

4

Solutions

Solution 4.1

4.1.1 The values of the signals are as follows:

	RegWrite	MemRead	ALUMux	MemWrite	ALUop	RegMux	Branch
a.	1	0	0 (Reg)	0	AND	1 (ALU)	0
b.	0	0	1 (Imm)	1	ADD	X	0

ALUMux is the control signal that controls the Mux at the ALU input, 0 (Reg) selects the output of the register file and 1 (Imm) selects the immediate from the instruction word as the second input to the ALU.

RegMux is the control signal that controls the Mux at the data input to the register file, 0 (ALU) selects the output of the ALU, and 1 (Mem) selects the output of memory.

A value of X is a “don’t care” (does not matter if signal is 0 or 1).

4.1.2 Resources performing a useful function for this instruction are:

a.	All except Data Memory and branch Add unit
b.	All except branch Add unit and write port of the Registers

4.1.3

	Outputs That Are Not Used	No Outputs
a.	Branch Add	Data Memory
b.	Branch Add, write port of Registers	None (all units produce outputs)

4.1.4 One long path for AND instruction is to read the instruction, read the registers, go through the ALU Mux, perform the ALU operation, and go through the Mux that controls the write data for Registers (I-Mem, Regs, Mux, ALU, and Mux). The other long path is similar, but goes through Control while registers are read (I-Mem, Control, Mux, ALU, Mux). There are other paths but they are shorter, such as the PC increment path (only Add and then Mux), the path to prevent branching (I-Mem to Control to Mux, so the Mux can use the Branch signal to select the PC + 4 input as the new value for PC), and the path that prevents a memory write (only I-Mem and then Control, etc.).

a.	Control is faster than registers, so the critical path is I-Mem, Regs, Mux, ALU, Mux.
b.	The two long paths are equal, so both are critical.

4.1.5 One long path is to read the instruction, read registers, use the Mux to select the immediate as the second ALU input, use ALU (compute address), access D-Mem, and use the Mux to select that as register data input, so we have I-Mem, Regs, Mux, ALU, D-Mem, Mux. The other long path is similar, but goes through Control instead of Regs (to generate the control signal for the ALU MUX). Other paths are shorter, and are similar to shorter paths described for 4.1.4.

a.	Control is faster than registers, so the critical path is I-Mem, Regs, Mux, ALU, Mux.
b.	The two long paths are equal, so both are critical.

4.1.6 This instruction has two kinds of long paths, those that determine the branch condition and those that compute the new PC. To determine the branch condition, we read the instruction, read registers or use the Control unit, then use the ALU Mux and then the ALU to compare the two values, then use the Zero output of the ALU to control the Mux that selects the new PC. As in 4.4.4 and 4.1.5:

a.	The first path (through Regs) is longer.
b.	The two long paths are equal, so both are critical.

To compute the PC, one path is to increment it by 4 (Add), add the offset (Add), and select that value as the new PC (Mux). The other path for computing the PC is to read the instruction (to get the offset) and use the branch Add unit and Mux. Both of the compute-PC paths are shorter than the critical path that determines the branch condition, because I-Mem is slower than the PC + 4 Add unit, and because ALU is slower than the branch Add.

Solution 4.2

4.2.1 Existing blocks that can be used for this instruction are:

a.	This instruction uses instruction memory, both existing read ports of Registers, the ALU (to compare Rs and Rt), and the write port of Registers.
b.	This instruction uses instruction memory, both register read ports, the ALU to add Rd and Rs together, data memory, and the write port in Registers.

4.2.2 New functional blocks needed for this instruction are:

a.	This instruction needs the Zero output of the ALU to be zero-extended to compute the value for Rd. Then we need to add this as another input to the Mux that selects the value to be written into Registers.
b.	None. This instruction can be implemented using existing blocks.

4.2.3 The new control signals are:

a.	We need a new control signal for the Mux that selects between values that can be written into Registers.
b.	None. This instruction can be implemented without adding new control signals. It only requires changes in the Control logic.

4.2.4 Clock cycle time is determined by the critical path, which for the given latencies happens to be to get the data value for the load instruction: I-Mem (read instruction), Regs (takes longer than Control), Mux (select ALU input), ALU, Data Memory, and Mux (select value from memory to be written into Registers). The latency of this path is $400\text{ps} + 200\text{ps} + 30\text{ps} + 120\text{ps} + 350\text{ps} + 30\text{ps} = 1130\text{ps}$.

New Clock Cycle Time	
a.	1430ps (1130ps + 300ps, ALU is on the critical path)
b.	1130ps. Control latency is now equal to Regs latency, so we have a second critical path (same as the existing one, but going through Control to generate the control signal for the Mux that selects second ALU input). This new critical path has the same latency as the existing one, so the clock cycle is unchanged.

4.2.5 The speedup comes from changes in clock cycle time and changes to the number of clock cycles we need for the program:

Benefit	
a.	We need 5% fewer cycles for a program, but cycle time is 1430 instead of 1130, so we have a speedup of $(1/0.95) \times (1130/1430) = 0.83$, which means we actually have a slowdown.
b.	Speedup is 1 (no change in number of cycles, no change in clock cycle time).

4.2.6 The cost is always the total cost of all components (not just those on the critical path, so the original processor has a cost of I-Mem, Regs, Control, ALU, D-Mem, 2 Add units, and 3 Mux units, for a total cost of $1000 + 200 + 500 + 100 + 2000 + 2 \times 30 + 3 \times 10 = 3890$.

We will compute cost relative to this baseline. The performance relative to this baseline is the speedup we computed in 4.2.5, and our cost/performance relative to the baseline is as follows:

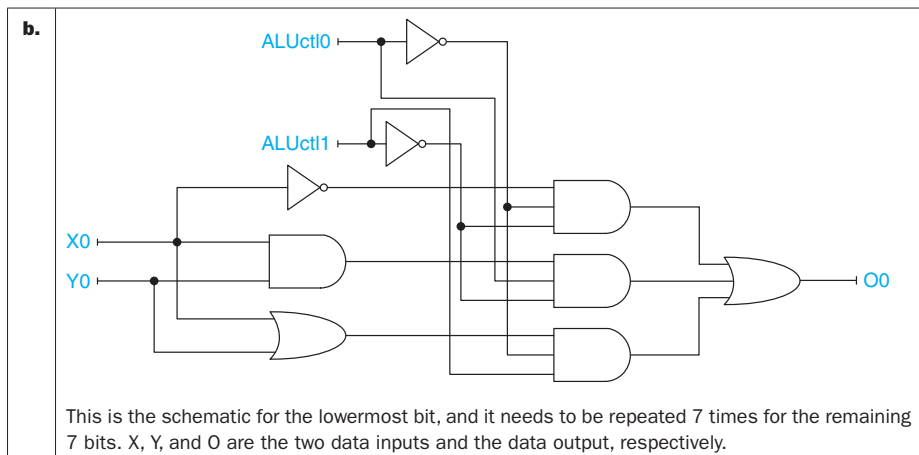
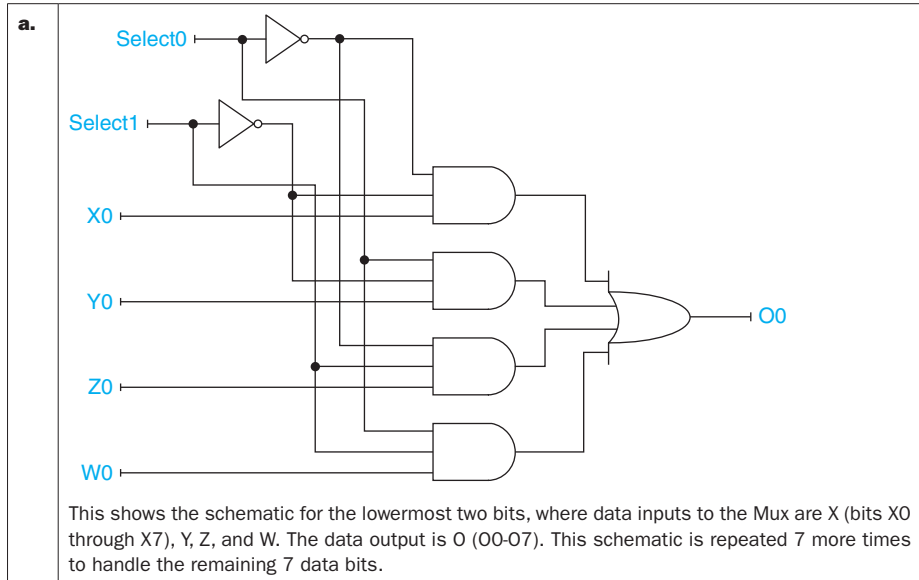
	New Cost	Relative Cost	Cost/Performance
a.	$3890 + 600 = 4490$	$4490/3890 = 1.15$	$1.15/0.83 = 1.39$. We are paying significantly more for significantly worse performance, so the cost/performance is a lot worse than with the unmodified processor.
b.	$3890 - 400 = 3490$	$3490/3890 = 0.9$	$0.9/1 = 0.9$. We are reducing cost and getting the same performance, so the cost/performance improves.

Solution 4.3

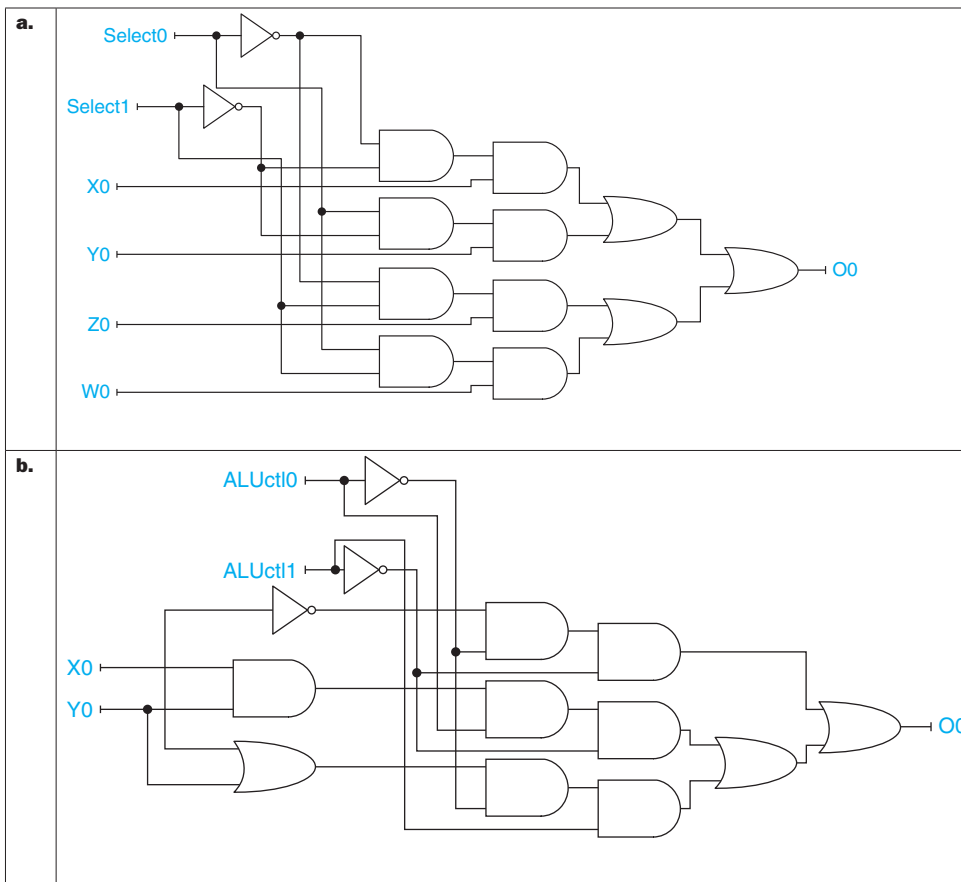
4.3.1

a.	Logic only.
b.	Logic only.

4.3.2



4.3.3



4.3.4 The latency of a path is the latency from an input (or a D-element output) to an output (or D-element input). The latency of the circuit is the latency of the path with the longest latency. Note that there are many correct ways to design the circuit in 4.3.2, and for each solution to 4.3.2 there is a different solution for this problem.

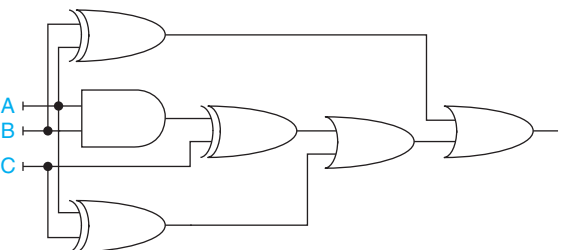
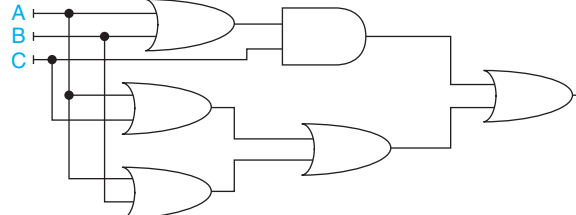
4.3.5 The cost of the implementation is simply the total cost of all its components. Note that there are many correct ways to design the circuit in 4.3.2, and for each solution to 4.3.2 there is a different solution for this problem.

4.3.6

a.	A three-input or a four-input gate has a lower latency than a cascade of two 2-input gates. This means that shorter overall latency is achieved by using 3- and 4-input gates (as in 4.3.2) rather than cascades of 2-input gates. The schematic shown for 4.3.2 turns out to already be optimal.
b.	Because multi-input AND and OR gates have the same latency as 2-input ones, we can use many-input gates to reduce the number of gates on the path from inputs to outputs. We can use De-Morgan's laws to convert sequences of gates into a circuit that only has NOT gates feeding into AND gates which feed into OR gates.

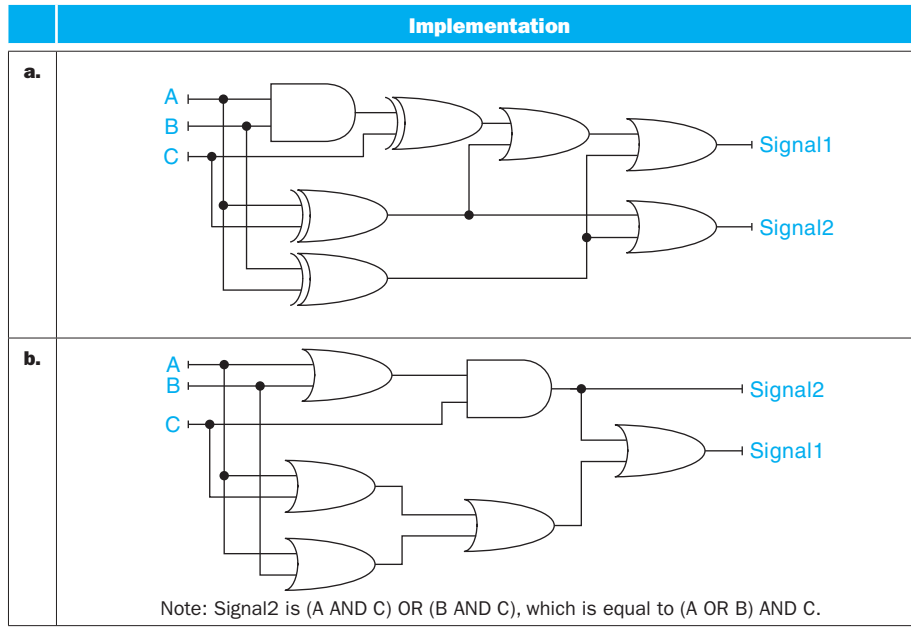
Solution 4.4

4.4.1 We show the implementation and also determine the latency (in gates) needed for 4.4.2.

	Implementation	Latency in Gates
a.		4
b.		4

4.4.2 See answer for 4.4.1 above.

4.4.3



4.4.4

a.	The critical path consists of AND, XOR, OR, and OR, for a total of 82ps.
b.	The critical path consists of three OR gates, for a total of 150ps.

