# BBM234 COMPUTER ORGANIZATION

# Course Info

- Instructors
  - Associate Professor Dr. Suleyman Tosun
    - E-mails: stosun@gmail.com
    - Office hour 1: Mon: 16:00-17:00
    - Office hour 2: Tue, 13:00-14:00

# Course Info

- Instructors
  - Assistant Professor Dr. Mehmet Köseoğlu
    - Office: Room 219
    - E-mail: [mkoseoglu@cs.hacettepe.edu.tr](mailto:mkoseoglu@cs.hacettepe.edu.tr)
    - Office hours:
      - Tuesday 12:30-13:30
      - Whenever I am in the office.

# Course Info

- Instructors
  - Dr. Hüseyin Temuçin
    - Office: Room 225
    - E-mail: htemucin @cs.hacettepe.edu.tr
    - Office hours:
      - Monday 13:30-14:30

# Lecture times

- Lectures:
    - Tuesday, 09:00-12:00
    - Section 1: D1
    - Section 2: D2
    - Section 3: D3

# Grading

- We will have two midterm and a final exams, two projects:
    - Two Midterms (50%: 25% each) – Weeks 7 and 11
    - Final (40 %)
    - Two Projects (10 %)
        - MIPS programming
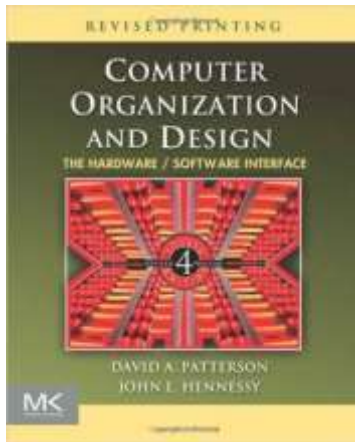        - Cache performance

# Academic Dishonesty

- Any student caught cheating on an exam or program will automatically fail the course and <u>will be referred to the department chair and/or dean</u>.

- When writing programs and doing homework, you may consult with me or the TA at any stage. You may seek help about the MIPS system or the editor from anyone at any time.

- To avoid cheating via collaboration, do not show your code and homework to any other classmates.

- If a classmate consults you for help after attempting to run his or her program, you may assist in determining why his or her code doesn't work, but <u>refrain from suggesting specific new code</u>.

- Do not lead your classmates into temptation: guard your print-outs. **We intend to use an automatic cheating-verification program**.

- Do not even take the risk!..

# Reference Books

*Digital Design and Computer Architecture*, D. Harris, S. Harris, Morgan Kaufmann, 2007.

- We will mention this book as "DDCA"

*Computer Organization and Design: The Hardware/Software Interface*, 4th Edition, John L. Hennessy & David A. Patterson, Morgan Kaufmann, 2011.
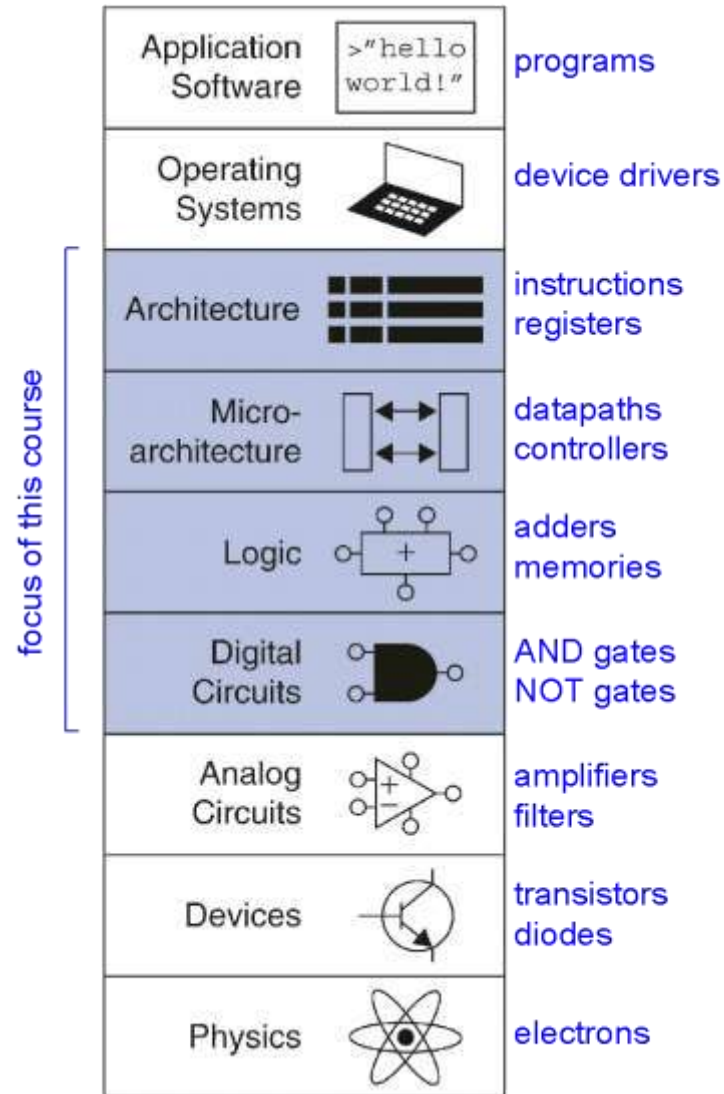
- We will mention this book as "COD"

# **Purpose of the Course**

- Understand what's under the hood of a computer
- Design and build a microprocessor
- Learn caches and their effects on performance
- Understand the virtual memory concept
- Learn how I/O systems are dealt with in hardware level (not in OS level)

# Where is this course in CS?

- This course is the follow up course of digital design course.

# What You Will Learn

- How programs are translated into the machine language
  - And how the hardware executes them
- The hardware/software interface
- What determines program performance
  - And how it can be improved
- How hardware designers improve performance
- What is parallel processing

# Tentative Schedule

| Week # | Topics | Book | Chapter |
|---|---|---|---|
| 1 | Introduction:Computer Abstraction and Technology | COD | 1 |
| 2 | Digital Building Blocks - Review | DDCA | 5 |
| 3-5 | Computer Architecture: MIPS Instruction Set Architecture | DDCA | 6 |
| 7 | Midterm 1 | | |
| 6,8 | Microarchitecture: Single-Cycle MIPS Processor | DDCA | 7 |
| 9-10 | Microarchitecture: Pipelined MIPS Processor | DDCA | 7 |
| 11 | Midterm 2 | | |
| 12-14 | Memory Systems: Memory System Performance Analysis, Caches, Virtual Memory, TLB | DDCA | 8 |
| | Final | | |

# Course Communication Platform

- We will use Piazza
- Join the class ASAP
  - Go to [www.piazza.com](www.piazza.com)
- **Signup Link:**
  - piazza.com/hacettepe.edu.tr/spring2019/bbm234
- Or - Click on Students Get Started button
- Type Hacettepe University
- Select Spring 2019
- Search BBM234 and
- Join the Class

# Chapter 1

# Computer Abstractions and Technology

# The Computer Revolution

- Progress in computer technology
  - Underpinned by Moore's Law
    - **Moore's law** is the observation that the number of transistors in a dense integrated circuit doubles approximately every two years
- Makes novel applications feasible
  - Computers in automobiles
  - Cell phones
  - Human genome project
  - World Wide Web
  - Search Engines
- Computers are pervasive

# Classes of Computers

- Personal computers
  - General purpose, variety of software
  - Subject to cost/performance tradeoff

- Server computers
  - Network based
  - High capacity, performance, reliability
  - Range from small servers to building sized

# Classes of Computers

- Supercomputers
    - High-end scientific and engineering calculations
    - Highest capability but represent a small fraction of the overall computer market

- Embedded computers
    - Hidden as components of systems
    - Stringent power/performance/cost constraints
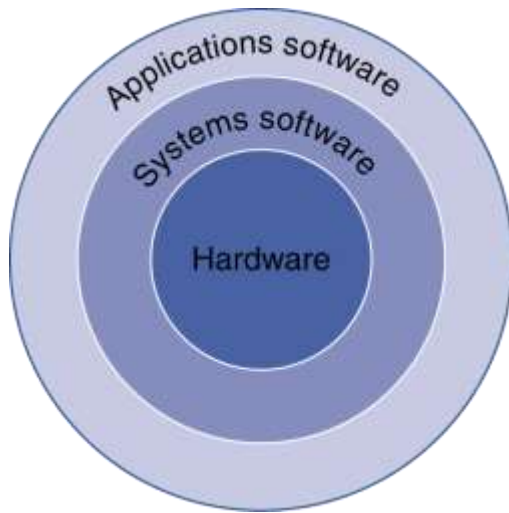
# Understanding Performance

- Algorithm
  - Determines number of operations executed
- Programming language, compiler, architecture
  - Determine number of machine instructions executed per operation
- Processor and memory system
  - Determine how fast instructions are executed
- I/O system (including OS)
  - Determines how fast I/O operations are executed

# Eight Great Ideas

- Design for **Moore's Law**

- Use **abstraction** to simplify design

- Make the **common case fast**

- Performance *via* **parallelism**

- Performance *via* **pipelining**

- Performance *via* **prediction**

- **Hierarchy** of memories

- **Dependability** *via* redundancy

MOORE'S LAW

ABSTRACTION

COMMON CASE FAST

PARALLELISM

PIPELINING

PREDICTION

HIERARCHY

DEPENDABILITY

# **Below Your Program**

- ## Application software

  - ### Written in high-level language

- ## System software

  - ### Compiler: translates HLL code to machine code

  - ### Operating System: service code

    - Handling input/output
    - Managing memory and storage
    - Scheduling tasks & sharing resources

- ## Hardware

  - ### Processor, memory, I/O controllers

# Levels of Program Code

- ## High-level language
  - Level of abstraction closer to problem domain
  - Provides for productivity and portability
- ## Assembly language
  - Textual representation of instructions
- ## Hardware representation
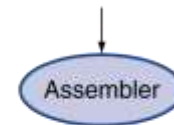  - Binary digits (bits)
  - Encoded instructions and data

High-level language program (in C)

```
swap(int v[], int k)
{int temp;
   temp = v[k];
   v[k] = v[k+1];
   v[k+1] = temp;
}
```

Compiler

Assembly language program (for MIPS)

```
swap:
   muli $2, $5,4
   add  $2, $4,$2
   lw   $15, 0($2)
   lw   $16, 4($2)
   sw   $16, 0($2)
   sw   $15, 4($2)
   jr   $31
```
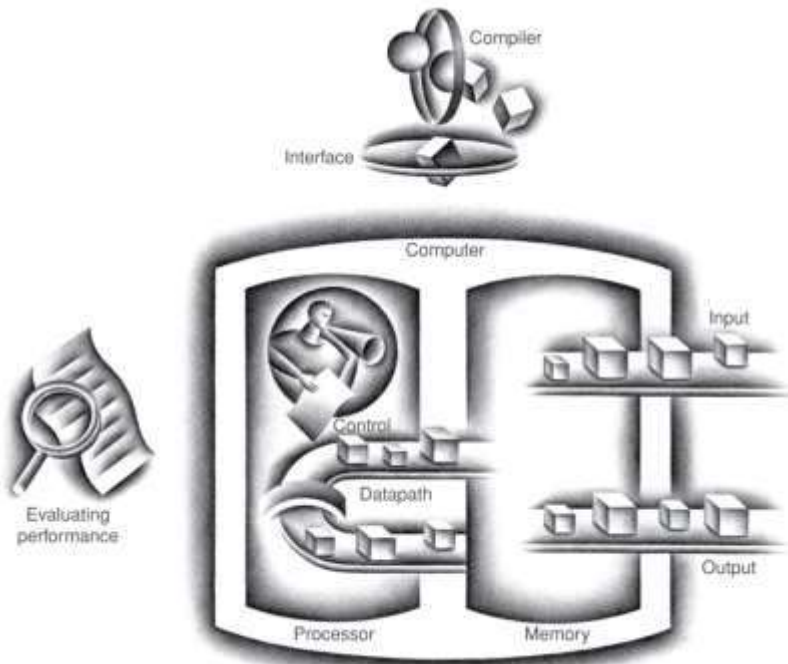
Assembler

Binary machine language program (for MIPS)

```
00000000101000010000000000011000
00000000000110000001100000100001
10001100011000100000000000000000
10001100111100100000000000000100
10101100111100100000000000000000
10101100011000100000000000000100
00000011111000000000000000001000
```

# Components of a Computer

The BIG Picture



- Same components for all kinds of computer
  - Desktop, server, embedded
- Input/output includes
  - User-interface devices
    - Display, keyboard, mouse
  - Storage devices
    - Hard disk, CD/DVD, flash
  - Network adapters
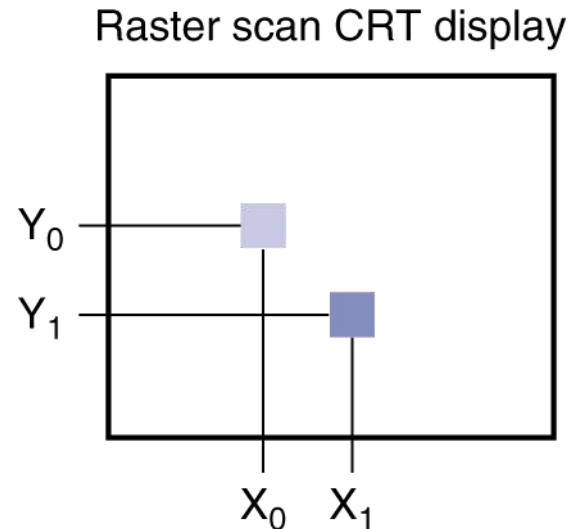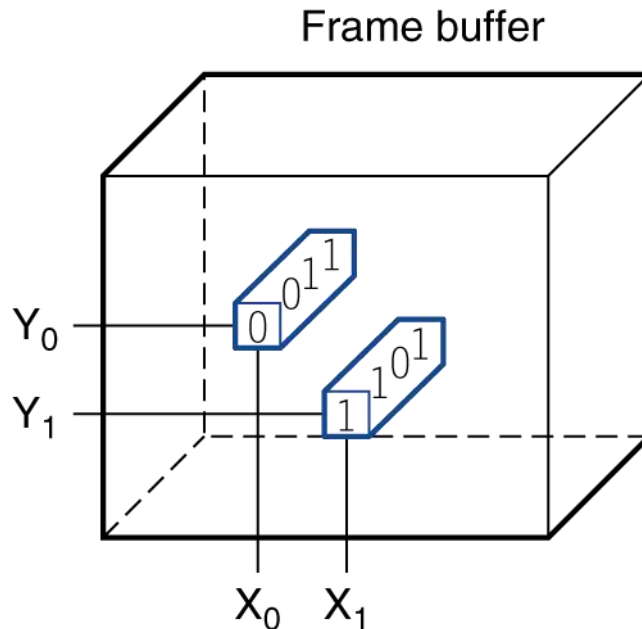    - For communicating with other computers

# Touchscreen

- PostPC device
- Supersedes keyboard and mouse
- Resistive and Capacitive types
  - Most tablets, smart phones use capacitive
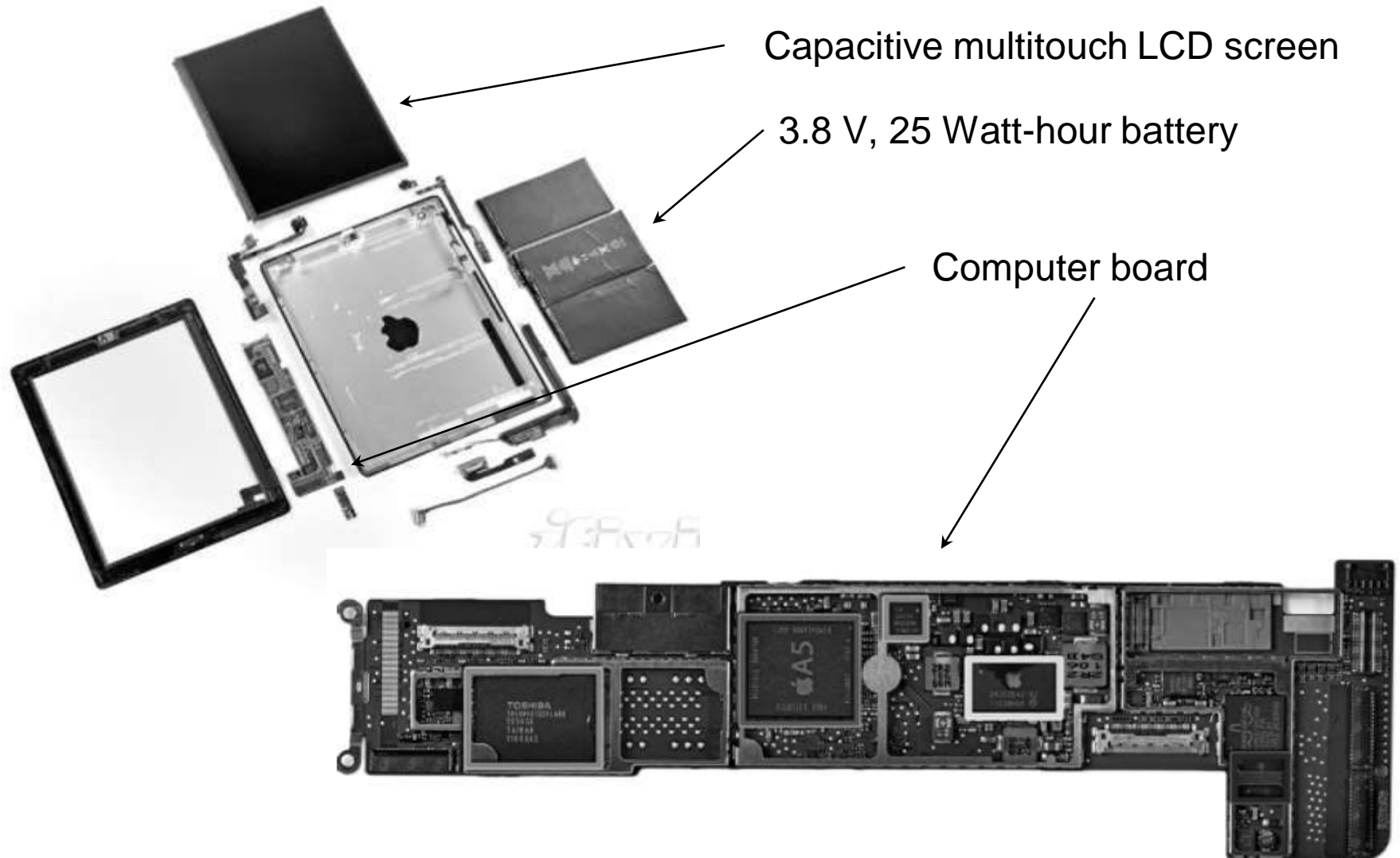  - Capacitive allows multiple touches simultaneously

# Through the Looking Glass

- LCD screen: picture elements (pixels)
  - Mirrors content of frame buffer memory

Frame buffer

Raster scan CRT display

# Opening the Box

Capacitive multitouch LCD screen

3.8 V, 25 Watt-hour battery

Computer board

# Inside the Processor (CPU)

- Datapath: performs operations on data
- Control: sequences datapath, memory, ...
- Cache memory
  - Small fast SRAM memory for immediate access to data

# Inside the Processor

- Apple A5

# Abstractions

- Abstraction helps us deal with complexity
  - Hide lower-level detail
- Instruction set architecture (ISA)
  - The hardware/software interface
- Application binary interface
  - The ISA plus system software interface
- Implementation
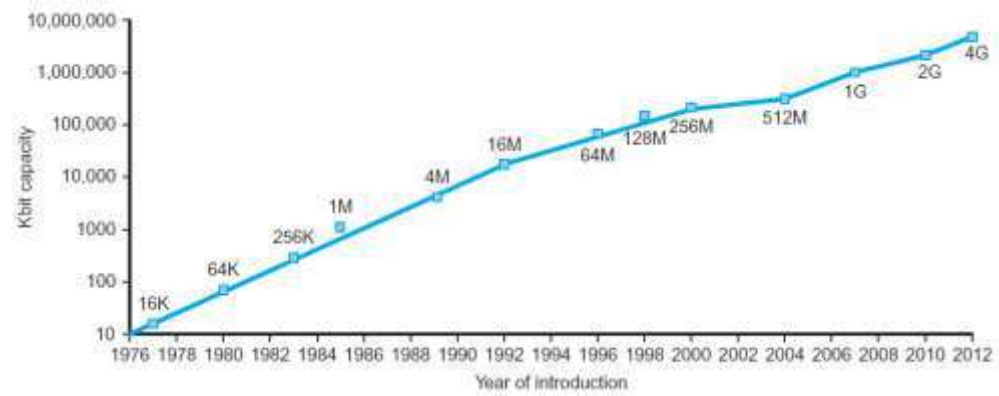  - The details underlying and interface

# A Safe Place for Data

- Volatile main memory
  - Loses instructions and data when power off
- Non-volatile secondary memory
  - Magnetic disk
  - Flash memory
  - Optical disk (CDROM, DVD)

# Technology Trends

- Electronics technology continues to evolve
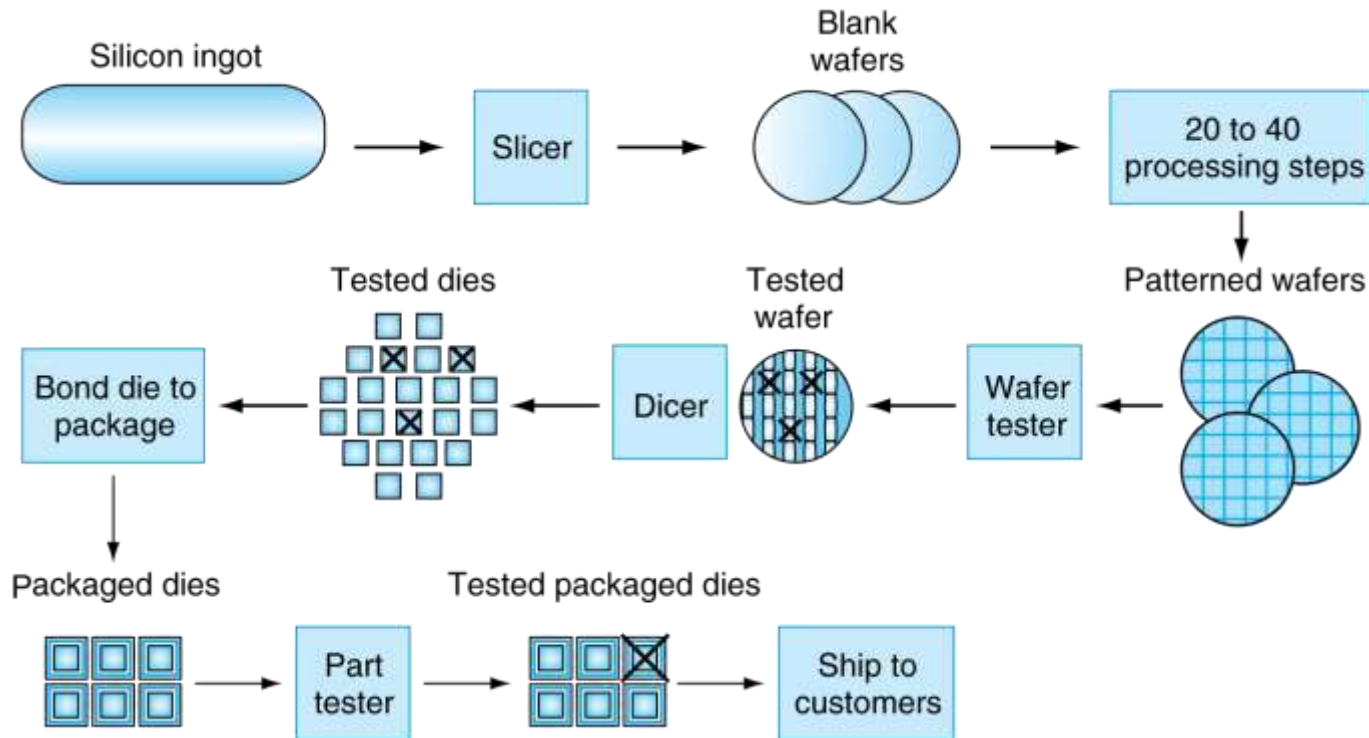    - Increased capacity and performance
    - Reduced cost



DRAM capacity

| Year | Technology | Relative performance/cost | |
|------|------------|---------------------------|---|
| 1951 | Vacuum tube | 1 | |
| 1965 | Transistor | 35 | |
| 1975 | Integrated circuit (IC) | 900 | |
| 1995 | Very large scale IC (VLSI) | 2,400,000 | |
| 2013 | Ultra large scale IC | 250,000,000,000 | |

