## Lecture 4: Performance measurement and Instruction Set Architectures

- Last Time
  - Introduction to performance
  - Computer benchmarks

- Today
  - Equations for performance analysis
  - Amdahl's Law
  - Introduction to ISAs

## Review: latency vs. throughput

- Pizza delivery example
  - Do you want your pizza hot?
    - Low latency
  - Or do you want your pizza to be inexpensive?
    - High throughput – lots of pizzas per hour
  - Two different delivery strategies for pizza company!

In this course:

We will focus primarily on latency (execution time for a single task)

## Improving Performance: Fundamentals

- Suppose we have a machine with two instructions
  - Instruction A executes in 100 cycles
  - Instruction B executes in 2 cycles

- We want better performance….
  - Which instruction do we improve?

## Speedup

- Make a change to an architecture
- Measure how much faster/slower it is

$$\text{Speedup} \ = \ \frac{\text{ExTime before change}}{\text{ExTime after change}}$$

## Speedup when we know details about the change

- Performance improvements depend on:
  - how good is enhancement (factor **S**)
  - how often is it used (fraction **p**)
- Speedup due to enhancement E:

$$\text{Speedup(E)} = \frac{\text{ExTime w/out E}}{\text{ExTime w/ E}} = \frac{\text{Perf w/ E}}{\text{Perf w/out E}}$$

$$ExTime_{new} = ExTime_{old} * \left[ (1-p) + \frac{p}{S} \right]$$

$$Speedup(E) = \frac{ExTime_{old}}{ExTime_{new}} = \frac{1}{(1-p) + \frac{p}{S}}$$

## Example

- FP instructions improved by 2x
- But….only 10% of instructions are FP

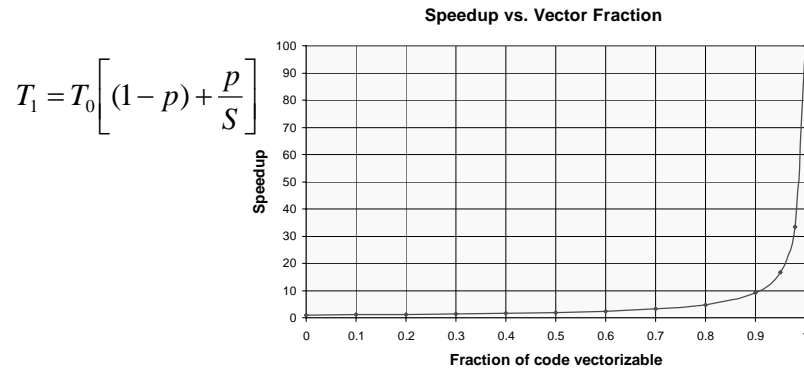$$ExTime_{new} = ExTime_{old} * \left( 0.9 + \frac{0.1}{2} \right) = 0.95 * ExTime_{old}$$

$$Speedup_{total} = \frac{1}{0.95} = 1.053$$

- Amdahl's Law:
  Speedup bounded by $\dfrac{1}{\text{fraction of time not enhanced}}$

## Amdahl's Law: Example 2

• Speed up vectorizable code by 100x

$$T_1 = T_0 \left[ (1-p) + \frac{p}{S} \right]$$

**Speedup vs. Vector Fraction**



(y-axis: **Speedup**, 0 to 100; x-axis: **Fraction of code vectorizable**, 0 to 1)

---

## Amdahl's Law: Summary message

• # Make the common case fast

• Examples:
  – All instructions require instruction fetch, only fraction require data
    $\Rightarrow$ optimize instruction access first

  – Data locality (spatial, temporal), small memories faster
    $\Rightarrow$ storage hierarchy: most frequent accesses to small, local memory

# CPU Performance Equation

- 3 components to execution time:

$$\text{CPU time} = \frac{\text{Seconds}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} * \frac{\text{Cycles}}{\text{Instruction}} * \frac{\text{Seconds}}{\text{Cycle}}$$

- Factors affecting CPU execution time:

|  | Inst. Count | CPI | Clock Rate |
|---|---|---|---|
| Program | X | | |
| Compiler | X | (X) | |
| Inst. Set | X | X | (X) |
| Organization | | X | X |
| MicroArch | | X | X |
| Technology | | | X |

- Consider all three elements when optimizing
- Workloads change!