4.4.5

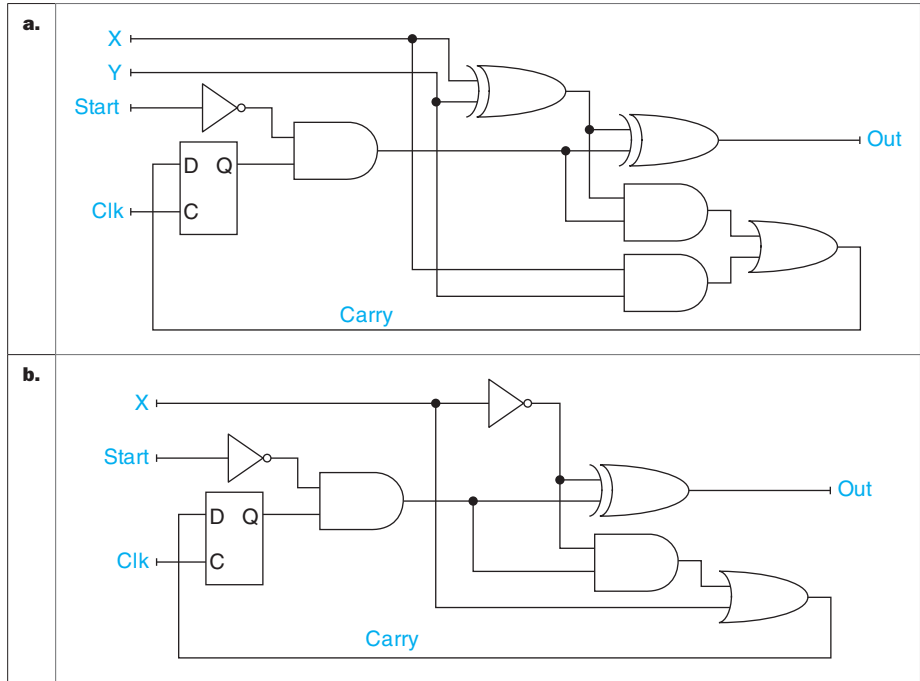
a.	The cost is 1 AND gate, 3 OR gates, and 3 XOR gates, for a total cost of 49.
b.	The cost is 1 AND gate and 5 OR gates, for a total cost of 18.

4.4.6 We already computed the cost of the combined circuit. Now we determine the cost of the separate circuits and the savings.

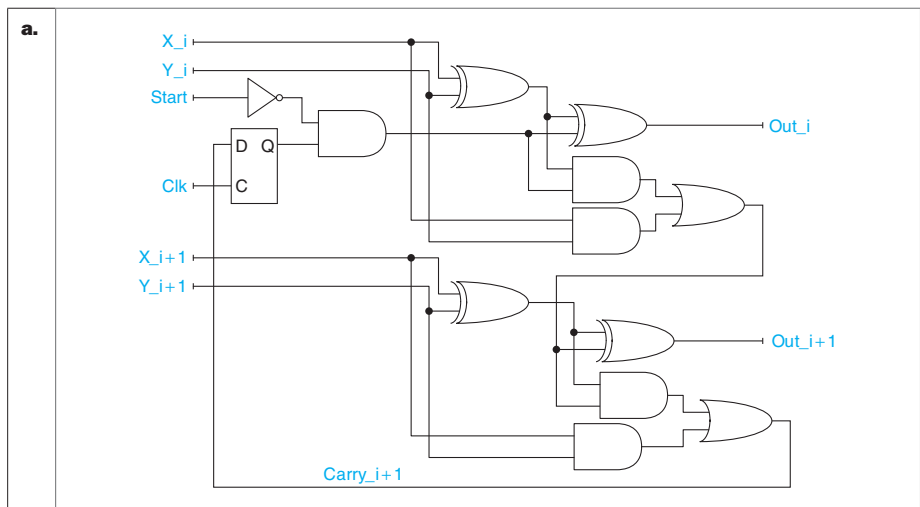
	Combined Cost	Separate Cost	Saved
a.	49	49 (no change)	0%
b.	18	27 (+ 2 AND and 1 OR gate)	$(27 - 18)/27 = 33\%$

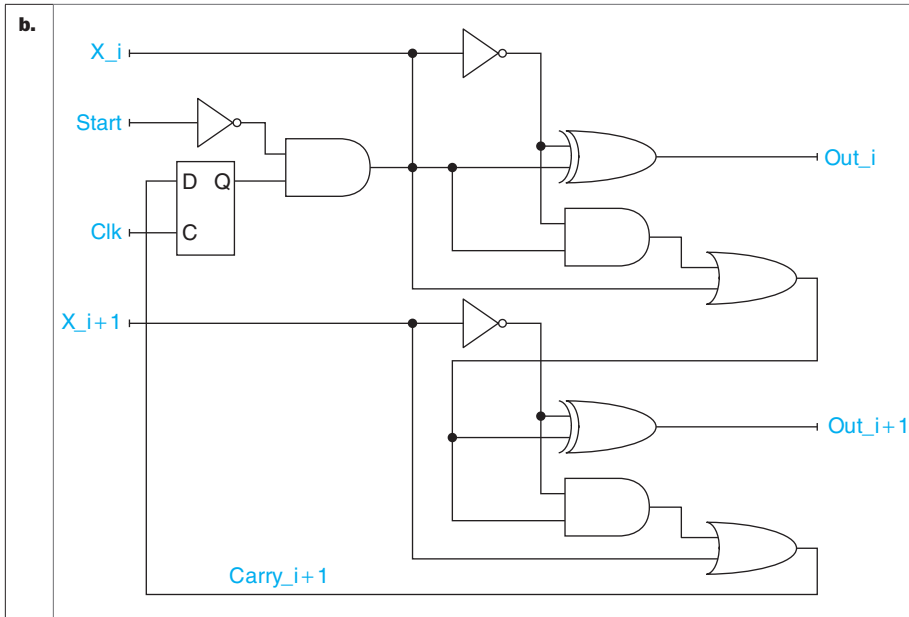
Solution 4.5

4.5.1



4.5.2





4.5.3

	Cycle Time	Operation Time
a.	76ps (NOT->AND->AND->OR->D)	$32 \times 76ps = 2432ps$
b.	500ps (NOT->AND->AND->OR->D)	$32 \times 500ps = 16000ps$

4.5.4

	Cycle Time	Speedup
a.	100ps (NOT->AND->AND->OR->AND->OR->D)	$(32 \times 76ps)/(16 \times 100ps) = 1.52$
b.	690ps (NOT->AND->AND->OR->AND->OR->D)	$(32 \times 500ps)/(16 \times 690ps) = 1.45$

4.5.5

	Circuit 1	Circuit 2
a.	40 (1 NOT, 3 AND, 1 OR, 2 XOR, 1 D)	64 (1 NOT, 5 AND, 2 OR, 4 XOR, 1 D)
b.	13 (2 NOT, 2 AND, 1 OR, 1 XOR, 1 D)	21 (3 NOT, 3 AND, 2 OR, 2 XOR, 1 D)

4.5.6

	Cost/Performance for Circuit 1	Cost/Performance for Circuit 2	Circuit 1 vs. Circuit 2
a.	$40 \times 32 \times 76 = 97280$	$64 \times 16 \times 100 = 102400$	Cost/performance of Circuit 2 is worse by about 5.3%
b.	$13 \times 32 \times 500 = 208000$	$21 \times 16 \times 690 = 231840$	Cost/performance of Circuit 2 is worse by about 11.5%

Solution 4.6

4.6.1 I-Mem takes longer than the Add unit, so the clock cycle time is equal to the latency of the I-Mem:

a.	200ps
b.	750ps

4.6.2 The critical path for this instruction is through the instruction memory, Sign-extend and Shift-left-2 to get the offset, Add unit to compute the new PC, and Mux to select that value instead of PC + 4. Note that the path through the other Add unit is shorter, because the latency of I-Mem is longer than the latency of the Add unit. We have:

a.	$200\text{ps} + 15\text{ps} + 10\text{ps} + 70\text{ps} + 20\text{ps} = 315\text{ps}$
b.	$750\text{ps} + 100\text{ps} + 0\text{ps} + 200\text{ps} + 50\text{ps} = 1100\text{ps}$

4.6.3 Conditional branches have the same long-latency path that computes the branch address as unconditional branches do. Additionally, they have a long-latency path that goes through Registers, Mux, and ALU to compute the PCSrc condition. The critical path is the longer of the two, and the path through PCSrc is longer for these latencies:

a.	$200\text{ps} + 90\text{ps} + 20\text{ps} + 90\text{ps} + 20\text{ps} = 420\text{ps}$
b.	$750\text{ps} + 300\text{ps} + 50\text{ps} + 250\text{ps} + 50\text{ps} = 1400\text{ps}$

4.6.4

a.	PC-relative branches.
b.	All instructions except unconditional jumps without a register operand (jal, j).

4.6.5

a.	PC-relative unconditional branch instructions. We saw in 4.6.3 that this is not on the critical path of conditional branches, and it is only needed for PC-relative branches. Note that MIPS does not have actual unconditional branches (BNE zero, zero, Label plays that role so there is no need for unconditional branch opcodes) so for MIPS the answer to this question is actually "None."
b.	All instructions except unconditional jumps without a register operand (jal, j).

4.6.6 Of the two instruction (BNE and ADD), BNE has a longer critical path so it determines the clock cycle time. Note that every path for ADD is shorter than or equal to the corresponding path for BNE, so changes in unit latency will not affect this. As a result, we focus on how the unit's latency affects the critical path of BNE:

a.	This unit is not on the critical path, so the only way for this unit to become critical is to increase its latency until the path for address computation through sign extend, shift left, and branch add becomes longer than the path for PCSrc through Registers, Mux, and ALU. The latency of Regs, Mux, and ALU is 200ps and the latency of Sign-extend, Shift-left-2, and Add is 95ps, so the latency of Shift-left-2 must be increased by 105ps or more for it to affect clock cycle time.
b.	This unit is already on the critical path of BNE, so changes in its latency affect the clock cycle time directly. Even if we speed this unit up to have zero latency, the path through Regs, Mux, and ALU will take 300ps and remain a critical path (because Sign-extend, Shift-left-2, and Add also take 300ps).

Solution 4.7

4.7.1 The longest-latency path for ALU operations is through I-Mem, Regs, Mux (to select ALU operand), ALU, and Mux (to select value for register write). Note that the only other path of interest is the PC-increment path through Add (PC + 4) and Mux, which is much shorter. So for the I-Mem, Regs, Mux, ALU, Mux path we have:

a.	$200\text{ps} + 90\text{ps} + 20\text{ps} + 90\text{ps} + 20\text{ps} = 420\text{ps}$
b.	$750\text{ps} + 300\text{ps} + 50\text{ps} + 250\text{ps} + 50\text{ps} = 1400\text{ps}$

4.7.2 The longest-latency path for LW is through I-Mem, Regs, Mux (to select ALU input), ALU, D-Mem, and Mux (to select what is written to register). The only other interesting paths are the PC-increment path (which is much shorter) and the path through Sign-extend unit in address computation instead of through Registers. However, Regs has a longer latency than Sign-extend, so for I-Mem, Regs, Mux, ALU, D-Mem, and Mux path we have:

a.	$200\text{ps} + 90\text{ps} + 20\text{ps} + 90\text{ps} + 250\text{ps} + 20\text{ps} = 670\text{ps}$
b.	$750\text{ps} + 300\text{ps} + 50\text{ps} + 250\text{ps} + 500\text{ps} + 50\text{ps} = 1900\text{ps}$

4.7.3 The answer is the same as in 4.7.2 because the LW instruction has the longest critical path. The longest path for SW is shorter by one Mux latency (no write to register), and the longest path for ADD or BNE is shorter by one D-Mem latency.

4.7.4 The data memory is used by LW and SW instructions, so the answer is:

a.	$25\% + 10\% = 35\%$
b.	$30\% + 20\% = 50\%$

4.7.5 The sign-extend circuit is actually computing a result in every cycle, but its output is ignored for ADD and NOT instructions. The input of the sign-extend circuit is needed for ADDI (to provide the immediate ALU operand), BEQ (to provide the PC-relative offset), and LW and SW (to provide the offset used in addressing memory) so the answer is:

a.	$20\% + 25\% + 25\% + 10\% = 80\%$
b.	$10\% + 10\% + 30\% + 20\% = 70\%$

4.7.6 The clock cycle time is determined by the critical path for the instruction that has the longest critical path. This is the LW instruction, and its critical path goes through I-Mem, Regs, Mux, ALU, D-Mem, and Mux so we have:

a.	D-Mem has the longest latency, so we reduce its latency from 250ps to 225ps, making the clock cycle 25ps shorter. The speedup achieved by reducing the clock cycle time is then $670\text{ps}/645\text{ps} = 1.039$.
b.	I-Mem has the longest latency, so we reduce its latency from 750ps to 675ps, making the clock cycle 75ps shorter. The speedup achieved by reducing the clock cycle time is then $1900\text{ps}/1825\text{ps} = 1.041$.

Solution 4.8

4.8.1 To test for a stuck-at-0 fault on a wire, we need an instruction that puts that wire to a value of 1 and has a different result if the value on the wire is stuck at zero:

a.	If this signal is stuck at zero, an instruction that writes to an odd-numbered register will end up writing to the even-numbered register. So if we place a value of zero in R30 and a value of 1 in R31, and then execute ADD R31, R30, R30 the value of R31 is supposed to be zero. If bit 0 of the Write Register input to the Registers unit is stuck at zero, the value is written to R30 instead and R31 will be 1.
b.	The MIPS architecture requires instructions to be word-aligned (lowest two bits of the instruction address are always zero). Because of this, we cannot execute an instruction that would set the specified signal to 1, so we cannot test for this stuck-at-0 fault.

4.8.2 The test for stuck-at-zero requires an instruction that sets the signal to 1 and the test for stuck-at-1 requires an instruction that sets the signal to 0. Because the signal cannot be both 0 and 1 in the same cycle, we cannot test the same signal simultaneously for stuck-at-0 and stuck-at-1 using only one instruction. The test for stuck-at-1 is analogous to the stuck-at-0 test:

a.	We can place a value of zero in R31 and a value of 1 in R30, then use ADD R30, R31, R31 which is supposed to place 0 in R30. If this signal is stuck-at-1, the write goes to R31 instead, so the value in R30 remains 1.
b.	If this signal is stuck-at-1, a branch instruction, such as BNE zero, zero, Label will result in a non-aligned PC (lowermost bit will be 1).

4.8.3

a.	We need to rewrite the program to use only odd-numbered registers.
b.	With this fault, every conditional branch results in a fetch of a misaligned instruction. This prevents any conditional changes in control flow, so the faulty processor is unusable.

4.8.4

a.	To set the MemRead signal to 1 (in order to test for stuck-at-0 fault), we need a load instruction. If MemRead is stuck-at-0, the memory does not get read and the value placed in the register is “random” (whatever happened to be at the output of the memory unit). Unfortunately, this “random” value can be the same as the one already in the register, so this test is not conclusive.
b.	To test for this fault, we need an instruction whose MemRead is 1, so it has to be a load. The instruction also needs to have RegDst set to 0, which is the case for loads. Finally, the instruction needs to have a different result if MemRead is set to 0. For a load, setting MemRead to zero would result in not reading memory at all, so the value placed in the register is “random” (whatever happened to be at the output of the memory unit). Unfortunately, this “random” value can be the same as the one already in the register, so this test is not conclusive.

4.8.5

a.	<p>If Jump is stuck at 0, the PC after a jump is not the jump address. Instead, the PC is either incremented (PC + 4) or computed as if this was a PC-relative branch. To test for this fault, we can place a jump instruction at a low address that jumps to a high address. If the Jump signal is stuck at 0, the PC after the jump will be much lower than it should be.</p> <p>To set the MemRead signal to 1 (in order to test for stuck-at-0 fault), we need a load instruction. If MemRead is stuck-at-0, the memory does not get read and the value placed in the register is “random” (whatever happened to be at the output of the memory unit). Unfortunately, this “random” value can be the same as the one already in the register, so this test is not conclusive.</p>
b.	To test for this fault, we need an instruction whose Jump is 1, so it has to be the jump instruction. However, for the jump instruction the RegDst signal is “don’t care” because it does not write to any registers, so the implementation may or may not allow us to set RegDst to 0 so we can test for this fault. As a result, we cannot reliably test for this fault.

4.8.6 Each single-instruction test “covers” all faults that, if present, result in different behavior for the test instruction. To test for as many of these faults as possible in a single instruction, we need an instruction that sets as many of these signals to a value that would be changed by a fault. Some signals cannot be tested using this single-instruction method, because the fault on a signal could still result in completely correct execution of all instructions that trigger the fault.