# Semiconductor Technology

- Silicon:  semiconductor
- Add materials to transform properties:
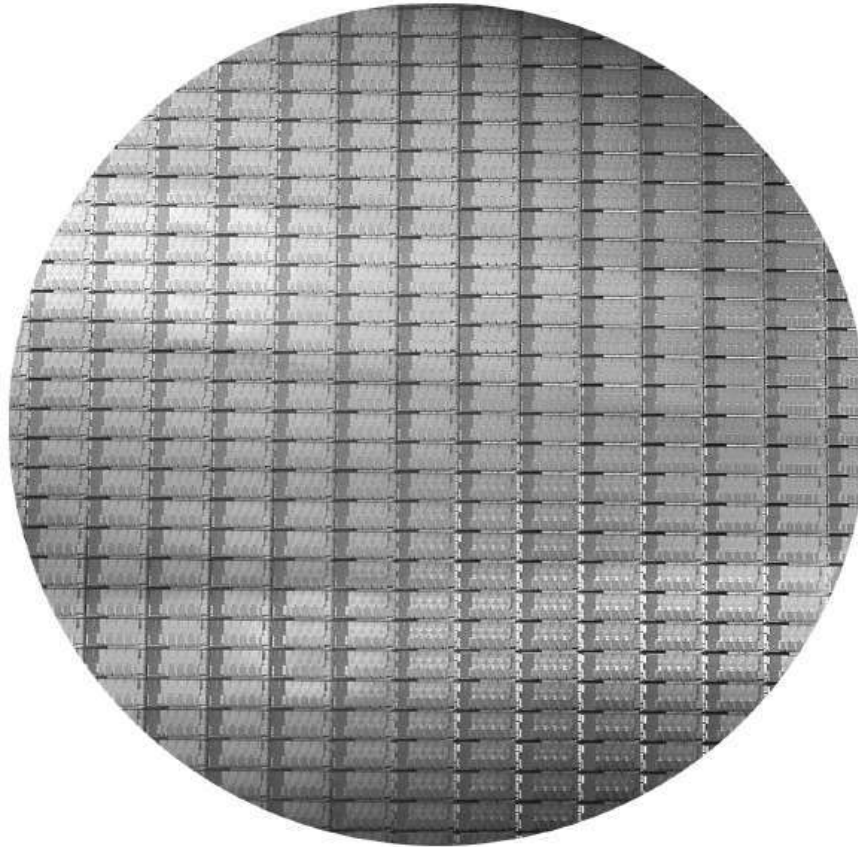  - Conductors
  - Insulators
  - Switch

# Manufacturing ICs



- Yield: proportion of working dies per wafer
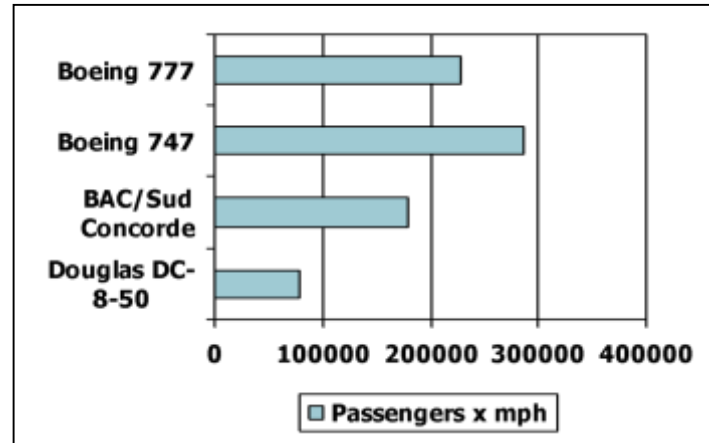
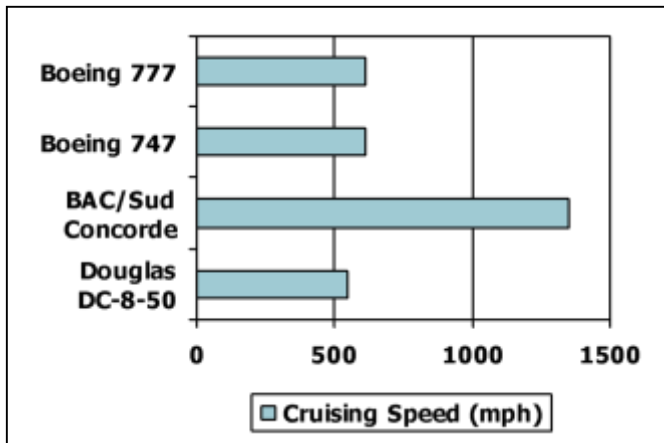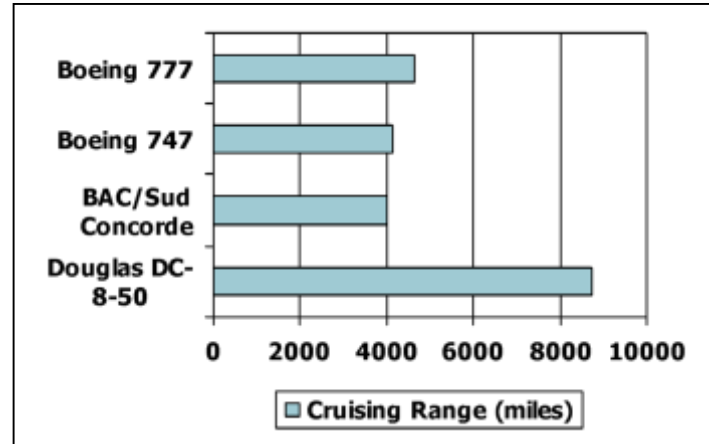# IC Manufacturing Steps

# Intel Core i7 Wafer



- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

# Defining Performance

■  Which airplane has the best performance?

# Response Time and Throughput

- Response time
    - How long it takes to do a task
- Throughput
    - Total work done per unit time
        - e.g., tasks/transactions/… per hour
- How are response time and throughput affected by
    - Replacing the processor with a faster version?
    - Adding more processors?
- We'll focus on response time for now…

# Relative Performance

- Define Performance = 1/Execution Time

- "X is $n$ time faster than Y"

$$\text{Performance}_X / \text{Performance}_Y$$
$$= \text{Execution time}_Y / \text{Execution time}_X = n$$

- Example: time taken to run a program
  - 10s on A, 15s on B
  - Execution Time$_B$ / Execution Time$_A$
    = 15s / 10s = 1.5
  - So A is 1.5 times faster than B

# Measuring Execution Time

- Elapsed time
  - Total response time, including all aspects
    - Processing, I/O, OS overhead, idle time
  - Determines system performance
- CPU time
  - Time spent processing a given job
    - Discounts I/O time, other jobs' shares
  - Comprises user CPU time and system CPU time
  - Different programs are affected differently by CPU and system performance

# CPU Clocking

- Operation of digital hardware governed by a constant-rate clock



- Clock period: duration of a clock cycle
  - e.g., 250ps = 0.25ns = $250 \times 10^{-12}$s
- Clock frequency (rate): cycles per second
  - e.g., 4.0GHz = 4000MHz = $4.0 \times 10^9$Hz

# CPU Time

$$\text{CPU Time} = \text{CPU Clock Cycles} \times \text{Clock Cycle Time}$$

$$= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}$$

- Performance improved by
  - Reducing number of clock cycles
  - Increasing clock rate
  - Hardware designer must often trade off clock rate against cycle count

# CPU Time Example

- Computer A: 2GHz clock, 10s CPU time
- Designing Computer B
  - Aim for 6s CPU time
  - Can do faster clock, but causes 1.2 × clock cycles
- How fast must Computer B clock be?

$$\text{Clock Rate}_B = \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6s}$$

$$\text{Clock Cycles}_A = \text{CPU Time}_A \times \text{Clock Rate}_A$$

$$= 10s \times 2\text{GHz} = 20 \times 10^9$$

$$\text{Clock Rate}_B = \frac{1.2 \times 20 \times 10^9}{6s} = \frac{24 \times 10^9}{6s} = 4\text{GHz}$$

# Instruction Count and CPI

$$\text{Clock Cycles} = \text{Instruction Count} \times \text{Cycles per Instruction}$$

$$\text{CPU Time} = \text{Instruction Count} \times \text{CPI} \times \text{Clock Cycle Time}$$

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

- Instruction Count for a program
  - Determined by program, ISA and compiler
- Average cycles per instruction
  - Determined by CPU hardware
  - If different instructions have different CPI
    - Average CPI affected by instruction mix

# CPI Example

- Computer A: Cycle Time = 250ps, CPI = 2.0
- Computer B: Cycle Time = 500ps, CPI = 1.2
- Same ISA
- Which is faster, and by how much?

$$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$$
$$= I \times 2.0 \times 250ps = I \times 500ps$$

A is faster…

$$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$$
$$= I \times 1.2 \times 500ps = I \times 600ps$$

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600ps}{I \times 500ps} = 1.2$$

…by this much

# CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^{n} (CPI_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$CPI = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^{n} \left( CPI_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}} \right)$$

Relative frequency

# CPI Example

- Alternative compiled code sequences using instructions in classes A, B, C

| Class | A | B | C |
|---|---|---|---|
| CPI for class | 1 | 2 | 3 |
| IC in sequence 1 | 2 | 1 | 2 |
| IC in sequence 2 | 4 | 1 | 1 |

- Sequence 1: IC = 5
  - Clock Cycles
    $= 2 \times 1 + 1 \times 2 + 2 \times 3$
    $= 10$
  - Avg. CPI = 10/5 = 2.0

- Sequence 2: IC = 6
  - Clock Cycles
    $= 4 \times 1 + 1 \times 2 + 1 \times 3$
    $= 9$
  - Avg. CPI = 9/6 = 1.5

# Performance Summary

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- Performance depends on
  - Algorithm: affects IC, possibly CPI
  - Programming language: affects IC, CPI
  - Compiler: affects IC, CPI
  - Instruction set architecture: affects IC, CPI, $T_c$

# Power Trends

- In CMOS IC technology

$$Power = Capacitive\ load \times Voltage^2 \times Frequency$$

| ×30 | | 5V → 1V | ×1000 |

# Reducing Power

- Suppose a new CPU has
    - 85% of capacitive load of old CPU
    - 15% voltage and 15% frequency reduction

$$\frac{P_{new}}{P_{old}} = \frac{C_{old} \times 0.85 \times (V_{old} \times 0.85)^2 \times F_{old} \times 0.85}{C_{old} \times V_{old}^2 \times F_{old}} = 0.85^4 = 0.52$$

- The power wall
    - We can't reduce voltage further
    - We can't remove more heat
- How else can we improve performance?

# Uniprocessor Performance

Constrained by power, instruction-level parallelism, memory latency

# Multiprocessors

- Multicore microprocessors
  - More than one processor per chip
- Requires explicitly parallel programming
  - Compare with instruction level parallelism
    - Hardware executes multiple instructions at once
    - Hidden from the programmer
  - Hard to do
    - Programming for performance
    - Load balancing
    - Optimizing communication and synchronization

# SPEC CPU Benchmark

- Programs used to measure performance
  - Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
  - Develops benchmarks for CPU, I/O, Web, …

- SPEC CPU2006
  - Elapsed time to execute a selection of programs
    - Negligible I/O, so focuses on CPU performance
  - Normalize relative to reference machine
  - Summarize as geometric mean of performance ratios
    - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^{n} \text{Execution time ratio}_i}$$

# CINT2006 for Intel Core i7 920

| Description | Name | Instruction Count x $10^9$ | CPI | Clock cycle time (seconds x $10^{-9}$) | Execution Time (seconds) | Reference Time (seconds) | SPECratio |
|---|---|---|---|---|---|---|---|
| Interpreted string processing | perl | 2252 | 0.60 | 0.376 | 508 | 9770 | 19.2 |
| Block-sorting compression | bzip2 | 2390 | 0.70 | 0.376 | 629 | 9650 | 15.4 |
| GNU C compiler | gcc | 794 | 1.20 | 0.376 | 358 | 8050 | 22.5 |
| Combinatorial optimization | mcf | 221 | 2.66 | 0.376 | 221 | 9120 | 41.2 |
| Go game (AI) | go | 1274 | 1.10 | 0.376 | 527 | 10490 | 19.9 |
| Search gene sequence | hmmer | 2616 | 0.60 | 0.376 | 590 | 9330 | 15.8 |
| Chess game (AI) | sjeng | 1948 | 0.80 | 0.376 | 586 | 12100 | 20.7 |
| Quantum computer simulation | libquantum | 659 | 0.44 | 0.376 | 109 | 20720 | 190.0 |
| Video compression | h264avc | 3793 | 0.50 | 0.376 | 713 | 22130 | 31.0 |
| Discrete event simulation library | omnetpp | 367 | 2.10 | 0.376 | 290 | 6250 | 21.5 |
| Games/path finding | astar | 1250 | 1.00 | 0.376 | 470 | 7020 | 14.9 |
| XML parsing | xalancbmk | 1045 | 0.70 | 0.376 | 275 | 6900 | 25.1 |
| Geometric mean | – | – | – | – | – | – | 25.7 |

# Concluding Remarks

- Cost/performance is improving
    - Due to underlying technology development
- Hierarchical layers of abstraction
    - In both hardware and software
- Instruction set architecture
    - The hardware/software interface
- Execution time: the best performance measure
- Power is a limiting factor
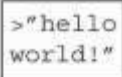    - Use parallelism to improve performance

# Chapter 4

*HARDWARE DESCRIPTION LANGUAGES*

*Digital Design and Computer Architecture*, 2nd Edition

David Money Harris and Sarah L. Harris

ELSEVIER

# Chapter 4 :: Topics

- **Introduction**
- **Combinational Logic**
- **Structural Modeling**
- **Sequential Logic**
- **More Combinational Logic**
- **Finite State Machines**
- **Parameterized Modules**
- **Testbenches**

| | |
|---|---|
| Application Software | >"hello world!" |
| Operating Systems | |
| Architecture | |
| Micro-architecture | |
| Logic | |
| Digital Circuits | |
| Analog Circuits | |
| Devices | |
| Physics | |

ELSEVIER

# Introduction

- ## Hardware description language (HDL):
  - specifies logic function only
  - Computer-aided design (CAD) tool produces or *synthesizes* the optimized gates
- ## Most commercial designs built using HDLs
- ## Two leading HDLs:
  - **SystemVerilog**
    - developed in 1984 by Gateway Design Automation
    - IEEE standard (1364) in 1995
    - Extended in 2005 (IEEE STD 1800-2009)
  - **VHDL 2008**
    - Developed in 1981 by the Department of Defense
    - IEEE standard (1076) in 1987
    - Updated in 2008 (IEEE STD 1076-2008)

# HDL to Gates

- **Simulation**
  - Inputs applied to circuit
  - Outputs checked for correctness
  - Millions of dollars saved by debugging in simulation instead of hardware

- **Synthesis**
  - Transforms HDL code into a *netlist* describing the hardware (i.e., a list of gates and the wires connecting them)

**IMPORTANT:**

When using an HDL, think of the **hardware** the HDL should produce

ELSEVIER

# SystemVerilog Modules



a
b
c

Verilog
Module

y

## Two types of Modules:

- **Behavioral:** describe what a module does

- **Structural:** describe how it is built from simpler modules

# Behavioral SystemVerilog

## SystemVerilog:

```
module example(input  logic a, b, c,
                output logic y);
  assign y = ~a & ~b & ~c | a & ~b & ~c | a & ~b &  c;
endmodule
```

## SystemVerilog:

```
module example(input  logic a, b, c,
                output logic y);
  assign y = ~a & ~b & ~c | a & ~b & ~c | a & ~b &  c;
endmodule
```

## SystemVerilog:

```
module example(input  logic a, b, c,
               output logic y);
   assign y = ~a & ~b & ~c | a & ~b & ~c | a & ~b &  c;
endmodule
```

## Synthesis:



un5_y

y

un8_y

# SystemVerilog Syntax

- ## Case sensitive

  - **Example:** `reset` and `Reset` are not the same signal.

- ## No names that start with numbers

  - **Example:** `2mux` is an invalid name

- ## Whitespace ignored

- ## Comments:

  - `// single line comment`
  - `/* multiline`

    `comment */`

# Structural Modeling - Hierarchy

```systemverilog
module and3(input  logic a, b, c,
            output logic y);
  assign y = a & b & c;
endmodule


module inv(input  logic a,
           output logic y);
  assign y = ~a;
endmodule


module nand3(input  logic a, b, c
             output logic y);
  logic n1;                          // internal signal

  and3 andgate(a, b, c, n1);   // instance of and3
  inv  inverter(n1, y);        // instance of inverter
endmodule
```

# Bitwise Operators

```
module gates(input  logic [3:0]  a, b,
             output logic [3:0] y1, y2, y3, y4, y5);
  /* Five different two-input logic
     gates acting on 4 bit busses */
  assign y1 = a & b;      // AND
  assign y2 = a | b;      // OR
  assign y3 = a ^ b;      // XOR
  assign y4 = ~(a & b); // NAND
  assign y5 = ~(a | b); // NOR
endmodule
```



//      single line comment

/*…*/  multiline comment

# Reduction Operators

```verilog
module and8(input  logic [7:0] a,
            output logic       y);
    assign y = &a;
    // &a is much easier to write than
    // assign y = a[7] & a[6] & a[5] & a[4] &
    //            a[3] & a[2] & a[1] & a[0];
endmodule
```

# Conditional Assignment

```
module mux2(input  logic [3:0] d0, d1,
            input  logic        s,
            output logic [3:0] y);
    assign y = s ? d1 : d0;
endmodule
```



? :    is also called a *ternary operator* because it operates on 3 inputs: `s`, `d1`, and `d0`.

# Internal Variables

```
module fulladder(input  logic a, b, cin,
                 output logic s, cout);
   logic p, g;    // internal nodes

   assign p = a ^ b;
   assign g = a & b;

   assign s = p ^ cin;
   assign cout = g | (p & cin);
endmodule
```

# Precedence

## Order of operations

**Highest**

| | |
|---|---|
| ~ | NOT |
| *, /, % | mult, div, mod |
| +, - | add,sub |
| <<, >> | shift |
| <<<, >>> | arithmetic shift |
| <, <=, >, >= | comparison |
| ==, != | equal, not equal |
| &, ~& | AND, NAND |
| ^, ~^ | XOR, XNOR |
| \|, ~\| | OR, NOR |
| ?: | ternary operator |

**Lowest**

ELSEVIER

# Numbers

## Format: N'Bvalue

**N** = number of bits, **B** = base

**N'B** is optional but recommended (default is decimal)

| Number | # Bits | Base | Decimal Equivalent | Stored |
|--------|--------|------|--------------------|--------|
| 3'b101 | 3 | binary | 5 | 101 |
| 'b11 | unsized | binary | 3 | 00…0011 |
| 8'b11 | 8 | binary | 3 | 00000011 |
| 8'b1010_1011 | 8 | binary | 171 | 10101011 |
| 3'd6 | 3 | decimal | 6 | 110 |
| 6'o42 | 6 | octal | 34 | 100010 |
| 8'hAB | 8 | hexadecimal | 171 | 10101011 |
| 42 | Unsized | decimal | 42 | 00…0101010 |

ELSEVIER

```
assign y = {a[2:1], {3{b[0]}}, a[0], 6'b100_010};
```

**// if y is a 12-bit signal, the above statement produces:**
```
y = a[2] a[1] b[0] b[0] b[0] a[0] 1 0 0 0 1 0
```

```
// underscores (_) are used for formatting only to make
   it easier to read. SystemVerilog ignores them.
```

## SystemVerilog:

```
module mux2_8(input  logic [7:0] d0, d1,
              input  logic       s,
              output logic [7:0] y);

  mux2 lsbmux(d0[3:0], d1[3:0], s, y[3:0]);
  mux2 msbmux(d0[7:4], d1[7:4], s, y[7:4]);
endmodule
```

# Z: Floating Output

## SystemVerilog:

```
module tristate(input  logic [3:0] a,
                input  logic       en,
                output logic [3:0] y);
   assign y = en ? a : 4'bz;
endmodule
```



en

a[3:0]   [3:0]   [3:0]          [3:0]   [3:0]   y[3:0]

y_1[3:0]

ELSEVIER

HARDWARE DESCRIPTION LANGUAGES

# Sequential Logic

- SystemVerilog uses **Idioms** to describe latches, flip-flops and FSMs

- Other coding styles may simulate correctly but produce incorrect hardware

# Always Statement

**General Structure:**

```
always @(sensitivity list)
    statement;
```

Whenever the event in `sensitivity list` occurs, `statement` is executed

# D Flip-Flop

```
module flop(input  logic       clk,
            input  logic [3:0] d,
            output logic [3:0] q);

   always_ff @(posedge clk)
     q <= d;                    // pronounced "q gets d"

endmodule
```



q[3:0]

# Resettable D Flip-Flop

```
module flopr(input   logic        clk,
             input   logic        reset,
             input   logic [3:0] d,
             output logic [3:0] q);

  // synchronous reset
  always_ff @(posedge clk)
    if (reset) q <= 4'b0;
    else       q <= d;

endmodule
```

# Resettable D Flip-Flop

```systemverilog
module flopr(input  logic          clk,
             input  logic          reset,
             input  logic [3:0] d,
             output logic [3:0] q);

   // asynchronous reset
   always_ff @(posedge clk, posedge reset)
     if (reset) q <= 4'b0;
     else       q <= d;

endmodule
```

# D Flip-Flop with Enable

```
module flopren(input  logic        clk,
               input  logic        reset,
               input  logic        en,
               input  logic [3:0]  d,
               output logic [3:0]  q);

  // asynchronous reset and enable
  always_ff @(posedge clk, posedge reset)
    if       (reset) q <= 4'b0;
    else if (en)     q <= d;

endmodule
```

# Latch

```
module latch(input  logic       clk,
             input  logic [3:0] d,
             output logic [3:0] q);

  always_latch
    if (clk) q <= d;

endmodule
```



**Warning**: We don't use latches in this text. But you might write code that inadvertently implies a latch. Check synthesized hardware – if it has latches in it, there's an error.

# Other Behavioral Statements

- Statements that must be inside `always` statements:
  - `if/else`
  - `case,casez`

# Combinational Logic using always

```
// combinational logic using an always statement
module gates(input  logic [3:0] a, b,
             output logic [3:0] y1, y2, y3, y4, y5);
  always_comb            // need begin/end because there is
    begin                // more than one statement in always
      y1 = a & b;     // AND
      y2 = a | b;     // OR
      y3 = a ^ b;     // XOR
      y4 = ~(a & b);  // NAND
      y5 = ~(a | b);  // NOR
    end
endmodule
```

**This hardware could be described with assign statements using fewer lines of code, so it's better to use assign statements in this case.**

# Combinational Logic using case

```
module sevenseg(input  logic [3:0] data,
                output logic [6:0] segments);
  always_comb
    case (data)
      //                          abc_defg
      0: segments =        7'b111_1110;
      1: segments =        7'b011_0000;
      2: segments =        7'b110_1101;
      3: segments =        7'b111_1001;
      4: segments =        7'b011_0011;
      5: segments =        7'b101_1011;
      6: segments =        7'b101_1111;
      7: segments =        7'b111_0000;
      8: segments =        7'b111_1111;
      9: segments =        7'b111_0011;
      default: segments = 7'b000_0000; // required
    endcase
endmodule
```

# Combinational Logic using `case`

- `case` statement implies combinational logic **only if** all possible input combinations described
- Remember to use **default** statement

ELSEVIER

# Combinational Logic using casez

```systemverilog
module priority_casez(input  logic [3:0] a,
                      output logic [3:0] y);

  always_comb
    casez(a)
      4'b1???: y = 4'b1000;   // ? = don't care
      4'b01??: y = 4'b0100;
      4'b001?: y = 4'b0010;
      4'b0001: y = 4'b0001;
      default: y = 4'b0000;
    endcase
endmodule
```

# Blocking vs. Nonblocking Assignment

- <= is **nonblocking** assignment
  - Occurs simultaneously with others

- = is **blocking** assignment
  - Occurs in order it appears in file

```
// Good synchronizer using
// nonblocking assignments
module syncgood(input  logic clk,
                input  logic d,
                output logic q);
  logic n1;
  always_ff @(posedge clk)
    begin
      n1 <= d;  // nonblocking
      q  <= n1; // nonblocking
    end
endmodule
```

```
// Bad synchronizer using
// blocking assignments
module syncbad(input logic  clk,
               input  logic d,
               output logic q);
  logic n1;
  always_ff @(posedge clk)
    begin
      n1 = d;  // blocking
      q  = n1; // blocking
    end
endmodule
```

# Parameterized Modules

**2:1 mux:**

```
module mux2
  #(parameter width = 8)  // name and default value
   (input  logic [width-1:0] d0, d1,
    input  logic             s,
    output logic [width-1:0] y);
  assign y = s ? d1 : d0;
endmodule
```

**Instance with 8-bit bus width (uses default):**

```
  mux2 mux1(d0, d1, s, out);
```

**Instance with 12-bit bus width:**

```
  mux2 #(12) lowmux(d0, d1, s, out);
```

# Testbenches

- HDL that tests another module: *device under test* (dut)

- Not synthesizeable

- Types:
  - Simple
  - Self-checking
  - Self-checking with testvectors

- Write SystemVerilog code to implement the following function in hardware:

$$y = \overline{\overline{b}\,\overline{c}} + a\overline{b}$$

- Name the module `sillyfunction`

# Testbench Example

- Write SystemVerilog code to implement the following function in hardware:

$$y = \overline{b}\,\overline{c} + a\overline{b}$$

```
module sillyfunction(input  logic a, b, c,
                     output logic y);
   assign y = ~b & ~c | a & ~b;
endmodule
```

# Simple Testbench

```
module testbench1();
  logic a, b, c;
  logic y;
  // instantiate device under test
  sillyfunction dut(a, b, c, y);
  // apply inputs one at a time
  initial begin
    a = 0; b = 0; c = 0; #10;
    c = 1; #10;
    b = 1; c = 0; #10;
    c = 1; #10;
    a = 1; b = 0; c = 0; #10;
    c = 1; #10;
    b = 1; c = 0; #10;
    c = 1; #10;
  end
endmodule
```

# Self-checking Testbench

```
module testbench2();
  logic  a, b, c;
  logic y;
  sillyfunction dut(a, b, c, y);  // instantiate dut
  initial begin // apply inputs, check results one at a time
    a = 0; b = 0; c = 0; #10;
    if (y !== 1) $display("000 failed.");
    c = 1; #10;
    if (y !== 0) $display("001 failed.");
    b = 1; c = 0; #10;
    if (y !== 0) $display("010 failed.");
    c = 1; #10;
    if (y !== 0) $display("011 failed.");
    a = 1; b = 0; c = 0; #10;
    if (y !== 1) $display("100 failed.");
    c = 1; #10;
    if (y !== 1) $display("101 failed.");
    b = 1; c = 0; #10;
    if (y !== 0) $display("110 failed.");
    c = 1; #10;
    if (y !== 0) $display("111 failed.");
  end
endmodule
```

# Testbench with Testvectors

- Testvector file: inputs and expected outputs

- Testbench:
    1. Generate clock for assigning inputs, reading outputs
    2. Read testvectors file into array
    3. Assign inputs, expected outputs
    4. Compare outputs with expected outputs and report errors

# Testbench with Testvectors

- ## Testbench clock:
  - assign inputs (on rising edge)
  - compare outputs with expected outputs (on falling edge).



