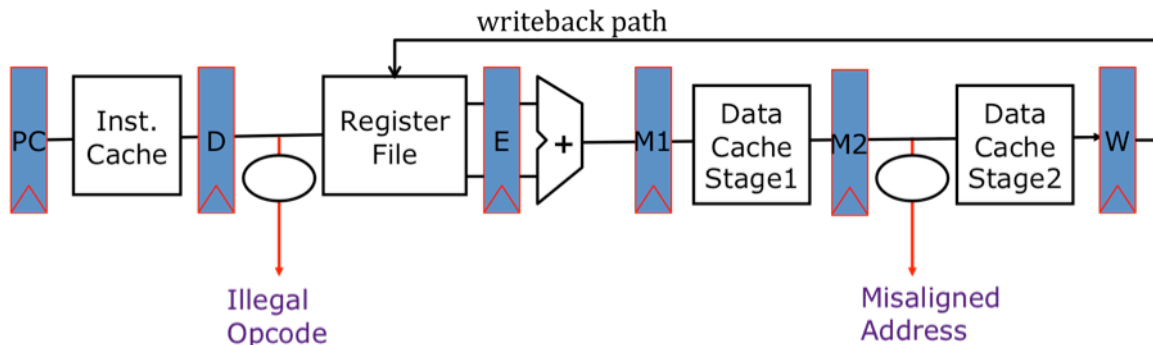


Figure 1-A. 6-stage pipeline. For clarity, bypass paths are not shown.

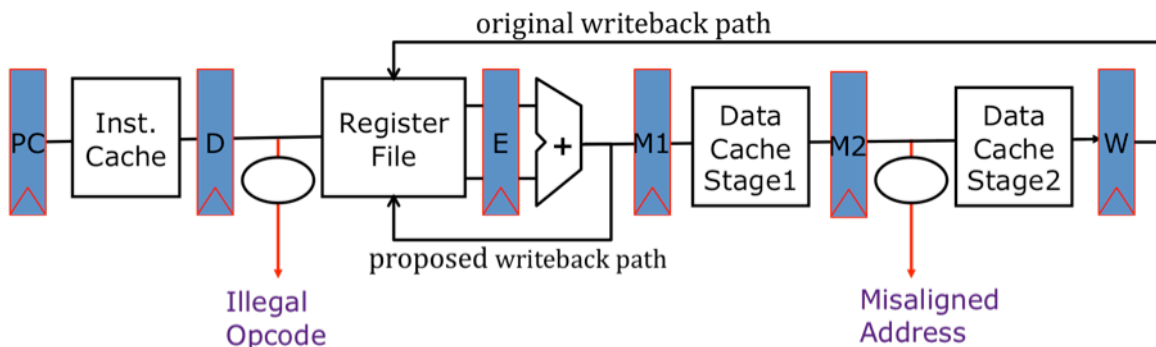


Figure 1-B. 6-stage pipeline with proposed additional write port.

Problem 3.A

Hazards: Second Write Port

The second write port allows some bypass paths to be removed without adding stalls in the decode stage. Explain how the second write port improves performance by eliminating such stalls and give a short code sequence that would have required an interlock to execute correctly with only a single write port and with the same bypass paths removed.

The second write port improves performance by resolving some RAW hazards earlier than they would be if ALU operations had to wait until writeback to provide their results to subsequent dependent instructions. It would help with the following instruction sequence:

```
add x1, x2, x3
add x4, x5, x6
add x7, x1, x9
```

The important insight is that the second write port cannot resolve data hazards for immediately back-to-back instructions. An arithmetic instruction in the EX stage writes back as it leaves the EX stage; therefore, the bypass path is necessary if the next instruction has a RAW dependency and is allowed to leave the ID stage.

Problem 3.B

Hazards: Bypasses Removed and New Hazards

After the second write port is added, which bypass paths can be removed in this new pipeline without introducing additional stalls? List each removed bypass individually. Are any new hazards added to the pipeline due to the earlier writeback of arithmetic instructions?

The bypass path from the end of M1 to the end of ID can be removed. (Credit was also given for the bypass path from the beginning of M2 to the beginning of EX, since these are equivalent.)

Additionally, ALU results no longer must be bypassed from the end of M2 or the end of WB, but these bypass paths are still used to forward load results to earlier stages.

There are multiple potential WAW hazards that must be appropriately addressed by the control logic. The two instructions writing at the same time must be appropriately prioritized. Also, if an arithmetic instruction is in M1 and a load with the same destination register is in M2, the write of the earlier load can clobber the result of the older instruction, leading to an incorrect architectural state. The control logic needs to be modified to handle these situations by suppressing the writes of older instructions when they conflict with the writes of newer instructions.

