

Exercises for  
EITF20 Computer Architecture  
HT1 2010

Anders Ardö  
Department of Electrical and Information Technology, EIT  
Lund University

December 4, 2013



**LUND INSTITUTE OF TECHNOLOGY**  
Lund University

# Contents

<b>1</b>	<b>Performance</b>	<b>3</b>
<b>2</b>	<b>ISA</b>	<b>7</b>
<b>3</b>	<b>Pipelining</b>	<b>9</b>
<b>4</b>	<b>Memory systems, Cache</b>	<b>19</b>
<b>5</b>	<b>Memory systems, Virtual Memory</b>	<b>24</b>
<b>6</b>	<b>Storage systems, I/O</b>	<b>26</b>
<b>7</b>	<b>Home Assignment - online quiz</b>	<b>29</b>
<b>8</b>	<b>Answers</b>	<b>30</b>
8.1	Performance . . . . .	30
8.2	ISA . . . . .	32
8.3	Pipelining I . . . . .	35
8.4	Memory systems, Cache I . . . . .	45
8.5	Memory systems, Virtual Memory . . . . .	51
8.6	Storage systems, I/O . . . . .	54

# 1 Performance

**Exercise 1.1** Describe the following items concerning computer architecture:

- a) “make the common case fast”
- b) “locality of reference”
- c) The “SPEC” benchmark series
- d) Amdahl’s Law

**Exercise 1.2** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 1.5

Component type	Product	Performance	Power
Processor	Sun Niagara 8-core	1.2 GHz	72-79W peak
	Intel Pentium 4	2 GHz	48.9-66W
DRAM	Kingston X64C3AD2 1 GB	184-pin	3.7W
	Kingston D2N3 1 GB	240-pin	2.3W
Hard drive	DiamondMax 16	5400 rpm	7.0W read/seek, 2.9 W idle
	DiamondMax Plus 9	7200 rpm	7.9W read/seek, 4.0 W idle

**Figure 1.23** Power consumption of several computer components.

[10/10/20] <1.6, 1.7> One critical factor in powering a server farm is cooling. If heat is not removed from the computer efficiently, the fans will blow hot air back onto the computer, not cold air. We will look at how different design decisions affect the necessary cooling, and thus the price, of a system. Use Figure 1.23 for your power calculations.

- a. [10] <1.6> A cooling door for a rack costs \$4000 and dissipates 14 KW (into the room; additional cost is required to get it out of the room). How many servers with a Sun Niagara 8-core processor, 1 GB 240-pin DRAM, and a single 5400 rpm hard drive can you cool with one cooling door?
- b. [10] <1.6, 1.8> You are considering providing fault tolerance for your hard drive. RAID 1 doubles the number of disks (see Chapter 6). Now how many systems can you place on a single rack with a single cooler?
- c. [20] <1.8> In a single rack, the MTTF of each processor is 4500 hours, of the hard drive is 9 million hours, and of the power supply is 30K hours. For a rack with 8 processors, what is the MTTF for the rack?

**Exercise 1.3** Amdahl observed: “a fairly obvious conclusion which can be drawn at this point is that the effort expended on achieving high parallel processing rates is wasted unless it is accompanied by achievements in sequential processing rates of very nearly the same magnitude”. This was then formulated as Amdahl’s Rule. This rule can be applied in various areas.

1. A common transformation required in graphics processors is square root. Implementations of floating-point (FP) square root vary significantly in performance. Suppose FP square root (FPSQR) is responsible for 20 % of the execution time of a critical graphics benchmark. One proposal is to enhance the FPSQR hardware and speed this up by a factor of 10. The alternative is just to make all FP instructions run faster; FP instructions are responsible for half of the execution time. By how much do the FP instructions have to be accelerated to achieve the same performance as achieved by inserting the specialized hardware?

2. Suppose that we want to enhance the processor used for Web serving. The new processor is 10 times faster on computation in the Web serving application than the original processor. Assuming that the original processor is busy with computation 40 % of the time and is waiting for I/O 60 % of the time, what is the overall speedup gained by incorporating the enhancement?

**Exercise 1.4** Use Amdahl's law to illustrate why it is important to keep a computer system balanced in terms of relative performance between for example I/O speed and raw CPU speed.

**Exercise 1.5** Three enhancements with the following speed-ups are proposed for a new architecture:

- Enhancement A: Speed-up = 30;
- Enhancement B: Speed-up = 20;
- Enhancement C: Speed-up = 15.

Only one enhancement is usable at a time.

1. How can Amdahl's Law be formulated to handle multiple enhancements?
2. If enhancements A and B are usable for 25% of the time, what fraction of the time must enhancement C be used to achieve an overall speed-up of 10?

Assume the enhancements can be used 25%, 35% and 10% of the time for enhancements A, B, and C respectively.

3. For what fraction of the reduced execution time is no enhancement is in use?

Assume, for some benchmark, the possible fraction of use is 15% for each of the enhancements A and B and 70% for enhancement C. We want to maximize performance.

4. If only one enhancement can be implemented, which should it be? If two enhancements can be implemented, which should be chosen?

**Exercise 1.6** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 1.12

[10/10/Discussion/10/20/Discussion] <1.7> Make the following calculations on the raw data in order to explore how different measures color the conclusions one can make. (Doing these exercises will be much easier using a spreadsheet.)

- a. [10] <1.8> Create a table similar to that shown in Figure 1.26, except express the results as normalized to the Pentium D for each benchmark.
- b. [10] <1.9> Calculate the arithmetic mean of the performance of each processor. Use both the original performance and your normalized performance calculated in part (a).
- c. [Discussion] <1.9> Given your answer from part (b), can you draw any conflicting conclusions about the relative performance of the different processors?
- d. [10] <1.9> Calculate the geometric mean of the normalized performance of the dual processors and the geometric mean of the normalized performance of the single processors for the Dhrystone benchmark.
- e. [20] <1.9> Plot a 2D scatter plot with the *x*-axis being Dhrystone and the *y*-axis being the memory benchmark.
- f. [Discussion] <1.9> Given your plot in part (e), in what area does a dual-processor gain in performance? Explain, given your knowledge of parallel processing and architecture, why these results are as they are.

Chip	# of cores	Clock frequency (MHz)	Memory performance	Dhrystone performance
Athlon 64 X2 4800+	2	2,400	3,423	20,718
Pentium EE 840	2	2,200	3,228	18,893
Pentium D 820	2	3,000	3,000	15,220
Athlon 64 X2 3800+	2	3,200	2,941	17,129
Pentium 4	1	2,800	2,731	7,621
Athlon 64 3000+	1	1,800	2,953	7,628
Pentium 4 570	1	2,800	3,501	11,210
Processor X	1	3,000	7,000	5,000

**Figure 1.26** Performance of several processors on two benchmarks.

**Exercise 1.7** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 1.13

[10/10/20] <1.9> Imagine that your company is trying to decide between a single-processor system and a dual-processor system. Figure 1.26 gives the performance on two sets of benchmarks—a memory benchmark and a processor benchmark. You know that your application will spend 40% of its time on memory-centric computations, and 60% of its time on processor-centric computations.

- a. [10] <1.9> Calculate the weighted execution time of the benchmarks.
- b. [10] <1.9> How much speedup do you anticipate getting if you move from using a Pentium 4 570 to an Athlon 64 X2 4800+ on a CPU-intensive application suite?
- c. [20] <1.9> At what ratio of memory to processor computation would the performance of the Pentium 4 570 be equal to the Pentium D 820?

**Exercise 1.8** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 1.14

[10/10/20/20] <1.10> Your company has just bought a new dual Pentium processor, and you have been tasked with optimizing your software for this processor. You will run two applications on this dual Pentium, but the resource requirements are not equal. The first application needs 80% of the resources, and the other only 20% of the resources.

- a. [10] <1.10> Given that 40% of the first application is parallelizable, how much speedup would you achieve with that application if run in isolation?
- b. [10] <1.10> Given that 99% of the second application is parallelizable, how much speedup would this application observe if run in isolation?
- c. [20] <1.10> Given that 40% of the first application is parallelizable, how much *overall system speedup* would you observe if you parallelized it?
- d. [20] <1.10> Given that 99% of the second application is parallelizable, how much overall system speedup would you get?

## 2 ISA

**Exercise 2.1** What is a load-store architecture and what are the advantages/disadvantages of such an architecture compared to other GPR (General Purpose Register) architectures?

**Exercise 2.2** What are the advantages/disadvantages of fixed-length and variable-length instruction encodings?

**Exercise 2.3** Assume an instruction set that uses a fixed 16-bit instruction length. Operand specifiers are 6 bits in length. There are 5 two operand instructions and 33 zero operand instructions. What is the maximum number of one-operand instructions that can be encoded using the fixed 16-bit instruction length?

**Exercise 2.4** Consider this high-level language code sequence of three statements:

A = B + C;

B = A + C;

D = A - B;

Use the technique of copy propagation (see figure B.20) to transform the code sequence to the point where no operand is a computed value. Note the instances in which the transformation has reduced the computational work of a statement and those cases where the work has increased.

What does this suggest about the technical challenge faced in trying to satisfy the desire for optimizing compilers?

**Exercise 2.5** A given processor has 32 registers, uses 16-bit immediates and has 142 instructions in its ISA. In a given program,

- 20 % of the instructions take 1 input register and have 1 output register.,
- 30 % have 2 input registers and 1 output register,
- 25 % have 1 input register, 1 output register and take an immediate input as well,
- and the remaining 25 % have one immediate input and 1 output register.

1. For each of the 4 types of instructions, how many bits are required? Assume that the ISA requires that all instructions be a multiple of 8 bits in length.

2. How much less memory does the program take up if variable-length instruction set encoding is used as opposed to fixed-length encoding?

**Exercise 2.6** Compute the effective CPI for MIPS using figure B.27. Suppose we have made the following measurements of average CPI for instructions:

Instruction	Average CPI
All ALU instructions	1.0
Load/Store	1.4
Conditional branches	
taken	2.0
not taken	1.5
Jumps	1.2

Assume that 60 % of the conditional branches are taken and that all instructions in the 'other' category in Fig B.27 are ALU instructions. Average the instructions frequencies of gcc and gcc to obtain the instruction mix.

**Exercise 2.7** Your task is to compare the memory efficiency of the following instruction set architectures:

- Accumulator -- All operations occur between a single register and a memory location. There are two accumulators of which one is selected by the opcode;
- Memory-memory – All instruction addresses reference only memory locations
- Stack – All operations occur on top of the stack. Push and pop are the only instructions that access memory; all others remove their operands from the stack and replace them with the result. The implementation uses a hardwired stack for only the top two stack entries, which keeps the processor circuit very small and low cost. Additional stack positions are kept in memory locations, and accesses to these stack positions require memory references.
- Load-store -- All operations occur in registers, and register-to register instructions have three register names per instruction.

To measure memory efficiency, make the following assumptions about all 4 instruction sets:

- All instructions are an integral number of bytes in length;
  - The opcode is always 1 byte (8 bits);
  - Memory accesses use direct, or absolute addressing.
  - The variables A, B, C, and D are initially in memory
- a. Invent your own assembly language mnemonics (fig B.2 provides a useful sample to generalize), and for each architecture write the best equivalent assembly language code for this high level language code:  
A = B + C;  
B = A + C;  
D = A - B;
- b. Assume the given code sequence is from a small, embedded computer application, such as a microwave oven controller that uses 16-bit memory addresses and data operands. If a load-store architecture is used, assume that it has 16 general-purpose registers. For each architecture of your choice answer the following questions:
1. How many instruction bytes are fetched?
  2. How many bytes of data are transferred from/to memory?
  3. Which architecture is the most efficient as measured in code size?
  4. Which architecture is most efficient as measured by total memory traffic (code + data)?



### 3 Pipelining

**Exercise 3.1** Consider following assembly-language program:

```

1: MOV   R3, R7
2: LD    R8, (R3)
3: ADD   R3, R3, 4
4: LOAD  R9, (R3)
5: BNE   R8, R9, L3

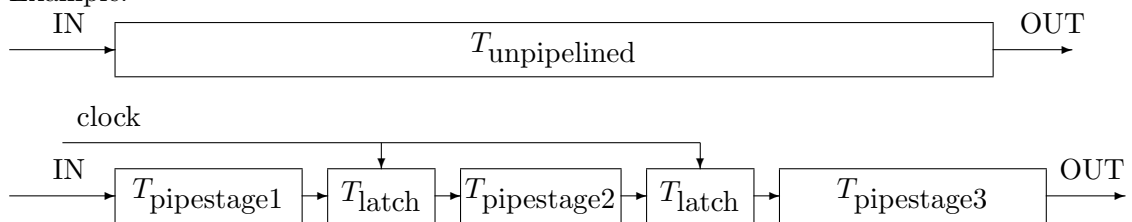
```

- This program includes WAW, RAW, and WAR dependencies. Show these?
- What is the difference between a dependency and hazard?
- What is the difference between a name dependency and a true dependency?
- Which of WAW, RAW, WAR are true dependencies and which are name dependencies?

**Exercise 3.2** a) Define an expression (Speedup= ) for pipeline speedup as a function of:

- $T_{\text{unpipelined}}$ : execution time for the non-pipelined unit.
- $\max(T_{\text{pipestage}})$ : the execution time of the slowest pipe-stage.
- $T_{\text{latch}}$ : overhead time (setup time) for the storage element between pipe-stages.