Solution 4.9

4.9.1

	Binary	Hexadecimal
a.	101011 10000 00100 0000000001100100	AA040064
b.	000000 00010 00011 00001 00000 101010	0043082A

4.9.2

	Read Register 1	Actually Read?	Read Register 2	Actually Read?
a.	16 (10000 _b)	Yes	4 (00100 _b)	Yes
b.	2 (00010 _b)	Yes	3 (00011 _b)	Yes

4.9.3

	Read Register 1	Register Actually Written?
a.	Either 4 (00100 _b) or 0 (don't know because RegDst is X)	No
b.	1 (00001 _b)	Yes

4.9.4

	Control Signal 1	Control Signal 2
a.	ALUSrc = 1	Branch = 0
b.	Jump = 0	RegDst = 1

4.9.5 We use I31 through I26 to denote individual bits of Instruction[31:26], which is the input to the Control unit:

a.	ALUSrc = I31
b.	Jump = (NOT I31) AND I27

4.9.6 If possible, we try to reuse some or all of the logic needed for one signal to help us compute the other signal at a lower cost:

a.	ALUSrc = I31 Branch = I28
b.	RegDst = NOT I31 Jump = RegDst AND I27

Solution 4.10

To solve the problems in this exercise, it helps to first determine the latencies of different paths inside the processor. Assuming zero latency for the Control unit, the critical path is the path to get the data for a load instruction, so we have I-Mem, Mux, Regs, Mux, ALU, D-Mem, and Mux on this path.

4.10.1 The Control unit can begin generating MemWrite only after I-Mem is read. It must finish generating this signal before the end of the clock cycle. Note that MemWrite is actually a write-enable signal for D-Mem flip-flops, and the actual write is triggered by the edge of the clock signal, so MemWrite need not arrive before that time. So the Control unit must generate the MemWrite in one clock cycle, minus the I-Mem access time:

	Critical Path	Maximum Time to Generate MemWrite
a.	$200\text{ps} + 20\text{ps} + 90\text{ps} + 20\text{ps} + 90\text{ps} + 250\text{ps} + 20\text{ps} = 690\text{ps}$	$690\text{ps} - 200\text{ps} = 490\text{ps}$
b.	$750\text{ps} + 50\text{ps} + 300\text{ps} + 50\text{ps} + 250\text{ps} + 500\text{ps} + 50\text{ps} = 1950\text{ps}$	$1950\text{ps} - 750\text{ps} = 1200\text{ps}$

4.10.2 All control signals start to be generated after I-Mem read is complete. The most slack a signal can have is until the end of the cycle, and MemWrite and RegWrite are both needed only at the end of the cycle, so they have the most slack. The time to generate both signals without increasing the critical path is the one computed in 4.10.1.

4.10.3 MemWrite and RegWrite are only needed by the end of the cycle. RegDst, Jump, and MemtoReg are needed one Mux latency before the end of the cycle, so they are more critical than MemWrite and RegWrite. Branch is needed two Mux latencies before the end of the cycle, so it is more critical than these. MemRead is needed one D-Mem plus one Mux latency before the end of the cycle, and D-Mem has more latency than a Mux, so MemRead is more critical than Branch. ALUOp must get to ALU control in time to allow one ALU Ctrl, one ALU, one D-Mem, and one Mux latency before the end of the cycle. This is clearly more critical than MemRead. Finally, ALUSrc must get to the pre-ALU Mux in time, one Mux, one ALU, one D-Mem, and one Mux latency before the end of the cycle. Again, this is more critical than MemRead. Between ALUOp and ALUSrc, ALUOp is more critical than ALUSrc if ALU control has more latency than a Mux. If ALUOp is the most critical, it must be generated one ALU Ctrl latency before the critical-path signals can go through Mux, Regs, and Mux. If the ALUSrc signal is the most critical, it must be generated while the critical path goes through Mux and Regs. We have:

	The Most Critical Control Signal Is	Time to Generate It without Affecting the Clock Cycle Time
a.	ALUOp (30ps > 20ps)	$20\text{ps} + 90\text{ps} + 20\text{ps} - 30\text{ps} = 100\text{ps}$
b.	ALUOp (70ps > 50ps)	$50\text{ps} + 300\text{ps} + 50\text{ps} - 70\text{ps} = 330\text{ps}$

For the next three problems, it helps to compute for each signal how much time we have to generate it before it starts affecting the critical path. We already did this for RegDst and RegWrite in 4.10.1, and in 4.10.3 we described how to do it for the remaining control signals. We have:

	RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
a.	470ps	470ps	450ps	220ps	470ps	100ps	490ps	110ps	490ps
b.	1150ps	1150ps	1100ps	650ps	1150ps	330ps	1200ps	350ps	1200ps

The difference between the allowed time and the actual time to generate the signal is called “slack.” For this problem, the allowed time will be the maximum time the signal can take without affecting clock cycle time. If slack is positive, the signal arrives before it is actually needed and it does not affect clock cycle time. If the slack is positive, the signal is late and the clock cycle time must be adjusted. We now compute the slack for each signal:

	RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
a.	-30ps	-30ps	0ps	20ps	20ps	-100ps	-10ps	10ps	-10ps
b.	50ps	150ps	0ps	-150ps	-50ps	30ps	-100ps	-50ps	0ps

4.10.4 With this in mind, the clock cycle time is what we computed in 4.10.1, plus the absolute value of the most negative slack. We have:

	Control Signal with the Most Negative Slack Is	Clock Cycle Time with Ideal Control Unit (from 4.10.1)	Actual Clock Cycle Time with These Signal Latencies
a.	ALUOp (-100ps)	690ps	790ps
b.	MemRead (-150ps)	1950ps	2100ps

4.10.5 It only makes sense to pay to speed up signals with negative slack, because improvements to signals with positive slack cost us without improving performance. Furthermore, for each signal with negative slack, we need to speed it up only until we eliminate all its negative slack, so we have:

	Signals with Negative Slack	Per-Processor Cost to Eliminate All Negative Slack
a.	RegWrite (-10ps) RegDst and Jump (-30ps) ALUOp (-100ps)	170ps at \$1/5ps = \$34
b.	MemtoReg and ALUSrc (-50ps) MemWrite (-100ps) MemRead (-150ps)	350ps at \$1/5ps = \$70

4.10.6 The signal with the most negative slack determines the new clock cycle time. The new clock cycle time increases the slack of all signals until there is no remaining negative slack. To minimize cost, we can then slow down signals that end up having some (positive) slack. Overall, the cost is minimized by slowing signals down by:

	RegDst	Jump	Branch	MemRead	MemtoReg	ALUOp	MemWrite	ALUSrc	RegWrite
a.	70ps	70ps	100ps	120ps	120ps	0ps	90ps	110ps	90ps
b.	200ps	300ps	150ps	0ps	100ps	180ps	50ps	100ps	150ps

Solution 4.11

4.11.1

	Sign-Extend	Jump's Shift-Left-2
a.	000000000000000000000000000010100	00011000100000000000001010000
b.	000000000000000000000000100000101010	0010000010000010000010101000

4.11.2

	ALUOp[1-0]	Instruction[5-0]
a.	00	010100
b.	10	101010

4.11.3

	New PC	Path
a.	PC + 4	PC to Add (PC + 4) to branch Mux to jump Mux to PC
b.	PC + 4	PC to Add (PC + 4) to branch Mux to jump Mux to PC

4.11.4

	WrReg Mux	ALU Mux	Mem/ALU Mux	Branch Mux	Jump Mux
a.	2 or 0 (RegDst is X)	20	X	PC + 4	PC + 4
b.	1	-128	0	PC + 4	PC + 4

4.11.5

	ALU	Add (PC + 4)	Add (Branch)
a.	-3 and 20	PC and 4	PC + 4 and 20×4
b.	-32 and -128	PC and 4	PC + 4 and 2090×4

4.11.6

	Read Register 1	Read Register 2	Write Register	Write Data	RegWrite
a.	3	2	X (2 or 0)	X	0
b.	4	2	1	0	1

Solution 4.12**4.12.1**

	Pipelined	Single-Cycle
a.	350ps	1250ps
b.	220ps	950ps

4.12.2

	Pipelined	Single-Cycle
a.	1750ps	1250ps
b.	1100ps	950ps

4.12.3

	Stage to Split	New Clock Cycle Time
a.	ID	300ps
b.	EX	210ps

4.12.4

a.	35%
b.	30%

4.12.5

a.	65%
b.	70%

4.12.6 We already computed clock cycle times for pipelined and single-cycle organizations in 4.12.1, and the multi-cycle organization has the same clock cycle time as the pipelined organization. We will compute execution times relative to the pipelined organization. In single-cycle, every instruction takes one (long) clock cycle. In pipelined, a long-running program with no pipeline stalls completes one instruction in every cycle. Finally, a multi-cycle organization completes an LW in

5 cycles, an SW in 4 cycles (no WB), an ALU instruction in 4 cycles (no MEM), and a BEQ in 4 cycles (no WB). So we have the speedup of pipeline:

	Multi-Cycle Execution Time Is X Times Pipelined Execution Time, where X is	Single-Cycle Execution Time Is X Times Pipelined Execution Time, Where X Is
a.	$0.20 \times 5 + 0.80 \times 4 = 4.20$	$1250\text{ps}/350\text{ps} = 3.57$
b.	$0.15 \times 5 + 0.85 \times 4 = 4.15$	$950\text{ps}/220\text{ps} = 4.32$

Solution 4.13

4.13.1

	Instruction Sequence	Dependencies
a.	I1: SW R16, -100(R6) I2: LW R4, 8(R16) I3: ADD R5, R4, R4	RAW on R4 from I2 to I3
b.	I1: OR R1, R2, R3 I2: OR R2, R1, R4 I3: OR R1, R1, R2	RAW on R1 from I1 to I2 and I3 RAW on R2 from I2 to I3 WAR on R2 from I1 to I2 WAR on R1 from I2 to I3 WAW on R1 from I1 to I3

4.13.2 In the basic five-stage pipeline WAR and WAW dependences do not cause any hazards. Without forwarding, any RAW dependence between an instruction and the next two instructions (if register read happens in the second half of the clock cycle and the register write happens in the first half). The code that eliminates these hazards by inserting NOP instructions is:

	Instruction Sequence	
a.	SW R16, -100(R6) LW R4, 8(R16) NOP NOP ADD R5, R4, R4	Delay I3 to avoid RAW hazard on R4 from I2
b.	OR R1, R2, R3 NOP NOP OR R2, R1, R4 NOP NOP OR R1, R1, R2	Delay I2 to avoid RAW hazard on R1 from I1 Delay I3 to avoid RAW hazard on R2 from I2

4.13.3 With full forwarding, an ALU instruction can forward a value to the EX stage of the next instruction without a hazard. However, a load cannot forward to

the EX stage of the next instruction (but can't to the instruction after that). The code that eliminates these hazards by inserting NOP instructions is:

Instruction Sequence		
a.	SW R16, -100(R6) LW R4, 8(R16) NOP ADD R5, R4, R4	Delay I3 to avoid RAW hazard on R4 from I2 Value for R4 is forwarded from I2 now
b.	OR R1, R2, R3 OR R2, R1, R4 OR R1, R1, R2	No RAW hazard on R1 from I1 (forwarded) No RAW hazard on R2 from I2 (forwarded)

4.13.4 The total execution time is the clock cycle time times the number of cycles. Without any stalls, a three-instruction sequence executes in 7 cycles (5 to complete the first instruction, then one per instruction). The execution without forwarding must add a stall for every NOP we had in 4.13.2, and execution forwarding must add a stall cycle for every NOP we had in 4.13.3. Overall, we get:

	No Forwarding	With Forwarding	Speedup Due to Forwarding
a.	$(7 + 2) \times 250\text{ps} = 2250\text{ps}$	$(7 + 1) \times 300\text{ps} = 2400\text{ps}$	0.94 (This is really a slowdown)
b.	$(7 + 4) \times 180\text{ps} = 1980\text{ps}$	$7 \times 240\text{ps} = 1680\text{ps}$	1.18

4.13.5 With ALU-ALU-only forwarding, an ALU instruction can forward to the next instruction, but not to the second-next instruction (because that would be forwarding from MEM to EX). A load cannot forward at all, because it determines the data value in MEM stage, when it is too late for ALU-ALU forwarding. We have:

Instruction Sequence		
a.	SW R16, -100(R6) LW R4, 8(R16) ADD R5, R4, R4	ALU-ALU forwarding of R4 from I2
b.	OR R1, R2, R3 OR R2, R1, R4 OR R1, R1, R2	ALU-ALU forwarding of R1 from I1 ALU-ALU forwarding of R2 from I2

4.13.6

	No Forwarding	With ALU-ALU Forwarding Only	Speedup with ALU-ALU Forwarding
a.	$(7 + 2) \times 250\text{ps} = 2250\text{ps}$	$7 \times 290\text{ps} = 2030\text{ps}$	1.11
b.	$(7 + 4) \times 180\text{ps} = 1980\text{ps}$	$7 \times 210\text{ps} = 1470\text{ps}$	1.35

Solution 4.14

4.14.1 In the pipelined execution shown below, *** represents a stall when an instruction cannot be fetched because a load or store instruction is using the memory in that cycle. Cycles are represented from left to right, and for each instruction we show the pipeline stage it is in during that cycle:

	Instruction	Pipeline Stage	Cycles
a.	SW R16,12(R6)	IF ID EX MEM WB	11
	LW R16,8(R6)	IF ED EX MEM WB	
	BEQ R5,R4,Lb1	IF ID EX MEM WB	
	ADD R5,R1,R4	*** ** IF ID EX MEM WB	
	SLT R5,R15,R4	IF ID EX MEM WB	
b.	SW R2,0(R3)	IF ID EX MEM WB	9
	OR R1,R2,R3	IF ED EX MEM WB	
	BEQ R2,R0,Lb1	IF ID EX MEM WB	
	ADD R1,R4,R3	*** IF ID EX MEM WB	

We cannot add NOPs to the code to eliminate this hazard—NOPs need to be fetched just like any other instructions, so this hazard must be addressed with a hardware hazard detection unit in the processor.

4.14.2 This change only saves one cycle in an entire execution without data hazards (such as the one given). This cycle is saved because the last instruction finishes one cycle earlier (one less stage to go through). If there were data hazards from loads to other instructions, the change would help eliminate some stall cycles.

	Instructions Executed	Cycles with 5 Stages	Cycles with 4 Stages	Speedup
a.	5	$4 + 5 = 9$	$3 + 5 = 8$	$9/8 = 1.13$
b.	4	$4 + 4 = 8$	$3 + 4 = 7$	$8/7 = 1.14$

4.14.3 Stall-on-branch delays the fetch of the next instruction until the branch is executed. When branches execute in the EXE stage, each branch causes two stall cycles. When branches execute in the ID stage, each branch only causes one stall cycle. Without branch stalls (e.g., with perfect branch prediction) there are no stalls, and the execution time is 4 plus the number of executed instructions. We have:

	Instructions Executed	Branches Executed	Cycles with Branch in EXE	Cycles with Branch in ID	Speedup
a.	5	1	$4 + 5 + 1 \times 2 = 11$	$4 + 5 + 1 \times 1 = 10$	$11/10 = 1.10$
b.	4	1	$4 + 4 + 1 \times 2 = 10$	$4 + 4 + 1 \times 1 = 9$	$10/9 = 1.11$

4.14.4 The number of cycles for the (normal) 5-stage and the (combined EX/MEM) 4-stage pipeline is already computed in 4.14.2. The clock cycle time is equal to the latency of the longest-latency stage. Combining EX and MEM stages affects clock time only if the combined EX/MEM stage becomes the longest-latency stage:

	Cycle Time with 5 Stages	Cycle Time with 4 Stages	Speedup
a.	200ps (IF)	210ps (MEM + 20ps)	$(9 \times 200)/(8 \times 210) = 1.07$
b.	200ps (ID, EX, MEM)	220ps (MEM + 20ps)	$(8 \times 200)/(7 \times 220) = 1.04$

4.14.5

	New ID Latency	New EX Latency	New Cycle Time	Old Cycle Time	Speedup
a.	180ps	140ps	200ps (IF)	200ps (IF)	$(11 \times 200)/(10 \times 200) = 1.10$
b.	300ps	190ps	300ps (ID)	200ps (ID, EX, MEM)	$(10 \times 200)/(9 \times 300) = 0.74$

4.14.6 The cycle time remains unchanged: a 20ps reduction in EX latency has no effect on clock cycle time because EX is not the longest-latency stage. The change does affect execution time because it adds one additional stall cycle to each branch. Because the clock cycle time does not improve but the number of cycles increases, the speedup from this change will be below 1 (a slowdown). In 4.14.3 we already computed the number of cycles when branch is in EX stage. We have:

	Cycles with Branch in EX	Execution Time (Branch in EX)	Cycles with Branch in MEM	Execution Time (Branch in MEM)	Speedup
a.	$4 + 5 + 1 \times 2 = 11$	$11 \times 200\text{ps} = 2200\text{ps}$	$4 + 5 + 1 \times 3 = 12$	$12 \times 200\text{ps} = 2400\text{ps}$	0.92
b.	$4 + 4 + 1 \times 2 = 10$	$10 \times 200\text{ps} = 2000\text{ps}$	$4 + 4 + 1 \times 3 = 11$	$11 \times 200\text{ps} = 2200\text{ps}$	0.91

Solution 4.15

4.15.1

a.	This instruction behaves like a normal load until the end of the MEM stage. After that, it behaves like an ADD, so we need another stage after MEM to compute the result, and we need additional wiring to get the value of Rt to this stage.
b.	This instruction behaves like a load until the end of the MEM stage. After that, we need another stage to compare the value against Rt. We also need to add an input to the PC Mux that takes the value of Rd, and the Mux select signal must now include the result of the new comparison. We also need an extra read port in Registers because the instruction needs three registers to be read.