---

# Cycles Per Instruction (CPI)

- Depends on the instruction

$$CPI_i = \text{Execution time of instruction } i * \text{Clock Rate}$$

- Average cycles per instruction

$$CPI = \sum_{i=1}^{n} CPI_i * F_i \quad \text{where } F_i = \frac{IC_i}{IC_{tot}}$$

- Example:

| Op | Freq | Cycles | CPI(i) | %time |
|---|---|---|---|---|
| ALU | 50% | 1 | 0.5 | 33% |
| Load | 20% | 2 | 0.4 | 27% |
| Store | 10% | 2 | 0.2 | 13% |
| Branch | 20% | 2 | 0.4 | 27% |
| | | CPI(total) | 1.5 | |

## Comparing and Summarizing Performance

- Fair way to summarize performance?
- Capture in a single number?

- Example:  Which of the following machines is best?

| | Computer A | Computer B | Computer C |
|---|---|---|---|
| Program 1 | 1 | 10 | 20 |
| Program 2 | 1000 | 100 | 20 |
| Total Time | 1001 | 110 | 40 |

---

## Means

Arithmetic mean $\quad \dfrac{1}{n}\sum_{i=1}^{n} T_i \quad$ Can be weighted:  $a_i T_i$
Represents total execution time

Harmonic mean $\quad \dfrac{n}{\sum_{i=1}^{n} \dfrac{1}{R_i}} \quad R_i = 1/T_i$

Geometric mean $\quad \left( \prod_{i=1}^{n} \dfrac{T_i}{T_{ri}} \right)^{\frac{1}{n}} \quad$ Good for mean of ratios,
where the ratio is with respect to a reference.

## Platform Performance Trends



Log of Performance (y-axis) vs. Year (x-axis)

Supercomputers
Mainframes
Minicomputers
Microprocessors

1970  1975  1980  1985  1990  1995

---

## Is Speed the Last Word in Performance?

- Depends on the application!
- Cost
  - Not just processor, but other components (ie. memory)
- Power consumption
  - Trade power for performance in many applications
- Capacity
  - Many database applications are I/O bound and disk bandwidth is the precious commodity

## Performance Measurement Summary

- Best benchmarks are real programs
  - Spec, TPC, Doom3

- Pitfalls still exist
  - Whole system measurement
  - Workload may not match user's

- Key concepts
  - Throughput and Latency
  - Equations: CPI, Amdahl's Law, …

- Next time
  - Instruction set architectures (ISA)
  - Read P&H 2.1 – 2.6

---

## Instruction set architectures

# LC-3 Instruction Set

|       | 15 14 13 12 | 11 10 9 | 8 7 6 | 5 | 4 3 | 2 1 0 |
|-------|-------------|---------|-------|---|-----|-------|
| ADD+  | 0001 | DR | SR1 | 0 | 00 | SR2 |
| ADD+  | 0001 | DR | SR1 | 1 | imm5 | |
| AND+  | 0101 | DR | SR1 | 0 | 00 | SR2 |
| AND+  | 0101 | DR | SR1 | 1 | imm5 | |
| BR    | 0000 | n z p | PCoffset9 | | | |
| JMP   | 1100 | 000 | BaseR | 000000 | | |
| JSR   | 0100 | 1 | PCoffset11 | | | |
| JSRR  | 0100 | 0 00 | BaseR | 000000 | | |
| LD+   | 0010 | DR | PCoffset9 | | | |
| LDI+  | 1010 | DR | PCoffset9 | | | |

⋮

---

# LC-3 ISA – Continued

- Memory
  - 16-bit data, 16-bit addresses
- Addressing modes
  - register, immediate
  - PC+offset, base+offset, indirect PC+offset
- 8 registers
  - R7 updated on JSR, JSRR
- State
  - PC, registers, memory, condition code registers
- 16-bit instructions
  - Simple instruction set

# Write programs using the ISA



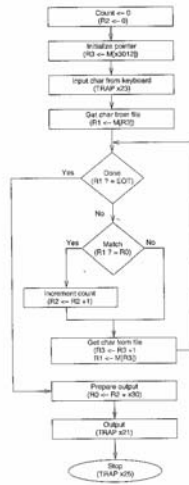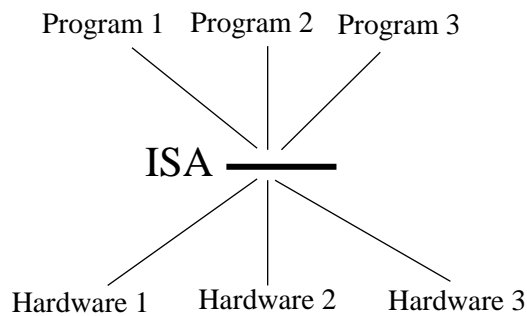Figure 5.17  A machine language program that implements the algorithm of Figure 5.16