Example:



- Define an expression for pipeline speedup as a function of (use only these):
  - noofstages: number of pipe stages (assume equal pipe stage execution time).
  - branch\_freq (bf): relative frequency of branches in the program.
  - branch\_penalty (bp): number of clock cycles lost due to a branch.
- Give an example of a piece of assembly code that contains WAW, RAW and WAR hazards and identify them. (Use for example the assembly instruction `ADDD Rx, Ry, Rz` which stores  $(Ry+Rz)$  in  $Rx$ )

**Exercise 3.3** Use the following code fragment:

```

Loop: LD    F0, 0(R2)
      LD    F4, 0(R3)
      MULD  F0, F0, F4
      ADDD  F2, F0, F2
      DADDUI R2, R2, #8
      DADDUI R3, R3, #8
      DSUBU R5, R4, R2
      BNEZ  R5, Loop

```

Assume that the initial values of R2 is 0 and R4 is 792.

For this exercise assume the standard five-stage integer pipeline and the MIPS FP pipeline as describe in section A.5. If structural hazards are due to write-back contention, assume the earliest instruction gets priority and other instructions are stalled.

- a) Show the timing of this instruction sequence for the MIPS FP pipeline without any forwarding or bypassing hardware but assuming a register read and a write in the same clock cycle “forwards” through the register file. Assume that the branch is handled by flushing the pipeline. If all memory references hit the cache, how many cycles does this loop take to execute?
- b) Show the timing of this instruction sequence for the MIPS FP pipeline with normal forwarding and bypassing hardware. Assume that the branch is handled by predicting it as not taken. If all memory references hit in the cache, how many cycles does this loop take to execute?

**Exercise 3.4** Suppose the branch frequencies (as percentage of all instructions) are as follows:

Conditional branches	15 %
Jumps and calls	1 %
Conditional branches	60 % are taken

We are examining a four-deep pipeline where the branch is resolved at the end of the second cycle for unconditional branches and at the end of the third cycle for conditional branches. Assuming that only the first pipe stage can always be done independent of whether the branch goes and ignoring other pipeline stalls, how much faster would the machine be without any branch hazards?

**Exercise 3.5** A reduced hardware implementation of the classic five stage RISC pipeline might use the EX stage hardware to perform a branch instruction comparison and then not actually deliver the branch target PC to the IF stage until the clock cycle in which the branch instruction reaches the MEM stage. Control hazard stalls can be reduced by resolving branch instructions in ID, but improving performance in one respect may reduce performance in other circumstances. How does determining branch outcome in the ID stage have the potential to increase data hazard stall cycles?

**Exercise 3.6** a) Show with an example using two assembler instructions only that exceptions may be generated in a pipeline in another order than the execution order. Assume a 5 stage pipeline with stages IF, ID, EXE, MEM, and WB.

- b) What is meant with “precise exceptions” and why is it important?

**Exercise 3.7** Consider an unpipelined processor. Assume that it has 1-ns clock cycle and that it uses 4 cycles for ALU operations and 5 cycles for branches and 4 cycles for memory operations. Assume that the relative frequencies of these operations are 50 %, 35 % and 15 % respectively. Suppose that due to clock skew and set up, pipelining the processor adds 0.15 ns of overhead to the clock. Ignoring any latency impact, how much speed up in the instruction execution rate will we gain from a pipeline?

**Exercise 3.8** Explain what a control hazard is and how to avoid it.

**Exercise 3.9** • Identify all data dependencies in the following code, assuming that we are using the 5-stage MIPS pipelined datapath. Which dependencies can be resolved via forwarding?

```
ADD R2,R5,R4
ADD R4,R2,R5
SW  R5,100(R2)
ADD R3,R2,R4
```

- Consider executing the following code on the 5-stage pipelined datapath: Which registers are read during the fifth clock cycle, and which registers are written at the end of the fifth clock cycle? Consider only the registers in the register file (i.e., R1, R2, R3, etc.)

```
ADD R1,R2,R3
ADD R4,R5,R6
ADD R7,R8,R9
ADD R10,R11,R12
ADD R13,R14,R15
```

**Exercise 3.10** Briefly give two ways in which loop unrolling can increase performance and one in which it can decrease performance.

**Exercise 3.11** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 2.2

		<b>Latencies beyond single cycle</b>
Loop:	LD F2,0(Rx)	Memory LD +3
I0:	MULTD F2,F0,F2	Memory SD +1
I1:	DIVD F8,F2,F0	Integer ADD, SUB +0
I2:	LD F4,0(Ry)	Branches +1
I3:	ADDD F4,F0,F4	ADDD +2
I4:	ADDD F10,F8,F2	MULTD +4
I5:	SD F4,0(Ry)	DIVD +10
I6:	ADDI Rx,Rx,#8	
I7:	ADDI Ry,Ry,#8	
I8:	SUB R20,R4,Rx	
I9:	BNZ R20,Loop	

**Figure 2.35** Code and latencies for Exercises 2.1 through 2.6.

[10] <1.8, 2.1, 2.2> Think about what latency numbers really mean—they indicate the number of cycles a given function requires to produce its output, nothing more. If the overall pipeline stalls for the latency cycles of each functional unit, then you are at least guaranteed that any pair of back-to-back instructions (a “producer” followed by a “consumer”) will execute correctly. But not all instruction pairs have a producer/consumer relationship. Sometimes two adjacent instructions have nothing to do with each other. How many cycles would the loop body in the code sequence in Figure 2.35 require if the pipeline detected true data dependences and only stalled on those, rather than blindly stalling everything just because one functional unit is busy? Show the code with <stall> inserted where necessary to accommodate stated latencies. (*Hint:* An instruction with latency “+2” needs 2 <stall> cycles to be inserted into the code sequence. Think of it this way: a 1-cycle instruction has latency 1 + 0, meaning zero extra wait states. So latency 1 + 1 implies 1 stall cycle; latency 1 +  $N$  has  $N$  extra stall cycles.)

**Exercise 3.12** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 2.7

[15] <2.1> Computers spend most of their time in loops, so multiple loop iterations are great places to speculatively find more work to keep CPU resources busy. Nothing is ever easy, though; the compiler emitted only one copy of that loop's code, so even though multiple iterations are handling distinct data, they will appear to use the same registers. To keep register usages multiple iterations from colliding, we rename their registers. Figure 2.36 shows example code that we would like our hardware to rename.

A compiler could have simply unrolled the loop and used different registers to avoid conflicts, but if we expect our hardware to unroll the loop, it must also do the register renaming. How? Assume your hardware has a pool of temporary registers (call them T registers, and assume there are 64 of them, T0 through T63) that it can substitute for those registers designated by the compiler. This rename hardware is indexed by the source register designation, and the value in the table is the T register of the last destination that targeted that register. (Think of these table values as producers, and the src registers are the consumers; it doesn't much matter where the producer puts its result as long as its consumers can find it.) Consider the code sequence in Figure 2.36. Every time you see a destination register in the code, substitute the next available T, beginning with T9. Then update all the src registers accordingly, so that true data dependences are maintained. Show the resulting code. (*Hint:* See Figure 2.37.)

---

```
Loop: LD      F2,0(Rx)
I0:  MULTD   F5,F0,F2
I1:  DIVD    F8,F0,F2
I2:  LD      F4,0(Ry)
I3:  ADDD    F6,F0,F4
I4:  ADDD    F10,F8,F2
I5:  SD      F4,0(Ry)
```

---

**Figure 2.36** Sample code for register renaming practice.

---

```
I0:  LD      T9,0(Rx)
I1:  MULTD   T10,F0,T9
. . .
```

---

**Figure 2.37** Expected output of register renaming.

**Exercise 3.13** Consider the following assembly program:

```

0      ADD      r3,r31,r2
1      LW       r6,0(r3)
2      ANDI    r7,r5,#3
3      ADD      r1,r6,r0
4      SRL     r7,r0,#8
5      OR      r2,r4,r7
6      SUB     r5,r3,r4
7      ADD     r15,r1,r10
8      LW      r6,0(r5)
9      SUB     r2,r1,r6
10     ANDI    r3,r7,#15

```

Assume the use of a four-stage pipeline: fetch (IF), decode/issue (DI), execute (EX) and write back (WB). Assume that all pipeline stages take one clock cycle except for the execute stage. For simple integer arithmetic and logical instructions, the execute stage takes one cycle, but for a load from memory, five cycles are needed in the execute stage. Suppose we have a simple scalar pipeline but allow some sort of out-of-order execution that results in the following table

	<b>Instruction</b>	<b>IF</b>	<b>DI</b>	<b>EX</b>	<b>WB</b>
	0 ADD r3,r31,r2	0	1	2	3
	1 LW r6,0(r3)	1	2	4 <sup>a</sup>	9
	2 ANDI r7,r5,#3	2	3	5 <sup>b</sup>	6
for the first seven instructions:	3 ADD r1,r6,r0	3	4	10	11
	4 SRL r7,r0,#8	4	5	6	7
	5 OR r2,r4,r7	5	6	8	10
	6 SUB r5,r3,r4	6	7	9	12 <sup>c</sup>

A number in the table indicates the clock cycle a certain instruction starts at a pipeline stage. There are a lot of implementation details that can be deduced from the execution table above.

- Explain why the first lw-instruction (instruction 1) cannot start in the execute stage until clock cycle 4.
- Explain why the first and-instruction (instruction 2) cannot start the execution stage until clock cycle 5.
- Explain why the first sub-instruction (instruction 6) cannot start the write back stage until clock cycle 12.
- Complete the table for the remaining instructions.
- Suppose instruction 2 was changed to: ANDI r6,r5,#3. What implications would that have on the design of the pipeline? How would the table look like?

**Exercise 3.14** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 2.11

[10/10/10] <2.3> Assume a five-stage single-pipeline microarchitecture (fetch, decode, execute, memory, write back) and the code in Figure 2.41. All ops are 1 cycle except LW and SW, which are 1 + 2 cycles, and branches, which are 1 + 1 cycles. There is no forwarding. Show the phases of each instruction per clock cycle for one iteration of the loop.

- [10] <2.3> How many clock cycles per loop iteration are lost to branch overhead?
- [10] <2.3> Assume a static branch predictor, capable of recognizing a backwards branch in the decode stage. Now how many clock cycles are wasted or branch overhead?
- [10] <2.3> Assume a dynamic branch predictor. How many cycles are lost or a correct prediction?

**Exercise 3.15** Schedule the following code using Tomasulo's algorithm assuming the hardware has three Load-units with a two-cycle execution latency, three Add/Sub units with 2 cycles execution latency, and two Mult/Div units where Mult has an execution latency of 10 cycles and Div 40 cycles. Assume the first instruction (LD F6,34(R2)) is issued in cycle 1.

```
LD    F6,34(R2)
LD    F2,34(R3)
MULTD F0,F2,F4
SUBD  F8,F6,F2
DIVD  F10,F0,F6
ADDD  F6,F8,F2
```

1. In which clock cycle (numbered 0,1,2,...) does the second LD instruction complete?
2. In which clock cycle does the MULTD instruction complete?
3. In which clock cycle does the ADDD instruction complete?

**Exercise 3.16** Assume a processor with a standard five-stage pipeline (IF, ID, EX, MEM, WB) and a branch prediction unit (a branch history table) in the ID-stage. Branch resolution is performed in the EX-stage. There are four cases for conditional branches:

- The branch is not taken and correctly predicted as not taken (NT/PNT)
- The branch is not taken and predicted as taken (NT/PT)
- The branch is taken and predicted as not taken (T/PNT)
- The branch is taken and correctly predicted as taken (T/PT)

Suppose that the branch penalties with this design are:

- NT/PNT: 0 cycles

- T/PT: 1 cycle
  - NT/PT, T/PNT: 2 cycles
- a) Describe how the CPU Performance Equation can be modified to take the performance of the branch prediction unit into account. Define the information you need to know to assess the performance.
  - b) Use the answer in the assignment above to calculate the average CPI for the processor assuming a base CPI of 1.2. Assume 20% conditional branches (disregard from other branches) and that 65% of these are taken on average. Assume further that the branch prediction unit mispredicts 12% of the conditional branches.
  - c) In order to increase the clock frequency from 500 MHz to 600 MHz, a designer splits the IF-stage into two stages, IF1 and IF2. This makes it easier for the instruction cache to deliver instructions in time. This also affects the branch penalties for the branch prediction unit as follows:
    - NT/PNT: 0 cycles
    - T/PT: 2 cycles
    - NT/PT, T/PNT: 3 cycles

How much faster is this new processor than the previous that runs on 500 MHz?

- d) Propose a solution to reduce the average branch penalty even further.

**Exercise 3.17** A processor with dynamic scheduling and issue bound operand fetch (operands are fetched during instruction issue) has 3 execution units – one LOAD/STORE unit, one ADD/SUB unit and one MUL/DIV unit. It has a reservation station with 1 slot per execution unit and a single register file. Starting with the following instruction sequence in the instruction fetch buffer and empty reservation stations, for each instruction find the cycle in which it will be issued and the cycle in which it will write result.

```

LOAD R6, 34(R12)
LOAD R2, 45(R13)
MUL R0, R2, R4
SUB R8, R2, R6
DIV R10, R0, R6
ADD R6, R8, R2

```

Assumptions: out of order issue, out of order execution, 4 stage pipeline (instruction fetch, decode and issue, execute, and write back), no data forwarding, and sufficient instruction window size.