- Testbench clock also used as clock for synchronous sequential circuits

# Testvectors File

- File: `example.tv`

- contains vectors of abc_yexpected

```
000_1
001_0
010_0
011_0
100_1
101_1
110_0
111_0
```

# 1. Generate Clock

```systemverilog
module testbench3();
  logic         clk, reset;
  logic         a, b, c, yexpected;
  logic         y;
  logic [31:0] vectornum, errors;    // bookkeeping variables
  logic [3:0]  testvectors[10000:0]; // array of testvectors

  // instantiate device under test
  sillyfunction dut(a, b, c, y);

  // generate clock
  always     // no sensitivity list, so it always executes
    begin
      clk = 1; #5; clk = 0; #5;
    end
```

# 2. Read Testvectors into Array

```verilog
// at start of test, load vectors and pulse reset

initial
  begin
    $readmemb("example.tv", testvectors);
    vectornum = 0; errors = 0;
    reset = 1; #27; reset = 0;
  end



// Note: $readmemh reads testvector files written in
// hexadecimal
```

# 3. Assign Inputs & Expected Outputs

```verilog
// apply test vectors on rising edge of clk
always @(posedge clk)
  begin
    #1; {a, b, c, yexpected} = testvectors[vectornum];
  end
```

# 4. Compare with Expected Outputs

```verilog
// check results on falling edge of clk
   always @(negedge clk)
    if (~reset) begin // skip during reset
      if (y !== yexpected) begin
        $display("Error: inputs = %b", {a, b, c});
        $display("  outputs = %b (%b expected)",y,yexpected);
        errors = errors + 1;
      end

// Note: to print in hexadecimal, use %h. For example,
//        $display("Error: inputs = %h", {a, b, c});
```

# 4. Compare with Expected Outputs

```verilog
// increment array index and read next testvector
    vectornum = vectornum + 1;
    if (testvectors[vectornum] === 4'bx) begin
        $display("%d tests completed with %d errors",
            vectornum, errors);
      $finish;
    end
  end
endmodule


// === and !== can compare values that are 1, 0, x, or z.
```

# Chapter 5

*Digital Design and Computer Architecture*, 2nd Edition

David Money Harris and Sarah L. Harris

# Chapter 5 :: Topics

- **Introduction**
- **Arithmetic Circuits**
- **Number Systems**
- **Sequential Building Blocks**
- **Memory Arrays**
- **Logic Arrays**

# Introduction

- **Digital building blocks:**
  - Gates, multiplexers, decoders, registers, arithmetic circuits, counters, memory arrays, logic arrays

- **Building blocks demonstrate hierarchy, modularity, and regularity:**
  - Hierarchy of simpler components
  - Well-defined interfaces and functions
  - Regular structure easily extends to different sizes

- **Will use these building blocks in Chapter 7 to build microprocessor**

# 1-Bit Adders

**Half Adder**



| A | B | $C_{out}$ | S |
|---|---|-----------|---|
| 0 | 0 | | |
| 0 | 1 | | |
| 1 | 0 | | |
| 1 | 1 | | |

S   =
$C_{out}$  =

**Full Adder**



| $C_{in}$ | A | B | $C_{out}$ | S |
|----------|---|---|-----------|---|
| 0 | 0 | 0 | | |
| 0 | 0 | 1 | | |
| 0 | 1 | 0 | | |
| 0 | 1 | 1 | | |
| 1 | 0 | 0 | | |
| 1 | 0 | 1 | | |
| 1 | 1 | 0 | | |
| 1 | 1 | 1 | | |

S   =
$C_{out}$ =

# 1-Bit Adders

**Half Adder**



| A | B | $C_{out}$ | S |
|---|---|-----------|---|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |

$S$ =

$C_{out}$ =

**Full Adder**



| $C_{in}$ | A | B | $C_{out}$ | S |
|----------|---|---|-----------|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

$S$ =

$C_{out}$ =

# 1-Bit Adders

**Half Adder**



| A | B | $C_{out}$ | S |
|---|---|-----------|---|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |

$$S = A \oplus B$$
$$C_{out} = AB$$

**Full Adder**



| $C_{in}$ | A | B | $C_{out}$ | S |
|----------|---|---|-----------|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

$$S = A \oplus B \oplus C_{in}$$
$$C_{out} = AB + AC_{in} + BC_{in}$$

# Multibit Adders (CPAs)

- Types of carry propagate adders (CPAs):
  - Ripple-carry            (slow)
  - Carry-lookahead         (fast)
  - Prefix                  (faster)
- Carry-lookahead and prefix adders faster for large adders but require more hardware

**Symbol**

# Subtracter

**Symbol**

**Implementation**

# Comparator: Equality

## Symbol



A        B

$/4$      $/4$

=

Equal

## Implementation



$A_3$
$B_3$

$A_2$
$B_2$

$A_1$
$B_1$

$A_0$
$B_0$

Equal

# Comparator: Less Than



A

B

N

N

-

N

[N-1]

A < B

ELSEVIER

# Arithmetic Logic Unit (ALU)

A      B

N     N

ALU

N

Y

$\frac{}{3}$ F

| $F_{2:0}$ | Function |
|-----------|----------|
| 000 | A & B |
| 001 | A \| B |
| 010 | A + B |
| 011 | not used |
| 100 | A & ~B |
| 101 | A \| ~B |
| 110 | A - B |
| 111 | SLT |

ELSEVIER

# ALU Design



| $F_{2:0}$ | Function |
|-----------|----------|
| 000 | A & B |
| 001 | A \| B |
| 010 | A + B |
| 011 | not used |
| 100 | A & ~B |
| 101 | A \| ~B |
| 110 | A - B |
| 111 | SLT |

ELSEVIER

# Set Less Than (SLT) Example



- Configure 32-bit ALU for SLT operation: $A = 25$ and $B = 32$

# Set Less Than (SLT) Example



- Configure 32-bit ALU for SLT operation: $A = 25$ and $B = 32$
  - $A < B$, so $Y$ should be 32-bit representation of 1 (0x00000001)
  - $F_{2:0} = 111$
    - $F_2 = 1$ (adder acts as subtracter), so 25 - 32 = -7
    - -7 has 1 in the most significant bit ($S_{31} = 1$)
    - $F_{1:0} = 11$ multiplexer selects $Y = S_{31}$ (zero extended) = 0x00000001.

# Shifters

- **Logical shifter:** shifts value to left or right and fills empty spaces with 0's
  - Ex: 11001 >> 2 =
  - Ex: 11001 << 2 =

- **Arithmetic shifter:** same as logical shifter, but on right shift, fills empty spaces with the old most significant bit (msb).
  - Ex: 11001 >>> 2 =
  - Ex: 11001 <<< 2 =

- **Rotator:** rotates bits in a circle, such that bits shifted off one end are shifted into the other end
  - Ex: 11001 ROR 2 =
  - Ex: 11001 ROL 2 =

ELSEVIER

# Shifters

- **Logical shifter:**
  - Ex: 11001 >> 2 = 00110
  - Ex: 11001 << 2 = 00100

- **Arithmetic shifter:**
  - Ex: 11001 >>> 2 = **11**110
  - Ex: 11001 <<< 2 = 00100

- **Rotator:**
  - Ex: 11001 ROR 2 = 01110
  - Ex: 11001 ROL 2 = 00111

ELSEVIER

# Shifter Design

# Shifters as Multipliers, Dividers

- ## $A << N = A \times 2^N$

    - **Example:** $00001 << 2 = 00100$  ($1 \times 2^2 = 4$)
    - **Example:** $11101 << 2 = 10100$  ($-3 \times 2^2 = -12$)

- ## $A >>> N = A \div 2^N$

    - **Example:** $01000 >>> 2 = 00010$  ($8 \div 2^2 = 2$)
    - **Example:** $10000 >>> 2 = 11100$  ($-16 \div 2^2 = -4$)

ELSEVIER

# Number Systems

- Numbers we can represent using binary representations
  - **Positive numbers**
    - Unsigned binary
  - **Negative numbers**
    - Two's complement
    - Sign/magnitude numbers

- What about **fractions**?

ELSEVIER

# Counters

- Increments on each clock edge
- Used to cycle through numbers. For example,
  - 000, 001, 010, 011, 100, 101, 110, 111, 000, 001…
- Example uses:
  - Digital clock displays
  - Program counter: keeps track of current instruction executing

**Symbol**          **Implementation**

# Shift Registers

- Shift a new bit in on each clock edge
- Shift a bit out on each clock edge
- *Serial-to-parallel converter*: converts serial input ($S_{in}$) to parallel output ($Q_{0:N-1}$)

**Symbol:**

**Implementation:**

# Shift Register with Parallel Load

- When *Load* = 1, acts as a normal *N*-bit register

- When *Load* = 0, acts as a shift register

- Now can act as a *serial-to-parallel converter* ($S_{in}$ to $Q_{0:N-1}$) or a *parallel-to-serial converter* ($D_{0:N-1}$ to $S_{out}$)

# Memory Arrays

- Efficiently store large amounts of data

- 3 common types:
  - Dynamic random access memory (DRAM)
  - Static random access memory (SRAM)
  - Read only memory (ROM)

- $M$-bit data value read/ written at each unique $N$-bit address

Address — $N$ → **Array**

$M$ Data

ELSEVIER

# Memory Arrays

- 2-dimensional array of bit cells

- Each bit cell stores one bit

- $N$ address bits and $M$ data bits:

  - $2^N$ rows and $M$ columns

  - **Depth:** number of rows (number of words)

  - **Width:** number of columns (size of word)

  - **Array size:** depth × width = $2^N \times M$

Address $\xrightarrow{\;N\;}$ **Array**

Data ↕ $M$

Address $\xrightarrow{\;2\;}$ **Array**

Data ↕ 3

Address   Data

| | | | |
|---|---|---|---|
| 11 | 0 | 1 | 0 |
| 10 | 1 | 0 | 0 |
| 01 | 1 | 1 | 0 |
| 00 | 0 | 1 | 1 |

depth

width

ELSEVIER

DIGITAL BUILDING BLOCKS

# Memory Array Example

- $2^2 \times 3$-bit array

- Number of words: 4

- Word size: 3-bits

- For example, the 3-bit word stored at address 10 is 100

# Memory Arrays



Address —10⟋→ **1024-word x 32-bit Array** ↕ 32 Data

# Memory Array Bit Cells



(a)  (b)

# Memory Array Bit Cells



bitline

wordline

stored bit

bitline = **0**

wordline = 1

stored bit = 0

bitline = **1**

wordline = 1

stored bit = 1

(a)

bitline = **Z**

wordline = 0

stored bit = 0

bitline = **Z**

wordline = 0

stored bit = 1

(b)

ELSEVIER

DIGITAL BUILDING BLOCKS

# Memory Array

- **Wordline:**
  - like an enable
  - single row in memory array read/written
  - corresponds to unique address
  - only one wordline HIGH at once

# Types of Memory

- Random access memory (RAM): **volatile**
- Read only memory (ROM): **nonvolatile**

# RAM: Random Access Memory

- **Volatile:** loses its data when power off
- Read and written quickly
- Main memory in your computer is RAM (DRAM)

Historically called *random* access memory because any data word accessed as easily as any other (in contrast to sequential access memories such as a tape recorder)

# ROM: Read Only Memory

- **Nonvolatile:** retains data when power off

- Read quickly, but writing is impossible or slow

- Flash memory in cameras, thumb drives, and digital cameras are all ROMs

Historically called *read only* memory because ROMs were written at manufacturing time or by burning fuses. Once ROM was configured, it could not be written again. This is no longer the case for Flash memory and other types of ROMs.

# Types of RAM

- **DRAM** (Dynamic random access memory)
- **SRAM** (Static random access memory)
- Differ in how they store data:
  - DRAM uses a capacitor
  - SRAM uses cross-coupled inverters

# DRAM

- Data bits stored on capacitor

- *Dynamic* because the value needs to be refreshed (rewritten) periodically and after read:
  - Charge leakage from the capacitor degrades the value
  - Reading destroys the stored value

# DRAM

# SRAM

bitline

wordline

stored
bit

$\overline{\text{bitline}}$

bitline

wordline

# Bistable Circuit

- Fundamental building block of other state elements
- Two outputs: $Q$, $\overline{Q}$
- No inputs

# Memory Arrays Review

# ROM: Dot Notation

# Multi-ported Memories

- Port: address/data pair
- 3-ported memory
  - 2 read ports (A1/RD1, A2/RD2)
  - 1 write port (A3/WD3, WE3 enables writing)
- **Register file:** small multi-ported memory

# Write port

# Read ports

# Logic Arrays

- **PLAs** (Programmable logic arrays)
  - AND array followed by OR array
  - Combinational logic only
  - Fixed internal connections

- **FPGAs** (Field programmable gate arrays)
  - Array of Logic Elements (LEs)
  - Combinational and sequential logic
  - Programmable internal connections

ELSEVIER

# Chapter 6

*Digital Design and Computer Architecture*, 2$^{nd}$ Edition

David Money Harris and Sarah L. Harris

# Chapter 6 :: Topics

- **Introduction**

- **Assembly Language**

- **Machine Language**

- **Programming**

- **Addressing Modes**

- **Lights, Camera, Action: Compiling, Assembling, & Loading**

- **Odds and Ends**

# Introduction

- Jumping up a few levels of abstraction

- **Architecture:** programmer's view of computer
  - Defined by instructions & operand locations

- **Microarchitecture:** how to implement an architecture in hardware (covered in Chapter 7)

| Layer | Examples |
|---|---|
| Application Software | programs |
| Operating Systems | device drivers |
| Architecture | instructions registers |
| Micro-architecture | datapaths controllers |
| Logic | adders memories |
| Digital Circuits | AND gates NOT gates |
| Analog Circuits | amplifiers filters |
| Devices | transistors diodes |
| Physics | electrons |

ELSEVIER

# Assembly Language

- **Instructions:** commands in a computer's language
  - **Assembly language:** human-readable format of instructions
  - **Machine language:** computer-readable format (1's and 0's)
- **MIPS** architecture:
  - Developed by John Hennessy and his colleagues at Stanford and in the 1980's.
  - Used in many commercial systems, including Silicon Graphics, Nintendo, and Cisco

Once you've learned one architecture, it's easy to learn others

ELSEVIER

# John Hennessy

- President of Stanford University
- Professor of Electrical Engineering and Computer Science at Stanford since 1977
- Coinvented the Reduced Instruction Set Computer (RISC) with David Patterson
- Developed the MIPS architecture at Stanford in 1984 and cofounded MIPS Computer Systems
- As of 2004, over 300 million MIPS microprocessors have been sold

# Architecture Design Principles

Underlying design principles, as articulated by Hennessy and Patterson:

1. **Simplicity favors regularity**
2. **Make the common case fast**
3. **Smaller is faster**
4. **Good design demands good compromises**

# Instructions: Addition

**C Code**

```
a = b + c;
```

**MIPS assembly code**

```
add a, b, c
```

- **add:** mnemonic indicates operation to perform
- **b, c:** source operands (on which the operation is performed)
- **a:** destination operand (to which the result is written)

# Instructions: Subtraction

- Similar to addition - only mnemonic changes

| C Code | MIPS assembly code |
|--------|--------------------|
| `a = b - c;` | `sub a, b, c` |

- **sub:** mnemonic
- **b, c:** source operands
- **a:** destination operand

ELSEVIER

# Design Principle 1

## Simplicity favors regularity

- Consistent instruction format
- Same number of operands (two sources and one destination)
- easier to encode and handle in hardware

# Multiple Instructions

- More complex code is handled by multiple MIPS instructions.

**C Code**

```
a = b + c - d;
```

**MIPS assembly code**

```
add t, b, c  # t = b + c
sub a, t, d  # a = t - d
```

# Design Principle 2

## Make the common case fast

- MIPS includes only simple, commonly used instructions

- Hardware to decode and execute instructions can be simple, small, and fast

- More complex instructions (that are less common) performed using multiple simple instructions

- MIPS is a *reduced instruction set computer* **(RISC)**, with a small number of simple instructions

- Other architectures, such as Intel's x86, are *complex instruction set computers* **(CISC)**

ELSEVIER

# Operands

- Operand location: physical location in computer
  - Registers
  - Memory
  - Constants (also called *immediates*)

# Operands: Registers

- MIPS has 32 32-bit registers
- Registers are faster than memory
- MIPS called "32-bit architecture" because it operates on 32-bit data

# Design Principle 3

## Smaller is Faster

- MIPS includes only a small number of registers

# MIPS Register Set

| Name | Register Number | Usage |
|------|-----------------|-------|
| `$0` | 0 | the constant value 0 |
| `$at` | 1 | assembler temporary |
| `$v0-$v1` | 2-3 | Function return values |
| `$a0-$a3` | 4-7 | Function arguments |
| `$t0-$t7` | 8-15 | temporaries |
| `$s0-$s7` | 16-23 | saved variables |
| `$t8-$t9` | 24-25 | more temporaries |
| `$k0-$k1` | 26-27 | OS temporaries |
| `$gp` | 28 | global pointer |
| `$sp` | 29 | stack pointer |
| `$fp` | 30 | frame pointer |
| `$ra` | 31 | Function return address |

# Operands: Registers

- Registers:
  - $ before name
  - Example: $0, "register zero", "dollar zero"
- Registers used for specific purposes:
  - $0 always holds the constant value 0.
  - the *saved registers*, $s0-$s7, used to hold variables
  - the *temporary registers*, $t0 - $t9, used to hold intermediate values during a larger computation
  - Discuss others later

# Instructions with Registers

- Revisit add instruction

**C Code**

a = b + c

**MIPS assembly code**
```
# $s0 = a, $s1 = b, $s2 = c
add $s0, $s1, $s2
```

# Operands: Memory

- Too much data to fit in only 32 registers
- Store more data in memory
- Memory is large, but slow
- Commonly used variables kept in registers

ELSEVIER

# Word-Addressable Memory

- Each 32-bit data word has a unique address

| Word Address | Data | |
|:---:|:---:|:---|
| ⋮ | ⋮ | ⋮ |
| 00000003 | 4 0 F 3 0 7 8 8 | Word 3 |
| 00000002 | 0 1 E E 2 8 4 2 | Word 2 |
| 00000001 | F 2 F 1 A C 0 7 | Word 1 |
| 00000000 | A B C D E F 7 8 | Word 0 |

ELSEVIER

# Reading Word-Addressable Memory

- Memory read called *load*

- **Mnemonic:** *load word* (`lw`)

- **Format:**

  `lw $s0, 5($t1)`

- **Address calculation:**
  - add *base address* (`$t1`) to the *offset* (5)
  - address = (`$t1` + 5)

- **Result:**
  - `$s0` holds the value at address (`$t1` + 5)

  **Any register** may be used as base address

# Reading Word-Addressable Memory

- **Example:** read a word of data at memory address 1 into $s3
  - address = ($0 + 1) = 1
  - $s3 = 0xF2F1AC07 after load

**Assembly code**

```
lw $s3, 1($0)   # read memory word 1 into $s3
```

| Word Address | Data | |
|---|---|---|
| ⋮ | ⋮ | ⋮ |
| 00000003 | 4 0 F 3 0 7 8 8 | Word 3 |
| 00000002 | 0 1 E E 2 8 4 2 | Word 2 |
| 00000001 | F 2 F 1 A C 0 7 | Word 1 |
| 00000000 | A B C D E F 7 8 | Word 0 |

ELSEVIER

# Writing Word-Addressable Memory

- Memory write are called *store*
- **Mnemonic:** *store word* (`sw`)

# Writing Word-Addressable Memory

- **Example:** Write (store) the value in $t4 into memory address 7
  - add the base address ($0) to the offset (0x7)
  - address: ($0 + 0x7) = 7

  Offset can be written in decimal (default) or hexadecimal

  **Assembly code**

  ```
  sw $t4, 0x7($0)  # write the value in $t4
                   # to memory word 7
  ```

  Word Address          Data

  ⋮                     ⋮              ⋮

  | 00000003 | 4 0 F 3 0 7 8 8 | Word 3 |
  | 00000002 | 0 1 E E 2 8 4 2 | Word 2 |
  | 00000001 | F 2 F 1 A C 0 7 | Word 1 |
  | 00000000 | A B C D E F 7 8 | Word 0 |

# Byte-Addressable Memory

- Each data byte has unique address
- Load/store words or single bytes: load byte (`lb`) and store byte (`sb`)
- 32-bit word = 4 bytes, so word address increments by 4

| Word Address | Data | |
|---|---|---|
| ⋮ | ⋮ | ⋮ |
| 0000000C | 4 0 F 3 0 7 8 8 | Word 3 |
| 00000008 | 0 1 E E 2 8 4 2 | Word 2 |
| 00000004 | F 2 F 1 A C 0 7 | Word 1 |
| 00000000 | A B C D E F 7 8 | Word 0 |

width = 4 bytes

ELSEVIER

# Reading Byte-Addressable Memory

- The address of a memory word must now be multiplied by 4. For example,
  - the address of memory word 2 is $2 \times 4 = 8$
  - the address of memory word 10 is $10 \times 4 = 40$ (0x28)

- **MIPS is byte-addressed, not word-addressed**

ELSEVIER

# Reading Byte-Addressable Memory

- **Example:** Load a word of data at memory address 4 into $s3.

- $s3 holds the value 0xF2F1AC07 after load

**MIPS assembly code**

```
lw $s3, 4($0)  # read word at address 4 into $s3
```

| Word Address | | | Data | | | | |
|---|---|---|---|---|---|---|---|
| ⋮ | | | ⋮ | | | ⋮ | |
| 0000000C | 4 0 | F 3 | 0 7 | 8 8 | Word 3 |
| 00000008 | 0 1 | E E | 2 8 | 4 2 | Word 2 |
| 00000004 | F 2 | F 1 | A C | 0 7 | Word 1 |
| 00000000 | A B | C D | E F | 7 8 | Word 0 |

width = 4 bytes

ELSEVIER

# Writing Byte-Addressable Memory

- **Example:** stores the value held in $t7 into memory address 0x2C (44)

**MIPS assembly code**

```
sw $t7, 44($0)   # write $t7 into address 44
```

| Word Address | Data | |
|---|---|---|
| ⋮ | ⋮ | ⋮ |
| 0000000C | 4 0  F 3  0 7  8 8 | Word 3 |
| 00000008 | 0 1  E E  2 8  4 2 | Word 2 |
| 00000004 | F 2  F 1  A C  0 7 | Word 1 |
| 00000000 | A B  C D  E F  7 8 | Word 0 |

width = 4 bytes

# Big-Endian & Little-Endian Memory

- How to number bytes within a word?
- **Little-endian:** byte numbers start at the little (least significant) end
- **Big-endian:** byte numbers start at the big (most significant) end
- **Word address** is the **same** for big- or little-endian

| Big-Endian | | Little-Endian |
|:---:|:---:|:---:|
| Byte Address | Word Address | Byte Address |

| C | D | E | F | | C | | F | E | D | C |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 | 9 | A | B | | 8 | | B | A | 9 | 8 |
| 4 | 5 | 6 | 7 | | 4 | | 7 | 6 | 5 | 4 |
| 0 | 1 | 2 | 3 | | 0 | | 3 | 2 | 1 | 0 |

MSB                LSB                MSB                LSB

ELSEVIER

# Big-Endian & Little-Endian Memory

- Jonathan Swift's *Gulliver's Travels*: the Little-Endians broke their eggs on the little end of the egg and the Big-Endians broke their eggs on the big end

- It doesn't really matter which addressing type used – except when the two systems need to share data!

**Big-Endian**

| Byte Address |
| --- |
| ⋮ |

| C | D | E | F |
|---|---|---|---|
| 8 | 9 | A | B |
| 4 | 5 | 6 | 7 |
| 0 | 1 | 2 | 3 |

MSB        LSB

**Word Address**

| ⋮ |
| --- |
| C |
| 8 |
| 4 |
| 0 |

**Little-Endian**

| Byte Address |
| --- |
| ⋮ |

| F | E | D | C |
|---|---|---|---|
| B | A | 9 | 8 |
| 7 | 6 | 5 | 4 |
| 3 | 2 | 1 | 0 |

MSB        LSB

ELSEVIER

# Big-Endian & Little-Endian Example

- Suppose `$t0` initially contains 0x23456789

- After following code runs on big-endian system, what value is `$s0`?

- In a little-endian system?

```
sw $t0, 0($0)
lb $s0, 1($0)
```

# Big-Endian & Little-Endian Example

- Suppose `$t0` initially contains 0x23456789

- After following code runs on big-endian system, what value is `$s0`?

- In a little-endian system?

```
sw $t0, 0($0)
lb $s0, 1($0)
```

- Big-endian:    0x00000045

- Little-endian: 0x00000067

### Big-Endian                    Little-Endian

|              |     |     |     |     | Word    |     |     |     |     |              |
|--------------|-----|-----|-----|-----|---------|-----|-----|-----|-----|--------------|
| Byte Address | 0   | 1   | 2   | 3   | Address | 3   | 2   | 1   | 0   | Byte Address |
| Data Value   | 23  | 45  | 67  | 89  | 0       | 23  | 45  | 67  | 89  | Data Value   |
|              | MSB |     |     | LSB |         | MSB |     |     | LSB |              |

ELSEVIER

# Design Principle 4

**Good design demands good compromises**

- Multiple instruction formats allow flexibility
  - `add`, `sub`: use 3 register operands
  - `lw`, `sw`: use 2 register operands and a constant

- Number of instruction formats kept small
  - to adhere to design principles 1 and 3 (simplicity favors regularity and smaller is faster).