4.15.2

a.	We need to add a control signal that selects what the new stage does (just pass the value from memory through, or add the register value to it).
b.	We need a control signal similar to the existing “Branch” signal to control whether or not the new comparison is allowed to affect the PC. We also need to add one bit to the control signal that selects whether the target address is PC + 4 + Offs or the register value.

4.15.3

a.	The addition of a new stage either adds new forwarding paths (from the new stage to EX) or (if there is no forwarding) makes a stall due to a data hazard one cycle longer. Additionally, this instruction produces its result only at the end of the new stage, so even with forwarding it introduces a data hazard that requires a two-cycle stall if the ADDM instruction is immediately followed by a data-dependent instruction.
b.	The addition of a new stage either adds new forwarding paths (from the new stage to EX) or (if there is no forwarding) makes a stall due to a data hazard one cycle longer. The instruction itself creates a control hazard that leaves the next PC unknown until the BEQM instruction leaves the new stage, which is two cycles longer than for a normal BEQ.

4.15.4

a.	LW Rd, Offs(Rs) ADD Rd, Rt, Rd	E.g., ADDM can be used when trying to compute a sum of array elements.
b.	LW Rtmp, Offs(Rs) BNE Rtmp, Rt, Skip JR Rd Skip:	E.g., BEQM can be used when trying to determine if an array has an element with a specific value.

4.15.5 The instruction can be translated into simple MIPS-like micro-operations (see 4.15.4 for a possible translation). These micro-operations can then be executed by the processor with a “normal” pipeline.

4.15.6 We will compute the execution time for every replacement interval. The old execution time is simply the number of instructions in the replacement interval (CPI of 1). The new execution time is the number of instructions after we made the replacement, plus the number of added stall cycles. The new number of instructions is the number of instructions in the original replacement interval, plus the new instruction, minus the number of instructions it replaces:

	New Execution Time	Old Execution Time	Speedup
a.	$30 - (2 - 1) + 2 = 31$	30	0.97
b.	$40 - (3 - 1) + 1 = 39$	40	1.03

Solution 4.16

4.16.1 For every instruction, the IF/ID register keeps the PC + 4 and the instruction word itself. The ID/EX register keeps all control signals for the EX, MEM, and WB stages, PC + 4, the two values read from Registers, the sign-extended lowermost 16 bits of the instruction word, and Rd and Rt fields of the instruction word (even for instructions whose format does not use these fields). The EX/MEM register keeps control signals for the MEM and WB stages, the PC + 4 + Offset (where Offset is the sign-extended lowermost 16 bits of the instructions, even for instructions that have no offset field), the ALU result and the value of its Zero output, the value that was read from the second register in the ID stage (even for instructions that never need this value), and the number of the destination register (even for instructions that need no register writes; for these instructions the number of the destination register is simply a “random” choice between Rd or Rt). The MEM/WB register keeps the WB control signals, the value read from memory (or a “random” value if there was no memory read), the ALU result, and the number of the destination register.

4.16.2

	Need to be Read	Actually Read
a.	R6, R16	R6, R16
b.	R1, R0	R1, R0

4.16.3

	EX	MEM
a.	-100 + R6	Write value to memory
b.	R1 OR R0	Nothing

4.16.4

	Loop	
a.	2: LW R2,16(R2) 2: SLT R1,R2,R4 2: BEQ R1,R9,Loop 3: ADD R1,R2,R1 3: LW R2,0(R1) 3: LW R2,16(R2) 3: SLT R1,R2,R4 3: BEQ R1,R9,Loop	WB EX MEM WB ID EX MEM WB IF ID EX MEM WB IF ID EX MEM WB IF ID *** EX MEM IF *** ID *** IF ***

b.	LW R1,0(R1)	WB
	LW R1,0(R1)	EX MEM WB
	BEQ R1,R0,Loop	ID *** EX MEM WB
	LW R1,0(R1)	IF *** ID EX MEM WB
	AND R1,R1,R2	IF ID *** EX MEM WB
	LW R1,0(R1)	IF *** ID EX MEM
	LW R1,0(R1)	IF ID ***
	BEQ R1,R0,Loop	IF ***

4.16.5 In a particular clock cycle, a pipeline stage is not doing useful work if it is stalled or if the instruction going through that stage is not doing any useful work there. In the pipeline execution diagram from 4.16.4, a stage is stalled if its name is not shown for a particular cycle, and stages in which the particular instruction is not doing useful work are marked in red. Note that a BEQ instruction is doing useful work in the MEM stage, because it is determining the correct value of the next instruction's PC in that stage. We have:

	Cycles per Loop Iteration	Cycles in Which All Stages Do Useful Work	% of Cycles in Which All Stages Do Useful Work
a.	7	1	14%
b.	8	2	0%

4.16.6 The address of that first instruction of the third iteration (PC + 4 for the BEQ from the previous iteration) and the instruction word of the BEQ from the previous iteration.

Solution 4.17

4.17.1 Of all these instructions, the value produced by this adder is actually used only by a BEQ instruction when the branch is taken. We have:

a.	18% (60% of 30%)
b.	6% (60% of 10%)

4.17.2 Of these instructions, only ADD needs all three register ports (reads two registers and write one). BEQ and SW does not write any register, and LW only uses one register value. We have:

a.	40%
b.	60%

4.17.3 Of these instructions, only LW and SW use the data memory. We have:

a.	30% (25% + 5%)
b.	30% (20% + 10%)

4.17.4 The clock cycle time of a single-cycle is the sum of all latencies for the logic of all five stages. The clock cycle time of a pipelined datapath is the maximum latency of the five stage logic latencies, plus the latency of a pipeline register that keeps the results of each stage for the next stage. We have:

	Single-Cycle	Pipelined	Speedup
a.	760ps	215ps	3.53
b.	850ps	215ps	3.95

4.17.5 The latency of the pipelined datapath is unchanged (the maximum stage latency does not change). The clock cycle time of the single-cycle datapath is the sum of logic latencies for the four stages (IF, ID, WB, and the combined EX + MEM stage). We have:

	Single-Cycle	Pipelined
a.	610ps	215ps
b.	650ps	215ps

4.17.6 The clock cycle time of the two pipelines (5-stage and 4-stage) as explained for 4.17.5. The number of instructions increases for the 4-stage pipeline, so the speedup is below 1 (there is a slowdown):

	Instructions with 5-Stage	Instructions with 4-Stage	Speedup
a.	$1.00 \times I$	$1.00 \times I + 0.5 \times (0.25 + 0.05) \times I = 1.150 \times I$	0.87
b.	$1.00 \times I$	$1.00 \times I + 0.5 \times (0.20 + 0.10) \times I = 1.150 \times I$	0.87

Solution 4.18

4.18.1 No signals are asserted in IF and ID stages. For the remaining three stages we have:

	EX	MEM	WB
a.	ALUSrc = 1, ALUOp = 00, RegDst = 0	Branch = 0, MemWrite = 0, MemRead = 1	MemtoReg = 0, RegWrite = 1
b.	ALUSrc = 0, ALUOp = 10, RegDst = 1	Branch = 0, MemWrite = 0, MemRead = 0	MemtoReg = 1, RegWrite = 1

4.18.2 One clock cycle.

4.18.3 The PCSrc signal is 0 for this instruction. The reason against generating the PCSrc signal in the EX stage is that the AND must be done after the ALU computes its Zero output. If the EX stage is the longest-latency stage and the ALU output is on

its critical path, the additional latency of an AND gate would increase the clock cycle time of the processor. The reason in favor of generating this signal in the EX stage is that the correct next-PC for a conditional branch can be computed one cycle earlier, so we can avoid one stall cycle when we have a control hazard.

4.18.4

	Control Signal 1	Control Signal 2
a.	Generated in ID, used in EX	Generated in MEM, used in MEM
b.	Generated in ID, used in MEM	Generated in ID, used in WB

4.18.5

a.	None. PCSRc is only 1 for a taken branch, and ALUsrc is 0 for PC-relative branches.
b.	None. Branch is only 1 for conditional branches, and conditional branches do not write registers.

4.18.6 Signal 2 goes back through the pipeline. It affects execution of instructions that execute after the one for which the signal is generated, so it is not a time-travel paradox.

Solution 4.19

4.19.1 Dependences to the 1st next instruction result in 2 stall cycles, and the stall is also 2 cycles if the dependence is to both the 1st and 2nd next instruction. Dependences to only the 2nd next instruction result in one stall cycle. We have:

	CPI	Stall Cycles
a.	$1 + 0.35 \times 2 + 0.15 \times 1 = 1.85$	46% (0.85/1.85)
b.	$1 + 0.35 \times 2 + 0.25 \times 1 = 1.95$	49% (0.95/1.95)

4.19.2 With full forwarding, the only RAW data dependences that cause stalls are those from the MEM stage of one instruction to the 1st next instruction. Even these dependences cause only one stall cycle, so we have:

	CPI	Stall Cycles
a.	$1 + 0.20 = 1.20$	17% (0.20/1.20)
b.	$1 + 0.10 = 1.1$	13% (0.15/1.15)

4.19.3 With forwarding only from the EX/MEM register, EX to 1st dependences can be satisfied without stalls but any other dependences (even when together with EX to 1st) incur a one-cycle stall. With forwarding only from the MEM/WB register, EX to 2nd dependences incur no stalls. MEM to 1st dependences still incur a

one-cycle stall, and EX to 1st dependences now incur one stall cycle because we must wait for the instruction to complete the MEM stage to be able to forward to the next instruction. We compute stall cycles per instructions for each case as follows:

	EX/MEM	MEM/WB	Fewer Stall Cycles with
a.	$0.2 + 0.05 + 0.1 + 0.1 = 0.45$	$0.05 + 0.2 + 0.1 = 0.35$	MEM/WB
b.	$0.1 + 0.15 + 0.1 + 0.05 = 0.4$	$0.2 + 0.1 + 0.05 = 0.35$	MEM/WB

4.19.4 In 4.19.1 and 4.19.2 we have already computed the CPI without forwarding and with full forwarding. Now we compute time per instruction by taking into account the clock cycle time:

	Without Forwarding	With Forwarding	Speedup
a.	$1.85 \times 150\text{ps} = 277.5\text{ps}$	$1.20 \times 150\text{ps} = 180\text{ps}$	1.54
b.	$1.95 \times 300\text{ps} = 585\text{ps}$	$1.1 \times 350\text{ps} = 385\text{ps}$	1.52

4.19.5 We already computed the time per instruction for full forwarding in 4.19.4. Now we compute time per instruction with time-travel forwarding and the speedup over full forwarding:

	With Full Forwarding	Time-Travel Forwarding	Speedup
a.	$1.20 \times 150\text{ps} = 180\text{ps}$	$1 \times 250\text{ps} = 250\text{ps}$	0.72
b.	$1.1 \times 350\text{ps} = 385\text{ps}$	$1 \times 450\text{ps} = 450\text{ps}$	0.86

4.19.6

	EX/MEM	MEM/WB	Shorter Time per Instruction with
a.	$1.45 \times 150\text{ps} = 217.5$	$1.35 \times 150\text{ps} = 202.5\text{ps}$	MEM/WB
b.	$1.4 \times 330\text{ps} = 462$	$1.35 \times 320\text{ps} = 432\text{ps}$	MEM/WB

Solution 4.20

4.20.1

	Instruction Sequence	RAW	WAR	WAW
a.	I1: ADD R1, R2, R1 I2: LW R2, 0(R1) I3: LW R1, 4(R1) I4: OR R3, R1, R2	(R1) I1 to I2, I3 (R2) I2 to I4 (R1) I3 to I4	(R2) I1 to I2 (R1) I1, I2 to I3	(R1) I1 to I3

b.	I1: LW R1,0(R1) I2: AND R1,R1,R2 I3: LW R2,0(R1) I4: LW R1,0(R3)	(R1) I1 to I2 (R1) I2 to I3	(R1) I1 to I2 (R2) I2 to I3 (R1) I3 to I4	(R1) I1 to I2 (R1) I2 to I4
-----------	---	--------------------------------	---	--------------------------------

4.20.2 Only RAW dependences can become data hazards. With forwarding, only RAW dependences from a load to the very next instruction become hazards. Without forwarding, any RAW dependence from an instruction to one of the following 3 instructions becomes a hazard:

	Instruction Sequence	With Forwarding	Without Forwarding
a.	I1: ADD R1,R2,R1 I2: LW R2,0(R1) I3: LW R1,4(R1) I4: OR R3,R1,R2	(R1) I3 to I4	(R1) I1 to I2, I3 (R2) I2 to I4 (R1) I3 to I4
b.	I1: LW R1,0(R1) I2: AND R1,R1,R2 I3: LW R2,0(R1) I4: LW R1,0(R3)	(R1) I1 to I2	(R1) I1 to I2 (R1) I2 to I3

4.20.3 With forwarding, only RAW dependences from a load to the next two instructions become hazards because the load produces its data at the end of the second MEM stage. Without forwarding, any RAW dependence from an instruction to one of the following 4 instructions becomes a hazard:

	Instruction Sequence	With Forwarding	RAW
a.	I1: ADD R1,R2,R1 I2: LW R2,0(R1) I3: LW R1,4(R1) I4: OR R3,R1,R2	(R2) I2 to I4 (R1) I3 to I4	(R1) I1 to I2, I3 (R2) I2 to I4 (R1) I3 to I4
b.	I1: LW R1,0(R1) I2: AND R1,R1,R2 I3: LW R2,0(R1) I4: LW R1,0(R3)	(R1) I1 to I2	(R1) I1 to I2 (R1) I2 to I3

4.20.4

	Instruction Sequence	RAW
a.	I1: ADD R1, R2, R1 I2: LW R2, 0(R1) I3: LW R1, 4(R1) I4: OR R3, R1, R2	(R1) I1 to I2 (30 overrides -1)
b.	I1: LW R1, 0(R1) I2: AND R1, R1, R2 I3: LW R2, 0(R1) I4: LW R1, 0(R3)	(R1) I1 to I2 (0 overrides 4)

4.20.5 A register modification becomes “visible” to the EX stage of the following instructions only two cycles after the instruction that produces the register value leaves the EX stage. Our forwarding-assuming hazard detection unit only adds a one-cycle stall if the instruction that immediately follows a load is dependent on the load. We have:

	Instruction Sequence with Forwarding Stalls	Execution without Forwarding	Values after Execution
a.	I1: ADD R1, R2, R1 I2: LW R2, 0(R1) I3: LW R1, 4(R1) Stall I4: OR R3, R1, R2	R1 = 30 (Stall and after) R2 = 0 (I4 and after) R1 = 0 (after I4) R3 = 30 (after I4)	R0 = 0 R1 = 0 R2 = 0 R3 = 30
b.	I1: LW R1, 0(R1) Stall I2: AND R1, R1, R2 I3: LW R2, 0(R1) I4: LW R1, 0(R3)	R1 = 0 (I3 and after) R1 = 4 (after I4) R2 = 0 R1 = 0	R0 = 0 R1 = 0 R2 = 0 R3 = 3000