Execute cycles taken by different instructions are:

LOAD/STORE: 2

ADD/SUB: 1

MUL: 2

DIV: 4

**Exercise 3.18** This assignment focuses on processors with dynamic scheduling, speculative execution using a reorder buffer, and dynamic branch prediction.



```

LDI  R1, 0          # R1 = 0
LDI  R2, 8000       # R2 = 8000
LDI  R3, 0          # R3 = 0
LOOP:
LD   R4, 0(R1)
DMUL R5, R4, R4
DADD R5, R3, R5
SD   R5, 0(R1)      # new[i] = old[i-1] + old[i]*old[i]
DADDI R3, R4, 0
DADDI R1, R1, 8
BNE  R1, R2, LOOP

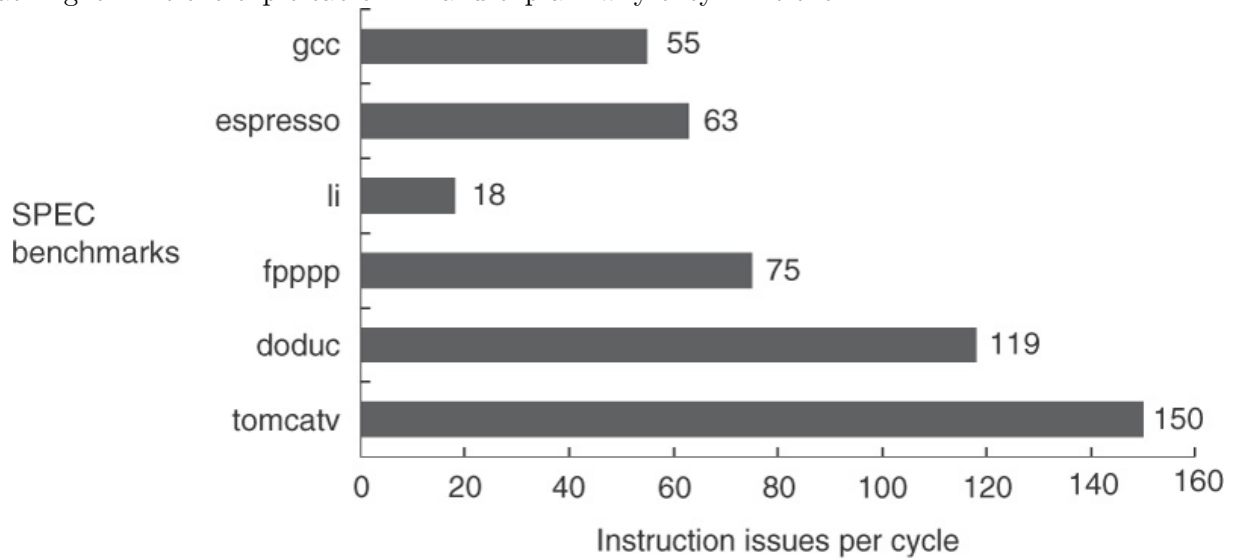
```

1. Explain how dynamic branch prediction works when using a 2-bit prediction scheme. Use the program above to demonstrate how the prediction works. Also, describe how the prediction works for a program containing two or more branches?
2. In what way does a reorder buffer help speculation? What are the key points when introducing support for speculation?
3. In a processor with dynamic scheduling according to Tomasulo's algorithm, hardware based speculation using a reorder buffer (ROB), and dynamic branch prediction, execution of each instruction follows a sequence of steps that is somewhat different depending on the type of instruction. Below are the steps carried out for an ALU instruction:
  - (a) Issue when reservation station and ROB entry is available
    - Read already available operands from registers and instruction
    - Send instruction to reservation station
    - Tag unavailable operands with ROB entry
    - Tag destination register with ROB entry
    - Write destination register to ROB entry
    - Mark ROB entry as busy
  2. Execute after issue
    - Wait for operand values on CDB (if not already available)
    - Compute result
  3. Write result when CDB and ROB available
    - Send result on CDB to reservation stations
    - Update ROB entry with result, and mark as ready
    - Free reservation station
  4. Commit when at head of ROB and ready
    - Update destination register with result from ROB entry
    - Untag destination register
    - Free ROB entry

What are the corresponding steps for handling a store instruction?

4. Explain how RAW hazards are resolved in the basic Tomasulo's algorithm.

**Exercise 3.19** As Figure 3.1, from the book Computer Architecture, A Quantitative Approach, below shows, the ILP available in many applications can be fairly high. Nevertheless, in practical processor implementations it can be hard to exploit this available ILP. List at least three factors that might limit the exploitable ILP and explain why they limit the ILP.



© 2007 Elsevier, Inc. All rights reserved.

## 4 Memory systems, Cache

**Exercise 4.1** Describe the concept “memory hierarchy”, and state why it is important. State the function of each part, normally used hardware components, and what problems they solve (if any).

**Exercise 4.2** Suppose that CPI for a given architecture (with a perfect memory system, using 32 bit addresses) is 1,5. We are considering the following cache systems:

- A 16 KB direct mapped “unified” cache using “write-back”. Miss-ratio = 2.9%. Does not affect the cycle length.
- A 16 KB 2-way set-associative “unified” cache using “write-back”. Miss-ratio = 2.2%. Increases the cycle length with a factor 1.2
- A 32 KB direct mapped “unified” cache using “write-back”. Miss-ratio = 2.0%. Increases the cycle length with a factor 1.25

Suppose a memory latency of 40 cycles, 4 bytes transferred per cycle and that 50% of the blocks are “dirty”. There are 32 bytes per block and 20% of the instructions are “data transfer” instructions. A “write buffer” is not used.

- a) Calculate effective CPI for the three cache systems.
- b) Recalculate the CPI using a system with a TLB with a 0.2 % miss-rate and a 20 cycle penalty. The caches are physically addressed.
- c) Which of the above cache-systems is best?
- d) Draw a block diagram of the set-associative cache. Show different address fields and describe how they are used to determine hit/miss.
- e) How does a TLB affect performance if the cache is virtually or physically addressed?

**Exercise 4.3** What are the design trade-offs between a large register file and a large data cache?

**Exercise 4.4** Suppose that it was possible to design a unified (common to instructions and data) first-level cache with two ports so that there would be no structural hazard in the the pipeline. What are the design trade-offs between using the unified cache compared to separate instruction and data cache? The total amount of cache memory should of course be the same in both cases.

**Exercise 4.5** What is a write-through cache? Is it faster/slower than a write-back cache with respect to the time it takes for writing.

**Exercise 4.6** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 5.1

Size	Direct	2-way LRU	4-way LRU	8-way LRU	Full LRU
1 KB	0.0863842--	0.0697167--	0.0634309--	0.0563450--	0.0533706--
2 KB	0.0571524--	0.0423833--	0.0360463--	0.0330364--	0.0305213--
4 KB	0.0370053--	0.0260286--	0.0222981--	0.0202763--	0.0190243--
8 KB	0.0247760--	0.0155691--	0.0129609--	0.0107753--	0.0083886--
16 KB	0.0159470--	0.0085658--	0.0063527--	0.0056438--	0.0050068--
32 KB	0.0110603--	0.0056101--	0.0039190--	0.0034628--	0.0030885--
64 KB	0.0066425--	0.0036625--	0.0009874--	0.0002666--	0.0000106--
128 KB	0.0035823--	0.0002341--	0.0000109--	0.0000058--	0.0000058--
256 KB	0.0026345--	0.0000092--	0.0000049--	0.0000051--	0.0000053--
512 KB	0.0014791--	0.0000065--	0.0000029--	0.0000029--	0.0000029--
1 MB	0.0000090--	0.0000058--	0.0000028--	0.0000028--	0.0000028--

**Figure 5.29** SPEC2000 data miss ratios (misses per 1000 instructions) [Cantin and Hill 2003].

- 5.1 [12/12/15/15] <5.2> The following questions investigate the impact of small and simple caches using CACTI, and assume a 90 nm (0.09  $\mu\text{m}$ ) technology.
- [12] <5.2> Compare the access times of 32 KB caches with 64-byte blocks and a single bank. What is the relative access times of two-way and four-way set associative caches in comparison to a direct-mapped organization?
  - [12] <5.2> Compare the access times of two-way set-associative caches with 64-byte blocks and a single bank. What is the relative access times of 32 KB and 64 KB caches in comparison to a 16 KB cache?
  - [15] <5.2> Does the access time for a typical level 1 cache organization increase with size roughly as the capacity in bytes  $B$ , the square root of  $B$ , or the log of  $B$ ?
  - [15] <5.2> Find the cache organization with the lowest average memory access time given the miss ratio table in Figure 5.29 and a cache access time budget of 0.90 ns. What is this organization, and does it have the lowest miss rate of all organizations for its capacity?

**Exercise 4.7** Let's try to show how you can make unfair benchmarks. Here are two machines with the same processor and main memory but different cache organizations. Assume that both processors run at 2 GHz, have a CPI of 1, and have a cache (read) miss time of 100 ns. Further, assume that writing a 32-bit word in main memory requires 100 ns (for the write-through cache). The caches are unified – they contain both instructions and data -, and each cache has a total capacity of 64 kB, not including tags and status bits. The cache on system A is two-way set associative and has 32-byte blocks. It is write-through and does not allocate a block on a write miss. The cache on system B is direct mapped and has 32-byte blocks. It is write-back and allocates a block on a write miss.

- Describe a program that makes system A run as fast as possible relative to system B's speed. How much faster is this program on system A as compared to system B.
- Describe a program that makes system B run as fast as possible relative to system A's speed. How much faster is this program on system B as compared to system A.

**Exercise 4.8** Assume having a two level memory hierarchy: a cache and a main memory, which are connected with a 32 bit wide bus. A hit in the cache can be executed within one clock cycle. At a cache miss an entire block must be replaced. This is done by sending the address (32 bits) to the memory which needs 4 clock cycles before it can send back a block by the bus. Every bus-transfer requires one clock cycle. The processor will need to wait until the entire block is in the cache. The following table shows the average miss ratio for different block sizes:

Block size (B, bytes)	Miss-ratio (M), %
4	4.5
16	2.4
32	1.6
64	1.0
128	0.64

- Which block size results in the best average memory-access time?
- If bus arbitration requires 2 clock cycles in average, which block size is the most optimal?
- Which block size is the best one if the bus between the cache and the main memory is widened to 64 bits (both with and without the cost of the bus-arbitration of assignment a and b)?

Assume a cache with a size of 256 bytes and a block size of 16 bytes. The blocks in the cache are numbered from 0 and upwards. Specify the sizes of the tag, index and byte fields when using 16 bits addressing, and in which block (or blocks) in the cache the address  $28E7_{16}$  is represented, for:

- A direct mapped cache.
- A 2-way set-associative cache.

**Exercise 4.9** a) What is a memory hierarchy? Describe the problem a memory hierarchy solves and explain why it works so well.

- A system has a 256 byte large cache. It is a set-associative cache with 16 sets. An address is 16 bits and the tag field in the address is 13 bits.
  - Specify the number of blocks in the cache.
  - Specify the size of the byte offset field.
  - Draw a sketch of how the cache is organized and describe what the different address fields are used for.
  - What tag would a byte with the address  $3333_{16}$  end up in, and which byte within the block is it?
- Explain and discuss two different ways to handle writes in the cache.

**Exercise 4.10** What is a replacement algorithm? Why is such an algorithm needed with cache memories? With which of the following strategies is a replacement algorithm needed:

- Direct mapping
- Set-associative mapping
- Fully associative mapping

**Exercise 4.11** Draw a schematics of how the following cache memory can be implemented: Total size 16kB, 4-way set-associative. Blocksize 16 byte, replacement algorithm LRU, uses write-back.

The schematics should among others show central MUXes, comparators and connections. It should clearly indicate how a physical memory address is translated to a cache position.

The greater detail shown the better.

**Exercise 4.12** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 5.4

[12/15] <5.2> Inspired by the usage of critical word first and early restart on level 1 cache misses, consider their use on level 2 cache misses. Assume a 1 MB L2 cache with 64-byte blocks and a refill path that is 16 bytes wide. Assume the L2 can be written with 16 bytes every 4 processor cycles, the time to receive the first 16-byte block from the memory controller is 100 cycles, each additional 16 B from main memory requires 16 cycles and data can be bypassed directly into the read port of the L2 cache. Ignore any cycles to transfer the miss request to the level 2 cache and the requested data to the level 1 cache.

- [12] <5.2> How many cycles would it take to service a level 2 cache miss with and without critical word first and early restart?
- [15] <5.2> Do you think critical word first and early restart would be more important for level 1 caches or level 2 caches, and what factors would contribute to their relative importance?

**Exercise 4.13** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 5.5

[10/12] <5.2> You are designing a write buffer between a write-through level 1 cache and a write-back level 2 cache. The level 2 cache write data bus is 16 bytes wide and can perform a write to an independent cache address every 4 processor cycles.

- [10] <5.2> How many bytes wide should each write buffer entry be?
- [12] <5.2> What speedup could be expected in the steady state by using a merging write buffer instead of a nonmerging buffer when zeroing memory by the execution of 32-bit stores if all other instructions could be issued in parallel with the stores and the blocks are present in the level 2 cache?

**Exercise 4.14** Three ways with hardware and/or software to decrease the time a program spends on (data) memory accesses are:

- nonblocking caches with “hit under miss”
- hardware prefetching with 4–8 stream buffers
- software prefetching with nonfaulting cache prefetch

Explain, for each of these methods, how it affects:

- a) miss rate
- b) memory bandwidth to the underlying memory
- c) number of executed instructions

**Exercise 4.15** In systems with a write-through L1 cache backed by a write-back L2 cache instead of main memory, a merging write buffer can be simplified.

- a) Explain how this can be done.
- b) Are there situations where having a full write buffer (instead of the simple version you have just proposed) could be helpful?

**Exercise 4.16** Shortly describe the three C’s model.

**Exercise 4.17** Explain where replacement policy fits into the three C’s model, and explain why this means that misses caused by a replacement policy are “ignored”- or more precisely cannot in general be definitively classified - by the three C’s model.

## 5 Memory systems, Virtual Memory

**Exercise 5.1** As caches increase in size, blocks often increase in size as well.

- a) If a large instruction cache has larger blocks, is there still a need for pre-fetching? Explain the interaction between pre-fetching and increased block size in instruction caches.
- b) Is there a need for data pre-fetch instructions when data blocks get larger?

**Exercise 5.2** Some memory systems handle TLB misses in software (as an exception), while others use hardware for TLB misses.