# Operands: Constants/Immediates

- `lw` and `sw` use constants or *immediates*
- *immediate*ly available from instruction
- 16-bit two's complement number
- `addi`: add immediate
- Subtract immediate (`subi`) necessary?

**C Code**

```
a = a + 4;
b = a – 12;
```

**MIPS assembly code**

```
# $s0 = a, $s1 = b
addi $s0, $s0, 4
addi $s1, $s0, -12
```

# Machine Language

- Binary representation of instructions

- Computers only understand 1's and 0's

- 32-bit instructions
  - Simplicity favors regularity: 32-bit data & instructions

- 3 instruction formats:
  - **R-Type:** register operands
  - **I-Type:** immediate operand
  - **J-Type:** for jumping (discuss later)

# R-Type

- *Register-type*
- 3 register operands:
    - `rs, rt`: source registers
    - `rd`: destination register
- Other fields:
    - `op`: the *operation code* or *opcode* (0 for R-type instructions)
    - `funct`: the *function*

        with opcode, tells computer what operation to perform
    - `shamt`: the *shift amount* for shift instructions, otherwise it's 0

## R-Type

| op | rs | rt | rd | shamt | funct |
|----|----|----|----|-------|-------|
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

© Digital Design and Computer Architecture, 2nd Edition, 2012

ELSEVIER

# R-Type Examples

## Assembly Code

```
add $s0, $s1, $s2

sub $t0, $t3, $t5
```

## Field Values

| op | rs | rt | rd | shamt | funct |
|----|----|----|----|-------|-------|
| 0 | 17 | 18 | 16 | 0 | 32 |
| 0 | 11 | 13 | 8 | 0 | 34 |
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

## Machine Code

| op | rs | rt | rd | shamt | funct | |
|----|----|----|----|-------|-------|--|
| 000000 | 10001 | 10010 | 10000 | 00000 | 100000 | (0x02328020) |
| 000000 | 01011 | 01101 | 01000 | 00000 | 100010 | (0x016D4022) |
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits | |

**Note** the order of registers in the assembly code:

```
add rd, rs, rt
```

ELSEVIER

# I-Type

- *Immediate-type*
- 3 operands:
  - `rs, rt`: register operands
  - `imm`: 16-bit two's complement immediate
- Other fields:
  - `op`: the opcode
  - Simplicity favors regularity: all instructions have opcode
  - Operation is completely determined by opcode

## I-Type

| op | rs | rt | imm |
|----|----|----|-----|
| 6 bits | 5 bits | 5 bits | 16 bits |

# I-Type Examples

## Assembly Code

```
addi $s0, $s1, 5

addi $t0, $s3, -12

lw   $t2, 32($0)

sw   $s1,  4($t1)
```

## Field Values

| op | rs | rt | imm |
|----|----|----|-----|
| 8 | 17 | 16 | 5 |
| 8 | 19 | 8 | -12 |
| 35 | 0 | 10 | 32 |
| 43 | 9 | 17 | 4 |

6 bits   5 bits   5 bits   16 bits

**Note** the differing order of registers in assembly and machine codes:

```
addi rt, rs, imm

lw   rt, imm(rs)

sw   rt, imm(rs)
```

## Machine Code

| op | rs | rt | imm | |
|----|----|----|-----|---|
| 001000 | 10001 | 10000 | 0000 0000 0000 0101 | (0x22300005) |
| 001000 | 10011 | 01000 | 1111 1111 1111 0100 | (0x2268FFF4) |
| 100011 | 00000 | 01010 | 0000 0000 0010 0000 | (0x8C0A0020) |
| 101011 | 01001 | 10001 | 0000 0000 0000 0100 | (0xAD310004) |

6 bits   5 bits   5 bits   16 bits

# Machine Language: J-Type

- *Jump-type*
- 26-bit address operand (`addr`)
- Used for jump instructions (`j`)

## J-Type

| op | addr |
|----|------|
| 6 bits | 26 bits |

# Review: Instruction Formats

## R-Type

| op | rs | rt | rd | shamt | funct |
|---|---|---|---|---|---|
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

## I-Type

| op | rs | rt | imm |
|---|---|---|---|
| 6 bits | 5 bits | 5 bits | 16 bits |

## J-Type

| op | addr |
|---|---|
| 6 bits | 26 bits |

# Power of the Stored Program

- 32-bit instructions & data stored in memory
- Sequence of instructions: only difference between two applications
- To run a new program:
  - No rewiring required
  - Simply store new program in memory
- Program Execution:
  - Processor *fetches* (reads) instructions from memory in sequence
  - Processor performs the specified operation

# The Stored Program

### Assembly Code

```
lw    $t2, 32($0)
add   $s0, $s1, $s2
addi  $t0, $s3, -12
sub   $t0, $t3, $t5
```

### Machine Code

```
0x8C0A0020
0x02328020
0x2268FFF4
0x016D4022
```

### Stored Program

| Address | Instructions |
|---------|--------------|
| ⋮ | ⋮ |
| 0040000C | 0 1 6 D 4 0 2 2 |
| 00400008 | 2 2 6 8 F F F 4 |
| 00400004 | 0 2 3 2 8 0 2 0 |
| 00400000 | 8 C 0 A 0 0 2 0 ← PC |
| ⋮ | ⋮ |

Main Memory

**Program Counter (PC):** keeps track of current instruction

ELSEVIER

# Interpreting Machine Code

- ## Start with opcode: tells how to parse rest
- ## If opcode all 0's
  - R-type instruction
  - Function bits tell operation
- ## Otherwise
  - opcode tells operation

### Machine Code

**(0x2237FFF1)**

| op | rs | rt | imm |
|---|---|---|---|
| 001000 | 10001 | 10111 | 1111 1111 1111 0001 |
| 2   2 | 3   7 | F   F | F   1 |

**(0x02F34022)**

| op | rs | rt | rd | shamt | funct |
|---|---|---|---|---|---|
| 000000 | 10111 | 10011 | 01000 | 00000 | 100010 |
| 0   2 | F   3 | 4   0 | 2   2 | | |

### Field Values

| op | rs | rt | imm |
|---|---|---|---|
| 8 | 17 | 23 | -15 |

| op | rs | rt | rd | shamt | funct |
|---|---|---|---|---|---|
| 0 | 23 | 19 | 8 | 0 | 34 |

### Assembly Code

`addi $s7, $s1, -15`

`sub $t0, $s7, $s3`

ELSEVIER

# Programming

- High-level languages:
  - e.g., C, Java, Python
  - Written at higher level of abstraction
- Common high-level software constructs:
  - if/else statements
  - for loops
  - while loops
  - arrays
  - function calls

# Ada Lovelace, 1815-1852

- Wrote the first computer program

- Her program calculated the Bernoulli numbers on Charles Babbage's Analytical Engine

- She was the only legitimate child of the poet Lord Byron

# Logical Instructions

- ## **and, or, xor, nor**
  - and: useful for **masking** bits
    - Masking all but the least significant byte of a value:
      0xF234012F AND 0x000000FF = 0x0000002F
  - or: useful for **combining** bit fields
    - Combine 0xF2340000 with 0x000012BC:
      0xF2340000 OR 0x000012BC = 0xF23412BC
  - nor: useful for **inverting** bits:
    - A NOR $0 = NOT A

- ## **andi, ori, xori**
  - 16-bit immediate is zero-extended (*not* sign-extended)
  - nori not needed

ELSEVIER

# Logical Instructions Example 1

## Source Registers

| $s1 | 1111 | 1111 | 1111 | 1111 | 0000 | 0000 | 0000 | 0000 |
|------|------|------|------|------|------|------|------|------|
| $s2 | 0100 | 0110 | 1010 | 0001 | 1111 | 0000 | 1011 | 0111 |

## Assembly Code

```
and $s3, $s1, $s2
or  $s4, $s1, $s2
xor $s5, $s1, $s2
nor $s6, $s1, $s2
```

## Result

| $s3 | | | | | | | | |
|------|------|------|------|------|------|------|------|------|
| $s4 | | | | | | | | |
| $s5 | | | | | | | | |
| $s6 | | | | | | | | |

ELSEVIER

# Logical Instructions Example 1

## Source Registers

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $s1 | 1111 | 1111 | 1111 | 1111 | 0000 | 0000 | 0000 | 0000 |
| $s2 | 0100 | 0110 | 1010 | 0001 | 1111 | 0000 | 1011 | 0111 |

## Assembly Code

```
and $s3, $s1, $s2
or  $s4, $s1, $s2
xor $s5, $s1, $s2
nor $s6, $s1, $s2
```

## Result

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $s3 | 0100 | 0110 | 1010 | 0001 | 0000 | 0000 | 0000 | 0000 |
| $s4 | 1111 | 1111 | 1111 | 1111 | 1111 | 0000 | 1011 | 0111 |
| $s5 | 1011 | 1001 | 0101 | 1110 | 1111 | 0000 | 1011 | 0111 |
| $s6 | 0000 | 0000 | 0000 | 0000 | 0000 | 1111 | 0100 | 1000 |

ELSEVIER

# Logical Instructions Example 2

**Source Values**

| $s1 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 1111 | 1111 |
|------|------|------|------|------|------|------|------|------|

| imm | 0000 | 0000 | 0000 | 0000 | 1111 | 1010 | 0011 | 0100 |
|------|------|------|------|------|------|------|------|------|

← zero-extended →

**Assembly Code**

```
andi $s2, $s1, 0xFA34
ori  $s3, $s1, 0xFA34
xori $s4, $s1, 0xFA34
```

**Result**

| $s2 | | | | | | | | |
|------|--|--|--|--|--|--|--|--|
| $s3 | | | | | | | | |
| $s4 | | | | | | | | |

ELSEVIER

# Logical Instructions Example 2

**Source Values**

| $s1 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 1111 | 1111 |
|------|------|------|------|------|------|------|------|------|

| imm | 0000 | 0000 | 0000 | 0000 | 1111 | 1010 | 0011 | 0100 |
|-----|------|------|------|------|------|------|------|------|

← zero-extended →

**Assembly Code**

```
andi $s2, $s1, 0xFA34
ori  $s3, $s1, 0xFA34
xori $s4, $s1, 0xFA34
```

**Result**

| $s2 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0011 | 0100 |
|------|------|------|------|------|------|------|------|------|
| $s3 | 0000 | 0000 | 0000 | 0000 | 1111 | 1010 | 1111 | 1111 |
| $s4 | 0000 | 0000 | 0000 | 0000 | 1111 | 1010 | 1100 | 1011 |

# Shift Instructions

- `sll:` shift left logical
  - **Example:** `sll $t0, $t1, 5  # $t0 <= $t1 << 5`
- `srl:` shift right logical
  - **Example:** `srl $t0, $t1, 5  # $t0 <= $t1 >> 5`
- `sra:` shift right arithmetic
  - **Example:** `sra $t0, $t1, 5  # $t0 <= $t1 >>> 5`

# Variable Shift Instructions

- `sllv:` shift left logical variable
  - **Example:** `sllv $t0, $t1, $t2 # $t0 <= $t1 << $t2`
- `srlv:` shift right logical variable
  - **Example:** `srlv $t0, $t1, $t2 # $t0 <= $t1 >> $t2`
- `srav:` shift right arithmetic variable
  - **Example**: `srav $t0, $t1, $t2 # $t0 <= $t1 >>> $t2`

# Shift Instructions

## Assembly Code

sll $t0, $s1, 2

srl $s2, $s1, 2

sra $s3, $s1, 2

## Field Values

| op | rs | rt | rd | shamt | funct |
|---|---|---|---|---|---|
| 0 | 0 | 17 | 8 | 2 | 0 |
| 0 | 0 | 17 | 18 | 2 | 2 |
| 0 | 0 | 17 | 19 | 2 | 3 |
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

## Machine Code

| op | rs | rt | rd | shamt | funct | |
|---|---|---|---|---|---|---|
| 000000 | 00000 | 10001 | 01000 | 00010 | 000000 | (0x00114080) |
| 000000 | 00000 | 10001 | 10010 | 00010 | 000010 | (0x00119082) |
| 000000 | 00000 | 10001 | 10011 | 00010 | 000011 | (0x00119883) |
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits | |

ELSEVIER

# Generating Constants

- ## 16-bit constants using `addi`:

**C Code**
```
// int is a 32-bit signed word
int a = 0x4f3c;
```

**MIPS assembly code**
```
# $s0 = a
addi $s0, $0, 0x4f3c
```

- ## 32-bit constants using load upper immediate (`lui`) and `ori`:

**C Code**
```
int a = 0xFEDC8765;
```

**MIPS assembly code**
```
# $s0 = a
lui $s0, 0xFEDC
ori $s0, $s0, 0x8765
```

# Multiplication, Division

- Special registers: `lo, hi`
- $32 \times 32$ multiplication, 64 bit result
  - `mult $s0, $s1`
  - Result in `{hi,lo}`
- 32-bit division, 32-bit quotient, remainder
  - `div $s0, $s1`
  - Quotient in `lo`
  - Remainder in `hi`
- Moves from `lo/hi` special registers
  - `mflo $s2`
  - `mfhi $s3`

# Branching

- Execute instructions out of sequence

- Types of branches:

  - **Conditional**

    - branch if equal (`beq`)
    - branch if not equal (`bne`)

  - **Unconditional**

    - jump (`j`)
    - jump register (`jr`)
    - jump and link (`jal`)

# Review: The Stored Program

Assembly Code

```
lw   $t2, 32($0)
add  $s0, $s1, $s2
addi $t0, $s3, -12
sub  $t0, $t3, $t5
```

Machine Code

```
0x8C0A0020
0x02328020
0x2268FFF4
0x016D4022
```

Stored Program

| Address | Instructions |
|---------|--------------|
| ⋮ | ⋮ |
| 0040000C | 0 1 6 D 4 0 2 2 |
| 00400008 | 2 2 6 8 F F F 4 |
| 00400004 | 0 2 3 2 8 0 2 0 |
| 00400000 | 8 C 0 A 0 0 2 0 ← PC |
| ⋮ | ⋮ |

Main Memory

# Conditional Branching (beq)

## # MIPS assembly

```
addi $s0, $0, 4          # $s0 = 0 + 4 = 4
addi $s1, $0, 1          # $s1 = 0 + 1 = 1
sll  $s1, $s1, 2         # $s1 = 1 << 2 = 4
beq  $s0, $s1, target    # branch is taken
addi $s1, $s1, 1         # not executed
sub  $s1, $s1, $s0       # not executed

target:                  # label
add  $s1, $s1, $s0       # $s1 = 4 + 4 = 8
```

**Labels** indicate instruction location. They can't be reserved words and must be followed by colon (:)

ELSEVIER

# The Branch Not Taken (bne)

## # MIPS assembly

```
addi      $s0, $0, 4              # $s0 = 0 + 4 = 4
addi      $s1, $0, 1              # $s1 = 0 + 1 = 1
sll       $s1, $s1, 2            # $s1 = 1 << 2 = 4
bne       $s0, $s1, target      # branch not taken
addi      $s1, $s1, 1            # $s1 = 4 + 1 = 5
sub       $s1, $s1, $s0         # $s1 = 5 - 4 = 1

target:
add       $s1, $s1, $s0         # $s1 = 1 + 4 = 5
```

# Unconditional Branching (`j`)

## # MIPS assembly

```
addi $s0, $0, 4                # $s0 = 4
addi $s1, $0, 1                # $s1 = 1
j         target               # jump to target
sra       $s1, $s1, 2          # not executed
addi      $s1, $s1, 1          # not executed
sub       $s1, $s1, $s0        # not executed

target:
add       $s1, $s1, $s0        # $s1 = 1 + 4 = 5
```

# Unconditional Branching (`jr`)

## # MIPS assembly

| | |
|---|---|
| **0x00002000** | `addi $s0, $0, 0x2010` |
| **0x00002004** | `jr   $s0` |
| **0x00002008** | `addi $s1, $0, 1` |
| **0x0000200C** | `sra  $s1, $s1, 2` |
| **0x00002010** | `lw   $s3, 44($s1)` |

`jr` is an **R-type** instruction.

# High-Level Code Constructs

- `if` statements

- `if/else` statements

- `while` loops

- `for` loops

# If Statement

**C Code**

```
if (i == j)
  f = g + h;

f = f – i;
```

**MIPS assembly code**

```
# $s0 = f, $s1 = g, $s2 = h
# $s3 = i, $s4 = j
```

# If Statement

**C Code**

```
if (i == j)
   f = g + h;

f = f - i;
```

**MIPS assembly code**

```
# $s0 = f, $s1 = g, $s2 = h
# $s3 = i, $s4 = j
      bne $s3, $s4, L1
      add $s0, $s1, $s2

L1: sub $s0, $s0, $s3
```

Assembly tests opposite case (`i != j`) of high-level code (`i == j`)

# If/Else Statement

**C Code**

**MIPS assembly code**

```
if (i == j)
  f = g + h;
else
  f = f - i;
```

# If/Else Statement

**C Code**

```
if (i == j)
  f = g + h;
else
  f = f - i;
```

**MIPS assembly code**

```
# $s0 = f, $s1 = g, $s2 = h
# $s3 = i, $s4 = j
        bne $s3, $s4, L1
        add $s0, $s1, $s2
        j   done
L1:    sub $s0, $s0, $s3
done:
```

# While Loops

**C Code**

```
// determines the power
// of x such that 2ˣ = 128
int pow = 1;
int x   = 0;

while (pow != 128) {
  pow = pow * 2;
  x = x + 1;
}
```

**MIPS assembly code**

Assembly tests for the opposite case (`pow == 128`) of the C code (`pow != 128`).

# While Loops

**C Code**

```
// determines the power
// of x such that 2ˣ = 128
int pow = 1;
int x   = 0;

while (pow != 128) {
  pow = pow * 2;
  x = x + 1;
}
```

**MIPS assembly code**

```
# $s0 = pow, $s1 = x

        addi $s0, $0, 1
        add  $s1, $0, $0
        addi $t0, $0, 128
while: beq  $s0, $t0, done
        sll  $s0, $s0, 1
        addi $s1, $s1, 1
        j    while
done:
```

> **Assembly tests for the opposite case (`pow == 128`) of the C code (`pow != 128`).**

# For Loops

```
for (initialization; condition; loop operation)
  statement
```

- **initialization:** executes before the loop begins
- **condition:** is tested at the beginning of each iteration
- **loop operation:** executes at the end of each iteration
- **statement:** executes each time the condition is met

# For Loops

**High-level code**

```
// add the numbers from 0 to 9
int sum = 0;
int i;

for (i=0; i!=10; i = i+1) {
  sum = sum + i;
}
```

**MIPS assembly code**

```
# $s0 = i, $s1 = sum
```

# For Loops

**C Code**

```
// add the numbers from 0 to 9
int sum = 0;
int i;

for (i=0; i!=10; i = i+1) {
  sum = sum + i;
}
```

**MIPS assembly code**

# For Loops

**C Code**

```
// add the numbers from 0 to 9
int sum = 0;
int i;

for (i=0; i!=10; i = i+1) {
  sum = sum + i;
}
```

**MIPS assembly code**

```
# $s0 = i, $s1 = sum
        addi $s1, $0, 0
        add  $s0, $0, $0
        addi $t0, $0, 10
for:    beq  $s0, $t0, done
        add  $s1, $s1, $s0
        addi $s0, $s0, 1
        j    for
done:
```

# Less Than Comparison

**C Code**
```
// add the powers of 2 from 1
// to 100
int sum = 0;
int i;

for (i=1; i < 101; i = i*2) {
  sum = sum + i;
}
```

**MIPS assembly code**

# Less Than Comparison

## C Code

```
// add the powers of 2 from 1
// to 100
int sum = 0;
int i;

for (i=1; i < 101; i = i*2) {
  sum = sum + i;
}
```

## MIPS assembly code

```
# $s0 = i, $s1 = sum
        addi $s1, $0, 0
        addi $s0, $0, 1
        addi $t0, $0, 101
loop:   slt  $t1, $s0, $t0
        beq  $t1, $0, done
        add  $s1, $s1, $s0
        sll  $s0, $s0, 1
        j    loop
done:
```

**$t1 = 1   if   i < 101**

# Arrays

- Access large amounts of similar data
- **Index**: access each element
- **Size**: number of elements

# Arrays

- 5-element array
- **Base address** = 0x12348000 (address of first element, `array[0]`)
- First step in accessing an array: load base address into a register

| | |
|---|---|
| 0x12340010 | array[4] |
| 0x1234800C | array[3] |
| 0x12348008 | array[2] |
| 0x12348004 | array[1] |
| 0x12348000 | array[0] |

ELSEVIER

# Accessing Arrays

```
// C Code
    int array[5];
    array[0] = array[0] * 2;
    array[1] = array[1] * 2;
```

# Accessing Arrays

```
// C Code
  int array[5];
  array[0] = array[0] * 2;
  array[1] = array[1] * 2;


# MIPS assembly code
# array base address = $s0
  lui  $s0, 0x1234           # 0x1234 in upper half of $S0
  ori  $s0, $s0, 0x8000      # 0x8000 in lower half of $s0

  lw   $t1, 0($s0)           # $t1 = array[0]
  sll  $t1, $t1, 1           # $t1 = $t1 * 2
  sw   $t1, 0($s0)           # array[0] = $t1

  lw   $t1, 4($s0)           # $t1 = array[1]
  sll  $t1, $t1, 1           # $t1 = $t1 * 2
  sw   $t1, 4($s0)           # array[1] = $t1
```

# Arrays using For Loops

```
// C Code
   int array[1000];
   int i;


   for (i=0; i < 1000; i = i + 1)
       array[i] = array[i] * 8;
```

```
# MIPS assembly code
# $s0 = array base address, $s1 = i
```

# Arrays Using For Loops

```
# MIPS assembly code
# $s0 = array base address, $s1 = i
# initialization code
  lui  $s0, 0x23B8        # $s0 = 0x23B80000
  ori  $s0, $s0, 0xF000   # $s0 = 0x23B8F000
  addi $s1, $0, 0         # i = 0
  addi $t2, $0, 1000      # $t2 = 1000

loop:
  slt  $t0, $s1, $t2      # i < 1000?
  beq  $t0, $0, done      # if not then done
  sll  $t0, $s1, 2        # $t0 = i * 4 (byte offset)
  add  $t0, $t0, $s0      # address of array[i]
  lw   $t1, 0($t0)        # $t1 = array[i]
  sll  $t1, $t1, 3        # $t1 = array[i] * 8
  sw   $t1, 0($t0)        # array[i] = array[i] * 8
  addi $s1, $s1, 1        # i = i + 1
  j    loop               # repeat
done:
```

# ASCII Code

- *American Standard Code for Information Interchange*

- Each text character has unique byte value
  - For example, S = 0x53, a = 0x61, A = 0x41
  - Lower-case and upper-case differ by 0x20 (32)

ELSEVIER

# Cast of Characters

| # | Char | # | Char | # | Char | # | Char | # | Char | # | Char |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 20 | space | 30 | 0 | 40 | @ | 50 | P | 60 | ' | 70 | p |
| 21 | ! | 31 | 1 | 41 | A | 51 | Q | 61 | a | 71 | q |
| 22 | " | 32 | 2 | 42 | B | 52 | R | 62 | b | 72 | r |
| 23 | # | 33 | 3 | 43 | C | 53 | S | 63 | c | 73 | s |
| 24 | $ | 34 | 4 | 44 | D | 54 | T | 64 | d | 74 | t |
| 25 | % | 35 | 5 | 45 | E | 55 | U | 65 | e | 75 | u |
| 26 | & | 36 | 6 | 46 | F | 56 | V | 66 | f | 76 | v |
| 27 | ' | 37 | 7 | 47 | G | 57 | W | 67 | g | 77 | w |
| 28 | ( | 38 | 8 | 48 | H | 58 | X | 68 | h | 78 | x |
| 29 | ) | 39 | 9 | 49 | I | 59 | Y | 69 | i | 79 | y |
| 2A | * | 3A | : | 4A | J | 5A | Z | 6A | j | 7A | z |
| 2B | + | 3B | ; | 4B | K | 5B | [ | 6B | k | 7B | { |
| 2C | , | 3C | < | 4C | L | 5C | \ | 6C | l | 7C | \| |
| 2D | - | 3D | = | 4D | M | 5D | ] | 6D | m | 7D | } |
| 2E | . | 3E | > | 4E | N | 5E | ^ | 6E | n | 7E | ~ |
| 2F | / | 3F | ? | 4F | O | 5F | _ | 6F | o | | |

ELSEVIER

# Function Calls

- **Caller:** calling function (in this case, `main`)
- **Callee:** called function (in this case, `sum`)

**C Code**

```c
void main()
{
  int y;
  y = sum(42, 7);
  ...
}

int sum(int a, int b)
{
  return (a + b);
}
```

# Function Conventions

- **Caller:**

  - passes **arguments** to callee
  - jumps to callee

- **Callee:**

  - **performs** the function
  - **returns** result to caller
  - **returns** to point of call
  - **must not overwrite** registers or memory needed by caller

ELSEVIER

# MIPS Function Conventions

- **Call Function:** jump and link (`jal`)
- **Return** from function: jump register (`jr`)
- **Arguments**: `$a0 - $a3`
- **Return value**: `$v0`

# Function Calls

**C Code**

```
int main() {
    simple();
    a = b + c;
}

void simple() {
    return;
}
```

**MIPS assembly code**

```
0x00400200 main: jal  simple
0x00400204       add  $s0, $s1, $s2
...

0x00401020 simple: jr $ra
```

**void** means that **simple** doesn't return a value

# Function Calls

## C Code

```
int main() {
    simple();
    a = b + c;
}

void simple() {
    return;
}
```

## MIPS assembly code

```
0x00400200 main:   jal   simple
0x00400204         add   $s0, $s1, $s2
...

0x00401020 simple: jr $ra
```