Problem 3.C

Precise Exceptions

Without further modifications, this pipeline may not support precise exceptions. Briefly explain why and provide a minimal code sequence that will result in an imprecise exception.

Illegal address exceptions are not detected until the start of the M2 stage. Since writebacks can occur at the end of the EX stage, it is possible for an arithmetic instruction following a memory access to an illegal address to have written its value back before the exception is detected, resulting in an imprecise exception. For example:

```
lw x1, -1(x0) # address -1 is misaligned
add x2, x3, x4 # x2 will be overwritten, but last instruction
has faulted!
```

Problem 3.D

Precise Exceptions: Implemented using a Interlock

Describe how precise exceptions can be implemented by adding a new interlock. Provide a minimal code sequence that would engage this interlock. Qualitatively, what is the performance impact of this solution?

Stall any ALU op in the ID stage if the instruction in the EX stage is a load or a store. The instruction sequence above engages this interlock. Loads and stores account for about 1/3 of dynamic instructions. Assuming that the instruction following a load or store is an arithmetic instruction 2/3 of the time, and ignoring the existing load-use delay, this solution will increase the CPI by $(1/3)*(2/3) = 2/9$. However, only a qualitative explanation was necessary for credit.

Problem 3.E

Precise Exceptions: Implemented using an Extra Read Port

Suppose you are additionally given the budget to add a new register file *read* port. Propose an alternative solution to implement precise exceptions in this pipeline without requiring any new interlocks.

In addition to writing an arithmetic instruction's destination register in the EX stage, also read its previous value and carry it down the pipeline. If an early writeback occurs before a preceding exception was detected, then the old value of *rd* is preserved in the M1 pipeline register and can be restored to the register file, maintaining precise state.

Note: It is better to read the previous value **as late as possible**, otherwise this read of *rd* might need an extra bypass path for the following instruction sequence:

```
ld x1, 0(x8)
ld x2, -1(x8) # misaligned
addi x1, x1, 4
```

This also depends on the interlocks used to resolve the WAW hazard mentioned in 3.B.

Problem 4: CISC vs RISC

For each of the following questions, select either *CISC* or *RISC*, depending on which ISA you feel would be best suited for the situation described. Also, briefly *explain your reasoning*.

Problem 4.A

Lack of Good Compilers I

Assume that compiler technology is poor, and therefore your users are far more apt to write all of their code in assembly. A _____ ISA would be best appreciated by these programmers.

CISC

RISC

CISC ISAs provided more complex, higher-level instructions such as string manipulation instructions and special addressing modes convenient for indexing tables (say for your company's payroll application). Two example CISC instructions: "DBcc: Test Condition, Decrement, and Branch" and "CMP2: Compare Register against Upper and Lower Bounds". This made life easy if you stared at assembly all day and could not hide behind convenient software abstractions/subroutines!

OR

CISC

RISC

A streamlined RISC ISA is far simpler for an assembly programmer to fully understand and reason about than all the idiosyncrasies that CISC ISAs tend to have, such as the variety of complex instructions for narrow use cases and the myriad addressing modes.

Problem 4.B

Lack of Good Compilers II

You desire to make compilers better at targeting your *yet-to-be-designed* machine. Therefore, you choose a _____ ISA, as it would be easiest for a compiler to target, thus allowing your users to write code in higher-level languages like C and Fortran and raise their productivity.

CISC

RISC

Compilers had difficulty targeting CISC ISAs in part because the complicated instructions have many difficult and hard to analyze side-effects. A load-store/register-register RISC ISA which limits side-effects to a single register or memory location per instruction is relatively easy for a compiler to understand, analyze, and schedule code for.

Problem 4.C

Fast Logic, Slow Memory

Assume that CPU logic is fast, *very* fast, while instruction fetch accesses are at least 10x slower (suppose you are the lead architect of the “709”). Which ISA style do you choose as a best match for the hardware’s limitations?

CISC

RISC

When instruction fetch takes 10x longer than a CPU logic operation, you are going to want to push as much compute as you can into each instruction! Certain especially complex CISC instructions can encode tens, even hundreds of equivalent RISC instructions. For example, a CISC instruction which performs a single expensive, multi-cycle string routine in hardware would be considerably faster than even an optimized RISC implementation that would need a loop with a series of loads, stores, and arithmetic instructions in the loop body.

Problem 4.D

Higher Performance(?)