- a) What are the trade-offs between these methods for handling TLB misses?
- b) Will TLB miss handling in software always be slower than TLB misses in hardware? Explain!
- c) Are there page table structures that would be difficult to handle in hardware, but possible in software? Are there any such structures that would be difficult for software to handle but easy for hardware to manage?

**Exercise 5.3** The difficulty of building a memory system to keep pace with faster CPUs is underscored by the fact that the raw material for main memory is the same as that found in the cheapest computer. The performance difference is rather based on the arrangement.

- a) List the four measures in answer to questions on block placement, identification, replacement and writing strategy that rule the hierarchical construction for virtual memory.
- b) What is the main purpose of the Translation-Lookaside Buffer within the memory hierarchy? Give an appropriate set of construction rules and explain why.
- c) Fill the following table with characteristic (typical) entries:

	TLB	1st-level cache	2nd-level cache	Virtual memory
Block size (in bytes)				
Block placement				
Overall size				

**Exercise 5.4** Designing caches for out-of-order (OOO) superscalar CPUs is difficult for several reasons. Clearly, the cache will need to be non-blocking and may need to cope with several outstanding misses. However, the access pattern for OOO superscalar processors differs from that generated by in-order execution.

What are the differences, and how might they affect cache design for OOO processors?

**Exercise 5.5** Consider the following three hypothetical, but not atypical, processors, which we run with the SPEC gcc benchmark

1. A simple MIPS two-issue static pipe running at a clock rate of 4 GHz and achieving a pipeline CPI of 0.8. This processor has a cache system that yields 0.005 misses per instruction.
2. A deeply pipelined version of a two-issue MIPS processor with slightly smaller caches and a 5 GHz clock rate. The pipeline CPI of the processor is 1.0, and the smaller caches yield 0.0055 misses per instruction on average.



3. A speculative MIPS, superscalar with a 64-entry window. It achieves one-half of the ideal issue rate measured for this window size (9 instruction issues per cycle). This processor has the smallest caches, which leads to 0.01 misses per instruction, but hides 25 scheduling. This processor has a 2.5 GHz clock.

Assume that the main memory time (which sets the miss penalty) is 50 ns. Determine the relative performance of these three processors.

**Exercise 5.6** a) Give three arguments for larger pages in virtual memory, and one against.

- b) Describe the concepts 'page', 'page fault', 'virtual address', 'physical address', 'TLB', and 'memory mapping' and how they are related.
- c) How much memory does the page table, indexed by the virtual page number, take for a system using 32 bit virtual addresses, 4 KB pages, and 4 bytes per page table entry. The system has 512 MB of physical memory.
- d) In order to save memory sometimes inverted page tables are used. Briefly describe how they are structured. How much memory would inverted page tables take for the above system.

**Exercise 5.7** A) Describe two cache memory optimization techniques that may improve hit performance (latency and throughput). For each technique, specify how it affects hit time and fetch bandwidth.

- B) Describe two cache memory optimization techniques that may reduce miss rate, and define the miss type (compulsory, capacity, conflict) that is primarily affected by each technique.
- C) Describe two cache memory optimization techniques that may reduce miss penalty.

## 6 Storage systems, I/O

For Hennessy/Patterson exercises 6.8 - 6.14.

For these exercises, assume that you have built a RAID system with six disks, plus a sufficient number of hot spares. Assume each disk is the 37 GB SCSI disk shown in Figure 6.3; assume each disk can sequentially read data at a peak of 142 MB/sec and sequentially write data at a peak of 85 MB/sec. Assume that the disks are connected to an Ultra320 SCSI bus that can transfer a total of 320 MB/sec. You can assume that each disk failure is independent and ignore other potential failures in the system. For the reconstruction process, you can assume that the overhead for any XOR computation or memory copying is negligible. During online reconstruction, assume that the reconstruction process is limited to use a total bandwidth of 10 MB/sec from the RAID system.

**Exercise 6.1** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.8

- 6.8 [10] <6.2> Assume that you have a RAID 4 system with six disks. Draw a simple diagram showing the layout of blocks across disks for this RAID system.

**Exercise 6.2** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.9

- 6.9 [10] <6.2, 6.4> When a single disk fails, the RAID 4 system will perform reconstruction. What is the expected time until a reconstruction is needed?

**Exercise 6.3** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.10

- 6.10 [10/10/10] <6.2, 6.4> Assume that reconstruction of the RAID 4 array begins at time  $t$ .
- [10] <6.2, 6.4> What read and write operations are required to perform the reconstruction?
  - [10] <6.2, 6.4> For offline reconstruction, when will the reconstruction process be complete?
  - [10] <6.2, 6.4> For online reconstruction, when will the reconstruction process be complete?

**Exercise 6.4** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.14

- 6.14 [10] <6.2, 6.4> RAID 6 is used to tolerate up to two simultaneous disk failures. Assume that you have a RAID 6 system based on row-diagonal parity, or RAID-DP; your six-disk RAID-DP system is based on RAID 4, with  $p = 5$ , as shown in Figure 6.5. If data disk 0 and data disk 3 fail, how can those disks be reconstructed? Show the sequence of steps that are required to compute the missing blocks in the first four stripes.

For Hennessy/Patterson exercises 6.19 - 6.22.

Here are your building blocks:

- A 10,000 MIPS CPU costing \$1000. Its MTTF is 1,000,000 hours.
- A 1000 MB/sec I/O bus with room for 20 Ultra320 SCSI buses and controllers.
- Ultra320 SCSI buses that can transfer 320 MB/sec and support up to 15 disks per bus (these are also called SCSI strings). The SCSI cable MTTF is 1,000,000 hours.
- An Ultra320 SCSI controller that is capable of 50,000 IOPS, costs \$250, and has an MTTF of 500,000 hours.
- A \$2000 enclosure supplying power and cooling to up to eight disks. The enclosure MTTF is 1,000,000 hours, the fan MTTF is 200,000 hours, and the power supply MTTF is 200,000 hours.
- The SCSI disks described in Figure 6.3.
- Replacing any failed component requires 24 hours.

You may make the following assumptions about your workload:

- The operating system requires 70,000 CPU instructions for each disk I/O.
- The workload consists of many concurrent, random I/Os, with an average size of 16 KB.

All of your constructed systems must have the following properties:

- You have a monetary budget of \$28,000.
- You must provide at least 1 TB of capacity.

**Exercise 6.5** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.19

- 6.19 [10] <6.2> You will begin by designing an I/O subsystem that is optimized only for capacity and performance (and not reliability), specifically IOPS. Discuss the RAID level and block size that will deliver the best performance.

**Exercise 6.6** Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.21

- 6.21 [10] <6.2, 6.4, 6.7> You will now redesign your system to optimize for reliability, by creating a RAID 10 or RAID 01 array. Your storage system should be robust not only to disk failures, but to controller, cable, power supply, and fan failures as well; specifically, a single component failure should not prohibit accessing both replicas of a pair. Draw a diagram illustrating how blocks are allocated across disks in the RAID 10 and RAID 01 configurations. Is RAID 10 or RAID 01 more appropriate in this environment?

Exercise 6.7 Hennessy/Patterson, Computer Architecture, 4th ed., exercise 6.22

- 6.22 [20/20/20/20/20] <6.2, 6.4, 6.7> Optimizing your RAID 10 or RAID 01 array only for reliability (but keeping within your capacity and monetary constraints), what is your RAID configuration?
- [20] <6.2, 6.4, 6.7> What is the overall MTTF of the components in your system?
  - [20] <6.2, 6.4, 6.7> What is the MTDL of your system?
  - [20] <6.2, 6.4, 6.7> What is the usable capacity of this system?
  - [20] <6.2, 6.4, 6.7> How much does your system cost?
  - [20] <6.2, 6.4, 6.7> Assuming a write-only workload, how many IOPS can you expect to deliver?

## 7 Home Assignment - online quiz

### OPTIONAL!

**However - An approved quiz will give you 2 extra points on the exam.**

Take the quiz available for Computer Architecture EITF20, at <http://moodle.eit.lth.se/>  
It will be open during weeks 6 and 7. You have to log in to be able to see it.

You log in with your STIL username/password.

If you have a problem contact the course coordinator, Anders Ardö or the course assistant.  
You can take the quiz any number of times during the time mentioned above.

When you have logged in, choose 'Computer Architecture EITF20' and click on the quiz 'Test Quiz ...'.

Then you can start answering questions. After all questions are answered you can send in your answers by clicking on 'Submit all and finish'. You will get a feedback saying how many correct answers you have. Both questions and numeric values in the quiz are selected randomly each time you try the quiz. Redo the test until you have at least 90 % correct in order to be approved.

## 8 Answers

### 8.1 Performance

1.1 a) Invest in resources that are used often!

b) A computer uses data (and instructions) that are often close in the address space (spacial locality), and also close in time (temporal locality). Makes caches feasible.

c) SPEC = “Standard Performance Evaluation Corporation”, a series of typical integer and floating point programs used to characterize the performance of a computer. Exists from several years: SPEC89, SPEC92, SPEC95, SPEC2000, ...

d) “The performance improvement to be gained from using some faster mode of execution is limited by the fraction of the time the faster mode can be used”.

Suppose that enhancement E accelerates a fraction F of the task by a factor S, and the remainder of the task is unaffected, then  $T_{exe}(withE) = T_{exe}(withoutE) * [(1 - F) + F/S]$ .

1.2 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

1.3 1. For the specialized hardware we find as speedup

$$FPSQR : S = 1/[(1-0.2) + 0.2/10] = 1/0.82 = 1.22$$

To achieve the same performance we have to make the floating-point instructions about 60% faster:

$$FP : S = 1/[(1-0.5) + 0.5/1.6] = 1/0.8125 = 1.23$$

1.3 With a speedup of 10 for a fraction of 0.4, the overall speedup becomes

$$S = 1/[0.6 + 0.4/10] = 1/0.64 = 1.56$$

1.4 Amdahl's law: system speedup limited by the slowest component:

- Assume 10% I/O
- CPU speedup = 10  $\implies$  System speedup = 5
- CPU speedup = 100  $\implies$  System speedup = 10

*I/O will more and more become a bottleneck!*

1.5 Amdahl's Law can be generalized to handle multiple enhancements. If only one enhancement can be used at a time during program execution, then for enhancements A,B,C,...,i

$$Speedup = \left[ 1 - \sum_i FE_i + \sum_i \frac{FE_i}{SE_i} \right]^{-1}$$

where  $FE_i$  is the fraction of time that enhancement  $i$  can be used and  $SE_i$  is the speedup of enhancement  $i$ . For a single enhancement the equation reduces to the familiar form of Amdahl's Law.

**1.5 1.** With three enhancements we have

$$Speedup = \left[ 1 - (FE_A + FE_B + FE_C) + \left( \frac{FE_A}{SE_A} + \frac{FE_B}{SE_B} + \frac{FE_C}{SE_C} \right) \right]^{-1}$$

**1.5 2.** Substituting in the known quantities gives

$$10 = \left[ 1 - (0.25 + 0.25 + FE_C) + \left( \frac{0.25}{30} + \frac{0.25}{20} + \frac{FE_C}{15} \right) \right]^{-1}$$

Solving yields  $FE_C = 0.45$ . Thus, the third enhancement (C) must be usable 45% of the time.

**1.5 3.** Let  $T_e$  and  $TNE_e$  denote execution time with enhancement and the time during enhanced execution in which no enhancements are in use, respectively. Let  $T_{original}$  and  $FNE_{original}$  stand for execution time without enhancements and the fraction of that time that cannot be enhanced. Finally, let  $FNE_e$  represent the fraction of the reduced (enhanced) execution time for which no enhancement is in use. By definition

$$FNE_e = TNE_e / T_e$$

. Because the time spent executing code that cannot be enhanced is the same whether enhancements are in use or not, and by Amdahl's Law, we have

$$\frac{TNE_e}{T_e} = \frac{FNE_{original} * T_{original}}{T_{original} / Speedup}$$

Canceling factors and substituting equivalent expressions for  $FNE_{original}$  and Speedup yields

$$\frac{FNE_{original} * T_{original}}{T_{original} / Speedup} = \frac{1 - \sum_i FE_i}{1 - \sum_i FE_i + \sum_i \frac{FE_i}{SE_i}}$$

Substituting with known quantities,

$$FNE_e = \frac{1 - (0.25 + 0.35 + 0.10)}{1 - (0.25 + 0.35 + 0.10) + \left( \frac{0.25}{30} + \frac{0.35}{20} + \frac{0.10}{15} \right)} = \frac{0.3}{0.3325} = 90\%$$

**1.5 4.** Let the speedup when implementing only enhancement  $i$  be  $Speedup_i$ , and let  $Speedup_{ij}$  denote the speedup when employing enhancement  $i$  and  $j$ .

$$Speedup_A = \left( 1 - 0.15 + \frac{0.15}{30} \right)^{-1} = 1.17$$

$$Speedup_B = \left( 1 - 0.15 + \frac{0.15}{20} \right)^{-1} = 1.17$$

$$Speedup_C = \left( 1 - 0.70 + \frac{0.70}{15} \right)^{-1} = 2.88$$

Thus, if only one enhancement can be implemented, enhancement 3 offers much greater speedup.

$$Speedup_{AB} = \left( 1 - (0.15 + 0.15) + \left( \frac{0.15}{30} + \frac{0.15}{20} \right) \right)^{-1} = 1.40$$

$$Speedup_{AC} = \left( 1 - (0.15 + 0.70) + \left( \frac{0.15}{30} + \frac{0.70}{15} \right) \right)^{-1} = 4.96$$

$$Speedup_{BC} = \left( 1 - (0.15 + 0.70) + \left( \frac{0.15}{20} + \frac{0.70}{15} \right) \right)^{-1} = 4.90$$

Thus, if only a pair of enhancements can be implemented, enhancements A and C offer the greatest speedup.