**jal:** jumps to `simple`
  $$\$ra = PC + 4 = 0x00400204$$

**jr $ra:** jumps to address in `$ra` (0x00400204)

# Input Arguments & Return Value

## MIPS conventions:

- Argument values: `$a0 - $a3`
- Return value: `$v0`

# Input Arguments & Return Value

**C Code**

```c
int main()
{
  int y;
  ...
  y = diffofsums(2, 3, 4, 5);   // 4 arguments
  ...
}

int diffofsums(int f, int g, int h, int i)
{
  int result;
  result = (f + g) - (h + i);
  return result;                    // return value
}
```

# Input Arguments & Return Value

**MIPS assembly code**

```
# $s0 = y

main:
  ...
  addi $a0, $0, 2     # argument 0 = 2
  addi $a1, $0, 3     # argument 1 = 3
  addi $a2, $0, 4     # argument 2 = 4
  addi $a3, $0, 5     # argument 3 = 5
  jal  diffofsums     # call Function
  add  $s0, $v0, $0   # y = returned value
  ...

# $s0 = result
diffofsums:
  add $t0, $a0, $a1   # $t0 = f + g
  add $t1, $a2, $a3   # $t1 = h + i
  sub $s0, $t0, $t1   # result = (f + g) - (h + i)
  add $v0, $s0, $0    # put return value in $v0
  jr  $ra             # return to caller
```

# Input Arguments & Return Value

**MIPS assembly code**

```
# $s0 = result
diffofsums:
    add $t0, $a0, $a1   # $t0 = f + g
    add $t1, $a2, $a3   # $t1 = h + i
    sub $s0, $t0, $t1   # result = (f + g) - (h + i)
    add $v0, $s0, $0    # put return value in $v0
    jr  $ra             # return to caller
```

- diffofsums overwrote 3 registers: $t0, $t1, $s0
- diffofsums can use *stack* to temporarily store registers

# The Stack

- Memory used to temporarily save variables

- Like stack of dishes, last-in-first-out (LIFO) queue

- *Expands*: uses more memory when more space needed

- *Contracts*: uses less memory when the space is no longer needed

# The Stack

- Grows down (from higher to lower memory addresses)

- Stack pointer: $sp points to top of the stack

| Address | Data | | Address | Data |
|---|---|---|---|---|
| 7FFFFFFC | 12345678 | ← $sp | 7FFFFFFC | 12345678 |
| 7FFFFFF8 | | | 7FFFFFF8 | AABBCCDD |
| 7FFFFFF4 | | | 7FFFFFF4 | 11223344 | ← $sp |
| 7FFFFFF0 | | | 7FFFFFF0 | |
| ⋮ | ⋮ | | ⋮ | ⋮ |

# How Functions use the Stack

- Called functions must have no unintended side effects

- But `diffofsums` overwrites 3 registers: $t0, $t1, $s0

```
# MIPS assembly
# $s0 = result
diffofsums:
   add $t0, $a0, $a1   # $t0 = f + g
   add $t1, $a2, $a3   # $t1 = h + i
   sub $s0, $t0, $t1   # result = (f + g) - (h + i)
   add $v0, $s0, $0    # put return value in $v0
   jr  $ra             # return to caller
```

# Storing Register Values on the Stack

```
# $s0 = result
diffofsums:
  addi $sp, $sp, -12    # make space on stack
                        # to store 3 registers

  sw   $s0, 8($sp)      # save $s0 on stack
  sw   $t0, 4($sp)      # save $t0 on stack
  sw   $t1, 0($sp)      # save $t1 on stack
  add  $t0, $a0, $a1    # $t0 = f + g
  add  $t1, $a2, $a3    # $t1 = h + i
  sub  $s0, $t0, $t1    # result = (f + g) - (h + i)
  add  $v0, $s0, $0     # put return value in $v0
  lw   $t1, 0($sp)      # restore $t1 from stack
  lw   $t0, 4($sp)      # restore $t0 from stack
  lw   $s0, 8($sp)      # restore $s0 from stack
  addi $sp, $sp, 12     # deallocate stack space
  jr   $ra              # return to caller
```

# The stack during `diffofsums` Call

| Address | Data |
|---------|------|
| FC | ? |  ← $sp
| F8 |  |
| F4 |  |
| F0 |  |
| ⋮ | ⋮ |

(a)

| Address | Data |
|---------|------|
| FC | ? |
| F8 | $s0 |
| F4 | $t0 |
| F0 | $t1 |  ← $sp
| ⋮ | ⋮ |

stack frame

(b)

| Address | Data |
|---------|------|
| FC | ? |  ← $sp
| F8 |  |
| F4 |  |
| F0 |  |
| ⋮ | ⋮ |

(c)

ELSEVIER

# Registers

| Preserved *Callee-Saved* | Nonpreserved *Caller-Saved* |
|:---:|:---:|
| `$s0-$s7` | `$t0-$t9` |
| `$ra` | `$a0-$a3` |
| `$sp` | `$v0-$v1` |
| **stack above `$sp`** | **stack below `$sp`** |

# Multiple Function Calls

```
proc1:
    addi $sp, $sp, -4    # make space on stack
    sw   $ra, 0($sp)     # save $ra on stack
    jal  proc2
    ...
    lw   $ra, 0($sp)     # restore $s0 from stack
    addi $sp, $sp, 4     # deallocate stack space
    jr  $ra              # return to caller
```

# Storing Saved Registers on the Stack

```
# $s0 = result
diffofsums:
    addi $sp, $sp, -4   # make space on stack to
                        # store one register

    sw  $s0, 0($sp)     # save $s0 on stack
                        # no need to save $t0 or $t1

    add $t0, $a0, $a1   # $t0 = f + g
    add $t1, $a2, $a3   # $t1 = h + i
    sub $s0, $t0, $t1   # result = (f + g) - (h + i)
    add $v0, $s0, $0    # put return value in $v0
    lw  $s0, 0($sp)     # restore $s0 from stack
    addi $sp, $sp, 4    # deallocate stack space
    jr  $ra             # return to caller
```
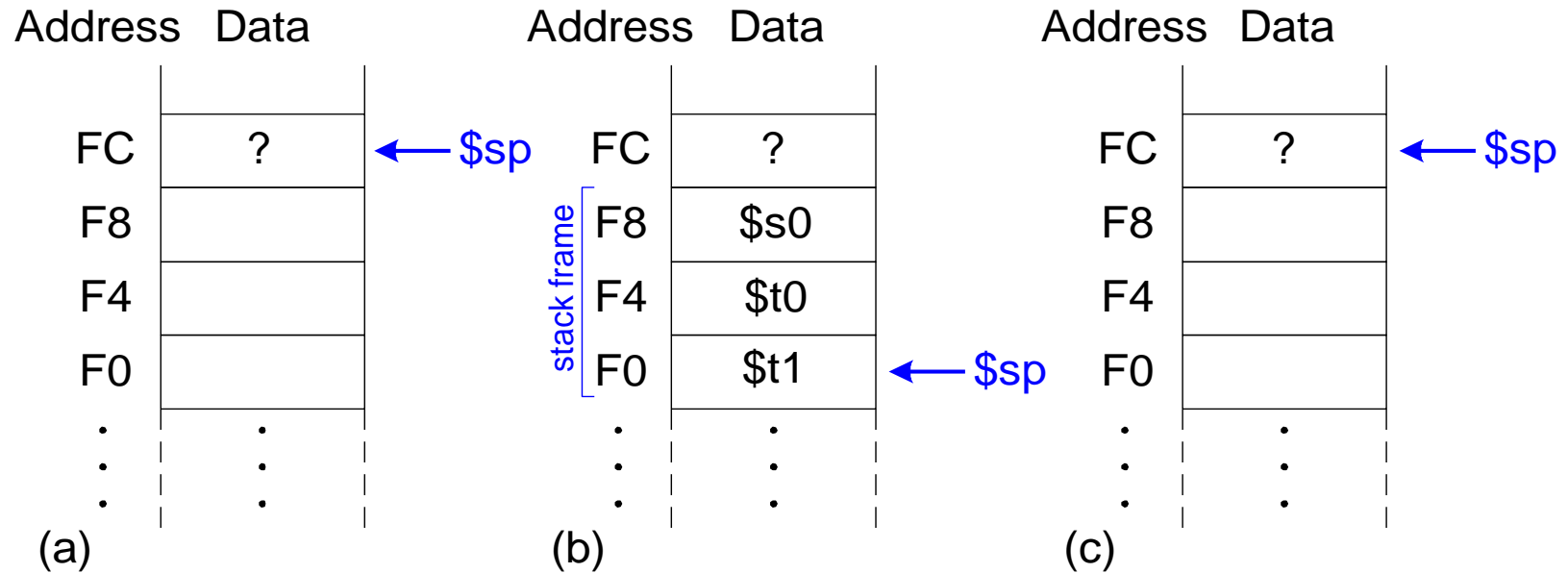
# Recursive Function Call

**High-level code**

```
int factorial(int n) {
  if (n <= 1)
    return 1;
  else
    return (n * factorial(n-1));
}
```

# Recursive Function Call

## MIPS assembly code

```
0x90  factorial: addi $sp, $sp, -8   # make room
0x94             sw   $a0, 4($sp)    # store $a0
0x98             sw   $ra, 0($sp)    # store $ra
0x9C             addi $t0, $0, 2
0xA0             slt  $t0, $a0, $t0  # a <= 1 ?
0xA4             beq  $t0, $0, else  # no: go to else
0xA8             addi $v0, $0, 1     # yes: return 1
0xAC             addi $sp, $sp, 8    # restore $sp
0xB0             jr   $ra            # return
0xB4      else:  addi $a0, $a0, -1   # n = n - 1
0xB8             jal  factorial      # recursive call
0xBC             lw   $ra, 0($sp)    # restore $ra
0xC0             lw   $a0, 4($sp)    # restore $a0
0xC4             addi $sp, $sp, 8    # restore $sp
0xC8             mul  $v0, $a0, $v0  # n * factorial(n-1)
0xCC             jr   $ra            # return
```

ELSEVIER

# Stack During Recursive Call

Address   Data

|  |  |
|---|---|
| FC |  | ← $sp |
| F8 |  |
| F4 |  |
| F0 |  |
| EC |  |
| E8 |  |
| E4 |  |
| E0 |  |
| DC |  |

Address   Data

|  |  |
|---|---|
| FC |  | ← $sp |
| F8 | $a0 (0x3) |
| F4 | $ra | ← $sp |
| F0 | $a0 (0x2) |
| EC | $ra (0xBC) | ← $sp |
| E8 | $a0 (0x1) |
| E4 | $ra (0xBC) | ← $sp |
| E0 |  |
| DC |  |

Address   Data

|  |  |
|---|---|
| FC |  | ← $sp   $v0 = 6 |
| F8 | $a0 (0x3) |
| F4 | $ra | ← $sp   $a0 = 3 / $v0 = 3 x 2 |
| F0 | $a0 (0x2) |
| EC | $ra (0xBC) | ← $sp   $a0 = 2 / $v0 = 2 x 1 |
| E8 | $a0 (0x1) |
| E4 | $ra (0xBC) | ← $sp   $a0 = 1 / $v0 = 1 x 1 |
| E0 |  |
| DC |  |

# Function Call Summary

- **Caller**
  - Put arguments in `$a0-$a3`
  - Save any needed registers (`$ra`, maybe `$t0-t9`)
  - `jal callee`
  - Restore registers
  - Look for result in `$v0`
- **Callee**
  - Save registers that might be disturbed (`$s0-$s7`)
  - Perform function
  - Put result in `$v0`
  - Restore registers
  - `jr $ra`

# Addressing Modes

## How do we address the operands?

- Register Only

- Immediate

- Base Addressing

- PC-Relative

- Pseudo Direct

# Addressing Modes

## Register Only

- Operands found in registers
  - **Example:** `add $s0, $t2, $t3`
  - **Example:** `sub $t8, $s1, $0`

## Immediate

- 16-bit immediate used as an operand
  - **Example:** `addi $s4, $t5, -73`
  - **Example:** `ori  $t3, $t7, 0xFF`

# Addressing Modes

## Base Addressing

- Address of operand is:

  ```
  base address + sign-extended immediate
  ```

  - **Example:** `lw  $s4, 72($0)`
    - address = `$0 + 72`

  - **Example:** `sw  $t2, -25($t1)`
    - address = `$t1 - 25`

# Addressing Modes

## PC-Relative Addressing

```
0x10                    beq    $t0, $0, else
0x14                    addi   $v0, $0, 1
0x18                    addi   $sp, $sp, i
0x1C                    jr           $ra
0x20        else:       addi   $a0, $a0, -1
0x24                    jal    factorial
```

### Assembly Code

```
beq $t0, $0, else

(beq $t0, $0, 3)
```

### Field Values

| op | rs | rt | imm | | |
|----|----|----|-----|---|---|
| 4 | 8 | 0 | 3 | | |
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

# Addressing Modes

## Pseudo-direct Addressing

```
0x0040005C            jal    sum
...
0x004000A0    sum:    add    $v0, $a0, $a1
```

JTA    0000 0000 0100 0000 0000 0000 1010 0000   (0x004000A0)

26-bit addr    0000 0000 0100 0000 0000 0000 1010 0000   (0x0100028)

0  1  0  0  0  2  8

### Field Values

| op | imm |
|---|---|
| 3 | 0x0100028 |
| 6 bits | 26 bits |

### Machine Code

| op | addr | |
|---|---|---|
| 000011 | 00 0001 0000 0000 0000 0010 1000 | (0x0C100028) |
| 6 bits | 26 bits | |

# Addressing Modes Summary

# How to Compile & Run a Program

High Level Code

↓

Compiler

↓

Assembly Code

↓

Assembler

↓

Object File

Object Files
Library Files

↓

Linker

↓

Executable

↓

Loader

↓

Memory

ELSEVIER

# Grace Hopper, 1906-1992

- Graduated from Yale University with a Ph.D. in mathematics

- Developed first compiler

- Helped develop the COBOL programming language

- Highly awarded naval officer

- Received World War II Victory Medal and National Defense Service Medal, among others

COMPUTER ARCHITECTURE

ELSEVIER

# What is Stored in Memory?

- Instructions (also called *text*)
- Data
  - Global/static: allocated before program begins
  - Dynamic: allocated within program

- How big is memory?
  - At most $2^{32} = 4$ gigabytes (4 GB)
  - From address 0x00000000 to 0xFFFFFFFF

# MIPS Memory Map

Address          Segment

0xFFFFFFFC

          Reserved

0x80000000
0x7FFFFFFC        Stack
                    ↓
          Dynamic Data
                    ↑
0x10010000        Heap
0x1000FFFC

          Static Data

0x10000000
0x0FFFFFFC

          Text

0x00400000
0x003FFFFC

          Reserved

0x00000000

ELSEVIER

# Example Program: C Code

```c
int f, g, y;   // global variables


int main(void)
{
  f = 2;
  g = 3;
  y = sum(f, g);

  return y;
}



int sum(int a, int b) {
  return (a + b);
}
```

# Example Program: MIPS Assembly

```c
int f, g, y;  // global

int main(void)
{

  f = 2;
  g = 3;

  y = sum(f, g);
  return y;
}

int sum(int a, int b) {
  return (a + b);
}
```

```
.data
f:
g:
y:
.text
main:
  addi $sp, $sp, -4    # stack frame
  sw   $ra, 0($sp)     # store $ra
  addi $a0, $0, 2      # $a0 = 2
  sw   $a0, f          # f = 2
  addi $a1, $0, 3      # $a1 = 3
  sw   $a1, g          # g = 3
  jal  sum             # call sum
  sw   $v0, y          # y = sum()
  lw   $ra, 0($sp)     # restore $ra
  addi $sp, $sp, 4     # restore $sp
  jr   $ra             # return to OS
sum:
  add  $v0, $a0, $a1   # $v0 = a + b
  jr   $ra             # return
```

ELSEVIER

# Example Program: Symbol Table

| Symbol | Address |
|--------|---------|
| f | 0x10000000 |
| g | 0x10000004 |
| y | 0x10000008 |
| main | 0x00400000 |
| sum | 0x0040002C |

# Example Program: Executable

| Executable file header | Text Size | Data Size |
|---|---|---|
| | 0x34 (52 bytes) | 0xC (12 bytes) |
| **Text segment** | **Address** | **Instruction** | |
| | 0x00400000 | 0x23BDFFFC | addi $sp, $sp, -4 |
| | 0x00400004 | 0xAFBF0000 | sw   $ra, 0 ($sp) |
| | 0x00400008 | 0x20040002 | addi $a0, $0, 2 |
| | 0x0040000C | 0xAF848000 | sw   $a0, 0x8000 ($gp) |
| | 0x00400010 | 0x20050003 | addi $a1, $0, 3 |
| | 0x00400014 | 0xAF858004 | sw   $a1, 0x8004 ($gp) |
| | 0x00400018 | 0x0C10000B | jal    0x0040002C |
| | 0x0040001C | 0xAF828008 | sw   $v0, 0x8008 ($gp) |
| | 0x00400020 | 0x8FBF0000 | lw    $ra, 0 ($sp) |
| | 0x00400024 | 0x23BD0004 | addi $sp, $sp, -4 |
| | 0x00400028 | 0x03E00008 | jr    $ra |
| | 0x0040002C | 0x00851020 | add  $v0, $a0, $a1 |
| | 0x00400030 | 0x03E00008 | jr    $ra |
| **Data segment** | **Address** | **Data** | |
| | 0x10000000 | f | |
| | 0x10000004 | g | |
| | 0x10000008 | y | |

ELSEVIER

# Example Program: In Memory

Address     Memory

| | |
|---|---|
| | Reserved |
| 0x7FFFFFFC | Stack ↓ |
| | ↑ |
| 0x10010000 | Heap |
| | . . . |
| | y |
| | g |
| 0x10000000 | f |
| | . . . |
| | 0x03E00008 |
| | 0x00851020 |
| | 0x03E00008 |
| | 0x23BD0004 |
| | 0x8FBF0000 |
| | 0xAF828008 |
| | 0x0C10000B |
| | 0xAF858004 |
| | 0x20050003 |
| | 0xAF848000 |
| | 0x20040002 |
| | 0xAFBF0000 |
| 0x00400000 | 0x23BDFFFC |
| | Reserved |

← $sp = 0x7FFFFFFC

← $gp = 0x10008000

← PC = 0x00400000

ELSEVIER

# Odds & Ends

- Pseudoinstructions
- Exceptions
- Signed and unsigned instructions
- Floating-point instructions

# Pseudoinstructions

| Pseudoinstruction | MIPS Instructions |
|---|---|
| li $s0, 0x1234AA77 | lui $s0, 0x1234<br>ori $s0, 0xAA77 |
| clear $t0 | add $t0, $0, $0 |
| move $s1, $s2 | add $s2, $s1, $0 |
| nop | sll $0, $0, 0 |

# Exceptions

- Unscheduled function call to *exception handler*

- Caused by:
  - Hardware, also called an *interrupt*, e.g., keyboard
  - Software, also called *traps*, e.g., undefined instruction

- When exception occurs, the processor:
  - Records the cause of the exception
  - Jumps to exception handler (at instruction address 0x80000180)
  - Returns to program

# Exception Registers

- Not part of register file
  - **Cause**: Records cause of exception
  - **EPC** (Exception PC): Records PC where exception occurred
- `EPC` and `Cause`: part of Coprocessor 0
- Move from Coprocessor 0
  - `mfc0 $t0, EPC`
  - Moves contents of `EPC` into `$t0`

# Exception Causes

| Exception | Cause |
|-----------|-------|
| Hardware Interrupt | 0x00000000 |
| System Call | 0x00000020 |
| Breakpoint / Divide by 0 | 0x00000024 |
| Undefined Instruction | 0x00000028 |
| Arithmetic Overflow | 0x00000030 |

ELSEVIER

# Exception Flow

- Processor saves cause and exception PC in `Cause` and `EPC`

- Processor jumps to exception handler (0x80000180)

- Exception handler:
  - Saves registers on stack
  - Reads `Cause` register

    ```
    mfc0 $t0, Cause
    ```
  - Handles exception
  - Restores registers
  - Returns to program

    ```
    mfc0 $k0, EPC
    jr $k0
    ```

# Signed & Unsigned Instructions

- Addition and subtraction
- Multiplication and division
- Set less than

# Addition & Subtraction

- **Signed:** `add, addi, sub`
  - Same operation as unsigned versions
  - But processor takes exception on overflow
- **Unsigned:** `addu, addiu, subu`
  - Doesn't take exception on overflow

**Note:** `addiu` sign-extends the immediate

ELSEVIER

# Multiplication & Division

- **Signed:** `mult, div`
- **Unsigned:** `multu, divu`

# Set Less Than

- **Signed:** `slt, slti`
- **Unsigned:** `sltu, sltiu`

**Note:** `sltiu` sign-extends the immediate before comparing it to the register

# Loads

- **Signed:**
  - Sign-extends to create 32-bit value to load into register
  - Load halfword: `lh`
  - Load byte: `lb`

- **Unsigned:**
  - Zero-extends to create 32-bit value
  - Load halfword unsigned: `lhu`
  - Load byte: `lbu`

# Floating-Point Instructions

- Floating-point coprocessor (Coprocessor 1)
- 32 32-bit floating-point registers ($f0-$f31)
- Double-precision values held in two floating point registers
  - e.g., $f0 and $f1, $f2 and $f3, etc.
  - Double-precision floating point registers: $f0, $f2, $f4, etc.

ELSEVIER

# Floating-Point Instructions

| Name | Register Number | Usage |
|------|-----------------|-------|
| `$fv0 - $fv1` | 0, 2 | return values |
| `$ft0 - $ft3` | 4, 6, 8, 10 | temporary variables |
| `$fa0 - $fa1` | 12, 14 | Function arguments |
| `$ft4 - $ft8` | 16, 18 | temporary variables |
| `$fs0 - $fs5` | 20, 22, 24, 26, 28, 30 | saved variables |

# F-Type Instruction Format

- $\texttt{Opcode} = 17\ (010001_2)$

- Single-precision:
  - $\texttt{cop} = 16\ (010000_2)$
  - add.s, sub.s, div.s, neg.s, abs.s, etc.

- Double-precision:
  - $\texttt{cop} = 17\ (010001_2)$
  - add.d, sub.d, div.d, neg.d, abs.d, etc.