Figure 5.16  An algorithm to count occurrences of a character

Lecture 4

---

# ISA is an interface (abstraction layer)

Program 1    Program 2    Program 3

ISA ——

Hardware 1    Hardware 2    Hardware 3

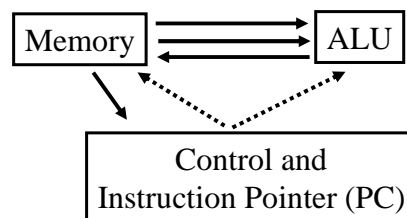Lecture 4

## Instruction Set Architecture is a Contract

- **Contract** between programmer and the hardware
  - Defines visible state of the system
  - Defines how state changes in response to instructions

- Programmer:
  - ISA is model of how a program will execute
- Hardware Designer:
  - ISA is formal definition of the correct way to execute a program

- ISA specification
  - The binary encodings of the instruction set
  - How instructions modify state of the machine
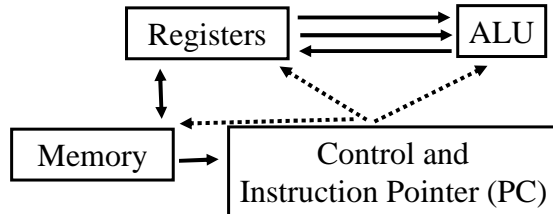
## ISA includes a model of the machine

A very simple model....

## A more typical ISA machine model

Registers → ALU

Memory → Control and Instruction Pointer (PC)

---

## ISA Basics

Instruction formats
Instruction types
Addressing modes

instruction

| Op | Mode | Ra | Rb |
| --- | --- | --- | --- |

Mem

Regs

Before State

Data types
Operations
Interrupts/Events

Mem

Regs

After State

Machine state
Memory organization
Register organization

## Architecture vs. Implementation

- Architecture: defines what a computer system does in response to a program and a set of data
  - Programmer visible elements of computer system

- Implementation: defines how a computer does it
  - Sequence of steps to complete operations
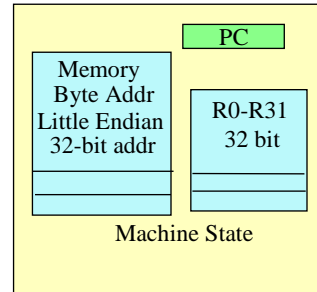  - Time to execute each operation
  - Hidden "bookkeeping" functions

## Examples

- Architecture or Implementation?

  - Number of GP registers
  - Width of memory bus
  - Binary representation of the instruction
    `sub r4,r2,#27`
  - Number of cycles to execute FP instruction
  - How condition code bits are set on a move instruction
  - Size of the instruction cache
  - Type of FP format
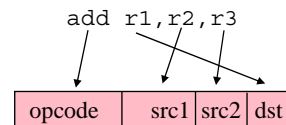
## Machine State

- Registers
  - Size/Type
    - Program Counter (= IP)
    - accumulators
    - index registers
    - general registers
    - control registers
- Memory
  - Visible hierarchy (if any)
  - Addressibility
    - byte, word, bit
    - byte order (endian-ness)
    - maximum size
  - protection/relocation

| PC |
| --- |

| Memory<br>Byte Addr<br>Little Endian<br>32-bit addr | R0-R31<br>32 bit |
| --- | --- |

Machine State

---

## Components of Instructions

- Operations (opcodes)
- Number of operands
- Operand specifiers

```
add r1,r2,r3
```

| opcode | | src1 | src2 | dst |
| --- | --- | --- | --- | --- |

- Instruction encodings
- Instruction classes
  - ALU ops (add, sub, shift)
  - Branch (beq, bne, etc.)
  - Memory (ld/st)

## Operand Number

- No Operands       HALT
  NOP

- 1 operand       NOT R4       R4 $\Leftarrow$ R4
  JMP _L1

- 2 operands       ADD R1, R2       R1 $\Leftarrow$ R1 + R2
  LDI R3, #1234

- 3 operands       ADD R1, R2, R3       R1 $\Leftarrow$ R2 + R3

- > 3 operands       MADD R4,R1,R2,R3    R4 $\Leftarrow$ R1+(R2*R3)

---

## Effect of Operand Number

```
E = (C+D)*(C-D)
```
                             **Assign**
                             **C $\Rightarrow$ r1**
                             **D $\Rightarrow$ r2**
                             **E $\Rightarrow$ r3**

**3 operand machine**            **2 operand machine**

```
add  r3,r1,r2          mov  r3,r1
sub  r4,r1,r2          add  r3,r2
mult r3,r4,r3          sub  r2,r1
                       mult r3,r2
```

# Summary

- ISA definition
  - system state (general/special registers, memory)
  - the effect of each operation on the system state

- Next Time
  - Homework #1 is due – at start of class
  - Addressing modes
  - Data types
  - Common instruction types
  - Case studies: MIPS + others
- Reading assignment – Appendix A.1-A.6 (on CD!)