Selecting the fastest enhancement(s) may not yield the highest speedup. As Amdahl's Law states, an enhancement contributes to speedup only for the fraction of time it can be used.

**1.6** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**1.7** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**1.8** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

## 8.2 ISA

**2.1** A load-store architecture is one in which only load and store instructions can access the memory system. In other GPR architectures, some or all of the other instructions may read their operands from or write their results to the memory system.

The primary advantage of non-load-store architectures is the reduced number of instructions required to implement a program and lower pressure on the register file.

The advantage of load-store architectures is that limiting the set of instructions that can access the memory system makes the micro-architecture simpler, which often allows implementation at a higher clock rate.

Depending on whether the clock rate increase or the decrease in number of instructions is more significant, either approach can result in greater performance.

**2.2** Variable-length instruction encodings reduce the amount of memory that programs take up, since each instruction takes only as much space as it requires. Instructions in a fixed-length encoding all take up as much storage as the longest instruction in the ISA, meaning that there is some number of wasted bits in the encoding of instructions that take fewer operands, don't allow immediate constants, etc.

However, variable-length instruction sets require more complex instruction decode logic than fixed-length instructions sets, and they make it harder to calculate the address of the next instruction in memory. Therefore, processors with fixed-length instruction sets can often be implemented at higher clock rate than processors with variable-length instruction sets.

**2.3** Use 2 bits for instruction type code (4 different types), two-operand instructions need to use two type-codes (we need to encode 5 instructions and they use 12 bits for operand specifiers leaving us with  $16 - 2 - 12 = 2$  bits for instruction code), one-operand 1 type code, and zero-operand 1 type code. One-operand instructions need a 6 bit operand specifier  $\Rightarrow 16 - 2 - 6 = 8$  bits  $\Rightarrow 256$  one-operand instructions

**2.4** Take the code sequence one line at a time.



1. $A = B + C ;$	The operands here are given, not computed by the code, so copy propagation will not transform this statement.
2. $B = A + C ;$	Here A is a computed value, so transform the code by substituting $A = B + C$ to get
$B = B + C + C ;$	No operand is computed
3. $D = A - B ;$	Both operands are computed so substitute for both to get
$D = (B + C) - (B + C + C) ;$	Simplify algebraically to get
$D = - C ;$	This is a given, not computed, operand

Copy propagation has increased the work in statement 2 from one addition to two. It has changed the work in statement 3 from subtraction to negation, possibly a savings.

The above suggests that writing optimizing compilers means incorporating sophisticated trade-off analysis capability to control any optimizing steps, if the best results are to be achieved.

**2.5** 1. 142 instructions  $\Rightarrow$  **8** bits ( $128 < 142 < 256$ ); 32 registers  $\Rightarrow$  **5** bits; **16** bit immediates

- 1 reg in, 1 reg out:  $8 + 5 + 5 = 18$  bits  $\Rightarrow$  24 bits
- 2 reg in, 1 reg out:  $8 + 5 + 5 + 5 = 23$  bits  $\Rightarrow$  24 bits
- 1 reg in, 1 reg out, 1 imm:  $8 + 5 + 5 + 16 = 34$  bits  $\Rightarrow$  40 bits
- 1 imm in, 1 reg out:  $8 + 16 + 5 = 29$  bits  $\Rightarrow$  32 bits

2. Since the largest instruction type requires 40-bit instructions, the fixed-length encoding will have 40 bits per instruction. Each instruction type in the variable encoding will use the number of bits from 1:  $0.2 * 24 + 0.3 * 24 + 0.25 * 40 + 0.25 * 32 = 30$  bits on average, ie 25 % less space.

**2.6** The first challenge of this exercise is to obtain the instruction mix. The instruction frequencies in Figure B.27 must add to 100.0, although gap and gcc add to 100.2 and 99.5 %, respectively, because of rounding errors. Because each total must in reality be 100.0, we should not attempt to scale the per instruction average frequencies by the shown totals of 100.2 and 99.5. However, in computing the average frequencies to one significant digit to the right of the decimal point, we should be careful to use an unbiased rounding scheme so that the total of the averaged frequencies is kept as close to 100 as possible. One such scheme is called round to even, which makes the least significant digit always even. For example, 0.15 rounds to 0.2, but 0.25 also rounds to 0.2. For a summation of terms, round to even will not accumulate an error as would, for example, the scheme rounding up where 0.15 rounds to 0.2 and 0.25 rounds to 0.3.

The gap and gcc the average instruction frequencies are shown below.

Instruction	Average frequency gap, gcc	Category
load	25.8	load/store
store	11.8	load/store
add	20.0	ALU
sub	2.0	ALU
mul	0.8	ALU
compare	4.4	ALU
load imm	3.6	ALU
cond branch	10.7	Cond branch
cond move	0.5	ALU
jump	0.8	jump
call	1.1	jump
return	1.1	jump
shift	2.4	ALU
and	4.4	ALU
or	8.2	ALU
xor	2.0	ALU
other logical	0.2	ALU

The exercise statement gives CPI information in terms of four major instruction categories, with two subcategories for conditional branches. To compute the average CPI we need to aggregate the instruction frequencies to match these categories. This is the second challenge, because it is easy to miscategorize instructions. The four main categories are ALU, load/store, conditional branch, and jumps. ALU instructions are any that take operands from the set of registers and return a result to that set of registers. Load/store instructions access memory. Conditional branch instructions must be able to set the program counter to a new value based on a condition. Jump-type instructions set the program counter to a new value no matter what.

With the above category definitions, the frequency of ALU instructions is the sum of the frequencies of the add, sub, mul, compare, load imm (remember, this instruction does not access memory, instead the value to load is encoded in a field within the instruction itself), cond move (implemented as an OR instruction between a controlling register and the register with the data to move to the destination register), shift, and, or, xor, and other logical for a total of 48.5 %. The frequency of load/store instructions is the sum of the frequencies of load and store for a total of 37.6 %. The frequency of conditional branches is 10.7 %. Finally, the frequency of jumps is the sum of the frequencies of the jump-type instructions, namely jump, call, and return, for a total of 3.0 %.

$$Effective\_CPI = \sum_{categories} Instruction\_category\_frequency * Clock\_cycles\_for\_category$$

$$= 0.485 * 1.0 + 0.367 * 1.4 + 0.107 * (0.6 * 2.0 + (1 - 0.6) * 1.5) + 0.03 * 1.2 = 1.24$$

## 2.7 a. Assembly programs:

Stack	Acc	Load-store	Mem-mem
Push B	Load B	Load R2,B	Add A,B,C
Push C	Add C	Load R3,C	Add B,A,C
Add	Store A	Add R1,R2,R3	Sub D,A,B
Push Top	Add C	Add R2,R1,R3	
Push Top	Store B	Sub R4,R1,R2	
Pop A	Sub A	Store R4,D	
Push C	Store D		
Add			
Push Top			
Pop B			
Sub			
Pop D			

b. Code sizes

Stack: 1 byte for opcode plus 2 bytes for memory addresses when needed

Accumulator: 1 byte for opcode plus 2 bytes for memory addresses

Load-store: 1 byte for opcode, half a byte for subject register, plus either 2 bytes for memory addresses or 1 byte for operand registers (round to whole bytes).

Mem-mem: 1 byte for opcode plus 6 bytes for memory addresses

	Stack	Acc	Load-store	Mem-mem
1. Instruction bytes fetched	$6 * 1 + 6 * 3$	$7 * 3$	$3 * 3 + 3 * 4$	$3 * 7$
2. Data memory transfers	$6 * 2$	$7 * 2$	$3 * 2$	$9 * 2$
3. Code size (= 1.)	24	<b>21</b>	<b>21</b>	<b>21</b>
4. Total memory traffic (= 1. + 2.)	36	34	<b>27</b>	39

### 8.3 Pipelining I

3.1 a) WAW: 1,3;  
RAW: 1,2; 1,3, 3,4; 2,5; 4,5  
WAR: 2,3;

b) A hazard is created when there is a dependency between instructions and they are close enough that the overlap caused by pipelining (or reordering) would change the order of access to the dependent operand.

c) In a true dependency information is transmitted between instructions, while this is not the case for name dependencies.

d) Name dependencies: WAR, WAW  
True dependencies: RAW

3.2 a) 
$$\text{Speedup} = \frac{T_{\text{unpipelined}}}{\max(T_{\text{pipestage}}) + T_{\text{latch}}}$$

b) 
$$\text{Speedup} = \frac{\text{no of stages}}{1 + bf * bp}$$

- c) 1: MOV R3, R7  
 2: LD R8,(R3)  
 3: ADDDI R3, R3, 4  
 4: LD R9, (R3)  
 5: BNE R8, R9, Loop

WAW: 1,3 RAW: 1,2; 1,3; 2,5; 3,4; 4,5 WAR: 2,3;

**3.3** The pipeline of Sections A.4 and A.5 resolves branches in ID and has multiple execution function units. More than one instruction may be in execution at the same time, but the exercise statement says write-back contention is possible and is handled by processing one instruction at a time in that stage.

Figure A.30 lists the latencies for the functional units; however, it is important to note that these data are for functional unit results forwarded to the EX stage. In particular, despite a latency of 1 for data memory access, it is still the case that for a cache hit the MEM stage completes memory access in one clock cycle.

Finally, examining the code reveals that the loop iterates 99 times. With this and analysis of iteration timing on the pipeline, we can determine loop execution time.

- a) Figure 1 shows the timing of instructions from the loop for the first version of pipeline hardware. There are several stall cycles shown in the timing diagram:
- Cycles 5–6: MUL.D stalls in ID to wait for F0 and F4 to be written back by the L.D instructions.
  - Cycles 8–15: ADD.D stalls in ID to wait for MUL.D to write back F0.
  - Cycle 19: DSUBU stalls in ID to wait for DADDUI to write back R2.
  - Cycles 21–22: BNEZ stalls in ID to wait for DSUBU to write back R5. Because the register file can read and write in the same cycle, the BNEZ can read the DSUBU result and resolve in the same cycle in which DSUBU writes that result.
  - Cycle 20: While not labeled a stall, because initially it does not appear to be a stall, the fetch made in this cycle will be discarded because this pipeline design handles the uncertainty of where to fetch after a branch by flushing the stages with instructions fetched after the branch. The pipeline begins processing after the branch with the correct fetch in cycle 23.

There are no structural hazard stalls due to write-back contention because processing instructions as soon as otherwise possible happens to use WB at most once in any clock cycle.

Figure 1 shows two instructions simultaneously in execution in clock cycles 17 and 18, but because different functional units are handling each instruction there is no structural hazard.

The first iteration ends with cycle 22, and the next iteration starts with cycle 23. Thus, each of the 99 loop iterations will take 22 cycles, so the total loop execution time is  $99 * 22 = 2178$  clock cycles.

- b) Figure 2 shows the timing of instructions from the loop for the second version of pipeline hardware.

There are several stall cycles shown in the timing diagram:

Instruction	Clock cycle																										
	1	2	3	4	5	6	7	8	...	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27			
LD F0,0(R2)	F	D	E	M	W																						
LD F4,0(R3)		F	D	E	M	W																					
MULD F0,F0,F4			F	D	s	s	E	E	...	E	M	W															
ADD F2,F0,F2				F	s	s	D	s	...	s	s	s	E	E	E	E	M	W									
DADDUI R2,R2,#8						F	s	...	s	s	s	D	E	M	W												
DADDUI R3,R3,#8												F	D	E	M	W											
DSUBU R5,R4,R2													F	D	s	E	M	W									
BNEZ R5,Loop														F	s	D	s	r									
LD F0,0(R2)																F	s	s	F	D	E	M	W				

Figure 1: Pipeline timing diagram for the pipeline without forwarding, branches that flush the pipeline, memory references that hit in cache, and FP latencies from Figure A.30. The abbreviations F, D, E, M, and W denote the fetch, decode, execute, memory access, and write-back stages, respectively. Pipeline stalls are indicated by s, branch resolution by r. One complete loop iteration plus the first instruction of the subsequent iteration is shown to make clear how the branch is handled. Because branch instructions complete (resolve) in the decode stage, use of the following stages by the BNEZ instruction is not depicted.

- Cycle 5: MUL.D stalls at ID waiting for L.D to forward F4 to EX from MEM. F0 reaches the register file by the first half of cycle 5 and thus is read by ID during this stall cycle.
- Cycles 7–12: ADD.D stalls at ID waiting for MUL.D to produce and forward the new value for F0.
- Cycle 16: DSUBU is stalled at ID to avoid contention with ADD.D for the WB stage. Note the complexity of pipeline state analysis that the ID stage must perform to ensure correct pipeline operation.
- Cycle 18: BNEZ stalls in ID to wait for DSUBU to produce and forward the new value for R5. While forwarding may deliver the needed value earlier in the clock cycle than can reading from the register file, and so in principle the branch could resolve earlier in the cycle, the next PC value cannot be used until the IF stage is ready, which will be with cycle 19.
- Cycle 17: Initially this cycle does not appear to be a stall because branches are predicted not taken and this fetch is from the fall-through location. However, for all but the last loop iteration this branch is mispredicted. Thus, the fetch in cycle 17 must be redone at the branch target, as shown in cycle 19.

Again, there are instances of two instructions in the execute stage simultaneously, but using different functional units.

The first iteration ends in cycle 19 when DSUBU writes back R5. The second iteration begins with the fetch of L.D F0, 0(R2) in cycle 19. Thus, all iterations, except the last, take 18 cycles. The last iteration completes in a total of 19 cycles. However, if there were code following the instructions of the loop, they would start after only 16 cycles of the last iteration because the branch is predicted correctly for the last iteration.

The total loop execution time is  $98 * 18 + 19 = 1783$  clock cycles.