4.20.6

	Instruction Sequence with Forwarding Stalls	Correct Execution	Sequence with NOPs
a.	I1: ADD R1, R2, R1 I2: LW R2, 0(R1) I3: LW R1, 4(R1) Stall I4: OR R3, R1, R2	I1: ADD R1, R2, R1 Stall Stall I2: LW R2, 0(R1) I3: LW R1, 4(R1) Stall Stall I4: OR R3, R1, R2	ADD R1, R2, R1 NOP NOP LW R2, 0(R1) LW R1, 4(R1) NOP NOP OR R3, R1, R2

b.	I1: LW R1,0(R1) Stall I2: AND R1,R1,R2 I3: LW R2,0(R1) I4: LW R1,0(R3)	I1: LW R1,0(R1) Stall Stall I2: AND R1,R1,R2 Stall Stall I3: LW R2,0(R1) I4: LW R1,0(R3)	LW R1,0(R1) NOP NOP AND R1,R1,R2 NOP NOP LW R2,0(R1) LW R1,0(R3)
-----------	--	---	---

Solution 4.21

4.21.1

a.	ADD R5,R2,R1 NOP NOP LW R3,4(R5) LW R2,0(R2) NOP OR R3,R5,R3 NOP NOP SW R3,0(R5)
b.	LW R2,0(R1) NOP NOP AND R1,R2,R1 LW R3,0(R2) NOP LW R1,0(R1) NOP NOP SW R1,0(R2)

4.21.2 We can move up an instruction by swapping its place with another instruction that has no dependences with it, so we can try to fill some NOP slots with such instructions. We can also use R7 to eliminate WAW or WAR dependences so we can have more instructions to move up.

a.	<pre> I1: ADD R5, R2, R1 I3: LW R2, 0(R2) NOP I2: LW R3, 4(R5) NOP NOP I4: OR R3, R5, R3 NOP NOP I5: SW R3, 0(R5) </pre>	<p>Moved up to fill NOP slot</p> <p>Had to add another NOP here, so there is no performance gain</p>
b.	<pre> I1: LW R2, 0(R1) NOP NOP I2: AND R1, R2, R1 I3: LW R3, 0(R2) NOP I4: LW R1, 0(R1) NOP NOP I5: SW R1, 0(R2) </pre>	<p>No improvement is possible. There is a chain of RAW dependences from I1 to I2 to I4 to I5, and each step in the chain has to be separated by two instructions.</p>

4.21.3 With forwarding, the hazard detection unit is still needed because it must insert a one-cycle stall whenever the load supplies a value to the instruction that immediately follows that load. Without the hazard detection unit, the instruction that depends on the immediately preceding load gets the stale value the register had before the load instruction.

a.	Code executes correctly (for both loads, there is no RAW dependence between the load and the next instruction).
b.	I1 gets the value of R2 from before I1, not from I1 as it should. Also, I5 gets the value of R1 from I1, not from I4 as it should.

4.21.4 The outputs of the hazard detection unit are PCWrite, IF/IDWrite, and ID/EXZero (which controls the Mux after the output of the Control unit). Note that IF/IDWrite is always equal to PCWrite, and ED/ExZero is always the opposite of PCWrite. As a result, we will only show the value of PCWrite for each cycle. The outputs of the forwarding unit are ALUin1 and ALUin2, which control Muxes which select the first and second input of the ALU. The three possible values for ALUin1 or ALUin2 are 0 (no forwarding), 1 (forward ALU output from previous instruction), or 2 (forward data value for second-previous instruction). We have:

Solution 4.22

4.22.1

	Executed Instructions	Pipeline Cycles													
		1	2	3	4	5	6	7	8	9	10	11	12	13	14
a.	LW R2,0(R2) BEQ R2,R0,Label (T) LW R2,0(R2) BEQ R2,R0,Label (NT) OR R2,R2,R3 SW R2,0(R5)	IF	ID	EX	MEM	WB									
			IF	ID	***	EX	MEM	WB							
				IF	***	ID	EX	MEB	WB						
						IF	ID	***	EX	MEM	WB				
									IF	ID	EX	MEB	WB		
										IF	ID	EX	MEB	WB	
b.	LW R2,0(R1) BEQ R2,R0,Label2 (NT) LW R3,0(R2) BEQ R3,R0,Label1 (T) BEQ R2,R0,Label2 (T) SW R1,0(R2)	IF	ID	EX	MEM	WB									
			IF	ID	***	EX	MEB	WB							
							IF	ID	EX	MEB	WB				
								IF	ID	EX	***	EX	MEB	WB	
									IF	ID	***	ID	EX	MEB	WB
										IF	ID	EX	MEB	WB	

4.22.2

	Executed Instructions	Pipeline Cycles													
		1	2	3	4	5	6	7	8	9	10	11	12	13	14
a.	LW R2,0(R2) BEQ R2,R0,Label (T) OR R2,R2,R3 LW R2,0(R2) BEQ R2,R0,Label (NT) OR R2,R2,R3 SW R2,0(R5)	IF	ID	EX	MEM	WB									
			IF	ID	***	EX	MEB	WB							
				IF	***	ID	EX	MEB	WB						
						IF	ID	***	EX	MEM	WB				
							IF	***	ID	EX	MEM	WB			
									IF	ID	EX	MEM	WB		
										IF	ID	EX	MEM	WB	
b.	LW R2,0(R1) BEQ R2,R0,Label2 (NT) LW R3,0(R2) BEQ R3,R0,Label1 (T) ADD R1,R3,R1 BEQ R2,R0,Label2 (T) LW R3,0(R2) SW R1,0(R2)	IF	ID	EX	MEM	WB									
			IF	ID	***	EX	MEM	WB							
				IF	***	ID	EX	MEB	WB						
							IF	ID	EX	MEM	WB				
								IF	ID	EX	MEM	WB			
									IF	ID	EX	MEM	WB		
										IF	ID	EX	MEM	WB	
											IF	ID	EX	MEM	WB

4.22.3

a.	Label: LW R2,0(R2) BEZ R2,Label ; Taken once, then not taken OR R2,R2,R3 SW R2,0(R5)
b.	Label1: LW R2,0(R1) Label1: BEZ R2,Label2 ; Not taken once, then taken LW R3,0(R2) BEZ R3,Label1 ; Taken ADD R1,R3,R1 Label2: SW R1,0(R2)

4.22.4 The hazard detection logic must detect situations when the branch depends on the result of the previous R-type instruction, or on the result of two previous loads. When the branch uses the values of its register operands in its ID stage, the R-type instruction's result is still being generated in the EX stage. Thus we must stall the processor and repeat the ID stage of the branch in the next cycle. Similarly, if the branch depends on a load that immediately precedes it, the result of the load is only generated two cycles after the branch enters the ID stage, so we must stall the branch for two cycles. Finally, if the branch depends on a load that is the second-previous instruction, the load is completing its MEM stage when the branch is in its ID stage, so we must stall the branch for one cycle. In all three cases, the hazard is a data hazard.

Note that in all three cases we assume that the values of preceding instructions are forwarded to the ID stage of the branch if possible.

4.22.5 For 4.22.1 we have already shown the pipeline execution diagram for the case when branches are executed in the EX stage. The following is the pipeline diagram when branches are executed in the ID stage, including new stalls due to data dependences described for 4.22.4:

	Executed Instructions	Pipeline Cycles														
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
a.	LW R2,0(R2) BEQ R2,R0,Label (T) LW R2,0(R2) BEQ R2,R0,Label (NT) OR R2,R2,R3 SW R2,0(R5)	IF	ID	EX	MEM	WB										
			IF	***	***	ID	EX	MEB	WB							
						IF	ID	EX	MEB	WB						
							IF	***	***	ID	EX	MEB	WB			
										IF	ID	EX	MEB	WB		
											IF	ID	EX	MEB	WB	
b.	LW R2,0(R1) BEQ R2,R0,Label2 (NT) LW R3,0(R2) BEQ R3,R0,Label1 (T) BEQ R2,R0,Label2 (T) SW R1,0(R2)	IF	ID	EX	MEM	WB										
			IF	***	***	ID	EX	MEM	WB							
						IF	ID	EX	MEB	WB						
							IF	***	***	ID	EX	MEB	WB			
										IF	ID	EX	MEM	WB		
											IF	ID	EX	MEB	WB	

Now the speedup can be computed as:

a.	$14/14 = 1$
b.	$14/15 = 0.93$

4.22.6 Branch instructions are now executed in the ID stage. If the branch instruction is using a register value produced by the immediately preceding instruction, as we described for 4.22.4 the branch must be stalled because the preceding

instruction is in the EX stage when the branch is already using the stale register values in the ID stage. If the branch in the ID stage depends on an R-type instruction that is in the MEM stage, we need forwarding to ensure correct execution of the branch. Similarly, if the branch in the ID stage depends on an R-type of load instruction in the WB stage, we need forwarding to ensure correct execution of the branch. Overall, we need another forwarding unit that takes the same inputs as the one that forwards to the EX stage. The new forwarding unit should control two Muxes placed right before the branch comparator. Each Mux selects between the value read from Registers, the ALU output from the EX/MEM pipeline register, and the data value from the MEM/WB pipeline register. The complexity of the new forwarding unit is the same as the complexity of the existing one.

Solution 4.23

4.23.1 Each branch that is not correctly predicted by the always-taken predictor will cause 3 stall cycles, so we have:

	Extra CPI
a.	$3 \times (1 - 0.45) \times 0.25 = 0.41$
b.	$3 \times (1 - 0.65) \times 0.08 = 0.08$

4.23.2 Each branch that is not correctly predicted by the always-not-taken predictor will cause 3 stall cycles, so we have:

	Extra CPI
a.	$3 \times (1 - 0.55) \times 0.25 = 0.34$
b.	$3 \times (1 - 0.35) \times 0.08 = 0.16$

4.23.3 Each branch that is not correctly predicted by the 2-bit predictor will cause 3 stall cycles, so we have:

	Extra CPI
a.	$3 \times (1 - 0.85) \times 0.25 = 0.113$
b.	$3 \times (1 - 0.98) \times 0.08 = 0.005$

4.23.4 Correctly predicted branches had CPI of 1 and now they become ALU instructions whose CPI is also 1. Incorrectly predicted instructions that are converted also become ALU instructions with a CPI of 1, so we have:

	CPI without Conversion	CPI with Conversion	Speedup from Conversion
a.	$1 + 3 \times (1 - 0.85) \times 0.25 = 1.113$	$1 + 3 \times (1 - 0.85) \times 0.25 \times 0.5 = 1.056$	$1.113/1.056 = 1.054$
b.	$1 + 3 \times (1 - 0.98) \times 0.08 = 1.005$	$1 + 3 \times (1 - 0.98) \times 0.08 \times 0.5 = 1.002$	$1.005/1.002 = 1.003$

4.23.5 Every converted branch instruction now takes an extra cycle to execute, so we have:

	CPI without Conversion	Cycles per Original Instruction with Conversion	Speedup from Conversion
a.	1.113	$1 + (1 + 3 \times (1 - 0.85)) \times 0.25 \times 0.5 = 1.181$	$1.113/1.181 = 0.94$
b.	1.015	$1 + (1 + 3 \times (1 - 0.98)) \times 0.08 \times 0.5 = 1.042$	$1.005/1.042 = 0.96$

4.23.6 Let the total number of branch instructions executed in the program be B. Then we have:

	Correctly Predicted	Correctly Predicted Non-Loop-Back	Accuracy on Non-Loop-Back Branches
a.	$B \times 0.85$	$B \times 0.05$	$(B \times 0.05)/(B \times 0.20) = 0.25$ (25%)
b.	$B \times 0.98$	$B \times 0.18$	$(B \times 0.18)/(B \times 0.20) = 0.90$ (90%)

Solution 4.24

4.24.1

	Always Taken	Always Not-taken
a.	$2/4 = 50\%$	$2/4 = 50\%$
b.	$3/5 = 60\%$	$2/5 = 40\%$

4.24.2

	Outcomes	Predictor Value at Time of Prediction	Correct or Incorrect	Accuracy
a.	T, T, NT, NT	0,1,2,1	I,I,I,C	25%
b.	T, NT, T, T	0,1,0,1	I,C,I,I	25%

4.24.3 The first few recurrences of this pattern do not have the same accuracy as the later ones because the predictor is still warming up. To determine the accuracy in the “steady state,” we must work through the branch predictions until the predictor values start repeating (i.e., until the predictor has the same value at the start of the current and the next recurrence of the pattern).

	Outcomes	Predictor Value at Time of Prediction	Correct or Incorrect (in Steady State)	Accuracy in Steady State
a.	T,T,NT,NT	1 st occurrence: 0,1,2,1 2 nd occurrence: 0,1,2,1	I,I,I,C	25%
b.	T, NT, T, T, NT	1 st occurrence: 0,1,0,1,2 2 nd occurrence: 1,2,1,2,3 3 rd occurrence: 2,3,2,3,3 4 th occurrence: 2,3,2,3,3	C,I,C,C,I	60%

4.24.4 The predictor should be an N-bit shift register, where N is the number of branch outcomes in the target pattern. The shift register should be initialized with the pattern itself (0 for NT, 1 for T), and the prediction is always the value in the leftmost bit of the shift register. The register should be shifted after each predicted branch.

4.24.5 Since the predictor's output is always the opposite of the actual outcome of the branch instruction, the accuracy is zero.

4.24.6 The predictor is the same as in 4.24.4, except that it should compare its prediction to the actual outcome and invert (logical NOT) all the bits in the shift register if the prediction is incorrect. This predictor still always perfectly predicts the given pattern. For the opposite pattern, the first prediction will be incorrect, so the predictor's state is inverted and after that the predictions are always correct. Overall, there is no warm-up period for the given pattern, and the warm-up period for the opposite pattern is only one branch.

Solution 4.25

4.25.1

	Instruction 1	Instruction 2
a.	Invalid target address (EX)	Invalid data address (MEM)
b.	Invalid target address (EX)	Invalid data address (MEM)

4.25.2 The Mux that selects the next PC must have inputs added to it. Each input is a constant address of an exception handler. The exception detectors must be added to the appropriate pipeline stage and the outputs of these detectors must be used to control the pre-PC Mux, and also to convert to NOPs instructions that are already in the pipeline behind the exception-triggering instruction.

4.25.3 Instructions are fetched normally until the exception is detected. When the exception is detected, all instructions that are in the pipeline after the first instruction must be converted to NOPs. As a result, the second instruction never completes and does not affect pipeline state. In the cycle that immediately follows the

cycle in which the exception is detected, the processor will fetch the first instruction of the exception handler.

4.25.4

	Handler Address
a.	0x1000E230
b.	0x678A0000

The first instruction word from the handler address is fetched in the cycle after the one in which the original exception is detected. When this instruction is decoded in the next cycle, the processor detects that the instruction is invalid. This exception is treated just like a normal exception—it converts the instruction being fetched in that cycle into an NOP and puts the address of the Invalid Instruction handler into the PC at the end of the cycle in which the Invalid Instruction exception is detected.

4.25.5 This approach requires us to fetch the address of the handler from memory. We must add the code of the exception to the address of the exception vector table, read the handler's address from memory, and jump to that address. One way of doing this is to handle it like a special instruction that computes the address in EX, loads the handler's address in MEM, and sets the PC in WB.

4.25.6 We need a special instruction that allows us to move a value from the (exception) Cause register to a general-purpose register. We must first save the general-purpose register (so we can restore it later), load the Cause register into it, add the address of the vector table to it, use the result as an address for a load that gets the address of the right exception handler from memory, and finally jump to that handler.

Solution 4.26

4.26.1 All exception-related signals are 0 in all stages, except the one in which the exception is detected. For that stage, we show values of Flush signals for various stages, and also the value of the signal that controls the Mux that supplies the PC value.

	Stage	Signals
a.	ID	IF.Flush = ID.Flush = 1, PCSel = Exc
b.	EX	IF.Flush = ID.Flush = EX.Flush = 1, PCSel = Exc

4.26.2 The signals stored in the ID/EX stage are needed to execute the instruction if there are no exceptions. Figure 4.66 does not show it, but exception conditions from various stages are also supplied as inputs to the Control unit. The signal that goes directly to EX is EX.Flush and it is based on these exception condition inputs, not on the opcode of the instruction that is in the ID stage. In particular, the

EX.Flush signal becomes 1 when the instruction in the EX stage triggers an exception and must be prevented from completing.