- 3 register operands:
  - fs, ft: source operands
  - fd: destination operands

**F-Type**

| op | cop | ft | fs | fd | funct |
|----|-----|-----|-----|-----|-------|
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

# Floating-Point Branches

- Set/clear condition flag: `fpcond`
  - Equality: `c.seq.s, c.seq.d`
  - Less than: `c.lt.s, c.lt.d`
  - Less than or equal: `c.le.s, c.le.d`
- Conditional branch
  - `bclf:` branches if `fpcond` is FALSE
  - `bclt:` branches if `fpcond` is TRUE
- Loads and stores
  - `lwc1: lwc1 $ft1, 42($s1)`
  - `swc1: swc1 $fs2, 17($sp)`

# Looking Ahead

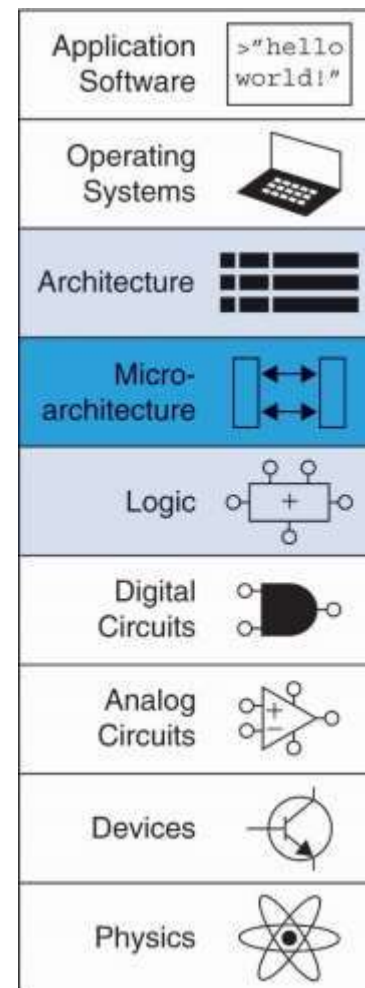**Microarchitecture** – building MIPS processor in hardware

**Bring colored pencils**

# Chapter 7

**MICROARCHITECTURE**

***Digital Design and Computer Architecture*, 2nd Edition**

David Money Harris and Sarah L. Harris

ELSEVIER

# Chapter 7 :: Topics

- **Introduction**

- **Performance Analysis**

- **Single-Cycle Processor**

- **Pipelined Processor**

# Introduction

- **Microarchitecture:** how to implement an architecture in hardware

- Processor:
  - **Datapath:** functional blocks
  - **Control:** control signals

| | |
|---|---|
| Application Software | programs |
| Operating Systems | device drivers |
| Architecture | instructions registers |
| Micro-architecture | datapaths controllers |
| Logic | adders memories |
| Digital Circuits | AND gates NOT gates |
| Analog Circuits | amplifiers filters |
| Devices | transistors diodes |
| Physics | electrons |

ELSEVIER

# Microarchitecture

- Multiple implementations for a single architecture:

  – **Single-cycle:** Each instruction executes in a single cycle

  – **Multicycle:** Each instruction is broken into series of shorter steps

  – **Pipelined:** Each instruction broken up into series of steps & multiple instructions execute at once

# Processor Performance

- ## Program execution time

**Execution Time = (#instructions)(cycles/instruction)(seconds/cycle)**

- ## Definitions:
  - CPI: Cycles/instruction
  - clock period: seconds/cycle
  - IPC: instructions/cycle = IPC
- ## Challenge is to satisfy constraints of:
  - Cost
  - Power
  - Performance

ELSEVIER

# MIPS Processor

- Consider subset of MIPS instructions:
    - R-type instructions: `and, or, add, sub, slt`
    - Memory instructions: `lw, sw`
    - Branch instructions: `beq`

# Architectural State

- Determines everything about a processor:
  - PC
  - 32 registers
  - Memory

# MIPS State Elements

# Single-Cycle MIPS Processor
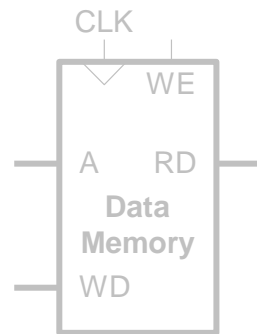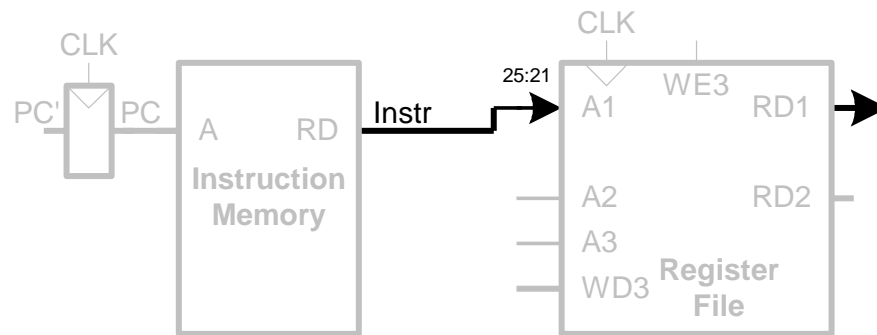
- Datapath

- Control

# Single-Cycle Datapath: `lw` fetch

**STEP 1:** Fetch instruction

**STEP 2:** Read source operands from RF

**STEP 3:** Sign-extend the immediate

# Single-Cycle Datapath: `lw` address

**STEP 4:** Compute the memory address

# Review: ALU



| $F_{2:0}$ | Function |
|-----------|----------|
| 000 | A & B |
| 001 | A \| B |
| 010 | A + B |
| 011 | not used |
| 100 | A & ~B |
| 101 | A \| ~B |
| 110 | A - B |
| 111 | SLT |

# Single-Cycle Datapath: `lw` Memory Read

- **STEP 5:** Read data from memory and write it back to register file

## STEP 6: Determine address of next instruction

## Write data in `rt` to memory

# Single-Cycle Datapath: R-Type

- Read from `rs` and `rt`
- Write *ALUResult* to register file
- Write to `rd` (instead of `rt`)

# Single-Cycle Datapath: `beq`

- Determine whether values in `rs` and `rt` are equal
- Calculate branch target address:

   BTA = (sign-extended immediate << 2) + (PC+4)

# Single-Cycle Processor

# Single-Cycle Control

| $F_{2:0}$ | Function |
|-----------|----------|
| 000 | A & B |
| 001 | A \| B |
| 010 | A + B |
| 011 | not used |
| 100 | A & ~B |
| 101 | A \| ~B |
| 110 | A - B |
| 111 | SLT |

# Control Unit: ALU Decoder

| ALUOp$_{1:0}$ | Meaning |
|---|---|
| 00 | Add |
| 01 | Subtract |
| 10 | Look at Funct |
| 11 | Not Used |

| ALUOp$_{1:0}$ | Funct | ALUControl$_{2:0}$ |
|---|---|---|
| 00 | X | 010 (Add) |
| 01 | X | 110 (Subtract) |
| 10 | 100000 (add) | 010 (Add) |
| 10 | 100010 (sub) | 110 (Subtract) |
| 10 | 100100 (and) | 000 (And) |
| 10 | 100101 (or) | 001 (Or) |
| 10 | 101010 (slt) | 111 (SLT) |

# Control Unit Main Decoder

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ |
|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | | | | | | | |
| lw | 100011 | | | | | | | |
| sw | 101011 | | | | | | | |
| beq | 000100 | | | | | | | |

# Control Unit: Main Decoder

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ |
|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | 1 | 1 | 0 | 0 | 0 | 0 | 10 |
| lw | 100011 | 1 | 0 | 1 | 0 | 0 | 0 | 00 |
| sw | 101011 | 0 | X | 1 | 0 | 1 | X | 00 |
| beq | 000100 | 0 | X | 0 | 1 | 0 | X | 01 |

ELSEVIER

# Extended Functionality: `addi`



**No change to datapath**

# Control Unit: `addi`

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ |
|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | 1 | 1 | 0 | 0 | 0 | 0 | 10 |
| `lw` | 100011 | 1 | 0 | 1 | 0 | 0 | 1 | 00 |
| `sw` | 101011 | 0 | X | 1 | 0 | 1 | X | 00 |
| `beq` | 000100 | 0 | X | 0 | 1 | 0 | X | 01 |
| **addi** | **001000** | | | | | | | |

# Control Unit: `addi`

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ |
|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | 1 | 1 | 0 | 0 | 0 | 0 | 10 |
| `lw` | 100011 | 1 | 0 | 1 | 0 | 0 | 1 | 00 |
| `sw` | 101011 | 0 | X | 1 | 0 | 1 | X | 00 |
| `beq` | 000100 | 0 | X | 0 | 1 | 0 | X | 01 |
| **`addi`** | **001000** | **1** | **0** | **1** | **0** | **0** | **0** | **00** |

ELSEVIER

# Control Unit: Main Decoder

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ | Jump |
|---|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | 1 | 1 | 0 | 0 | 0 | 0 | 10 | 0 |
| lw | 100011 | 1 | 0 | 1 | 0 | 0 | 1 | 00 | 0 |
| sw | 101011 | 0 | X | 1 | 0 | 1 | X | 00 | 0 |
| beq | 000100 | 0 | X | 0 | 1 | 0 | X | 01 | 0 |
| j | 000100 | | | | | | | | |

ELSEVIER

# Control Unit: Main Decoder

| Instruction | $Op_{5:0}$ | RegWrite | RegDst | AluSrc | Branch | MemWrite | MemtoReg | $ALUOp_{1:0}$ | Jump |
|---|---|---|---|---|---|---|---|---|---|
| R-type | 000000 | 1 | 1 | 0 | 0 | 0 | 0 | 10 | 0 |
| lw | 100011 | 1 | 0 | 1 | 0 | 0 | 1 | 00 | 0 |
| sw | 101011 | 0 | X | 1 | 0 | 1 | X | 00 | 0 |
| beq | 000100 | 0 | X | 0 | 1 | 0 | X | 01 | 0 |
| j | 000100 | 0 | X | X | X | 0 | X | XX | 1 |

# Review: Processor Performance

**Program Execution Time**

**= (#instructions)(cycles/instruction)(seconds/cycle)**

**= # instructions x CPI x $T_C$**

# Single-Cycle Performance



$T_C$ **limited by critical path (`lw`)**

# Single-Cycle Performance

- Single-cycle critical path:

$$T_c = t_{pcq\_PC} + t_{mem} + \max(t_{RFread}, t_{sext} + t_{mux}) + t_{ALU} + t_{mem} + t_{mux} + t_{RFsetup}$$

- Typically, limiting paths are:
  - memory, ALU, register file
  - $T_c = t_{pcq\_PC} + 2t_{mem} + t_{RFread} + t_{mux} + t_{ALU} + t_{RFsetup}$

# Single-Cycle Performance Example

| Element | Parameter | Delay (ps) |
|---|---|---|
| Register clock-to-Q | $t_{pcq\_PC}$ | 30 |
| Register setup | $t_{setup}$ | 20 |
| Multiplexer | $t_{mux}$ | 25 |
| ALU | $t_{ALU}$ | 200 |
| Memory read | $t_{mem}$ | 250 |
| Register file read | $t_{RFread}$ | 150 |
| Register file setup | $t_{RFsetup}$ | 20 |

$T_c = ?$

# Single-Cycle Performance Example

| Element | Parameter | Delay (ps) |
|---------|-----------|------------|
| Register clock-to-Q | $t_{pcq\_PC}$ | 30 |
| Register setup | $t_{setup}$ | 20 |
| Multiplexer | $t_{mux}$ | 25 |
| ALU | $t_{ALU}$ | 200 |
| Memory read | $t_{mem}$ | 250 |
| Register file read | $t_{RFread}$ | 150 |
| Register file setup | $t_{RFsetup}$ | 20 |

$$T_c = t_{pcq\_PC} + 2t_{mem} + t_{RFread} + t_{mux} + t_{ALU} + t_{RFsetup}$$
$$= [30 + 2(250) + 150 + 25 + 200 + 20] \text{ ps}$$
$$= 925 \text{ ps}$$

ELSEVIER

# Single-Cycle Performance Example

Program with 100 billion instructions:

**Execution Time** = # instructions x CPI x $T_C$

$\qquad$ = $(100 \times 10^9)(1)(925 \times 10^{-12}\,\text{s})$

$\qquad$ = **92.5 seconds**

ELSEVIER

# Multicycle MIPS Processor

- ## Single-cycle:

  + simple

  - cycle time limited by longest instruction (`lw`)

  - 2 adders/ALUs & 2 memories

- ## Multicycle:

  + higher clock speed

  + simpler instructions run faster

  + reuse expensive hardware on multiple cycles

  - sequencing overhead paid many times

- ## Same design steps: datapath & control

# Pipelined MIPS Processor

- Temporal parallelism
- Divide single-cycle processor into 5 stages:
  - Fetch
  - Decode
  - Execute
  - Memory
  - Writeback
- Add pipeline registers between stages

# Single-Cycle vs. Pipelined

## Single-Cycle

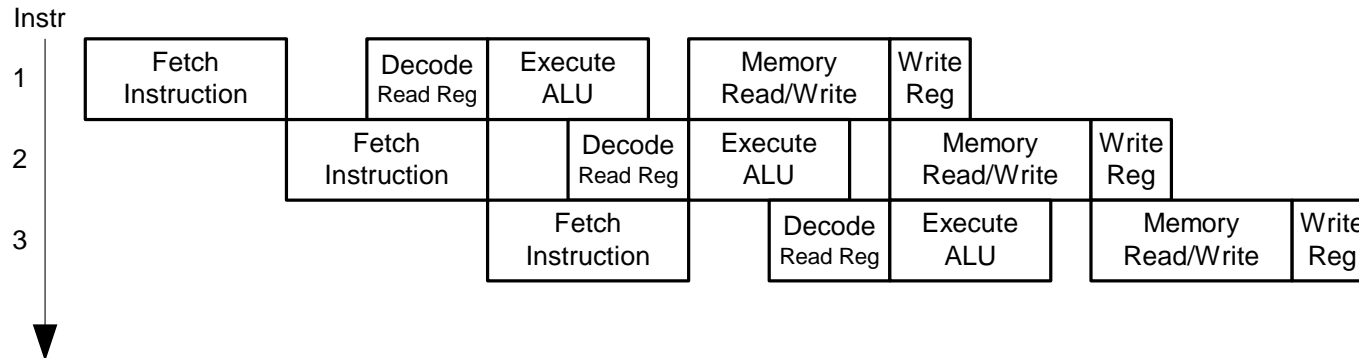| Instr | | | | | |
|-------|---|---|---|---|---|
| 1 | Fetch Instruction | Decode Read Reg | Execute ALU | Memory Read / Write | Write Reg |
| 2 | | | Fetch Instruction | Decode Read Reg | Execute ALU | Memory Read / Write | Write Reg |

Time (ps): 0 100 200 300 400 500 600 700 800 900 1000 1100 1200 1300 1400 1500 1600 1700 1800 1900

## Pipelined

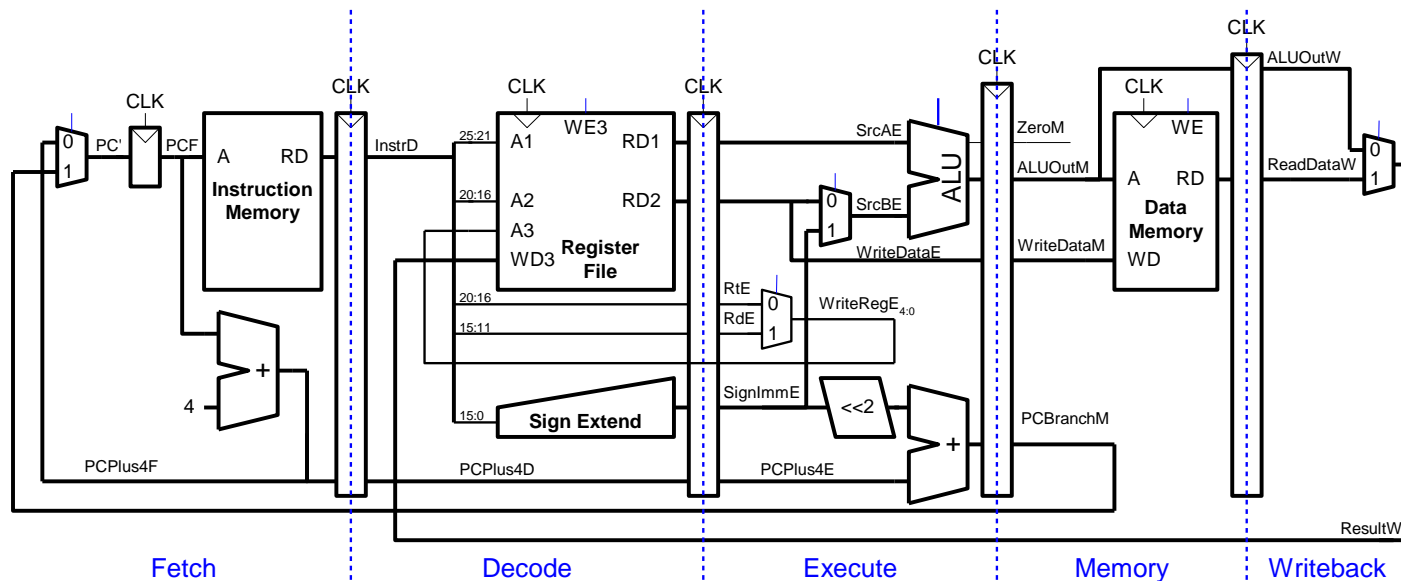| Instr | | | | | |
|-------|---|---|---|---|---|
| 1 | Fetch Instruction | Decode Read Reg | Execute ALU | Memory Read/Write | Write Reg |
| 2 | | Fetch Instruction | Decode Read Reg | Execute ALU | Memory Read/Write | Write Reg |
| 3 | | | Fetch Instruction | Decode Read Reg | Execute ALU | Memory Read/Write | Write Reg |

ELSEVIER

# Pipelined Processor Abstraction

# Single-Cycle & Pipelined Datapath

# Corrected Pipelined Datapath



***WriteReg* must arrive at same time as *Result***

# Pipelined Processor Control



- **Same control unit as single-cycle processor**
- **Control delayed to proper pipeline stage**

# Pipeline Hazards

- When an instruction depends on result from instruction that hasn't completed

- Types:

  - **Data hazard:** register value not yet written back to register file

  - **Control hazard:** next instruction not decided yet (caused by branches)

# Data Hazard

© *Digital Design and Computer Architecture, 2nd Edition, 2012*

# Handling Data Hazards

- Insert `nops` in code at compile time

- Rearrange code at compile time

- Forward data at run time

- Stall the processor at run time

ELSEVIER

# Compile-Time Hazard Elimination

- Insert enough `nops` for result to be ready
- Or move independent useful instructions forward

# Data Forwarding



add $s0, $s2, $s3
and $t0, $s0, $s1
or  $t1, $s4, $s0
sub $t2, $s0, $s5

# Data Forwarding

# Data Forwarding

- Forward to Execute stage from either:
  - Memory stage or
  - Writeback stage

- Forwarding logic for *ForwardAE*:

```
if        ((rsE != 0) AND (rsE == WriteRegM) AND RegWriteM)
    then    ForwardAE = 10
else if ((rsE != 0) AND (rsE == WriteRegW) AND RegWriteW)
    then    ForwardAE = 01
else        ForwardAE = 00
```

**Forwarding logic for *ForwardBE* same, but replace *rsE* with *rtE***

# Stalling



lw $s0, 40($0)

and $t0, $s0, $s1

or  $t1, $s4, $s0

sub $t2, $s0, $s5

# Stalling



lw $s0, 40($0)

and $t0, $s0, $s1

or  $t1, $s4, $s0

sub $t2, $s0, $s5

# Stalling Hardware

# Stalling Logic

```
lwstall =
   ((rsD==rtE) OR (rtD==rtE)) AND MemtoRegE


StallF = StallD = FlushE = lwstall
```

# Control Hazards

- **`beq:`**
  - branch not determined until 4<sup>th</sup> stage of pipeline
  - Instructions after branch fetched before branch occurs
  - These instructions must be flushed if branch happens

- **Branch misprediction penalty**
  - number of instruction flushed when branch is taken
  - May be reduced by determining branch earlier

# Control Hazards: Original Pipeline

# Control Hazards

# Early Branch Resolution



**Introduced another data hazard in Decode stage**

```
        1      2      3      4      5      6      7      8      9
                                                              Time (cycles)

20   beq $t1, $t2, 40      IM |lw| RF |$t1/$t2| - | DM | RF

24   and $t0, $s0, $s1          IM |and| RF |$s0/$s1| & | DM | RF         ] Flush
                                                                            this
28   or  $t1, $s4, $s0                                                   instruction

2C   sub $t2, $s0, $s5

30   ...
...
64   slt $t3, $s2, $s3               IM |slt| RF |$s2/$s3| slt | DM |$t3| RF
```

# Handling Data & Control Hazards

# Control Forwarding & Stalling Logic

- ## **Forwarding logic:**

  **ForwardAD** = ($rsD$ !=0) AND ($rsD$ == $WriteRegM$) AND $RegWriteM$
  **ForwardBD** = ($rtD$ !=0) AND ($rtD$ == WriteRegM) AND $RegWriteM$

- ## **Stalling logic:**

  **branchstall** = $BranchD$ AND $RegWriteE$ AND
          ($WriteRegE$ == $rsD$ OR $WriteRegE$ == $rtD$)
        OR
        $BranchD$ AND $MemtoRegM$ AND
          ($WriteRegM$ == $rsD$ OR $WriteRegM$ == $rtD$)

  $StallF$ = $StallD$ = $FlushE$ = $lwstall$ OR $branchstall$

ELSEVIER

# Branch Prediction

- Guess whether branch will be taken
  - Backward branches are usually taken (loops)
  - Consider history to improve guess

- Good prediction reduces fraction of branches requiring a flush

ELSEVIER

# Pipelined Performance Example

- SPECINT2000 benchmark:
    - 25% loads
    - 10% stores
    - 11% branches
    - 2% jumps
    - 52% R-type
- Suppose:
    - 40% of loads used by next instruction
    - 25% of branches mispredicted
    - All jumps flush next instruction
- **What is the average CPI?**

# Pipelined Performance Example

- SPECINT2000 benchmark:
  - 25% loads
  - 10% stores
  - 11% branches
  - 2% jumps
  - 52% R-type
- Suppose:
  - 40% of loads used by next instruction
  - 25% of branches mispredicted
  - All jumps flush next instruction
- **What is the average CPI?**
  - Load/Branch CPI = 1 when no stalling, 2 when stalling
  - $CPI_{lw} = 1(0.6) + 2(0.4) = 1.4$
  - $CPI_{beq} = 1(0.75) + 2(0.25) = 1.25$

  **Average CPI** $= (0.25)(1.4) + (0.1)(1) + (0.11)(1.25) + (0.02)(2) + (0.52)(1)$

  $= \textbf{1.15}$

# Pipelined Performance

- Pipelined processor critical path:

$$T_c = \max \{$$

$$t_{pcq} + t_{\text{mem}} + t_{\text{setup}}$$

$$2(t_{RFread} + t_{\text{mux}} + t_{\text{eq}} + t_{\text{AND}} + t_{\text{mux}} + t_{\text{setup}})$$

$$t_{pcq} + t_{\text{mux}} + t_{\text{mux}} + t_{\text{ALU}} + t_{\text{setup}}$$

$$t_{pcq} + t_{\text{memwrite}} + t_{\text{setup}}$$

$$2(t_{pcq} + t_{\text{mux}} + t_{\text{RFwrite}}) \}$$

# Pipelined Performance Example

| Element | Parameter | Delay (ps) |
|---|---|---|
| Register clock-to-Q | $t_{pcq\_PC}$ | 30 |
| Register setup | $t_{setup}$ | 20 |
| Multiplexer | $t_{mux}$ | 25 |
| ALU | $t_{ALU}$ | 200 |
| Memory read | $t_{mem}$ | 250 |
| Register file read | $t_{RFread}$ | 150 |
| Register file setup | $t_{RFsetup}$ | 20 |
| Equality comparator | $t_{eq}$ | 40 |
| AND gate | $t_{AND}$ | 15 |
| Memory write | $T_{memwrite}$ | 220 |
| Register file write | $t_{RFwrite}$ | 100 ps |

$$T_c = 2(t_{RFread} + t_{mux} + t_{eq} + t_{AND} + t_{mux} + t_{setup})$$
$$= 2[150 + 25 + 40 + 15 + 25 + 20] \text{ ps} = \textbf{550 ps}$$

# Pipelined Performance Example

Program with 100 billion instructions

**Execution Time** = (# instructions) × CPI × $T_c$

$$= (100 \times 10^9)(1.15)(550 \times 10^{-12})$$

**= 63 seconds**

# Processor Performance Comparison

| Processor | Execution Time (seconds) | Speedup (single-cycle as baseline) |
|---|---|---|
| **Single-cycle** | 92.5 | 1 |
| **Pipelined** | 63 | 1.47 |
| | | |

# Chapter 8

**MEMORY & I/O SYSTEMS**

*Digital Design and Computer Architecture*, 2nd Edition

David Money Harris and Sarah L. Harris

ELSEVIER

# Chapter 8 :: Topics

- **Introduction**
- **Memory System Performance Analysis**
- **Caches**
- **Virtual Memory**
- **Memory-Mapped I/O**
- **Summary**

# Introduction

- Computer performance depends on:
  - Processor performance
  - Memory system performance

## Memory Interface

# Processor-Memory Gap

In prior chapters, assumed access memory in 1 clock cycle – but hasn't been true since the 1980's

# Memory System Challenge

- Make memory system appear as fast as processor

- Use hierarchy of memories

- Ideal memory:
  - Fast
  - Cheap (inexpensive)
  - Large (capacity)

**But can only choose two!**

# Memory Hierarchy

| Technology | Price / GB | Access Time (ns) | Bandwidth (GB/s) |
|------------|-----------|------------------|------------------|
| SRAM | $10,000 | 1 | 25+ |
| DRAM | $10 | 10 - 50 | 10 |
| SSD | $1 | 100,000 | 0.5 |
| HDD | $0.1 | 10,000,000 | 0.1 |

Speed

Capacity

Cache

Main Memory

Virtual Memory

# Locality

Exploit locality to make memory accesses fast

- **Temporal Locality:**
  - Locality in time
  - If data used recently, likely to use it again soon
  - **How to exploit:** keep recently accessed data in higher levels of memory hierarchy

- **Spatial Locality:**
  - Locality in space
  - If data used recently, likely to use nearby data soon
  - **How to exploit:** when access data, bring nearby data into higher levels of memory hierarchy too

# Memory Performance

- **Hit:** data found in that level of memory hierarchy
- **Miss:** data not found (must go to next level)

| | |
|---|---|
| **Hit Rate** | = # hits / # memory accesses |
| | = 1 − Miss Rate |
| **Miss Rate** | = # misses / # memory accesses |
| | = 1 − Hit Rate |

- **Average memory access time (AMAT):** average time for processor to access data

$$\text{AMAT} = t_{\text{cache}} + MR_{\text{cache}}[t_{MM} + MR_{MM}(t_{VM})]$$

# Memory Performance Example 1

- A program has 2,000 loads and stores

- 1,250 of these data values in cache

- Rest supplied by other levels of memory hierarchy

- **What are the hit and miss rates for the cache?**

ELSEVIER

# Memory Performance Example 1

- A program has 2,000 loads and stores

- 1,250 of these data values in cache

- Rest supplied by other levels of memory hierarchy

- **What are the hit and miss rates for the cache?**

**Hit Rate** = 1250/2000 = **0.625**

**Miss Rate** = 750/2000 = **0.375** = 1 − Hit Rate

ELSEVIER

# Memory Performance Example 2

- Suppose processor has 2 levels of hierarchy: cache and main memory

- $t_{\text{cache}} = 1$ cycle, $t_{MM} = 100$ cycles

- **What is the AMAT of the program from Example 1?**

# Memory Performance Example 2

- Suppose processor has 2 levels of hierarchy: cache and main memory

- $t_{\text{cache}} = 1$ cycle, $t_{MM} = 100$ cycles

- **What is the AMAT of the program from Example 1?**

$$\text{AMAT} = t_{\text{cache}} + MR_{\text{cache}}(t_{MM})$$
$$= [1 + 0.375(100)] \text{ cycles}$$
$$= \textbf{38.5 cycles}$$

# Cache

- Highest level in memory hierarchy
- Fast (typically ~ 1 cycle access time)
- Ideally supplies most data to processor
- Usually holds most recently accessed data

# Cache Design Questions

- What data is held in the cache?

- How is data found?

- What data is replaced?