Instruction	Clock cycle																						
	1	2	3	4	5	6	7	.	.	.	12	13	14	15	16	17	18	19	20	21	22	23	
LD F0,0(R2)	F	D	E	M	W																		
LD F4,0(R3)		F	D	E	M	W																	
MULD F0,F0,F4			F	D	s	E	E	.	.	.	E	M	W										
ADDD F2,F0,F2				F	s	D	s	.	.	.	s	E	E	E	E	M	W						
DADDUI R2,R2,#8					F	s	.	.	.	s	D	E	M	W									
DADDUI R3,R3,#8											F	D	E	M	W								
DSUBU R5,R4,R2											F	D	s	E	M	W							
BNEZ R5,Loop												F	s	D	r								
LD F0,0(R2)															F	s	F	D	E	M	W		

Figure 2: Pipeline timing diagram for the pipeline with forwarding, branches handled by predicted-not-taken, memory references that hit in cache, and FP latencies from Figure A.30. The notation used is the same as in Figure 1.

**3.4** This exercise asks, “How much faster would the machine be . . . ,” which should make you immediately think speedup. In this case, we are interested in how the presence or absence of control hazards changes the pipeline speedup. Recall one of the expressions for the speedup from pipelining presented on page A-13

$$Pipeline\_speedup = \frac{1}{1 + Pipeline\_stalls} * Pipeline\_depth \tag{1}$$

where the only contributions to Pipeline stalls arise from control hazards because the exercise is only focused on such hazards. To solve this exercise, we will compute the speedup due to pipelining both with and without control hazards and then compare these two numbers.

For the “ideal” case where there are no control hazards, and thus stalls, Equation 1 yields

$$Pipeline\_speedup_{ideal} = \frac{1}{1 + 0} * 4 = 4 \tag{2}$$

where, from the exercise statement the pipeline depth is 4 and the number of stalls is 0 as there are no control hazards.

For the “real” case where there are control hazards, the pipeline depth is still 4, but the number of stalls is no longer 0 as it was in Equation 2. To determine the value of Pipeline stalls, which includes the effects of control hazards, we need three pieces of information. First, we must establish the “types” of control flow instructions we can encounter in a program. From the exercise statement, there are three types of control flow instructions: taken conditional branches, not-taken conditional branches, and jumps and calls. Second, we must evaluate the number of stall cycles caused by each type of control flow instruction. And third, we must find the frequency at which each type of control flow instruction occurs in code. Such values are given in the exercise statement.

To determine the second piece of information, the number of stall cycles created by each of the three types of control flow instructions, we examine how the pipeline behaves under the appropriate conditions. For the purposes of discussion, we will assume the four stages of the pipeline are Instruction Fetch, Instruction Decode, Execute, and Write Back (abbreviated IF, ID, EX, and WB, respectively). A specific structure is not necessary to solve the exercise; this structure was chosen simply to ease the following discussion.

First, let us consider how the pipeline handles a jump or call. Figure 3 illustrates the behavior of the pipeline during the execution of a jump or call. Because the first pipe stage can always be done independently of whether the control flow instruction goes or not, in cycle 2 the pipeline fetches the instruction following the jump or call (note that this is all we can do – IF must update the PC, and the next sequential address is the only address known at this point; however, this behavior will prove to be beneficial for conditional branches as we will see shortly). By the end of cycle 2, the jump or call resolves (recall that the exercise specifies that calls and jumps resolve at the end of the second stage), and the pipeline realizes that the fetch it issued in cycle 2 was to the wrong address (remember, the fetch in cycle 2 retrieves the instruction immediately following the control flow instruction rather than the target instruction), so the pipeline reissues the fetch of instruction  $i + 1$  in cycle 3. This causes a one-cycle stall in the pipeline since the fetches of instructions after  $i + 1$  occur one cycle later than they ideally could have.

Figure 4 illustrates how the pipeline stalls for two cycles when it encounters a taken conditional branch. As was the case for unconditional branches, the fetch issued in cycle 2 fetches the instruction after the branch rather than the instruction at the target of the branch. Therefore, when the branch finally resolves in cycle 3 (recall that the exercise specifies that conditional branches resolve at the end of the third stage), the pipeline realizes it must reissue the fetch for instruction  $i + 1$  in cycle 4, which creates the two-cycle penalty.

Figure 5 illustrates how the pipeline stalls for a single cycle when it encounters a not-taken conditional branch. For not-taken conditional branches, the fetch of instruction  $i + 1$  issued in cycle 2 actually obtains the correct instruction. This occurs because the pipeline fetches the next sequential instruction from the program by default—which happens to be the instruction that follows a not-taken branch. Once the conditional branch resolves in cycle 3, the pipeline determines it does not need to reissue the fetch of instruction  $i + 1$  and therefore can resume executing the instruction it fetched in cycle 2. Instruction  $i + 1$  cannot leave the IF stage until after the branch resolves because the exercise specifies the pipeline is only capable of using the IF stage while a branch is being resolved.

	Clock cycle					
Instruction	1	2	3	4	5	6
Jump or call	IF	ID	EX	WB		
$i+1$		IF	IF	ID	EX	...
$i+2$			stall	IF	ID	...
$i+3$				stall	IF	...

Figure 3: Effects of a jump or call Instruction on the pipeline.

	Clock cycle					
Instruction	1	2	3	4	5	6
Taken branch	IF	ID	EX	WB		
$i+1$		IF	stall	IF	ID	...
$i+2$			stall	stall	IF	...
$i+3$				stall	stall	...

Figure 4: Effects of a taken conditional branch on the pipeline.

Combining all of our information on control flow instruction type, stall cycles, and frequency leads us to Figure 6.

Instruction	Clock cycle					
	1	2	3	4	5	6
Not-taken branch	IF	ID	EX	WB		
i+1		IF	stall	ID	EX	...
i+2			stall	IF	ID	...
i+3				stall	IF	...

Figure 5: Effects of a not-taken conditional branch on the pipeline.

Note that this figure accounts for the taken/not-taken nature of conditional branches. With this information we can compute the stall cycles caused by control flow instructions:

$$Pipeline\_stalls_{real} = (1 * 1\%) + (2 * 9\%) + (1 * 6\%) = 0.25 \quad (3)$$

where each term is the product of a frequency and a penalty. We can now plug the appropriate value for  $Pipeline\_stalls_{real}$  into Equation 1 to arrive at the pipeline speedup in the “real” case:

$$Pipeline\_speedup_{real} = \frac{1}{1 + 0.25} * (4.0) = 3.2 \quad (4)$$

Finding the speedup of the ideal over the real pipelining speedups from Equations 2 and 4 leads us to the final answer:

$$Pipeline\_speedup_{without\_control\_hazards} = 4/3.2 = 1.25 \quad (5)$$

Control flow type	Frequency (per instruction)	Stalls (cycles)
Jumps and calls	1%	1
Conditional (taken)	15% * 60% = 9%	2
Conditional (not taken)	15% * 40% = 6%	1

Figure 6: Summary of the behavior of control flow instructions.

Thus, the presence of control hazards in the pipeline loses approximately 25% of the speedup you achieve without such hazards.

**3.5** If a branch outcome is to be determined earlier, then the branch must be able to read its operand equally early. Branch direction is controlled by a register value that may either be loaded or computed. If the branch register value comparison is performed in the EX stage, then forwarding can deliver a computed value produced by the immediately preceding instruction if that instruction needs only one cycle in EX. There is no data hazard stall for this case. Forwarding can deliver a loaded value without a data hazard stall if the load can perform memory access in a single cycle and if at least one instruction separates it from the branch.

If now the branch compare is done in the ID stage, the two forwarding cases just discussed will each result in one data hazard stall cycle because the branch will need its operand one cycle



before it exists in the pipeline. Often, instructions can be scheduled so that there are more instructions between the one producing the value and the dependent branch. With enough separation there is no data hazard stall.

So, resolving branches early reduces the control hazard stalls in a pipeline. However, without a sufficient combination of forwarding and scheduling, the savings in control hazard stalls will be offset by an increase in data hazard stalls.

**3.6** a) One such situation occurs if we have e.g. the instruction LW (Load Word) in the MEM stage performing a data memory read operation exhibiting a “page fault” (because the referenced data is not present in PM), and any other instruction (say e.g. ADD) issued after LW exhibiting a “page fault” in the IF stage (because the referenced instruction is not present in PM). In this case LW will generate an exception after the second instruction (ADD) generates an exception.

b) If the pipeline can be stopped so that the the instructions just before the faulting instruction are completed and those after it can be restarted from scratch, the pipeline is said to have precise exceptions.

It is important in order to be able to correctly implement virtual memory and IEEE arithmetic trap handlers.

**3.7** The average instruction execution time on an unpipelined processor is

$$clockcycle * Avg.CPI = 1ns * ((0.5 * 4) + (0.35 * 5) + (0.15 * 4)) = 4.35ns$$

The avg. instruction execution time on pipelined processor is =  $1ns + 0.15ns = 1.15ns$

So speed up =  $4.35/1.15 = 3.78$

**3.8** Control hazard: Instruction fetch depends on some in-flight instruction being executed. For example, the target address of a branch or jump is not immediately available after the branch / jump exits from fetch stage.

- design a hazard unit to detect the hazard and stall the pipeline
- branch prediction
- using delay slot in the instruction scheduling

**3.9** • – The second instruction is dependent upon the first (because of R2)  
 – The third instruction is dependent upon the first (because of R2)  
 – The fourth instruction is dependent upon the first (because of R2)  
 – The fourth instruction is dependent upon the second (because of R4)

All these dependencies can be solved by forwarding

- Registers R11 and R12 are read during the fifth clock cycle. Register R1 is written at the end of fifth clock cycle.

**3.10** Performance can be increased by:

- Fewer loop conditional evaluations.

- Fewer branches/jumps.
- Opportunities to reorder instructions across iterations.
- Opportunities to merge loads and stores across iterations.

Performance can be decreased by:

- Increased pressure on the I-cache.
- Large branch/jump distances may require slower instructions

**3.11** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**3.12** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**3.13** a) Because there is no forwarding from the WB-stage and the correct value of source register r3 is therefore not available until clock cycle 4.

b) Only one instruction can be issued in each clock cycle and since instruction 1 has to wait, instruction 2 also must wait one clock cycle.

c) The write back stage is occupied in clock cycles 10 and 11 by instructions that have been issued earlier.

d) Here is the continuation of the table:

	<b>Instruction</b>	<b>IF</b>	<b>DI</b>	<b>EX</b>	<b>WB</b>	<b>Comment</b>
0	ADD r3,r31,r2	0	1	2	3	
1	LW r6,0(r3)	1	2	4	9	
2	ANDI r7,r5,#3	2	3	5	6	
3	ADD r1,r6,r0	3	4	10	11	
4	SRL r7,r0,#8	4	5	6	7	
5	OR r2,r4,r7	5	6	8	10	
6	SUB r5,r3,r4	6	7	9	12	
7	ADD r15,r1,r10	7	8	12	13	Wait for register r1
8	LW r6,0(r5)	8	9	13	18	Wait for register r5
9	SUB r2,r1,r6	9	10	19	20	Wait for register r6
10	ANDI r3,r7,#15	10	11	14	15	Can execute out-of-order

e) The main implication is that instruction 2 is not allowed to execute out-of-order in relation to instruction 1. There is a name-dependence between these instructions and if instruction 2 completes before instruction 1 there will be a WAW-hazard. Instruction 2 is stalled in the DI stage and the table must be modified:

	<b>Instruction</b>	<b>IF</b>	<b>DI</b>	<b>EX</b>	<b>WB</b>
0	ADD r3,r31,r2	0	1	2	3
1	LW r6,0(r3)	1	2	4	9
2	<b>ANDI r6,r5,#3</b>	2	3	<b>10</b>	<b>11</b>
3	ADD r1,r6,r0	3	4	11	12
4	SRL r7,r0,#8	4	5	6	7
5	OR r2,r4,r7	5	6	8	11
6	SUB r5,r3,r4	6	7	10	13

**3.14** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**3.15** 1. 5

2. 16

3. 11

**3.16** a) The instruction count and the clock cycle time are not affected. Therefore the only modification is in the CPI count which will have an addition that comes from branch instructions. The addition depends on the relative frequency of conditional branch instructions,  $f_{branch}$ , the fraction of these that are taken,  $b_{taken}$ , and the fraction of the branches that are miss-predicted,  $b_{misspred}$ .

$$CPI = CPI_{base} + f_{branch} * ( b_{taken} * (b_{misspred} * 2 + (1 - b_{misspred}) * 1) + (1 - b_{taken}) * (b_{misspred} * 2 + (1 - b_{misspred}) * 0) )$$

b)  $f_{branch} = 0.2$ ,  $b_{taken} = 0.65$ ,  $b_{misspred} = 0.12$ .

$$\begin{aligned} CPI &= 1.2 + 0.2 * (0.65 * (0.12 * 2 + 0.88 * 1) + 0.35 * (0.12 * 2 + 0.88 * 0)) \\ &= 1.2 + 0.2 * (0.728 + 0.084) \\ &= 1.3624 \end{aligned}$$

c)

$$\begin{aligned} CPI_{new} &= 1.2 + 0.2 * (0.65 * (0.12 * 3 + 0.88 * 2) + 0.35 * (0.12 * 3 + 0.88 * 0)) \\ &= 1.2 + 0.2 * (1.378 + 0.126) \\ &= 1.5008 \end{aligned}$$

$$T_{exeold} = IC * CPI_{old} * 1/500$$

$$T_{exenew} = IC * CPI_{new} * 1/600$$

$$\begin{aligned} Speedup &= T_{exeold}/T_{exenew} \\ &= (CPI_{old} * 1/500)/(CPI_{new} * 1/600) \\ &= (1.3624/500)/(1.5008/600) \\ &= 1.089 \end{aligned}$$

The new processor with higher clock frequency is thus only 8.9% faster than the old even though the clock frequency has been increased by 20%. The sole reason for this is the branch hazards.

d) A branch target buffer can be used in the IF-stage. The branch penalty for correctly predicted branches will then be 0.