4.26.3 The disadvantage is that the exception handler begins executing one cycle later. Also, an exception condition normally checked in MEM cannot be delayed into WB, because at that time the instruction is updating registers and cannot be prevented from doing so.

4.26.4 When overflow is detected in EX, each exception results in a 3-cycle delay (IF, ID, and EX are cancelled). By moving overflow into MEM, we add one more cycle to this delay. To compute the speedup, we compute execution time per 100,000 instructions:

	Old Clock Cycle Time	New Clock Cycle Time	Old Time per 100,000 Instructions	New Time per 100,000 Instructions	Speedup
a.	250ps	220ps	250ps × 100,003	220ps × 100,004	1.13635
b.	200ps	175ps	200ps × 100,003	175ps × 100,004	1.14285

4.26.5 Exception control (Flush) signals are not really generated in the EX stage. They are generated by the Control unit, which is drawn as part of the ID stage, but we could have a separate “Exception Control” unit to generate Flush signals and this unit is not really a part of any stage.

4.26.6 Flush signals must be generated one Mux time before the end of the cycle. However, their generation can only begin after exception conditions are generated. For example, arithmetic overflow is only generated after the ALU operation in EX is complete, which is usually in the later part of the clock cycle. As a result, the Control unit actually has very little time to generate these signals, and they can easily be on the critical path that determines the clock cycle time.

Solution 4.27

4.27.1 When the invalid instruction (I3) is decoded, IF.Flush and ID.Flush signals are used to convert I3 and I4 into NOPs (marked with *). In the next cycle, in IF we fetch the first instruction of the exception handler, in ID we have an NOP (instead of I4, marked), in EX we have an NOP (instead of I3), and I1 and I2 still continue through the pipeline normally:

	Branch and Delay Slot	Pipeline
a.	I1: BEQ R5,R4,Label I2: SLT R5,R15,R4 I3: Invalid I4: Something I5: Handler	IF ID EX MEM WB IF ID EX MEM IF ID *EX IF *ID IF

b.	I1: BEQ R1,R0,Label	IF	ID	EX	MEM	WB
	I2: LW R1,0(R1)		IF	ID	EX	MEM
	I3: Invalid			IF	ID	*EX
	I4: Something				IF	*ID
	I5: Handler					IF

4.27.2 When I2 is in the MEM stage, it triggers an exception condition that results in converting I2 and I5 into NOPs (I3 and I4 are already NOPs by then). In the next cycle, we fetch I6, which is the first instruction of the exception handler for the exception triggered by I2.

	Branch and Delay Slot	Branch and Delay Slot				
a.	I1: BEQ R5,R4,Label	IF	ID	EX	MEM	WB
	I2: SLT R5,R15,R4		IF	ID	EX	MEM *WB
	I3: Invalid			IF	ID	*EX *ME
	I4: Something				IF	*ID *EX
	I5: Handler 1					IF *ID
	I6: Handler 2					IF
b.	I1: BEQ R1,R0,Label	IF	ID	EX	MEM	WB
	I2: LW R1,0(R1)		IF	ID	EX	MEM *WB
	I3: Invalid			IF	ID	*EX *ME
	I4: Something				IF	*ID *EX
	I5: Handler 1					IF *ID
	I6: Handler 2					IF

4.27.3 The EPC is the PC + 4 of the delay-slot instruction. As described in Section 4.9, the exception handler subtracts 4 from the EPC, so it gets the address of the instruction that generated the exception (I2, the delay-slot instruction). If the exception handler decides to resume execution of the application, it will jump to the I2. Unfortunately, this causes the program to continue as if the branch was not taken, even if it was taken.

4.27.4 The processor cancels the store instruction and other instructions (from the “Invalid instruction” exception handler) fetched after it, and then begins fetching instructions from the invalid data address handler. A major problem here is that the new exception sets the EPC to the instruction address in the “Invalid instruction” handler, overwriting the EPC value that was already there (address for continuing the program). If the invalid data address handler repairs the problem and attempts to continue the program, the “Invalid instruction” handler will be executed. However, if it manages to repair the problem and wants to continue the program, the EPC is incorrect (it was overwritten before it could be saved). This is the reason why exception handlers must be written carefully to avoid triggering exceptions themselves, at least until they have safely saved the EPC.

4.27.5 Not for store instructions. If we check for the address overflow in MEM, the store is already writing data to memory in that cycle and we can no longer “cancel” it. As a result, when the exception handler is called the memory is already

changed by the store instruction, and the handler cannot observe the state of the machine that existed before the store instruction.

4.27.6 We must add two comparators to the EX stage, one that compares the ALU result to WADDR, and another that compares the data value from Rt to WVAL. If one of these comparators detects equality and the instruction is a store, this triggers a “Watchpoint” exception. As discussed for 4.27.5, we cannot delay the comparisons until the MEM stage because at that time the store is already done and we need to stop the application at the point before the store happens.

Solution 4.28

4.28.1

a.	<pre> ADD R2, R0, R0 Again: BEQ R2, R8, End ADD R3, R2, R9 LW R4, 0(R3) SW R4, 1(R3) ADDI R2, R2, 2 BEQ R0, R0, Again End: </pre>
b.	<pre> ADD R5, R0, R0 Again: BEQ R5, R6, End ADD R10, R5, R1 LW R11, 0(R10) LW R10, 1(R10) SUB R10, R11, R10 ADD R11, R5, R2 SW R10, 0(R11) ADDI R5, R5, 2 BEW R0, R0, Again End: </pre>

4.28.2

	Instructions	Pipeline
a.	<pre> ADD R2, R0, R0 BEQ R2, R8, End ADD R3, R2, R9 LW R4, 0(R3) SW R4, 1(R3) ADDI R2, R2, 2 BEQ R0, R0, Again BEQ R2, R8, End ADD R3, R2, R9 LW R4, 0(R3) SW R4, 1(R3) ADDI R2, R2, 2 BEQ R0, R0, Again BEQ R2, R8, End </pre>	<pre> IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB </pre>

4.28.4

	Instructions	Pipeline
a.	ADD R2, R0, R0 ADD R3, R2, R9 BEQ R2, R8, End LW R4, 0(R3) ADDI R2, R2, 2 SW R4, 1(R3) BEQ R0, R0, Again ADD R3, R2, R9 BEQ R2, R8, End LW R4, 0(R3) ADDI R2, R2, 2 SW R4, 1(R3) BEQ R0, R0, Again ADD R3, R2, R9 BEQ R2, R8, End	<pre> IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB </pre>
b.	ADD R5, R0, R0 ADD R10, R5, R1 BEQ R5, R6, End LW R11, 0(R10) ADD R12, R5, R2 LW R10, 1(R10) ADDI R5, R5, 2 SUB R10, R11, R10 SW R10, 0(R12) BEQ R0, R0, Again ADD R10, R5, R1 BEQ R5, R6, End LW R11, 0(R10) ADD R12, R5, R2 LW R10, 1(R10) ADDI R5, R5, 2 SUB R10, R11, R10 SW R10, 0(R12) BEQ R0, R0, Again ADD R10, R5, R1 BEQ R5, R6, End	<pre> IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB </pre>

4.28.5

	CPI for 1-Issue	CPI for 2-Issue	Speedup
a.	1 (no data hazards)	0.83 (5 cycles for 6 instructions). In every iteration SW can execute in parallel with the next instruction.	1.20
b.	1.11 (10 cycles per 9 instructions). There is 1 stall cycle in each iteration due to a data hazard between the second LW and the next instruction (SUB).	1.06 (19 cycles per 18 instructions). Neither of the two LW instructions can execute in parallel with another instruction, and SUB stalls because it depends on the second LW. The SW instruction executes in parallel with ADDI in even-numbered iterations.	1.05

4.28.6

	CPI for 1-Issue	CPI for 2-Issue	Speedup
a.	1	0.67 (4 cycles for 6 instructions). In every iteration ADD and LW cannot execute in the same cycle because of a data dependence. The rest of the instructions can execute in pairs.	1.49
b.	1.11	0.83 (15 cycles per 18 instructions). In all iterations, SUB is stalled because it depends on the second LW. The only instructions that execute in odd-numbered iterations as a pair are ADDI and BEQ. In even-numbered iterations, only the two LW instructions cannot execute as a pair.	1.34

Solution 4.29

4.29.1 Note that all register read ports are active in every cycle, so 4 register reads (2 instructions with 2 reads each) happen in every cycle. We determine the number of cycles it takes to execute an iteration of the loop and the number of useful reads, then divide the two. The number of useful register reads for an instruction is the number of source register parameters minus the number of registers that are forwarded from prior instructions. We assume that register writes happen in the first half of the cycle and the register reads happen in the second half.

	Loop	Pipeline Stages	Useful Reads	% Useful
a.	ADD R2, R2, R3 BEQ R2, zero, Loop ADDI R1, R1, 4 LW R2, 0(R1) LW R3, 16(R1) ADD R2, R2, R1 ADD R2, R2, R3 BEQ R2, zero, Loop	ID EX ME WB ID ** EX ME WB IF ** ID EX ME WB IF ** ID ** EX ME WB IF ** ID EX ME WB IF ** ID EX ME WB IF ID EX ME WB IF ID ** EX ME WB	1 0 (R1 fw) 0 (R1 fw) 1 (R1, R2 fw) 0 (R2, R3 fw) 1 (R2 fw)	15% (3/(5×4))
b.	AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop LW R1, 0(R1) AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop	EX ME WB ID EX ME WB ID EX ME WB IF ID EX ME WB IF ID ** ** EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB	0 (R1 fw) 0 (R1, R2 fw) 1 1 (R1 fw)	12.5% (2/(4×4))

4.29.2 The utilization of read ports is lower with a wider-issue processor:

	Loop	Pipeline Stages	Useful Reads	% Useful
a.	ADD R2, R2, R3 BEQ R2, zero, Loop ADDI R1, R1, 4 LW R2, 0(R1) LW R3, 16(R1) ADD R2, R2, R1 ADD R2, R2, R3 BEQ R2, zero, Loop	ID ** ** EX ME WB ID ** ** ** EX ME WB IF ** ** ** ID EX ME WB IF ** ** ** ID ** EX ME WB IF ** ** ** ID ** EX ME WB IF ** ID ** EX ME WB IF ** ID ** ** EX ME WB IF ** ID ** ** ** EX ME WB	1 0 (R1 fw) 0 (R1 fw) 0 (R1,R2 fw) 0 (R2,R3 fw) 1 (R2 fw)	5.5% (2/(6 × 6))
b.	AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop LW R1, 0(R1) AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop LW R1, 0(R1) AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop LW R1, 0(R1) AND R1, R1, R2 LW R2, 0(R2) BEQ R1, zero, Loop	ID ** EX ME WB ID ** EX ME WB ID ** ** EX ME WB IF ** ** ID EX ME WB IF ** ** ID ** ** EX ME WB IF ** ** ID ** ** EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB IF ** ** ID ** ** EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB IF ** ** ID EX ME WB IF ID ** EX ME WB IF ID ** EX ME WB IF ID ** ** EX ME WB	0 (R1 fw) 0 (R1,R2 fw) 0 (R2 fw) 1 (R1 fw) 0 (R1 fw) 0 (R1,R2 fw) 1 1 (R1 fw) 0 (R1 fw) 0 (R1,R2 fw) 0 (R2 fw) 1 (R1 fw)	6.7% (4/(10 × 6))

4.29.3

	2 Ports Used	3 Ports Used
a.	1 cycle out of 6 (16.7%)	Never (0%)
b.	3 cycles out of 10 (30%)	1 cycle out of 10 (10%)

4.29.4

	Unrolled and Scheduled Loop	Comment
a.	NOP Loop: LW R2, 8(R1) LW R3, 24(R1) ADDI R1, R1, 8 ADD R2, R2, R1 ADD R2, R2, R3 BEQ R2, zero, Loop	We are able to complete one iteration of the unrolled loop every 4 cycles. Both loads and adds that come from the first original iteration of the unrolled loop can be eliminated (they are only used to compute R2 for BEQ, which is removed). We combine both ADDI instructions and then schedule the unrolled loop to execute in four cycles per (unrolled) iteration, which is optimal. Note the NOP before the loop, which is needed to ensure that BEQ always executes together with the first LW of the next iteration.

b.	<pre> NOP Loop: LW R1, 0(R1) LW R10, 0(R2) NOP AND R1, R1, R2 LW R1, 0(R1) AND R1, R1, R10 LW R2, 0(R10) BEQ R1, zero, Loop </pre>	<p>We are able to execute one iteration of the unrolled loop in 6 cycles, which is optimal. Note the NOP before the loop, which is needed to ensure that BEQ always executes together with the first LW of the next iteration.</p>
-----------	--	--

4.29.5 We determine the number of cycles needed to execute two iterations of the original loop (one iteration of the unrolled loop). Note that we cannot use CPI in our speedup computation because the two versions of the loop do not execute the same instructions.

	Original Loop	Unrolled Loop	Speedup
a.	$5 \times 2 = 10$	4	2.5
b.	$4 \times 2 = 8$	6	1.3

4.29.6 On a pipelined processor the number of cycles per iteration is easily computed by adding together the number of instructions and the number of stalls. The only stalls occur when an LW instruction is followed immediately with a RAW-dependent instruction, so we have:

	Original Loop	Unrolled Loop	Speedup
a.	$6 \times 2 = 12$	6	2
b.	$(4 + 1) \times 2 = 10$	9	1.1

Solution 4.30

4.30.1 Let p be the probability of having a mispredicted branch. Whenever we have an incorrectly predicted BEQ as the first of the two instructions in a cycle (the probability of this event is p), we waste one issue slot (half a cycle) and another two entire cycles. If the first instruction in a cycle is not a mispredicted BEQ but the second one is (the probability of this is $(1 - p) \times p$), we waste two cycles. Without these mispredictions, we would be able to execute two instructions per cycle. We have:

	CPI
a.	$0.5 + 0.05 \times 2.5 + 0.95 \times 0.05 \times 2 = 0.720$
b.	$0.5 + 0.01 \times 2.5 + 0.99 \times 0.01 \times 2 = 0.545$

4.30.2 Inability to predict a branch results in the same penalty as a mispredicted branch. We compute the CPI like in 4.30.1, but this time we also have a 2-cycle penalty if we have a correctly predicted branch in the first issue slot and another branch that would be correctly predicted in the second slot. We have:

	CPI with 2 Predicted Branches per Cycle	CPI with 1 Predicted Branch per Cycle	Speedup
a.	0.720	$0.5 + 0.05 \times 2.5 + 0.95 \times 0.05 \times 2 + 0.20 \times 0.20 \times 2 = 0.800$	1.11
b.	0.545	$0.5 + 0.01 \times 2.5 + 0.99 \times 0.01 \times 2 + 0.04 \times 0.04 \times 2 = 0.548$	1.01

4.30.3 We have a one-cycle penalty whenever we have a cycle with two instructions that both need a register write. Such instructions are ALU and LW instructions. Note that BEQ does not write registers, so stalls due to register writes and due to branch mispredictions are independent events. We have:

	CPI with 2 Register Writes per Cycle	CPI with 1 Register Write per Cycle	Speedup
a.	0.720	$0.5 + 0.05 \times 2.5 + 0.95 \times 0.05 \times 2 + 0.65 \times 0.65 \times 1 = 1.143$	1.59
b.	0.545	$0.5 + 0.01 \times 2.5 + 0.99 \times 0.01 \times 2 + 0.75 \times 0.75 \times 1 = 1.107$	2.03

4.30.4 We have already computed the CPI with the given branch prediction accuracy, and we know that the CPI with ideal branch prediction is 0.5, so:

	CPI with Given Branch Prediction	CPI with Perfect Branch Prediction	Speedup
a.	0.720	0.5	1.44
b.	0.545	0.5	1.09

4.30.5 The CPI with perfect branch prediction is now 0.25 (four instructions per cycle). A branch misprediction in the first issue slot of a cycle results in 2.75 penalty cycles (remaining issue slots in the same cycle plus 2 entire cycles), in the second issue slot 2.5 penalty cycles, in the third slot 2.25 penalty cycles, and in the last (fourth) slot 2 penalty cycles. We have:

	CPI with Given Branch Prediction	CPI with Perfect Branch Prediction	Speedup
a.	$0.25 + 0.05 \times 2.75 + 0.95 \times 0.05 \times 2.5 + 0.95^2 \times 0.05 \times 2.25 + 0.95^3 \times 0.05 \times 2 = 0.694$	0.25	2.77
b.	$0.25 + 0.01 \times 2.75 + 0.99 \times 0.01 \times 2.5 + 0.99^2 \times 0.01 \times 2.25 + 0.99^3 \times 0.01 \times 2 = 0.344$	0.25	1.37

The speedup from improved branch prediction is much larger in a 4-issue processor than in a 2-issue processor. In general, processors that issue more instructions per cycle gain more from improved branch prediction because each branch misprediction costs them more instruction execution opportunities (e.g., 4 per cycle in 4-issue vs. 2 per cycle in 2-issue).