**Focus on data loads, but stores follow same principles**

ELSEVIER

# What data is held in the cache?

- Ideally, cache anticipates needed data and puts it in cache

- But impossible to predict future

- Use past to predict future – temporal and spatial locality:

  - **Temporal locality:** copy newly accessed data into cache

  - **Spatial locality:** copy neighboring data into cache too

ELSEVIER

# Cache Terminology

- **Capacity ($C$):**
  - number of data bytes in cache
- **Block size ($b$):**
  - bytes of data brought into cache at once
- **Number of blocks ($B = C/b$):**
  - number of blocks in cache: $B = C/b$
- **Degree of associativity ($N$):**
  - number of blocks in a set
- **Number of sets ($S = B/N$):**
  - each memory address maps to exactly one cache set

ELSEVIER

# How is data found?

- Cache organized into $S$ sets

- Each memory address maps to exactly one set

- Caches categorized by # of blocks in a set:

  - **Direct mapped:** 1 block per set

  - *N*-**way set associative:** $N$ blocks per set

  - **Fully associative:** all cache blocks in 1 set

- Examine each organization for a cache with:

  - Capacity ($C$ = 8 words)

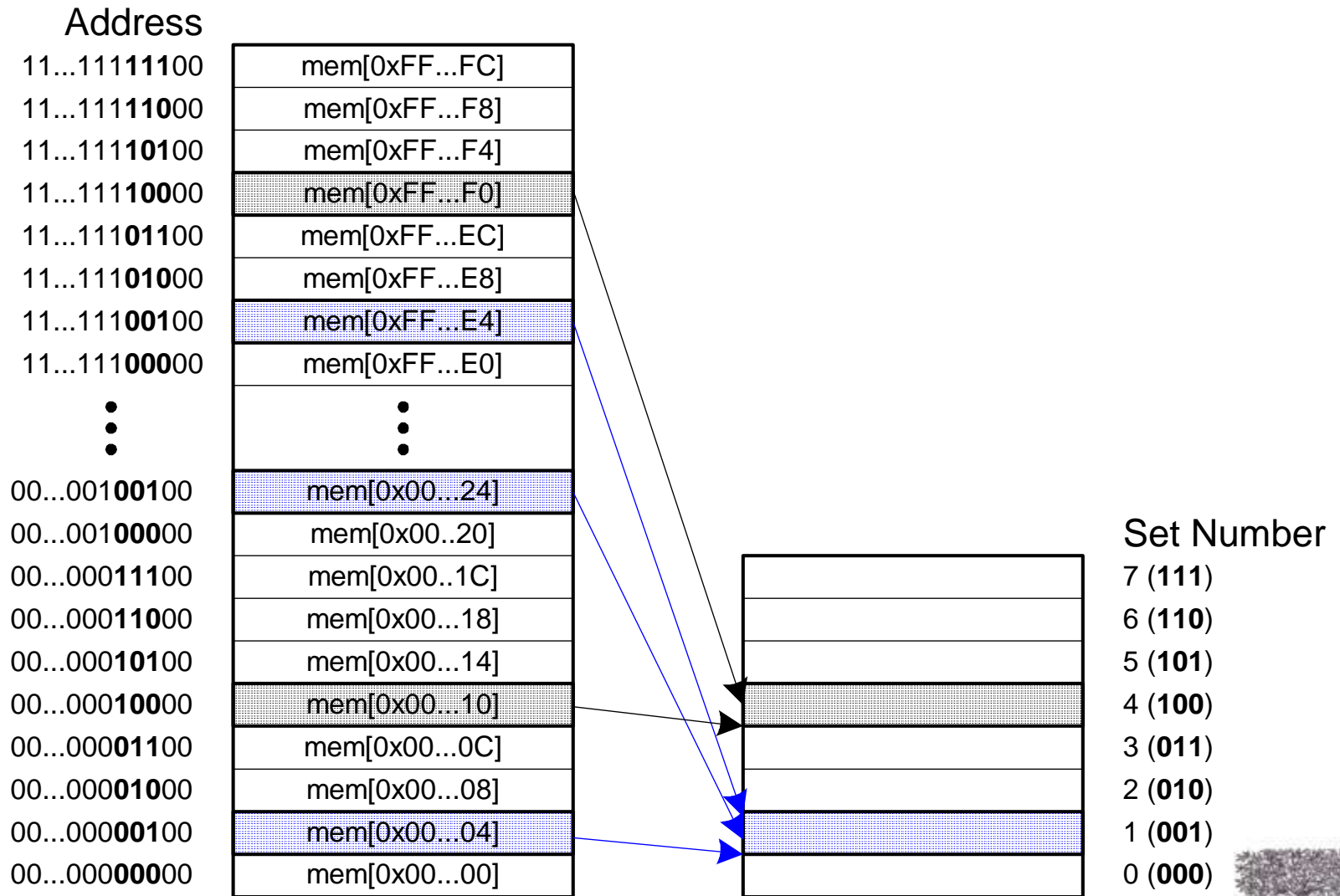  - Block size ($b$ = 1 word)

  - So, number of blocks ($B$ = 8)

ELSEVIER

# Example Cache Parameters

- $C$ = **8** words (capacity)

- $b$ = **1** word (block size)

- So, $B$ = **8** (# of blocks)

**Ridiculously small, but will illustrate organizations**

# Direct Mapped Cache

Address

| | |
|---|---|
| 11...1111**11**00 | mem[0xFF...FC] |
| 11...111**11**00 | mem[0xFF...F8] |
| 11...111**10**00 | mem[0xFF...F4] |
| 11...111**100**00 | mem[0xFF...F0] |
| 11...111**011**00 | mem[0xFF...EC] |
| 11...111**010**00 | mem[0xFF...E8] |
| 11...111**001**00 | mem[0xFF...E4] |
| 11...111**000**00 | mem[0xFF...E0] |

⋮ ⋮

| | |
|---|---|
| 00...001**001**00 | mem[0x00...24] |
| 00...001**000**00 | mem[0x00..20] |
| 00...000**111**00 | mem[0x00..1C] |
| 00...000**110**00 | mem[0x00...18] |
| 00...000**101**00 | mem[0x00...14] |
| 00...000**100**00 | mem[0x00...10] |
| 00...000**011**00 | mem[0x00...0C] |
| 00...000**010**00 | mem[0x00...08] |
| 00...000**001**00 | mem[0x00...04] |
| 00...000**000**00 | mem[0x00...00] |

$2^{30}$ Word Main Memory

Set Number

7 (**111**)
6 (**110**)
5 (**101**)
4 (**100**)
3 (**011**)
2 (**010**)
1 (**001**)
0 (**000**)

$2^{3}$ Word Cache

Chapter 8 <19>

ELSEVIER

# Direct Mapped Cache Hardware



8-entry x
(1+27+32)-bit
SRAM

# Direct Mapped Cache Performance

Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...00 | 001 | 00 |

3

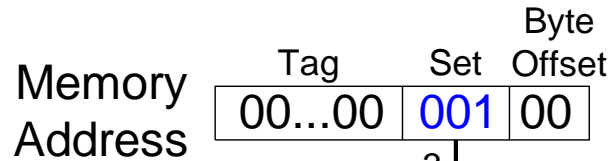| V | Tag | Data | |
|---|-----|------|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 1 | 00...00 | mem[0x00...0C] | Set 3 (011) |
| 1 | 00...00 | mem[0x00...08] | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] | Set 1 (001) |
| 0 | | | Set 0 (000) |

```
# MIPS assembly code

        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0xC($0)
        lw   $t3, 0x8($0)
        addi $t0, $t0, -1
        j    loop
done:
```

**Miss Rate = ?**

ELSEVIER

# Direct Mapped Cache Performance

Memory Address

| Tag | Set | Byte Offset |
|---|---|---|
| 00...00 | 001 | 00 |

3

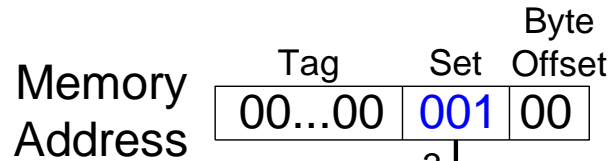| V | Tag | Data | |
|---|---|---|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 1 | 00...00 | mem[0x00...0C] | Set 3 (011) |
| 1 | 00...00 | mem[0x00...08] | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] | Set 1 (001) |
| 0 | | | Set 0 (000) |

```
# MIPS assembly code

       addi $t0, $0, 5
loop:  beq  $t0, $0, done
       lw   $t1, 0x4($0)
       lw   $t2, 0xC($0)
       lw   $t3, 0x8($0)
       addi $t0, $t0, -1
       j    loop
done:
```

**Miss Rate = 3/15**

**= 20%**

**Temporal Locality**

**Compulsory Misses**

© Digital Design and Computer Architecture, 2nd Edition, 2012

ELSEVIER

# Direct Mapped Cache: Conflict

Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...01 | 001 | 00 |

3

```
# MIPS assembly code

        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0x24($0)
        addi $t0, $t0, -1
        j    loop
done:
```

| V | Tag | Data | |
|---|-----|------|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 0 | | | Set 3 (011) |
| 0 | | | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] mem[0x00...24] | Set 1 (001) |
| 0 | | | Set 0 (000) |

**Miss Rate = ?**

ELSEVIER

# Direct Mapped Cache: Conflict

Memory Address

| Tag | Set | Byte Offset |
|-----|-----|-------------|
| 00...01 | 001 | 00 |

3

```
# MIPS assembly code

        addi  $t0, $0, 5
loop:   beq   $t0, $0, done
        lw    $t1, 0x4($0)
        lw    $t2, 0x24($0)
        addi  $t0, $t0, -1
        j     loop

done:
```

| V | Tag | Data | |
|---|-----|------|---|
| 0 | | | Set 7 (111) |
| 0 | | | Set 6 (110) |
| 0 | | | Set 5 (101) |
| 0 | | | Set 4 (100) |
| 0 | | | Set 3 (011) |
| 0 | | | Set 2 (010) |
| 1 | 00...00 | mem[0x00...04] mem[0x00...24] | Set 1 (001) |
| 0 | | | Set 0 (000) |

**Miss Rate = 10/10**

**= 100%**

**Conflict Misses**

ELSEVIER

# *N*-Way Set Associative Cache

# *N*-Way Set Associative Performance

```
# MIPS assembly code

        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0x24($0)
        addi $t0, $t0, -1
        j    loop
done:
```

**Miss Rate = ?**

|  | Way 1 | | | Way 0 | |  |
|---|---|---|---|---|---|---|
| V | Tag | Data | V | Tag | Data |  |
| 0 |  |  | 0 |  |  | Set 3 |
| 0 |  |  | 0 |  |  | Set 2 |
| 0 |  |  | 0 |  |  | Set 1 |
| 0 |  |  | 0 |  |  | Set 0 |

ELSEVIER

# *N*-Way Set Associative Performance

```
# MIPS assembly code


        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0x24($0)
        addi $t0, $t0, -1
        j    loop
done:
```

**Miss Rate = 2/10**

**= 20%**

**Associativity reduces conflict misses**

| Way 1 | | | Way 0 | | | |
|---|---|---|---|---|---|---|
| V | Tag | Data | V | Tag | Data | |
| 0 | | | 0 | | | Set 3 |
| 0 | | | 0 | | | Set 2 |
| 1 | 00...10 | mem[0x00...24] | 1 | 00...00 | mem[0x00...04] | Set 1 |
| 0 | | | 0 | | | Set 0 |

ELSEVIER

# Fully Associative Cache

| V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data | V | Tag | Data |
|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|---|-----|------|
|   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |   |     |      |

**Reduces conflict misses**

**Expensive to build**

# Spatial Locality?

- Increase block size:
  - Block size, **b = 4 words**
  - *C* = 8 words
  - Direct mapped (1 block per set)
  - Number of blocks, **B = 2** (*C/b* = 8/4 = 2)

# Cache with Larger Block Size

# Direct Mapped Cache Performance

```
        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0xC($0)
        lw   $t3, 0x8($0)
        addi $t0, $t0, -1
        j    loop
done:
```
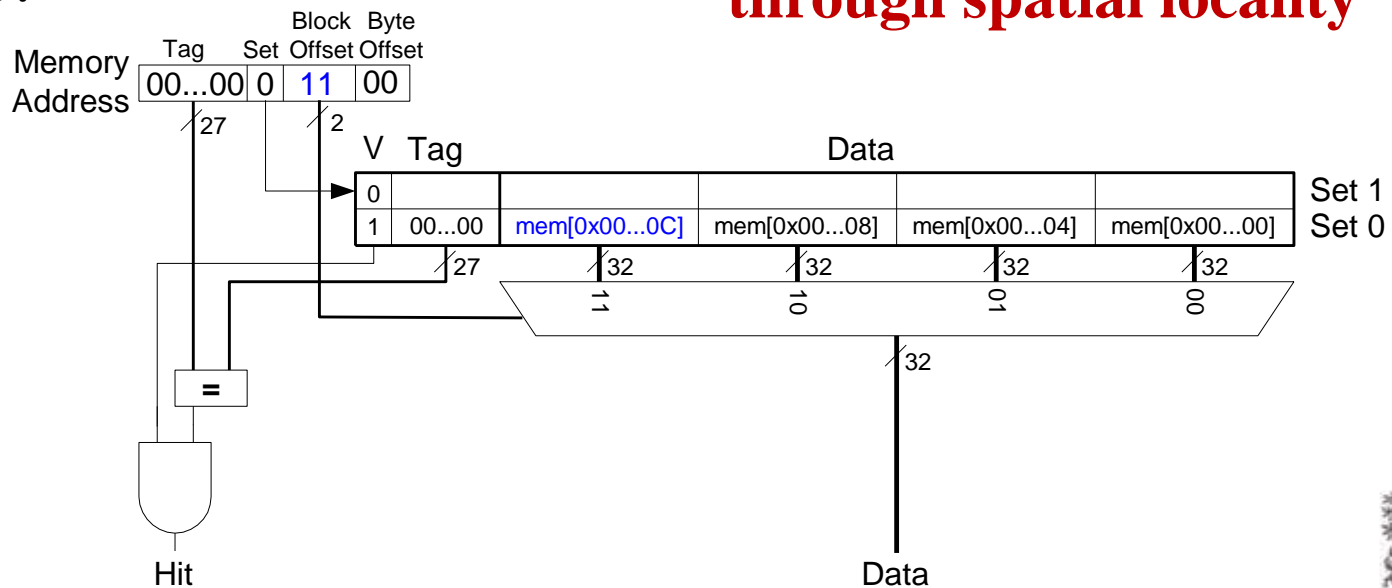
**Miss Rate = ?**

# Direct Mapped Cache Performance

```
        addi $t0, $0, 5
loop:   beq  $t0, $0, done
        lw   $t1, 0x4($0)
        lw   $t2, 0xC($0)
        lw   $t3, 0x8($0)
        addi $t0, $t0, -1
        j    loop
done:
```

**Miss Rate = 1/15**

**= 6.67%**

**Larger blocks**
**reduce compulsory misses**
**through spatial locality**

# Cache Organization Recap

- Capacity: $C$
- Block size: $b$
- Number of blocks in cache: $B = C/b$
- Number of blocks in a set: $N$
- Number of sets: $S = B/N$

| Organization | Number of Ways ($N$) | Number of Sets ($S = B/N$) |
|---|---|---|
| **Direct Mapped** | 1 | $B$ |
| **N-Way Set Associative** | $1 < N < B$ | $B / N$ |
| **Fully Associative** | $B$ | 1 |

# Capacity Misses

- Cache is too small to hold all data of interest at once
- If cache full: program accesses data X & evicts data Y
- *Capacity miss* when access Y again
- How to choose Y to minimize chance of needing it again?
- **Least recently used (LRU) replacement:** the least recently used block in a set evicted

ELSEVIER

# Types of Misses

- **Compulsory:** first time data accessed
- **Capacity:** cache too small to hold all data of interest
- **Conflict:** data of interest maps to same location in cache

**Miss penalty:** time it takes to retrieve a block from lower level of hierarchy

ELSEVIER

# MIPS assembly

```
lw $t0, 0x04($0)
lw $t1, 0x24($0)
lw $t2, 0x54($0)
```

| | | Way 1 | | | Way 0 | | |
|---|---|---|---|---|---|---|---|
| V | U | Tag | Data | V | Tag | Data | |
| 0 | 0 | | | 0 | | | Set 3 (11) |
| 0 | 0 | | | 0 | | | Set 2 (10) |
| 0 | 0 | | | 0 | | | Set 1 (01) |
| 0 | 0 | | | 0 | | | Set 0 (00) |

# LRU Replacement

```
# MIPS assembly

lw $t0, 0x04($0)
lw $t1, 0x24($0)
lw $t2, 0x54($0)
```

| V | U | Tag | Data | V | Tag | Data | |
|---|---|-----|------|---|-----|------|---|
| | Way 1 | | | | Way 0 | | |
| 0 | 0 | | | 0 | | | Set 3 (11) |
| 0 | 0 | | | 0 | | | Set 2 (10) |
| 1 | 0 | 00...010 | mem[0x00...24] | 1 | 00...000 | mem[0x00...04] | Set 1 (01) |
| 0 | 0 | | | 0 | | | Set 0 (00) |

(a)

| V | U | Tag | Data | V | Tag | Data | |
|---|---|-----|------|---|-----|------|---|
| | Way 1 | | | | Way 0 | | |
| 0 | 0 | | | 0 | | | Set 3 (11) |
| 0 | 0 | | | 0 | | | Set 2 (10) |
| 1 | 1 | 00...010 | mem[0x00...24] | 1 | 00...101 | mem[0x00...54] | Set 1 (01) |
| 0 | 0 | | | 0 | | | Set 0 (00) |

(b)

ELSEVIER

# Cache Summary

- **What data is held in the cache?**
  - Recently used data (temporal locality)
  - Nearby data (spatial locality)

- **How is data found?**
  - Set is determined by address of data
  - Word within block also determined by address
  - In associative caches, data could be in one of several ways

- **What data is replaced?**
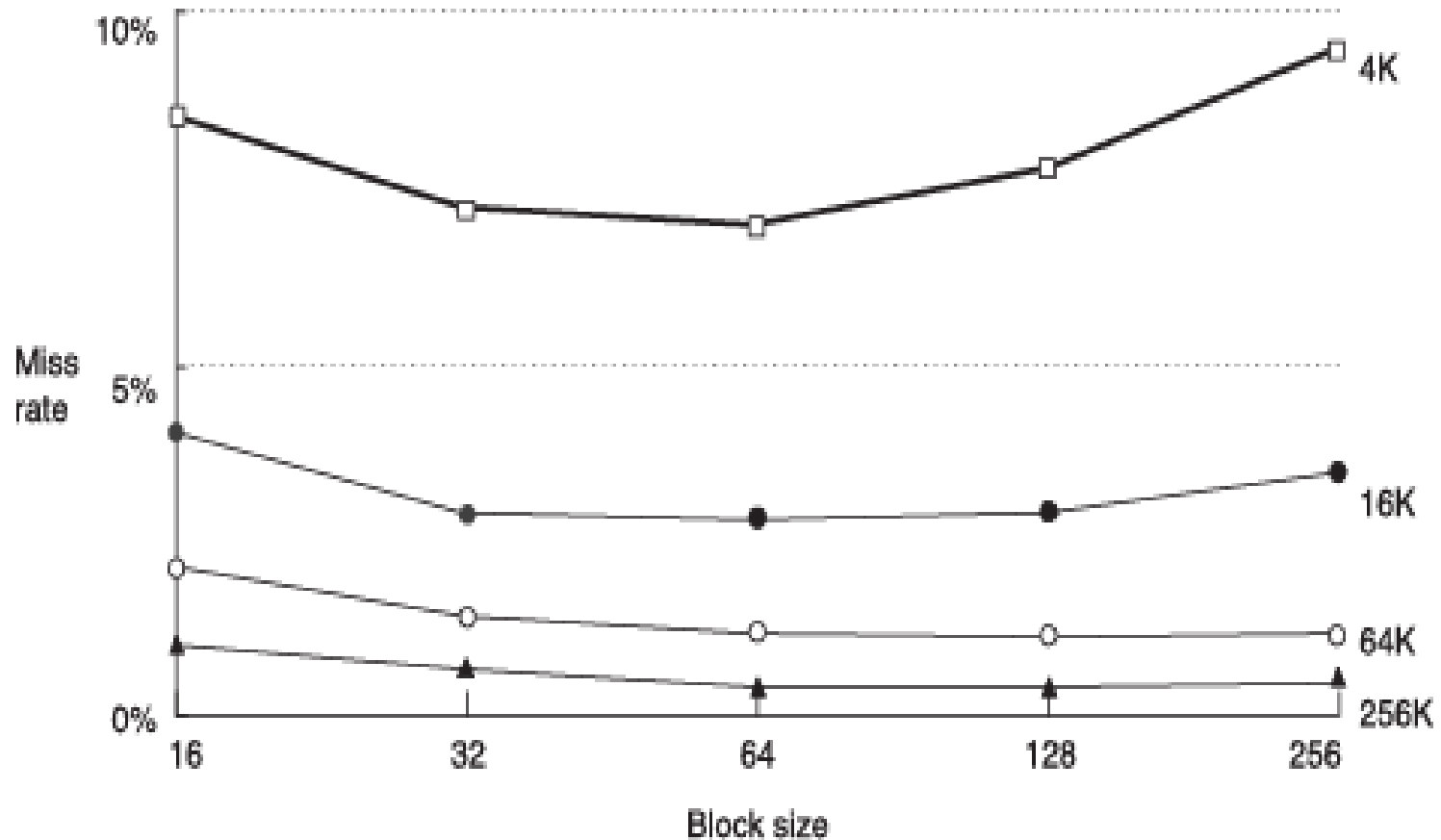  - Least-recently used way in the set

# Miss Rate Trends

- **Bigger caches reduce capacity misses**
- **Greater associativity reduces conflict misses**



**Adapted from Patterson & Hennessy,** *Computer Architecture: A Quantitative Approach,* **2011**

# Miss Rate Trends



- **Bigger blocks reduce compulsory misses**
- **Bigger blocks increase conflict misses**

# Multilevel Caches

- Larger caches have lower miss rates, longer access times

- Expand memory hierarchy to multiple levels of caches

- Level 1: small and fast (e.g. 16 KB, 1 cycle)

- Level 2: larger and slower (e.g. 256 KB, 2-6 cycles)

- Most modern PCs have L1, L2, and L3 cache

# Intel Pentium III Die



256 KB Level 2 Unified Cache

16 KB Level 1 Data Cache

16 KB Level 1 Instruction Cache

# Virtual Memory

- Gives the illusion of bigger memory

- Main memory (DRAM) acts as cache for hard disk

# Memory Hierarchy

| Technology | Price / GB | Access Time (ns) | Bandwidth (GB/s) |
|---|---|---|---|
| SRAM | $10,000 | 1 | 25+ |
| DRAM | $10 | 10 - 50 | 10 |
| SSD | $1 | 100,000 | 0.5 |
| HDD | $0.1 | 10,000,000 | 0.1 |

Speed ↑

Cache

Main Memory

Virtual Memory

Capacity →

- **Physical Memory:** DRAM (Main Memory)
- **Virtual Memory:** Hard drive
  – Slow, Large, Cheap

ELSEVIER

# Hard Disk

Magnetic Disks

Read/Write Head

**Takes milliseconds to *seek* correct location on disk**

# Virtual Memory

- **Virtual addresses**
  - Programs use virtual addresses
  - Entire virtual address space stored on a hard drive
  - Subset of virtual address data in DRAM
  - CPU translates virtual addresses into *physical addresses* (DRAM addresses)
  - Data not in DRAM fetched from hard drive

- **Memory Protection**
  - Each program has own virtual to physical mapping
  - Two programs can use same virtual address for different data
  - Programs don't need to be aware others are running
  - One program (or virus) can't corrupt memory used by another

ELSEVIER

# Cache/Virtual Memory Analogues

| Cache | Virtual Memory |
|---|---|
| Block | Page |
| Block Size | Page Size |
| Block Offset | Page Offset |
| Miss | Page Fault |
| Tag | Virtual Page Number |

**Physical memory acts as cache for virtual memory**

# Virtual Memory Definitions

- **Page size:** amount of memory transferred from hard disk to DRAM at once

- **Address translation:** determining physical address from virtual address

- **Page table:** lookup table used to translate virtual addresses to physical addresses

# Virtual & Physical Addresses

**Most accesses hit in physical memory**

**But programs have the large capacity of virtual memory**

# Address Translation



Virtual Address

30 29 28 ... 14 13 12 11 10 9 ... 2 1 0

| VPN | Page Offset |

19

Translation

12

15

| PPN | Page Offset |

26 25 24 ... 13 12 11 10 9 ... 2 1 0

Physical Address

# Virtual Memory Example

- **System:**

  – Virtual memory size: 2 GB = $2^{31}$ bytes

  – Physical memory size: 128 MB = $2^{27}$ bytes

  – Page size: 4 KB = $2^{12}$ bytes

ELSEVIER

# Virtual Memory Example

- ## System:

    - Virtual memory size: 2 GB = $2^{31}$ bytes

    - Physical memory size: 128 MB = $2^{27}$ bytes

    - Page size: 4 KB = $2^{12}$ bytes

- ## Organization:

    - Virtual address: **31** bits

    - Physical address: **27** bits

    - Page offset: **12** bits

    - # Virtual pages = $2^{31}/2^{12}$ = $\mathbf{2^{19}}$  (VPN = 19 bits)

    - # Physical pages = $2^{27}/2^{12}$ = $\mathbf{2^{15}}$ (PPN = 15 bits)

# Virtual Memory Example

- **19-bit virtual page numbers**
- **15-bit physical page numbers**

# Virtual Memory Example

What is the physical address
of virtual address **0x247C**?

# Virtual Memory Example

What is the physical address
of virtual address **0x247C?**

- VPN = **0x2**
- VPN 0x2 maps to PPN **0x7FFF**
- 12-bit page offset: **0x47C**
- Physical address = **0x7FFF**47C

# How to perform translation?

- **Page table**
  - Entry for each virtual page
  - Entry fields:
    - **Valid bit:** 1 if page in physical memory
    - **Physical page number:** where the page is located

# Page Table Example



VPN is index into page table

What is the physical address of virtual address **0x5F20**?