**3.17** The following chart shows the execution of the given instruction sequence cycle by cycle. The stages of instruction execution:

- F Instruction fetch
- D Decode and issue

- E1 Execute in LOAD/STORE unit
- E2 Execute in ADD/SUB unit
- E3 Execute in MUL/DIV unit
- W Write back into register file and reservation stations

Instruction		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
LOAD R6, 34(R12)	F	D	E1	E1	W											
LOAD R2, 45(R13)	F	r	r	r	D	E1	E1	W								
MUL R0, R2, R4	F	D	s	s	s	s	s	s	E3	E3	W					
SUB R8, R2, R6	F	D	s	s	s	s	s	s	E2	W						
DIV R10, R0, R6	F	r	r	r	r	r	r	r	r	r	D	E3	E3	E3	E3	W
ADD R6, R8, R2	F	r	r	r	r	r	r	r	r	D	E2	W				

Cycles in which an instruction is waiting for a reservation station are marked as 'r' and the cycles in which an instruction is waiting for one or more operands are marked as 's'. As seen in the time chart, the issue and write back cycles for various instructions are:

instruction	issue cycle	write back cycle
LOAD	1	4
LOAD	4	7
MUL	1	10
SUB	1	9
DIV	10	15
ADD	9	11

- 3.18** 1. See book pages 82-83 and Figure 2.4. For our example program we would have one entry corresponding to the BNE instruction in the end of the program. The prediction would evolve according to this table if we assume we start in state 00:

Execution of the BNE instruction	Prediction state before execution	Prediction	State after execution
First	00	Not taken (wrong)	01 (it was taken)
Second	01	Not taken (wrong)	11 (it was taken)
Third	11	Taken (correct)	11 (it was taken)
4th, 5th, ...	11	Taken (correct)	11 (it was taken)
1000th	11	Taken (wrong)	10 (was not taken)

For a program with two branches, we would get a 2-bit state entry for each branch as long as they do not share the same index in the branch-prediction buffer. Branches with different index (entries) are not correlated with each other.

- The key point for supporting speculation is to commit architectural visible changes in program order. It must be possible to cancel (flush) speculative results without any remaining visible changes in the registers or in the memory. Registers and memory should only be written in the commit stage. Before that, the reorder buffer holds the speculative (temporary) results.
- The corresponding steps for a store instruction is:
  - Issue when reservation station and ROB entry is available

- Read already available operands from registers and instruction
  - Send instruction to reservation station
  - Tag unavailable operands with ROB entry
  - Mark ROB entry as busy
- (b) Execute after issue
- Wait for operand values on CDB (if not already available)
  - Compute address and store it in ROB entry
- (c) Write result when CDB and ROB available
- Update ROB entry with source register value, and mark as ready
  - Free reservation station
- (d) Commit when at head of ROB and ready
- Write result (source register value) to memory at computed address
  - Free ROB entry

The important points are: (1) The actual write is done at the commit stage, (2) To carry out the write, the address needs to be computed (execute stage), and, (3) the source operand (the value to write) is needed before commit is possible.

4. RAW hazards are avoided by delaying instruction execution until the source operands are available. Instructions wait in a reservation station until source operands are available. Earlier instructions that write dependent values send their results directly to the waiting reservation stations.

**3.19** Some factors that limit the exploitable ILP (book Section 3.2):

Factor	Comment
Window size (the size of buffers to keep track of instructions)	At small window sizes, the processor simply cannot see the future instructions that could have been issued in parallel.
Branch, jump prediction	It is hard and expensive to reduce the miss prediction rate beyond a few percent. Misspeculations reduce the exploitable ILP.
Finite number of registers	Registers are needed to store temporary results. This can limit the amount of parallel work possible.
Imperfect alias analysis	If false dependencies through memory and name dependencies can be detected, more instructions can be issued in parallel. However, this is sometimes impractically complex.

## 8.4 Memory systems, Cache I

**4.1** In real life bigger memory is slower and faster memory is more expensive. We want to simultaneously increase the speed and decrease the cost. Speed is important because the widening performance gap between CPU and memory. Size is important since applications and data sets are growing bigger. Use several types of memory with varying speeds arranged in a hierarchy that is optimized with respect to the use of memory. Mapping functions provide address translations between levels.

Registers: internal ultra-fast memory for CPU; static register

Cache: speed up memory access; SRAM

Main memory: DRAM

VM: make memory larger, disk; safe sharing between processes of physical memory, protection, relocation

(archival storage, backup on tape)

**4.2** a) Block transfer time between cache and memory (penalty):  $40+32/4 = 48$  cycles.

Number of block transfers per instruction between cache and memory:

$$\begin{aligned} & \text{Cache accesses/instr} * (0.5 \text{ blocks writeback} + 1 \text{ block fetch}) * \text{Miss ratio} = \\ & (1+0.2) * (0.5+1) * \text{Miss ratio} = 1.8 * \text{Miss ratio} \end{aligned}$$

$$\text{CPI} = \text{baseCPI} + \text{MemoryStalls/instr} = \text{base CPI} + \text{BlockTransfers/instr} * \text{penalty}$$

$$\text{CPI 1} = 1.5 + 1.8 * 0.029 * 48 = 4.01$$

$$\text{CPI 2} = 1.5 + 1.8 * 0.022 * 48 = 3.40$$

$$\text{CPI 3} = 1.5 + 1.8 * 0.020 * 48 = 3.22$$

b) We must do a TLB access for each cache access since the caches are physically addressed. We then in all three cases get an extra CPI offset of:  $0.002 * 20 * 1.2 = 0.048$ .

c) Comparing execution times using CPU performance formula:

$$\text{EXE 1} = 4.01 * 1 * \text{CP} * \text{IC} = 4.01 * \text{CP} * \text{IC}$$

$$\text{EXE 1} = 3.40 * 1.2 * \text{CP} * \text{IC} = 4.08 * \text{CP} * \text{IC}$$

$$\text{EXE 1} = 3.22 * 1.25 * \text{CP} * \text{IC} = 4.025 * \text{CP} * \text{IC}$$

Cache 1 is the best.

d) Similar to fig C.5, page C-13 in Hennessy & Patterson.

e) In a virtually addressed cache the TLB is only accessed at cache misses. In a physically addressed cache TLB is accessed for each cache access.

**4.3** A large register file and a large data cache both serve the purpose of reducing memory traffic. From an implementation point of view, the same chip area can be used for either a large register file or a large data cache. With a larger register set, the instruction set must be changed so that it can address more register in the instructions. From a programming point of view, registers can be manipulated by program code, but the data cache is transparent to the user. In fact, the data cache is primarily involved in load/store operations. The addressing of a cache involves address translation and is more complicated than that of a register file.

**4.4** With a unified cache, the fraction of the cache devoted to instructions and data respectively may change from one program to another. This can be a benefit but also a problem for, e.g., a small loop that touches a lot of data. In a unified cache, instructions will probably be replaced by data as the loop executes. With two separate caches, the entire loop can fit in cache and performance will probably be improved. The design trade-offs involve choosing the correct division between instruction and data caches for separate cache memories, studying the miss rate for a unified cache, choosing a correct address mapping for a unified cache and a replacement policy.

**4.5** A cache write is called write-through when information is passed both to the block in the cache and to the block in the lower-level memory; when information is only written to the block, it is called write-back. Write-back is the fastest of the two as it occurs at the speed of the cache memory, while multiple writes within a block require only one write to the lower-level memory.

4.6 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

4.7 This exercise uses differences in cache organizations to illustrate how benchmarks can present a skewed perspective of system performance. Because the memory hierarchy heavily influences system performance, it is possible to develop code that runs poorly on a particular cache organization. This exercise should drive home not only an appreciation for the influence of cache organization on performance, but also an appreciation of how difficult it is for a single program to provide a reasonable summary of general system performance.

- Consider the MIPS code blurb in Figure 7. We make two assumptions in this code: First, the value of `r0` is zero; second, locations `f00` and `bar` both map onto the same set in both caches. For example, `f00` and `bar` could be `0x00000000` and `0x80000000` (these addresses are in hexadecimal), respectively, since both addresses reside in the set zero of either cache. On Cache A, this code only causes two compulsory misses to load the two instructions into the cache. After that, all accesses generated by the code hit the cache. For Cache B, all the accesses miss the cache because a direct-mapped cache can only store one block in each set, yet the program has two active blocks that map to the same set. The cache will “thrash” because when it generates an access to `f00`, the block containing `bar` is resident in the cache, and when it generates an access to `bar`, the block containing `f00` is resident in the cache. This is a good example of where a victim cache could eliminate the performance benefit of the associative cache. Keep in mind that in this example the information that Cache B misses on is always recently resident.

```
F00:      beqz    r0, bar   : branch iff r0 == 0
          .
          .
          .
bar:      beqz    r0, f00   ; branch iff r0 == 0
```

Figure 7: MIPS code that performs better on cache A.

With all access hits, Cache A allow the processor to maintain  $CPI = 1$ . Cache B misses each access at cost of 100 ns, or 200 clock cycles. This Cache B allows its processor to achieve  $CPI = 201$ . Cache A offers a speed-up of 201 over Cache B.

- Consider the code blurb shown in Figure 8. We make two assumptions: first, locations `baz` and `qux` and the location pointed to by `0(r1)` map to different sets within the caches and are all initially resident; second, `r0` is zero (as it always is for MIPS) This code illustrates the main thrust of a program that makes a system with Cache B outperform a system with Cache A, that is, one, that repeatedly writes to a location that is resident in both caches. Each time the `sw` executes on Cache A, the data stored are written to memory because Cache A is write through. For Cache B, the `sw` always finds the appropriate block in the cache (we assume the data at location `0(r1)` are resident in the cache) and updates only the cache block, as the cache is write back; the block is not written to main memory until it is replaced.

In the steady state, Cache B hits on every write, and, so, maintain  $CPI = 1$ . Cache A writes to memory on each store, consuming an extra 100 ns each time. Cache B allows the processor to complete one iteration in 2 clocks. With Cache A the processor needs 202 clocks per iteration. Cache B offers a speedup of 101 over Cache A.

```

Baz:      sw      0(r1), r0 ; store r0 to memory
Qux:      beqz    r0, baz  ;branch iff r0 == 0

```

Figure 8: MIPS code that performs better on cache B.

**4.8** Since a processor is clocked in discrete steps it's enough to count the extra cost you get on a miss. Here we have stated the miss cost as the number of clock cycles.

- a) Miss cost  $t_m(B) = M(B)/100 * (1 + 4 + B/4)$  which gives us  
 $t_m(32) = 0.016 * (5 + 32/4) = 0.2$  as the smallest value  $\Rightarrow B_{opt} = 32$
- b) Miss cost  $t_m(B) = M(B)/100 * (2 + 1 + 4 + B/4)$  which gives us  
 $t_m(64) = 0.010 * (2 + 1 + 64/4) = 0.23$  as the smallest value  $\Rightarrow B_{opt} = 64$
- c) I: Miss cost  $t_m(B) = M(B)/100 * (1 + 4 + \lceil B/8 \rceil)$  which gives us  
 $t_m(64) = 0.010 * (1 + 4 + 64/8) = 0.13$  as the smallest value  $\Rightarrow B_{opt} = 64$   
 II: Miss cost  $t_m(B) = M(B)/100 * (2 + 1 + 4 + \lceil B/8 \rceil)$  which gives us  
 $t_m(128) = 0.0064 * (2 + 1 + 4 + 128/8) = 0.1472$  as the smallest value  $\Rightarrow B_{opt} = 128$
- d) 16 byte in a block means that 4 bit in the byte offset field is needed. There are  $256/16 = 16$  blocks in the cache memory. This gives us that 4 bits are needed for the index field.

$$Block\_address = \lfloor address/block\_size \rfloor = \lfloor 28E7_{16}/10_{16} \rfloor = \lfloor 10471/16 \rfloor = 654$$

index = block address mod (number of blocks) =  $654 \bmod 16 = 14$

Address  $28E7_{16}$  is represented by block 14 in the cache memory.

- e) The block address is the same. We now have 8 sets ( 2 blocks in every set) which means that 3 bits is needed for the index field. The address is represented in  
 $set = 654 \bmod 8 = 6$ . This set includes block  $2*6$  and  $2*6 + 1$ .

**4.9** a) Modern computer systems needs large memories that at the same time isn't allowed to cost too much. Large memories are unfortunately also very expensive if they also shall be fast. The technique with memory hierarchy solves this problem by having a series of memories with small fast memories closest to the processor and gradually bigger and slower memories the further away from the processor you get. The level furthest away from the processor is big enough to have space for the total memory space you wish to address. The memory closest to the processor only have room for a fraction of the total addressable memory space. When the processor performs a memory reference (fetching of instructions or reading/ writing data) the first level in the memory hierarchy is checked to see if the wanted data is there. If is, then the memory referents is performed there. If it isn't there then the next level is checked and so forth. When you find what you are looking for in a memory hierarchy a block of data (or instructions) that contain this memory position will be moved to a higher level of the hierarchy. The size of the block you moves get smaller the higher up in the memory hierarchy you get.