4.30.6 With this pipeline, the penalty for a mispredicted branch is 20 cycles plus the fraction of a cycle due to discarding instructions that follow the branch in the same cycle. We have:

	CPI with Given Branch Prediction	CPI with Perfect Branch Prediction	Speedup
a.	$0.25 + 0.05 \times 20.75 + 0.95 \times 0.05 \times 20.5 + 0.95^2 \times 0.05 \times 20.25 + 0.95^3 \times 0.05 \times 20 = 4.032$	0.25	16.13
b.	$0.25 + 0.01 \times 20.75 + 0.99 \times 0.01 \times 20.5 + 0.99^2 \times 0.01 \times 20.25 + 0.99^3 \times 0.01 \times 20 = 1.053$	0.25	4.21

We observe huge speedups when branch prediction is improved in a processor with a very deep pipeline. In general, processors with deeper pipelines benefit more from improved branch prediction because these processors cancel more instructions (e.g., 20 stages worth of instructions in a 50-stage pipeline vs. 2 stages worth of instructions in a 5-stage pipeline) on each misprediction.

Solution 4.31

4.31.1 The number of cycles is equal to the number of instructions (one instruction is executed per cycle) plus one additional cycle for each data hazard which occurs when an LW instruction is immediately followed by a dependent instruction. We have:

	CPI
a.	$(12 + 3)/12 = 1.25$
b.	$(9 + 2)/9 = 1.22$

4.31.2 The number of cycles is equal to the number of instructions (one instruction is executed per cycle), plus the stall cycles due to data hazards. Data hazards occur when the memory address used by the instruction depends on the result of a previous instruction (EXE to ARD, 2 stall cycles) or the instruction after that (1 stall cycle), or when an instruction writes a value to memory and one of the next

two instructions reads a value from the same address (2 or 1 stall cycles). All other data dependences can use forwarding to avoid stalls. We have:

	Instructions	Stall Cycles	CPI
a.	I1: mov -4(esp), eax I2: mov -4(esp), edx I3: add (edi,eax,4),edx I4: mov edx, -4(esp) I5: mov -4(esp),eax I6: cmp 0, (edi,eax,4) I7: jne Label	1 (eax from I1) 2 (read from I4) 2 (eax from I6)	$(7 + 5)/7 = 1.71$
b.	I1: add 4, edx I2: mov (edx), eax I3: add 4(edx), eax I4: add 8(edx), eax I5: mov eax, -4(edx) I6: test edx, edx I7: jl Label	2 (edx from I2)	$(7 + 2)/7 = 1.29$

4.31.3 The number of instructions here is that from the x86 code, but the number of cycles per iteration is that from the MIPS code (we fetch x86 instructions, but after instructions are decoded we end up executing the MIPS version of the loop):

	CPI
a.	$15/7 = 2.14$
b.	$11/7 = 1.57$

4.31.4 Dynamic scheduling allows us to execute an independent “future” instruction when the one we should be executing stalls. We have:

	Instructions	Reordering	CPI
a.	I1: lw r2, -4(sp) I2: lw r3, -4(sp) I3: sll r2, r2, 2 I4: add r2, r2, r4 I5: lw r2, 0(r2) I6: add r3, r3, r2 I7: sw r3, -4(sp) I8: lw r2, -4(sp) I9: sll r2, r2, 2 I10: add r2, r2, r4 I11: lw r2, 0(r2) I12: bne r2, zero, Label	I6 stalls, and all subsequent instructions have dependences so this stall remains. I9 stalls, but we can do I2 from the next iteration instead. However, this makes I3 stall and we can't eliminate that stall. I12 stalls and all subsequent instructions that remain have dependences so this stall remains.	$(12 + 3)/12 = 1.25$

b.	<pre> I1: addi r4,r4,4 I2: lw r3,0(r4) I3: lw r2,4(r4) I4: add r2,r2,r3 I5: lw r3,8(r4) I6: add r2,r2,r3 I7: sw r2,-4(r4) I8: slt r1,r4,zero I9: bne r1,zero,Label </pre>	<p>I4 stalls, but we can do I8 instead.</p> <p>I6 stalls, and all remaining subsequent instructions have dependences so this stall remains.</p>	$(9 + 1)/9 = 1.11$
-----------	--	---	--------------------

4.31.5 We use t0, t1, etc. as names for new registers in our renaming. We have:

	Instructions	Stalls	CPI
a.	<pre> I1: lw t1,-4(sp) I2: lw t2,-4(sp) I3: sll t3,t1,2 I4: add t4,t3,r4 I5: lw t5,0(t4) I6: add r3,t2,t5 I7: sw r3,-4(sp) I8: lw t6,-4(sp) I9: sll t7,t6,2 I10: add t8,t7,r4 I11: lw r2,0(t8) I12: bne r2,zero,Label </pre>	<p>I6 stalls, and all subsequent instructions have dependences. Note that I8 reads what I7 wrote to memory, so these instructions are still dependent.</p> <p>I9 stalls, but we can do I1 from the next iteration instead.</p> <p>I12 stalls, but we can do I2 from the next iteration instead.</p>	$(12 + 1)/12 = 1.08$
b.	<pre> I1: addi r4,r4,4 I2: lw t1,0(r4) I3: lw t2,4(r4) I4: add t3,t2,t1 I5: lw r3,8(r4) I6: add r2,t3,r3 I7: sw r2,-4(r4) I8: slt r1,r4,zero I9: bne r1,zero,Label </pre>	<p>This loop can now execute without stalls. I4 would stall, but we can do I5 instead. After I5 we execute I4, so I6 no longer stalls.</p>	$9/9 = 1$

4.31.6 Note that now every time we execute an instruction it can be renamed differently. We have:

	Instructions	Reordering	CPI
a.	<pre> I1: lw t1,-4(sp) I2: lw t2,-4(sp) I3: sll t3,t1,2 I4: add t4,t3,r4 I5: lw t5,0(t4) I6: add t6,t2,t5 I7: sw t6,-4(sp) I8: lw t7,-4(sp) I9: sll t8,t7,2 I10: add t9,t8,r4 I11: lw t10,0(t9) I12: bne t10,zero,Label </pre>	<p>I6 stalls, and all subsequent instructions have dependences. Note that I8 reads what I7 wrote to memory, so these instructions are still dependent.</p> <p>I9 would stall, but we can do I1 from the next iteration instead.</p> <p>I12 would stall, but we can do I2 from the next iteration instead.</p>	$(12 + 1)/12 = 1.08$

b.	<pre> I1: addi t1,t1,4 I2: lw t2,0(t1) I3: lw t3,4(t1) I4: add t4,t3,t2 I5: lw t5,8(t1) I6: add t6,t4,t5 I7: sw t6,-4(t1) I8: slt t7,t1,zero I9: bne t7,zero,Label </pre>	<p>No stalls remain. I4 would stall, but we can do I5 instead. After I5 we execute I4, so I6 no longer stalls.</p> <p>In next iteration uses of r4 renamed to t3.</p>	9/9 = 1
-----------	---	---	---------

Solution 4.32

4.32.1 The expected number of mispredictions per instruction is the probability that a given instruction is a branch that is mispredicted. The number of instructions between mispredictions is one divided by the number of mispredictions per instruction. We get:

	Mispredictions per Instruction	Instructions between Mispredictions
a.	$0.25 \times (1 - 0.95)$	80
b.	$0.25 \times (1 - 0.99)$	400

4.32.2 The number of in-progress instructions is equal to the pipeline depth times the issue width. The number of in-progress branches can then be easily computed because we know what percentage of all instructions are branches. We have:

	In-progress Branches
a.	$15 \times 4 \times 0.25 = 15$
b.	$30 \times 4 \times 0.25 = 30$

4.32.3 We keep fetching from the wrong path until the branch outcome is known, fetching 4 instructions per cycle. If the branch outcome is known in stage N of the pipeline, all instructions are from the wrong path in $N - 1$ stages. In the Nth stage, all instructions after the branch are from the wrong path. Assuming that the branch is just as likely to be the 1st, 2nd, 3rd, or 4th instruction fetched in its cycle, we have on average 1.5 instructions from the wrong path in the Nth stage (3 if branch is 1st, 2 if branch is 2nd, 1 if branch is 3rd, and 0 if branch is last). We have:

	Wrong-path Instructions
a.	$(12 - 1) \times 4 \times 1.5 = 45.5$
b.	$(20 - 1) \times 4 \times 1.5 = 77.5$

4.32.4 We can compute the CPI for each processor, then compute the speedup. To compute the CPI, we note that we have determined the number of useful instructions between branch mispredictions (for 4.32.1) and the number of mis-fetched instructions per branch misprediction (for 4.32.3), and we know how many instructions in total are fetched per cycle (4 or 8). From that we can determine the number of cycles between branch mispredictions, and then the CPI (cycles per useful instruction). We have:

	4-Issue		8-Issue			Speedup
	Cycles	CPI	Mis-Fetched	Cycles	CPI	
a.	$(45.5 + 80)/4 = 31.4$	$31.4/80 = 0.392$	$(12 - 1) \times 8 \times 3.5 = 91.5$	$(91.5 + 80)/8 = 21.4$	$21.4/80 = 0.268$	1.46
b.	$(77.5 + 400)/4 = 119.4$	$119.4/400 = 0.298$	$(20 - 1) \times 8 \times 3.5 = 155.5$	$(155.5 + 400)/8 = 69.4$	$69.4/400 = 0.174$	1.72

4.32.5 When branches are executed one cycle earlier, there is one less cycle needed to execute instructions between two branch mispredictions. We have:

	“Normal” CPI	“Improved” CPI	Speedup
a.	$31.4/80 = 0.392$	$30.4/80 = 0.380$	1.033
b.	$119.4/400 = 0.298$	$118.4/400 = 0.296$	1.008

4.32.6

	“Normal” CPI	“Improved” CPI	Speedup
a.	$21.4/80 = 0.268$	$20.4/80 = 0.255$	1.049
b.	$69.4/400 = 0.174$	$68.4/400 = 0.171$	1.015

Speedups from this improvement are larger for the 8-issue processor than with the 4-issue processor. This is because the 8-issue processor needs fewer cycles to execute the same number of instructions, so the same 1-cycle improvement represents a large relative improvement (speedup).

Solution 4.33

4.33.1 We need two register reads for each instruction issued per cycle:

	Read Ports
a.	$2 \times 2 = 4$
b.	$8 \times 2 = 16$

4.33.2 We compute the time-per-instruction as CPI times the clock cycle time. For the 1-issue 5-stage processor we have a CPI of 1 and a clock cycle time of T.

For an N-issue K-stage processor we have a CPI of $1/N$ and a clock cycle of $T \times 5/K$. Overall, we get a speedup of:

Speedup	
a.	$15/5 \times 2 = 6$
b.	$30/5 \times 8 = 48$

4.33.3 We are unable to benefit from a wider issue width (CPI is 1), so we have:

Speedup	
a.	$15/5 = 3$
b.	$30/5 = 6$

4.33.4 We first compute the number of instructions executed between mispredicted branches. Then we compute the number of cycles needed to execute these instructions if there were no misprediction stalls, and the number of stall cycles due to a misprediction. Note that the number of cycles spent on a misprediction is the number of entire cycles (one less than the stage in which branches are executed) and a fraction of the cycle in which the mispredicted branch instruction is. The fraction of a cycle is determined by averaging over all possibilities. In an N-issue processor, we can have the branch as the first instruction of the cycle, in which case we waste $(N - 1)$ Nths of a cycle, or the branch can be the second instruction in the cycle, in which case we waste $(N - 2)$ Nths of a cycle, ..., or the branch can be the last instruction in the cycle, in which case none of that cycle is wasted. With all of this data we can compute what percentage of all cycles are misprediction stall cycles:

	Instructions between Branch Mispredictions	Cycles between Branch Mispredictions	Stall Cycles	% Stalls
a.	$1/(0.10 \times 0.04) = 250$	$250/2 = 125$	8.3	$8/(125 + 8.3) = 6\%$
b.	$1/(0.10 \times 0.02) = 500$	$500/8 = 62.5$	4.4	$4/(62.5 + 4.4) = 6\%$

4.33.5 We have already computed the number of stall cycles due to a branch misprediction, and we know how to compute the number of non-stall cycles between mispredictions (this is where the misprediction rate has an effect). We have:

	Stall Cycles between Mispredictions	Need # of Instructions between Mispredictions	Allowed Branch Misprediction Rate
a.	8.3	$8.3 \times 2/0.05 = 330$	$1/(330 \times 0.10) = 3.03\%$
b.	4.4	$4.4 \times 8/0.01 = 3550$	$1/(3550 \times 0.10) = 0.28\%$

The needed accuracy is 100% minus the allowed misprediction rate.

4.33.6 This problem is very similar to 4.33.5, except that we are aiming to have as many stall cycles as we have non-stall cycles. We get:

	Stall Cycles between Mispredictions	Need # of Instructions between Mispredictions	Allowed Branch Misprediction Rate
a.	8.3	$8.3 \times 2 = 16.5$	$1/(16.5 \times 0.10) = 60.1\%$
b.	4.4	$4.4 \times 8 = 35.5$	$1/(35.5 \times 0.10) = 28.2\%$

The needed accuracy is 100% minus the allowed misprediction rate.

Solution 4.34

4.34.1 We need an IF pipeline stage to fetch the instruction. Since we will only execute one kind of instruction, we do not need to decode the instruction but we still need to read registers. As a result, we will need an ID pipeline stage although it would be misnamed. After that, we have an EXE stage, but this stage is simpler because we know exactly which operation should be executed so there is no need for an ALU that supports different operations. Also, we need no Mux to select which values to use in the operation because we know exactly which value it will be. We have:

a.	In the ID stage we read two registers and we do not need a sign-extend unit. In the EXE stage we need an “And” unit whose inputs are the two register values read in the ID stage. After the EXE stage we have a WB stage which writes the result from the And unit into Rd (again, no Mux). Note that there is no MEM stage, so this is a 4-stage pipeline. Also note that the PC is always incremented by 4, so we do not need the other Add and Mux units that compute the new PC for branches and jumps.
b.	We read two registers in the ID stage, and we also need the sign-extend unit for the Offs field in the instruction word. In the EXE stage we need an Add unit whose inputs are the Rs register value and the sign-extended offset from the ID stage. After the EXE stage we use the output of the Add unit as a memory address in the MEM stage, and the value we read from Rt is used as a data value for a memory write. Note that there is no WB stage, so this is a 4-stage pipeline. Also note that the PC is always incremented by 4, so we do not need the other Add and Mux units that compute the new PC for branches and jumps.

4.34.2

a.	Assuming that the register write in WB happens in the first half of the cycle and the register reads in ID happen in the second half, we only need to forward the And result from the EX/WB pipeline register to the inputs of the And unit in the EXE stage of the next instruction (if that next instruction depends on the previous one). No hazard detection unit is needed because forwarding eliminates all hazards.
b.	There is no need for forwarding or hazard detection in this pipeline because there are no RAW data dependences between two store instructions.

4.34.3 We need to add some decoding logic to our ID stage. The decoding logic must simply check whether the opcode and funct filed (if there is a funct field)

match this instruction. If there is no match, we must put the address of the exception handler into the PC (this adds a Mux before the PC) and flush (convert to NOPs) the undefined instruction (write zeros to the ID/EX pipeline register) and the following instruction which has already been fetched (write zeros to the IF/ID pipeline register).