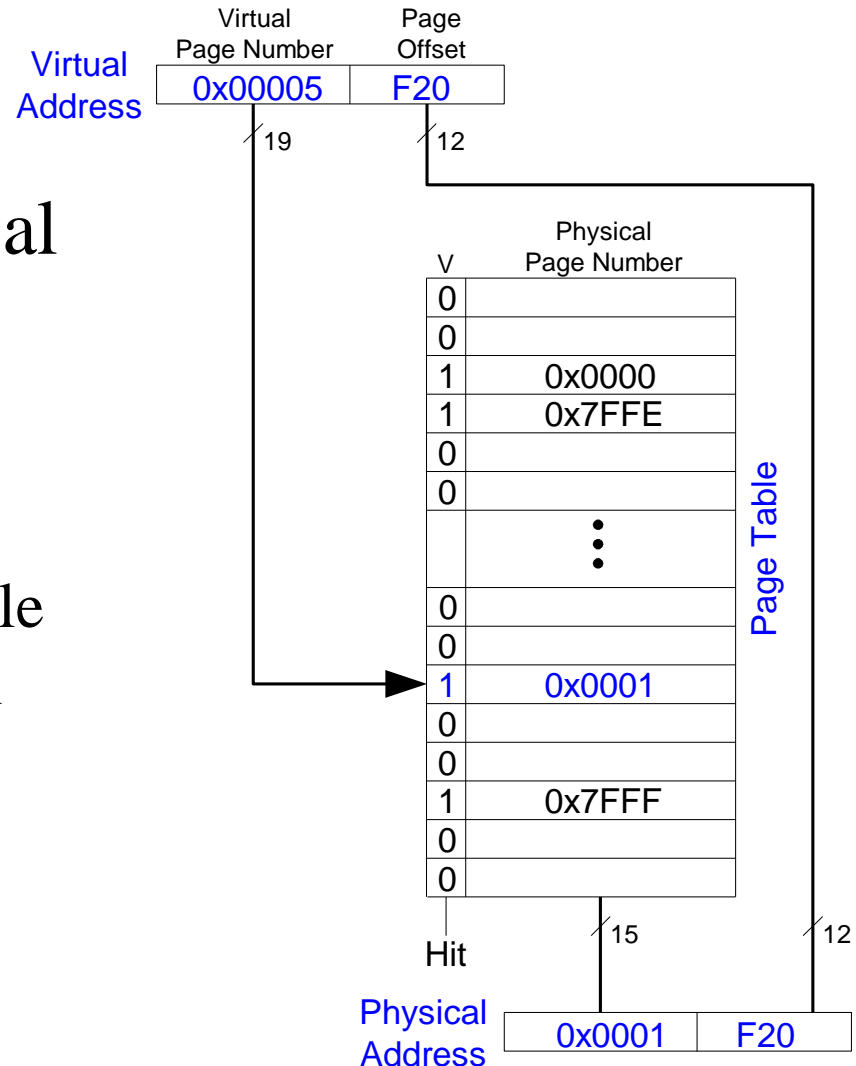| V | Physical Page Number |
|---|---|
| 0 | |
| 0 | |
| 1 | 0x0000 |
| 1 | 0x7FFE |
| 0 | |
| 0 | |
| | ⋮ |
| 0 | |
| 0 | |
| 1 | 0x0001 |
| 0 | |
| 0 | |
| 1 | 0x7FFF |
| 0 | |
| 0 | |

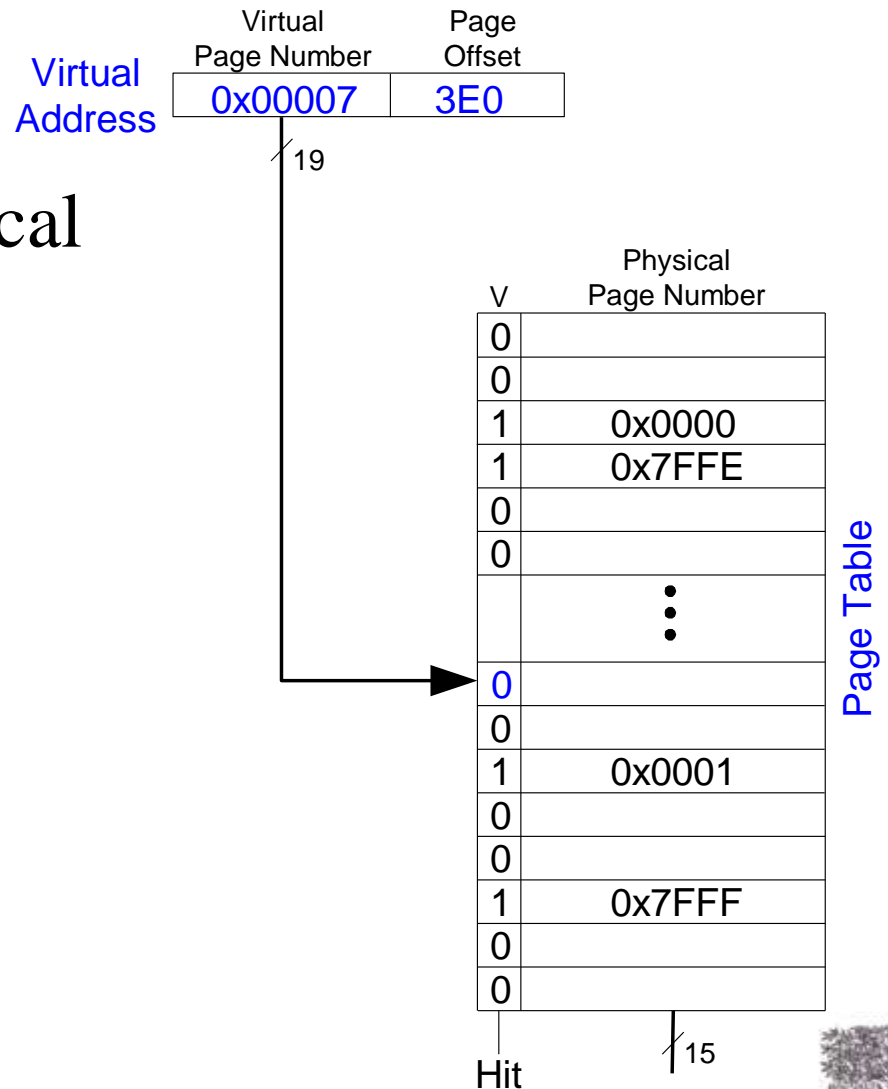Page Table

# Page Table Example 1

What is the physical address of virtual address **0x5F20**?

- – VPN = **5**
- – Entry 5 in page table VPN 5 => physical page **1**
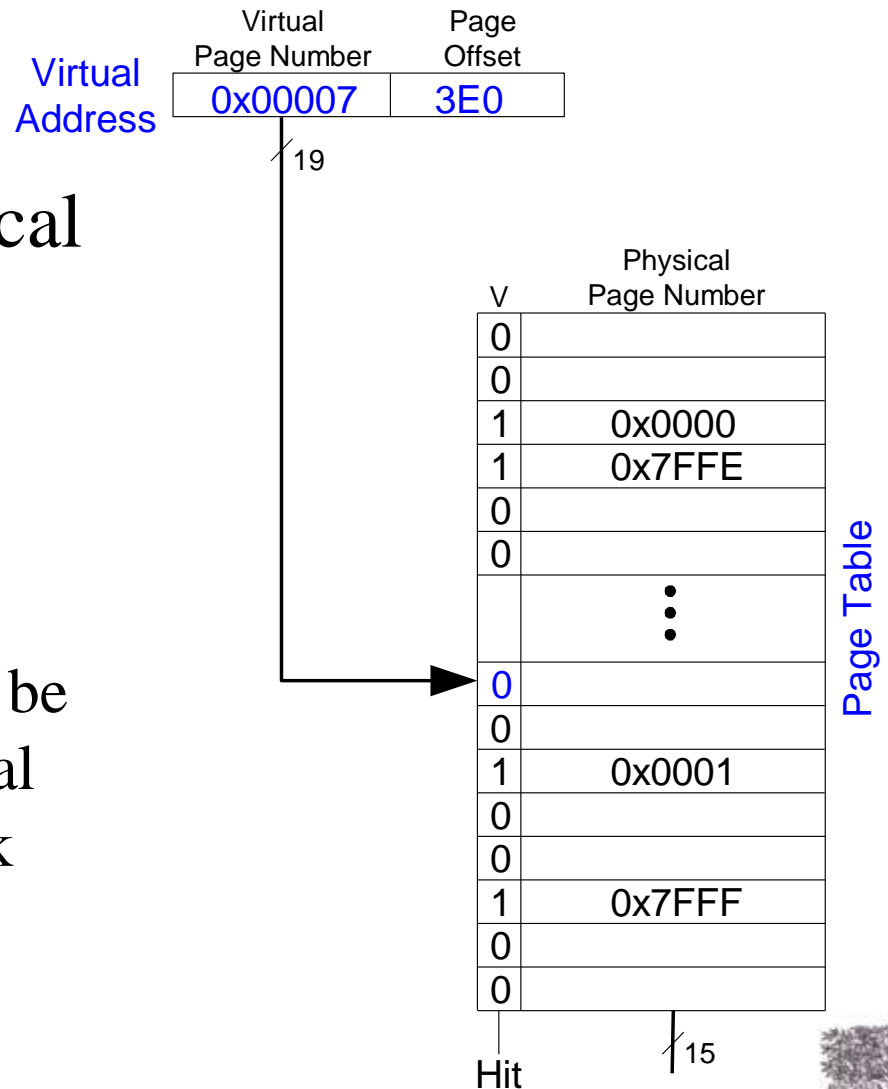- – Physical address: **0x1F20**

# Page Table Example 2

What is the physical address of virtual address **0x73E0**?

Virtual Address

| Virtual Page Number | Page Offset |
|---|---|
| 0x00007 | 3E0 |

19

| V | Physical Page Number |
|---|---|
| 0 | |
| 0 | |
| 1 | 0x0000 |
| 1 | 0x7FFE |
| 0 | |
| 0 | |
| ⋮ | |
| 0 | |
| 0 | |
| 1 | 0x0001 |
| 0 | |
| 0 | |
| 1 | 0x7FFF |
| 0 | |
| 0 | |

Page Table

15

Hit

ELSEVIER

# Page Table Example 2

What is the physical address of virtual address **0x73E0**?

- VPN = **7**
- Entry 7 is invalid
- Virtual page must be *paged* into physical memory from disk

Virtual Page Number

Page Offset

**Virtual Address**

| 0x00007 | 3E0 |

19

Physical Page Number

| V | |
|---|---|
| 0 | |
| 0 | |
| 1 | 0x0000 |
| 1 | 0x7FFE |
| 0 | |
| 0 | |
| ⋮ | |
| 0 | |
| 0 | |
| 1 | 0x0001 |
| 0 | |
| 0 | |
| 1 | 0x7FFF |
| 0 | |
| 0 | |

Page Table

Hit

15

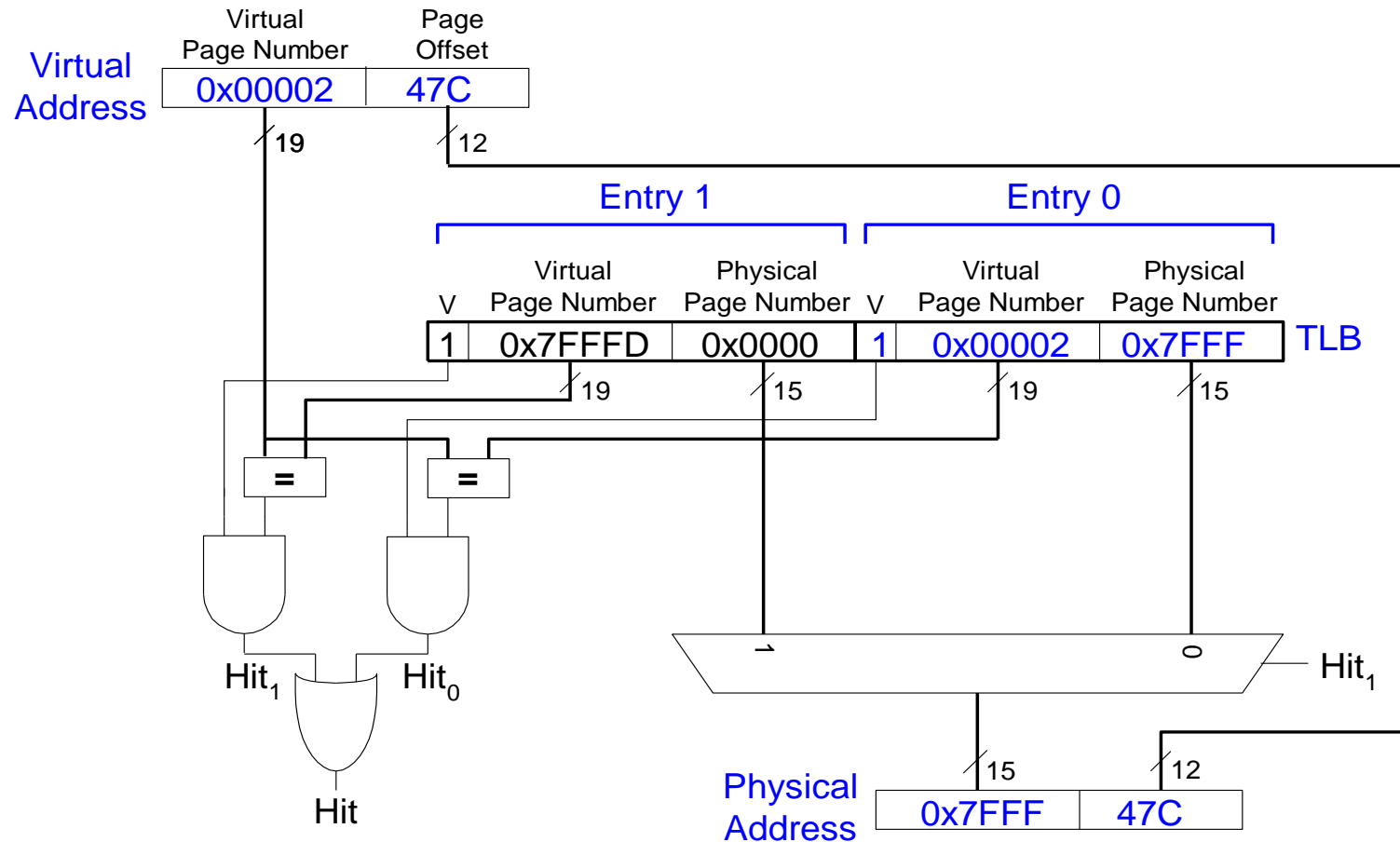ELSEVIER

# Page Table Challenges

- **Page table is large**
  - usually located in physical memory
- Load/store requires 2 main memory accesses:
  - one for translation (page table read)
  - one to access data (after translation)
- Cuts memory performance in half
  - *Unless we get clever...*

# Translation Lookaside Buffer (TLB)

- Small cache of most recent translations

- Reduces # of memory accesses for *most* loads/stores from 2 to 1

ELSEVIER

# TLB

- Page table accesses: high temporal locality
  - Large page size, so consecutive loads/stores likely to access same page
- TLB
  - Small: accessed in < 1 cycle
  - Typically 16 - 512 entries
  - Fully associative
  - > 99 % hit rates typical
  - Reduces # of memory accesses for most loads/stores from 2 to 1

# Example 2-Entry TLB

# Memory Protection

- Multiple processes (programs) run at once

- Each process has its own page table

- Each process can use entire virtual address space

- A process can only access physical pages mapped in its own page table

ELSEVIER

# Virtual Memory Summary

- Virtual memory increases **capacity**

- A subset of virtual pages in physical memory

- **Page table** maps virtual pages to physical pages – address translation

- A **TLB** speeds up address translation

- Different page tables for different programs provides **memory protection**
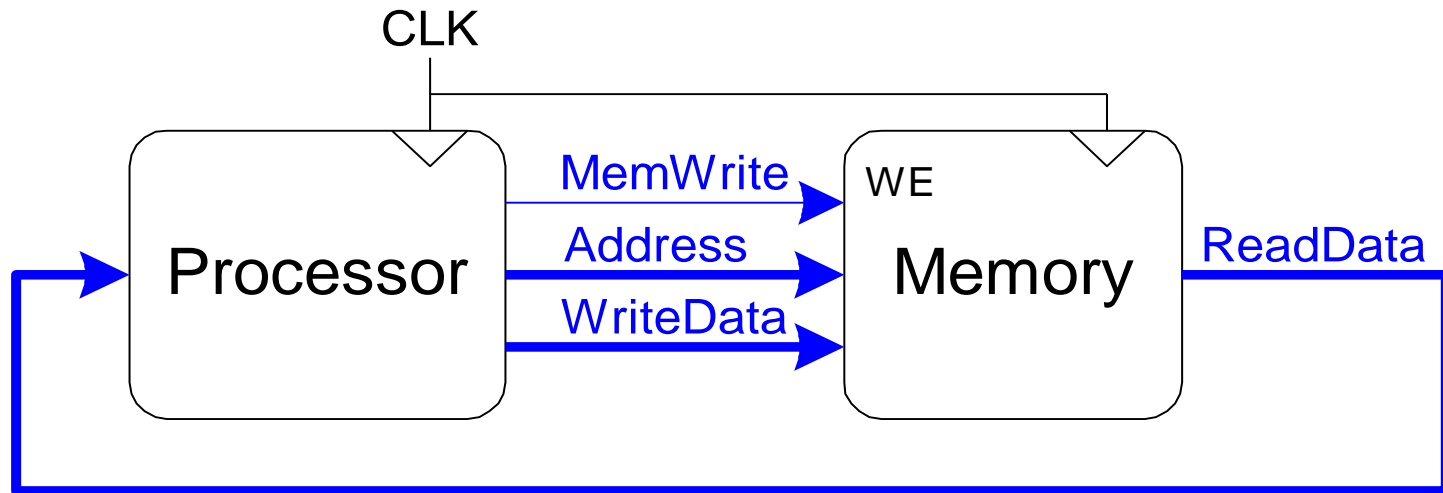
ELSEVIER

# Memory-Mapped I/O

- Processor accesses I/O devices just like memory (like keyboards, monitors, printers)
- Each I/O device assigned one or more address
- When that address is detected, data read/written to I/O device instead of memory
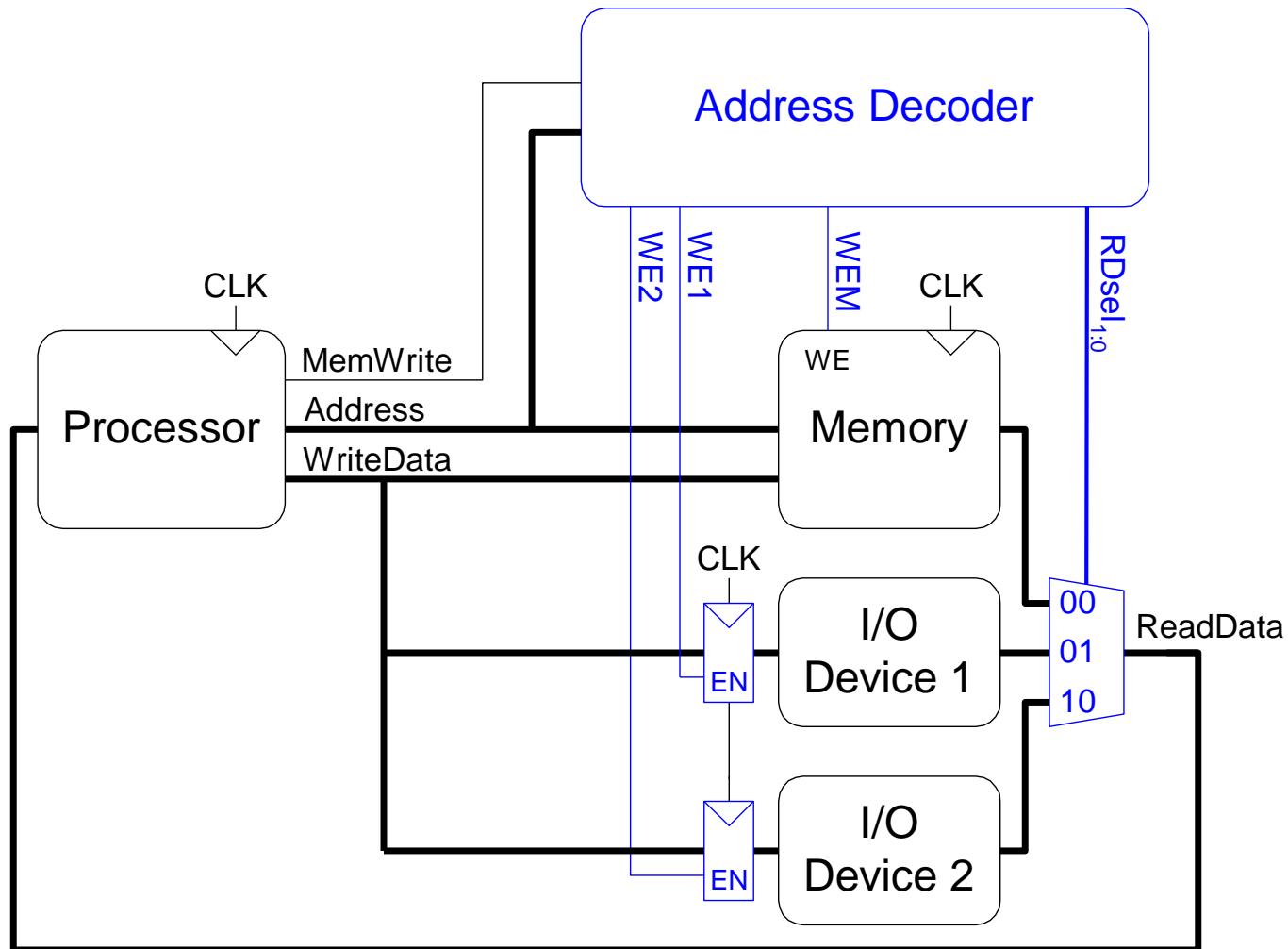- A portion of the address space dedicated to I/O devices

# Memory-Mapped I/O Hardware

- ## Address Decoder:
  - Looks at address to determine which device/memory communicates with the processor

- ## I/O Registers:
  - Hold values written to the I/O devices

- ## ReadData Multiplexer:
  - Selects between memory and I/O devices as source of data sent to the processor

# The Memory Interface
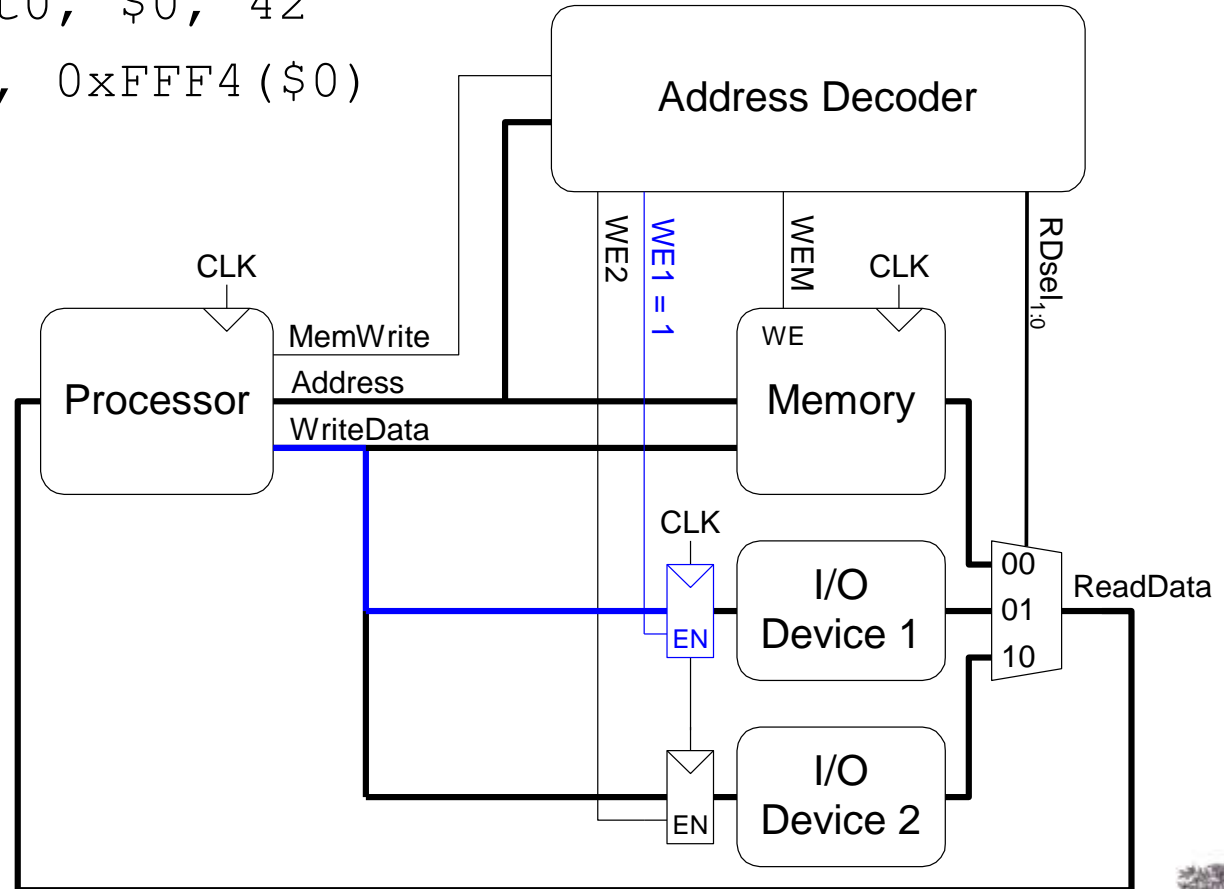
# Memory-Mapped I/O Hardware

# Memory-Mapped I/O Code

- Suppose I/O Device 1 is assigned the address 0xFFFFFFF4

    – **Write the value 42 to I/O Device 1**

    – **Read value from I/O Device 1 and place in $t3**

# Memory-Mapped I/O Code

- **Write the value 42 to I/O Device 1 (0xFFFFFFF4)**

```
addi $t0, $0, 42
sw $t0, 0xFFF4($0)
```

# Memory-Mapped I/O Code

- **Read the value from I/O Device 1 and place in $t3**

```
lw $t3, 0xFFF4($0)
```