Starting with a clean slate in the year 2021 (area/logic/memory is cheap), you think that a _____ ISA that would lend itself best to a very high performance processor (e.g., high frequency, highly pipelined).

CISC

RISC

Because RISC instructions tend to have simple, easy to analyze side-effects, they lend themselves more readily to pipelined micro-architectures which dynamically check for dependencies between instructions and interlock or bypass when dependencies arise. And because little work needs to be performed in each stage, the pipeline can be clocked at very high frequencies.

This advantage is evident in modern micro-architectures of old CISC ISAs: The frontend of the processor typically has a decoder which translates CISC instructions (e.g., x86 instructions) into RISC “micro-ops”, which a high-performance pipeline can then dynamically schedule for maximum performance.

For these CISC architectures such as x86 and IBM S/360, they are still around for legacy reasons. But if you had a chance at a clean slate, you would probably prefer a clean RISC implementation with a direct translation to the micro-architecture instead of using area and power on a CISC decoder front-end (not to mention the additional complexity forced on your memory system to handle the odd CISC addressing modes).

Problem 5: Iron Law of Processor Performance

Mark whether the following modifications will cause each of the *first three* categories to **increase**, **decrease**, or whether the modification will have **no effect**. Explain your reasoning.

For the final column “Overall Performance”, mark whether the following modifications **increase**, **decrease**, have **no effect**, or whether the modification will have an **ambiguous** effect. Explain your reasoning. If the modification has an **ambiguous** effect, describe the tradeoff in which it would be a beneficial modification or in which it would a detrimental modification (i.e., as an engineer would you suggest using the modification or not and why?).

		Instructions / Program	Cycles / Instruction	Seconds / Cycle	Overall Performance
a)	Adding a branch delay slot	Increase: NOPs must be inserted when the branch delay slot cannot be usefully filled	Decrease: Some control hazards are eliminated; also additional NOPs execute quickly because they have no data hazards.	No effect: will not meaningfully change the pipeline. ALSO ACCEPT: Decrease because no branch kill	Ambiguous: Depends on the program and how often the delay slot can be filled with useful work
b)	Adding a complex instruction	Decrease: if the added instruction can replace a sequence of instructions.	Increase: implementing the instruction can mean adding stages or making stages have more complex control logic.	Increase: more control logic and interlocks will often increase the critical path. ALSO ACCEPT: No effect	Ambiguous: if the program can take advantage of the new instruction, it can be worth the cost. This is a hard decision for an ISA designer to make!
c)	Reduce number of registers in the ISA	Increase: Values will more frequently be spilled to the stack, increasing number of loads and stores	Increase: more loads followed by dependent instructions will cause more stalls. Memory latency is hard to schedule around.	Decrease: fewer registers means shorter register file access time	Ambiguous: if the program uses few registers and thus spills rarely to memory, the faster reg. access times may win out. Also, your instructions may be able to be shorter, improving amongst other things code density a

d)	Improving memory access speed	No effect: since instructions make no assumption about memory speed.	Decrease: programs will spend less time stalled waiting for memory	Decrease: if memory access is on the critical path or memory was 1 cycle. ALSO ACCEPT: No effect: if memory is pipelined and just takes less cycles.	Improve: improving memory access time will increase performance of the whole system.
e)	Adding 16-bit versions of the most common instructions in RISC-V (normally 32 bits in length) to the ISA (i.e., make RISC-V a variable-length ISA)	No effect: The actual number of instructions is unchanged.	Decrease: since code size has shrunk, there will be fewer instruction cache (I\$) misses and less time spent waiting to fetch	Increase: decode becomes more complex with more formats, and instruction fetch has to deal with misalignment.	Ambiguous: the main advantage is smaller code size, which can improve I\$ hit rates and save on fetch energy (get more instructions per fetch). However, the more complex decode can offset these gains.
f)	For a given CISC ISA, changing the implementation of the micro-architecture from a microcoded engine to a RISC pipeline (with a CISC-to-RISC decoder on the frontend)	No effect: Since the ISA is not changing, the binary does not change, and thus there is no change to Inst/Program.	Decrease: Microcoded machines take several clock cycles to execute an instruction, while the RISC pipeline should have a CPI near 1 (thanks to pipelining).	No effect: the amount of work done in one pipeline stage and one microcode cycle are about the same. ALSO ACCEPT: Increase: the RISC pipeline introduces longer control paths and adds bypasses, which are likely to be on the critical path.	The decrease in CPI from pipeline far outweighs any critical path overhead of hardwired control logic.