4.34.4

a.	We need to replace the And unit in EXE with an ALU that supports either an Add or an And. The ALUOp signal to select between these operations must be supplied by the Control unit.
b.	The two operations are identical until the end of the EXE stage. After that, the ADDI operation must store the ALU output to the Rt register, so we must add the WB stage (SW did not need it). In fact, the work of the WB stage can be done in the MEM stage, so our pipeline remains a 4-stage pipeline. Our control logic must select whether we write the value of Rt to memory (for SW) or we write the ALU result to Rt (for ADDI).

4.34.5

a.	The same forwarding logic used for AND can be used for ADD, and we still need no hazard detection.
b.	Now we need forwarding because of ADDI instructions. Assuming that the register write in WB happens in the first half of the cycle and the register read in ID happens in the second half, we need to forward the Add result of an ADDI instruction from the EX/WB pipeline register to the first (register Rs) input of the Add unit in the EXE stage of the next instruction if that next instruction depends on the ADDI. We also need to forward that same Add result to replace the Rt value that will be stored into memory by the next SW instruction, if that instruction's Rt register is the same register as the Rt (result) register of the ADDI instruction. Fortunately, we still need no hazard detection.

4.34.6 The decoding logic must now check if the instruction matches either of the two instructions. After that, the exception handling is the same as for 4.34.3.

Solution 4.35

4.35.1 The worst case for control hazards is if the mispredicted branch instruction is the last one in its cycle and we have been fetching the maximum number of instructions in each cycle. Then the control hazard affects the remaining instructions in the branch's own pipeline stage and all instructions in stages between fetch and the branch execution stage. We have:

	Delay Slots Needed
a.	$10 \times 2 - 1 = 19$
b.	$15 \times 4 - 1 = 59$

4.35.2 If branches are executed in stage X, the number of stall cycles due to a misprediction is $(N - 1)$. These cycles are reduced by filling them with delay-slot instructions. We compute the number of execution (non-stall) cycles between mispredictions, and the speedup as follows:

	Non-stall Cycles between Mispredictions	Stall Cycles without Delay Slots	Stall Cycles with 4 Delay Slots	Speedup Due to Delay Slots
a.	$1/(0.25 \times (1 - 0.90) \times 2) = 20$	9	7	$(20 + 9)/(20 + 7) = 1.074$
b.	$1/(0.15 \times (1 - 0.96) \times 4) = 41.7$	14	13	$(41.7 + 14)/(41.7 + 13) = 1.018$

4.35.3 For 20% of the branches we add an extra instruction, for 30% of the branches we add two extra instructions, and for 40% of the branches we add three extra instructions. Overall, an average branch instruction is now accompanied by $0.20 + 0.30 \times 2 + 0.40 \times 3 = 2$ NOP instructions. Note that these NOPs are added for every branch, not just mispredicted ones. These NOP instructions add to the execution time of the program, so we have:

	Total Cycles between Mispredictions without Delay Slots	Stall Cycles with 4 Delay Slots	Extra Cycles Spent on NOPs	Speedup Due to Delay Slots
a.	$20 + 9 = 29$	7	$1 \times 20 \times 0.25 = 5$	$29/(20 + 7 + 5) = 0.906$
b.	$41.7 + 14 = 55.7$	13	$0.5 \times 41.7 \times 0.15 = 3.125$	$55.7/(41.7 + 13 + 3.125) = 0.963$

4.35.4

a.	<pre> add r2,zero,zero ; r1=0 Loop: beq r2,r3,End lb r10,1000(r2) ; Delay slot add r1,r1,r10 beq zero,zero,Loop addi r2,r2,1 ; Delay slot Exit: </pre>
b.	<pre> add r2,zero,zero ; r1=0 Loop: beq r2,r3,End lb r10,1000(r2) ; Delay slot lb r11,1001(r2) sub r12,r10,r11 add r1,r1,r12 beq zero,zero,Loop addi r2,r2,2 ; Delay slot Exit: </pre>

4.35.5

a.	<pre> add r2,zero,zero ; r1=0 Loop: beq r2,r3,End lb r10,1000(r2) ; Delay slot nop ; 2nd delay slot beq zero,zero,Loop add r1,r1,r10 ; Delay slot addi r2,r2,1 ; 2nd delay slot Exit: </pre>
b.	<pre> add r2,zero,zero ; r1=0 Loop: beq r2,r3,End lb r10,1000(r2) ; Delay slot lb r11,1001(r2) ; 2nd delay slot sub r12 r10,r11 beq zero,zero,Loop add r1,r1,r12 ; Delay slot addi r2,r2,2 ; 2nd delay slot Exit: </pre>

4.35.6 The maximum number of in-flight instructions is equal to the pipeline depth times the issue width. We have:

	Instructions in Flight	Instructions per Iteration	Iterations in Flight
a.	$15 \times 2 = 30$	5	$30/5 + 1 = 7$
b.	$25 \times 4 = 100$	7	$\text{roundUp}(100/7) + 1 = 16$

Note that an iteration is in flight when even one of its instructions is in flight. This is why we add one to the number we compute from the number of instructions in flight (instead of having an iteration entirely in flight, we can begin another one and still have the “trailing” one partially in flight) and round up.

Solution 4.36**4.36.1**

	Instruction	Translation
a.	SWINC Rt,Offset(Rs)	SW Rt,Offset(Rs) ADDI Rs,Rs,4
b.	SWI Rt,Rd(Rs)	ADD tmp,Rd,Rs SW Rt,0(tmp)

4.36.2 The ID stage of the pipeline would now have a lookup table and a micro-PC, where the opcode of the fetched instruction would be used to index into the lookup table. Micro-operations would then be placed into the ID/EX pipeline register, one per cycle, using the micro-PC to keep track of which micro-op is the next one to be output. In the cycle in which we are placing the last micro-op of an

instruction into the ID/EX register, we can allow the IF/ID register to accept the next instruction. Note that this results in executing up to one micro-op per cycle, but we are actually fetching instructions less often than that.

4.36.3

Instruction	
a.	We need to add an incrementer in the MEM stage. This incrementer would increment the value read from Rs while memory is being accessed. We also need to write this incremented value back into Rs.
b.	We can use the existing EX stage to perform this address calculation and then write to memory in the MEM stage. But we do need an additional (third) register read port because this instruction reads three registers in the ID stage, and we need to pass these three values to the EX stage.

4.36.4 Not often enough to justify the changes we need to make to the pipeline. Note that these changes slow down all the other instructions, so we are speeding up a relatively small fraction of the execution while slowing down everything else.

4.36.5 Each original ADDM instruction now results in executing two more instructions, and also adds a stall cycle (the ADD depends on the LW). As a result, each cycle in which we executed an ADDM instruction now adds three more cycles to the execution. We have:

Speedup from ADDM Translation	
a.	$1/(1 + 0.03 \times 3) = 0.92$
b.	$1/(1 + 0.05 \times 3) = 0.87$

4.36.6 Each translated ADDM adds the 3 stall cycles, but now half of the existing stalls are eliminated. We have:

Speedup from ADDM Translation	
a.	$1/(1 + 0.03 \times 3 - 0.12/2) = 0.97$
b.	$1/(1 + 0.05 \times 3 - 0.20/2) = 0.95$

Solution 4.37

4.37.1 All of the instructions use the instruction memory, the PC + 4 adder, the control unit (to decode the instruction), and the ALU. For the least utilized unit, we have:

a.	The result of the branch adder (add offset to PC + 4) is never used.
b.	The read port of the data memory is never used (no load instructions).

Note that the branch adder performs its operation in every cycle, but its result is actually used only when a branch is taken.

4.37.2 The read port is only used by LW and the write port by SW instructions. We have:

	Data Memory Read	Data Memory Write
a.	20% (1 out of 5)	20% (1 out of 5)
b.	0% (no LW)	25% (1 out of 4)

4.37.3 In the IF/ID pipeline register, we need 32 bits for the instruction word and 32 bits for PC + 4 for a total of 64 bits. In the ID/EX register, we need 32 bits for each of the two register values, the sign-extended offset/immediate value, and PC + 4 (for exception handling). We also need 5 bits for each of the three register fields from the instruction word (Rs,Rt,Rd), and 10 bits for all the control signals output by the Control unit. The total for the ID/EX register is 153 bits. In the EX/MEM register, we need 32 bits each for the value of register Rt and for the ALU result. We also need 5 bits for the number of the destination register and 4 bits for control signals. The total for the EX/MEM register is 73 bits. Finally, for the MEM/WB register we need 32 bits each for the ALU result and value from memory, 5 bits for the number of the destination register, and 2 bits for control signals. The total for MEM/WB is 71 bits. The grand total for all pipeline registers is 361 bits.

4.37.4 In the IF stage, the critical path is the I-Mem latency. In the ID stage, the critical path is the latency to read Regs. In the EXE stage, we have a Mux and then ALU latency. In the MEM stage we have the D-Mem latency, and in the WB stage we have a Mux latency and setup time to write Regs (which we assume is zero). For a single-cycle design, the clock cycle time is the sum of these per-stage latencies (for a load instruction). For a pipelined design, the clock cycle time is the longest of the per-stage latencies. To compare these clock cycle times, we compute a speedup based on clock cycle time alone (assuming the number of clock cycles is the same in single-cycle and pipelined designs, which is not true). We have:

	IF	ID	EX	MEM	WB	Single-Cycle	Pipelined	“Speedup”
a.	200ps	90ps	110ps	250ps	20ps	670ps	250ps	2.68
b.	750ps	300ps	300ps	500ps	50ps	1900ps	750ps	2.53

Note that this speedup is significantly lower than 5, which is the “ideal” speedup of 5-stage pipelining.

4.37.5 If we only support ADD instructions, we do not need the MUX in the WB stage, and we do not need the entire MEM stage. We still need Muxes before the ALU for forwarding. We have:

	IF	ID	EX	WB	Single-Cycle	Pipelined	“Speedup”
a.	200ps	90ps	110ps	0ps	400ps	200ps	2.00
b.	750ps	300ps	300ps	0ps	1135ps	750ps	1.80

Note that the “ideal” speedup from pipelining is now 4 (we removed the MEM stage), and the actual speedup is about half of that.

4.37.6 For the single-cycle design, we can reduce the clock cycle time by 1ps by reducing the latency of any component on the critical path by 1ps (if there is only one critical path). For a pipelined design, we must reduce latencies of all stages that have longer latencies than the target latency. We have:

	Single-Cycle	Needed Cycle Time for Pipelined	Cost for Pipelined
a.	$0.2 \times 670 = \$134$	$0.8 \times 250\text{ps} = 200\text{ps}$	\$50 (MEM)
b.	$0.2 \times 1900 = \$380$	$0.8 \times 750\text{ps} = 600\text{ps}$	\$150 (IF)

Note that the cost of improving the pipelined design by 20% is lower. This is because its clock cycle time is already lower, so a 20% improvement represents fewer picoseconds (and fewer dollars in our problem).

Solution 4.38

4.38.1 The energy for the two designs is the same: I-Mem is read, two registers are read, and a register is written. We have:

a.	$140\text{pJ} + 2 \times 70\text{ps} + 60\text{pJ} = 340\text{pJ}$
b.	$70\text{pJ} + 2 \times 40\text{pJ} + 40\text{pJ} = 190\text{pJ}$

4.38.2 The instruction memory is read for all instructions. Every instruction also results in two register reads (even if only one of those values is actually used). A load instruction results in a memory read and a register write, a store instruction results in a memory write, and all other instructions result in either no register write (e.g., BEQ) or a register write. Because the sum of memory read and register write energy is larger than memory write energy, the worst-case instruction is a load instruction. For the energy spent by a load, we have:

a.	$140\text{pJ} + 2 \times 70\text{pJ} + 60\text{pJ} + 140\text{pJ} = 480\text{pJ}$
b.	$70\text{pJ} + 2 \times 40\text{pJ} + 40\text{pJ} + 90\text{pJ} = 280\text{pJ}$

4.38.3 Instruction memory must be read for every instruction. However, we can avoid reading registers whose values are not going to be used. To do this, we must

add RegRead1 and RegRead2 control inputs to the Registers unit to enable or disable each register read. We must generate these control signals quickly to avoid lengthening the clock cycle time. With these new control signals, an LW instruction results in only one register read (we still must read the register used to generate the address), so we have:

	Energy before Change	Energy Saved by Change	% Savings
a.	$140\text{pJ} + 2 \times 70\text{pJ} + 60\text{pJ} + 140\text{pJ} = 480\text{pJ}$	70pJ	14.6%
b.	$70\text{pJ} + 2 \times 40\text{pJ} + 40\text{pJ} + 90\text{pJ} = 280\text{pJ}$	40pJ	14.3%

4.38.4 Before the change, the Control unit decodes the instruction while register reads are happening. After the change, the latencies of Control and Register Read cannot be overlapped. This increases the latency of the ID stage and could affect the processor's clock cycle time if the ID stage becomes the longest-latency stage. We have:

	Clock Cycle Time before Change	Clock Cycle Time after Change
a.	250ps (D-Mem in MEM stage)	No change ($150\text{ps} + 90\text{ps} < 250\text{ps}$)
b.	750ps (I-Mem in IF stage)	800ps (Ctl then Regs in ID stage)

4.38.5 If memory is read in every cycle, the value is either needed (for a load instruction), or it does not get past the WB Mux (or a non-load instruction that writes to a register), or it does not get written to any register (all other instructions, including stalls). This change does not affect clock cycle time because the clock cycle time must already allow enough time for memory to be read in the MEM stage. It does affect energy: a memory read occurs in every cycle instead of only in cycles when a load instruction is in the MEM stage.

4.38.6

	I-Mem Active Energy	I-Mem Latency	Clock Cycle Time	Total I-Mem Energy	Idle Energy %
a.	140pJ	200ps	250ps	$140\text{pJ} + 50\text{ps} \times 0.1 \times 140\text{pJ} / 200\text{ps} = 143.5\text{pJ}$	$3.5\text{pJ} / 143.5\text{pJ} = 2.44\%$
b.	70pJ	750ps	750ps	70pJ	0%

Solution 4.39

4.39.1 The number of instructions executed per second is equal to the number of instructions executed per cycle (IPC, which is $1/\text{CPI}$) times the number of cycles per second (clock frequency, which is $1/T$ where T is the clock cycle time). The IPC

is the percentage of cycle in which we complete an instruction (and not a stall), and the clock cycle time is the latency of the maximum-latency pipeline stage. We have:

	IPC	Clock Cycle Time	Clock Frequency	Instructions per Second
a.	0.75	350ps	2.86GHz	2.14×10^9
b.	0.80	220ps	4.55GHz	3.64×10^9

4.39.2 Power is equal to the product of energy per cycle times the clock frequency (cycles per second). The energy per cycle is the total of the energy expenditures in all five stages. We have:

	Clock Frequency	Energy per Cycle (in pJ)	Power (W)
a.	2.86GHz	$100 + 45 + 50 + 0.30 \times 150 + 0.45 \times 50 = 262.5$	0.75
b.	4.55GHz	$75 + 45 + 100 + 0.45 \times 100 + 0.50 \times 35 = 282.5$	1.28

4.39.3 The time that remains in the clock cycle after a circuit completes its work is often called slack. We determine the clock cycle time and then the slack for each pipeline stage:

	Clock Cycle Time	IF Slack	ID Slack	EX Slack	MEM Slack	WB Slack
a.	350ps	100ps	0ps	200ps	50ps	150ps
b.	220ps	20ps	50ps	0ps	10ps	70ps

4.39.4 All stages now have latencies equal to the clock cycle time. For each stage, we can compute the factor X for it by dividing the new latency (clock cycle time) by the original latency. We then compute the new per-cycle energy consumption for each stage by dividing its energy by its factor X. Finally, we re-compute the power dissipation:

	X for IF	X for ID	X for EX	X for MEM	X for WB	New Power (W)
a.	350/250	350/350	350/150	350/300	350/200	0.54
b.	220/200	220/170	220/220	220/210	220/150	1.17

4.39.5 This changes the clock cycle time to 1.1 of the original, which changes the factor X for each stage and the clock frequency. After that this problem is solved in the same way as 4.39.4. We get:

	X for IF	X for ID	X for EX	X for MEM	X for WB	New Power (W)
a.	385/250	385/350	385/150	385/300	385/200	0.45
b.	242/200	242/170	242/220	242/210	242/150	0.97

4.39.6 The X factor for each stage is the same as in 4.39.6, but this time in our power computation we divide the per-cycle energy of each stage by X^2 instead of x. We get:

	New Power (W)	Old Power (W)	Saved
a.	0.31	0.75	58.7%
b.	0.81	1.28	36.7%