The point of memory hierarchy is to have the data /instructions that the processor most often needs in the fast memory closest to the processor. You can afford fetching things that isn't used as often from the slower and bigger memory. The reason that it work so well is that computer programs often exhibits the principle of locality: Temporal locality:



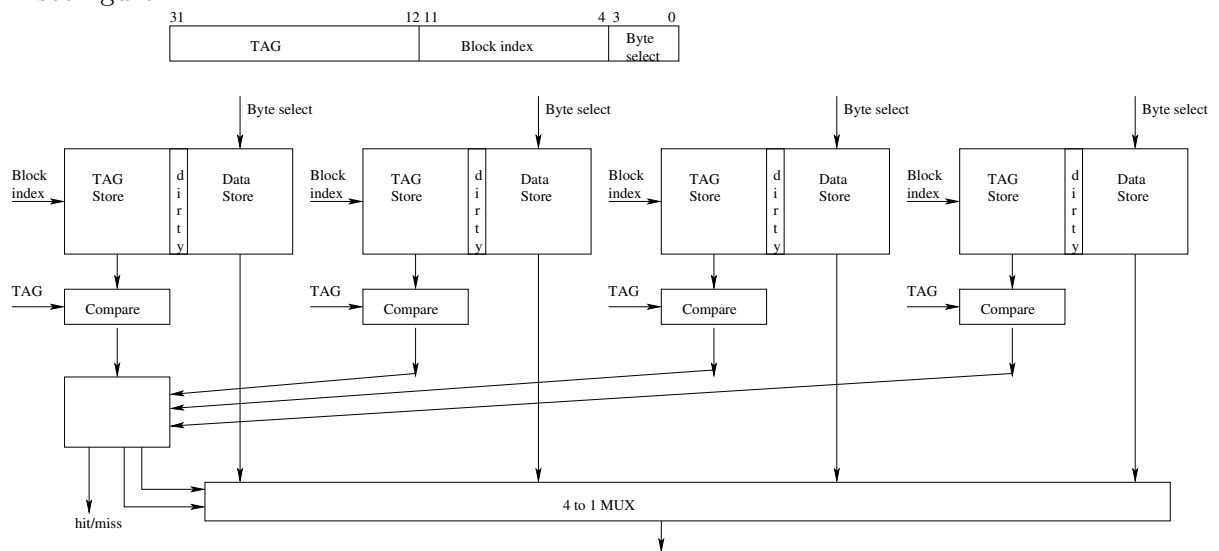
What has been used by the processor will with an high degree of likelihood be used again. This property leads to that it's enough that a higher lever of the memory hierarchy only contain a small part of the lower ones. Spatial locality: What is close (in the address space) to something the processor has accessed will with high degree of likelihood also be needed. By moving a block of data/instructions around what you need at a miss you don't have to fetch as many times.

- b) The byte offset is  $16-13=3$  bits. This means that the block-size is  $2^3 = 8$  bytes and the number of blocks is  $256/8 = 32$ . There are two blocks in each set which gives us that the cache is two-way set associative. Address  $3333_{16}$  have the tag =  $0666_{16}$ , and byte offset = 3.
- c) Write-through and Copy-back: Write-through: All the writes are also sent to the main memory even if they do hit in the cache. This means that with a block replacement you don't have to update the main memory at the cost of slower writes or more complicated hardware. Copy-back: Writes are only performed locally in the cache. For every block you then will also need a dirty bit that tells if the block in the cache has been written to. On a block change the block also needs to be written back to the main memory if the dirty bit is true.

**4.10** When a cache miss occurs, the controller must select a block to be replaced with the desired data. Three primary strategies for doing this are: Random, Least-recently used (LRU), and First-in first-out (FIFO).

A replacement algorithm is needed with set-associative and fully associative caches. For direct mapped caches there is no choice, the block to be replaced is uniquely determined by the address.

**4.11** see figure



**4.12** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**4.13** See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>

**4.14** a) miss rate:

- non-blocking: don't affect miss rate, the main thing happening is that the processor makes other useful things while the miss is handled
- hardware prefetching with stream buffers: a hit in the stream buffer cancels the cache request, ie the memory reference is not counted as a miss, which means that the miss rate will decrease.
- software prefetching: if correctly done miss rate will decrease

b) memory bandwidth:

- non-blocking: since the processor will have fewer stall cycles it will get a lower CPI and consequently the requirements on memory bandwidth will increase
- hardware prefetch and software prefetch: prefetching is a form of speculation which means that some of the memory traffic is unused which in turn might increase the need for memory bandwidth

c) number of executed instructions: will be unchanged for non-blocking and hardware prefetch. Software prefetch will add the prefetch-instructions, so number of executed instructions will increase.

**4.15** a) The merging write buffer links the CPU to the write-back L2 cache. Two CPU writes cannot merge if they are to different sets in L2. So, for each new entry into the buffer a quick check on only those address bits that determine the L2 set number need be performed at first. If there is no match in this 'screening' test, then the new entry is not merged. If there is a set number match, then all address bits can be checked for a definitive result.

b) As the associativity of L2 increases, the rate of false positive matches from the simplified check will increase, reducing performance.

**4.16** The three C's model sorts the causes for cache misses into three categories:

- Compulsory – The very first access can never be in cache and is therefore bound to generate a miss;
- Capacity – If the cache cannot contain all the blocks needed for a program, capacity misses may occur;
- Conflict – If the block placement strategy is set-associative or direct-mapped, conflict misses will occur because a block can be discarded and later retrieved if too many blocks map to its set.

**4.17** The three C's give insight into the cause of misses, but this simple model has its limits; it gives you insight into average behavior but may not explain an individual miss. For example, changing cache size changes conflict misses as well as capacity misses, since a larger cache spreads out references to more blocks. Thus, a miss might move from a capacity miss to a conflict miss as cache size changes. Note that the three C's also ignore replacement policy, since it is difficult to model and since, in general, it is less significant. In specific circumstances the replacement policy can actually lead to anomalous behavior, such as poorer miss rates for larger associativity, which is contradictory to the three C's model.

## 8.5 Memory systems, Virtual Memory

- 5.1** a) Program basic blocks are often short (less than 10 instructions). Even program run blocks, sequences of instructions executed between branches, are not very long. Pre-fetching obtains the next sequential block, but program execution does not continue to follow location PC, PC+4, PC+8, ..., for very long. So as blocks get larger the probability that a program will not execute all instructions in the block, but rather take a branch to another instruction address, increases. Pre-fetching instructions benefits performance when the program continues straight-line execution into the next block. So as instruction cache blocks increase in size, pre- fetching becomes less attractive.
- b) Data structures often comprise lengthy sequences of memory addresses. Program access of a data structure often takes the form of a sequential sweep. Large data blocks work well with such access patterns; pre- fetching is likely still of value due to the highly sequential access patterns. The efficiency of data pre-fetch can be enhanced through a suitable grouping of the data items taking the block limitations into account. This is especially noteworthy when the data-structure exceeds the cache size. Under such circumstances, it will become of critical importance to limit the amount of out-of-cache block references.
- 5.2** a) We can expect software to be slower due to the overhead of a context switch to the handler code, but the sophistication of the replacement algorithm can be higher for software and a wider variety of virtual memory organizations can be readily accommodated. Hardware should be faster, but less flexible.
- b) Yes, the speed of handling TLB misses in software is slower than the hardware solution, see answer in a).  
Factors other than whether miss handling is done in software or hardware can quickly dominate handling time. Is the page table itself paged? Can software implement a more efficient page table search algorithm than hardware? What about hardware TLB entry pre-fetching?
- c) Software solution is easier than handling in hardware when page table structures change dynamically on physical properties, e.g. grow or shrink in size; hardware solution is easier when table contents change dynamically, because static software scheduling is difficult.
- 5.3** a) – As miss penalty tends to be severe, one usually decides for a complex placement strategy. Usually one takes for full association.
- To reduce address translation time, a cache is added to remember the most likely translations, the Translation Lookaside Buffer.
  - Almost all operating systems rely on a replacement of the least-recently used (LRU) block indicated by a reference bit, which is logically set whenever a page is addressed.
  - Since the cost of an unnecessary access to the next-lower level is high, one usually includes a dirty bit. It allows blocks to be written to lower memory only if they have been altered since reading.
- b) The main purpose of the TLB is to accelerate the address translation for reading/writing virtual memory. A TLB entry holds a portion of the virtual address, a physical page frame number, a protection field, a valid bit, a use bit and a dirty bit. The latter two is not always used. The size of the page table is inversely proportional to the page size;

choosing a large page size allows larger caches with fast cache hit times with a small TLB. A small page size conserves storage, limiting the amount of internal fragmentation. Their combined effect can be seen in process start-up time, where a large page size lengthens invocation time but shortens page renewal times. Hence, the balance goes for large pages in large computers and vice-versa.

	TLB	1st-level cache	2nd-level cache	Virtual memory
Block size (in bytes)	4-32	16-256	1-4k	4096-65,536
c) Block placement	Full associative	2/4-way set associative	8/16-way set associative	Direct mapped
Overall size	32-8,192b	1 MB	2-16MB	32 MB – 1 TB

**5.4** Out-of-order (OOO) execution will change both the timing of and sequence of cache access with respect to that of in-order execution. Some specific differences and their effect on what cache design is most desirable are explored in the following.

Because OOO reduces data hazard stalls, the pace of cache access, both to instructions and data, will be higher than if execution were in order. Thus, the pipeline demand for available cache bandwidth is higher with OOO. This affects cache design in areas such as block size, write policy, and pre-fetching.

Block size has a strong effect on the delivered bandwidth between the cache and the next lower level in the memory hierarchy. A write-through write policy requires more bandwidth to the next lower memory level than does write back, generally, and use of a dirty bit further reduces the bandwidth demand of a write-back policy. Pre-fetching increases the bandwidth demand. Each of these cache design parameters – block size, write policy, and pre-fetching – is in competition with the pipeline for cache bandwidth, and OOO increases the competition. Cache design should adapt for this shift in bandwidth demand toward the pipeline.

Cache accesses for data and, because of exception, instructions occur during execution. OOO execution will change the sequence of these accesses and may also change their pacing.

A change in sequence will interact with the cache replacement policy. Thus, a particular cache and replacement policy that performs well on a chosen application when execution of the superscalar pipeline is in order may perform differently – even quite differently – when execution is OOO.

If there are multiple functional units for memory access, then OOO execution may allow bunching multiple accesses into the same clock cycle. Thus, the instantaneous or peak memory access bandwidth from the execution portion of the superscalar can be higher with OOO.

Imprecise exceptions are another cause of change in the sequence of memory accesses from that of in-order execution. With OOO some instructions from earlier in the program order may not have made their memory accesses, if any, at the time of the exception. Such accesses may become interleaved with instruction and data accesses of the exception-handling code. This increases the opportunity for capacity and conflict misses. So a cache design with size and/or associativity to deliver lower numbers of capacity and conflict misses may be needed to meet the demands of OOO.

**5.5** First, we use the miss penalty and miss rate information to compute the contribution to CPI from cache misses for each configuration. We do this with the formula:

$$\text{Cache CPI} = \text{Misses per instruction} * \text{Miss Penalty}$$

We need to compute the miss penalties for each system:

$$\text{Miss Penalty} = \frac{\text{Memory Access Time}}{\text{Clock Cycle}}$$

The clock cycle times for the processors are 250 ps, 200 ps, and 400 ps, respectively. Hence, the miss penalties are

$$1 : \frac{50 \text{ ns}}{250 \text{ ps}} = 200 \text{ cycles}$$

$$2 : \frac{50 \text{ ns}}{200 \text{ ps}} = 250 \text{ cycles}$$

$$3 : \frac{0.75 * 50 \text{ ns}}{400 \text{ ps}} = 94 \text{ cycles}$$

Applying this for each cache:

$$CPI_1 = 0.005 * 200 = 1.0$$

$$CPI_2 = 0.0055 * 250 = 1.4$$

$$CPI_3 = 0.01 * 94 = 0.94$$

We know the pipeline CPI contribution for everything but processor 3; its pipeline CPI is given by

$$\text{Pipeline CPI} = 1/\text{Issue rate} = \frac{1}{9 * 0.5} = 1/4.5 = 0.22$$

Now we find the CPI for each processor by adding the pipeline and cache CPI contributions:

$$1 : 0.8 + 1.0 = 1.8$$

$$2 : 1.0 + 1.4 = 2.4$$

$$3 : 0.22 + 0.94 = 1.16$$

Since this is the same architecture, we can compare instruction execution rates in millions of instructions per second (MIPS) to determine relative performance CR / CPI as

$$1 : \frac{4000 \text{ MHz}}{1.8} = 2222 \text{ MIPS}$$

$$2 : \frac{5000 \text{ MHz}}{2.4} = 2083 \text{ MIPS}$$

$$3 : \frac{2500 \text{ MHz}}{1.16} = 2155 \text{ MIPS}$$

In this example, the simple two-issue static superscalar looks best. In practice, performance depends on both the CPI and clock rate assumption.

**5.6** a) For:

- Size of the page table is inversely proportional to the page size
- Larger page sizes allow larger caches using a virtually indexed, physically tagged direct mapped cache.
- The number of TLB entries are restricted, larger page size means more memory mapped efficiently.

Against: Larger pages lead more wasted storage due to internal fragmentation.

- b) In a virtual memory system: Virtual address is a logical address space for a process. This is translated by a combination of hardware and software into a physical address which access main memory. This process is called memory mapping. The virtual address space is divided into *pages* (blocks of memory). Page fault: an access to a page which is not in physical memory. TLB, Translation Lookaside Buffer is a cache of address translations.
- c) Page table takes  $\frac{2^{32}}{2^{12}} * 4 = 2^{22} = 4$  Mbyte
- d) An inverted page table is like a fully associative cache where each page table entry contains the physical address and, as tag, the virtual address. It takes  $\frac{2^{29}}{2^{12}} * (4 + 4) = 2^{20} = 1$  Mbyte

5.7 See lecture slides, HP Ch. 5 and App C.

## 8.6 Storage systems, I/O

- 6.1 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.2 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.3 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.4 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.5 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.6 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>
- 6.7 See 'Case Study Solutions' at <http://www.elsevierdirect.com/companion.jsp?ISBN=9780123704900>