

CSC 411

Computer Organization (Spring 2022)

Lecture 4: Performance

Prof. Marco Alvarez, University of Rhode Island

SI prefixes

Factor	Name	Symbol	Factor	Name	Symbol
10^{24}	yotta	Y	10^{-1}	deci	d
10^{21}	zetta	Z	10^{-2}	centi	c
10^{18}	exa	E	10^{-3}	milli	m
10^{15}	peta	P	10^{-6}	micro	μ
10^{12}	tera	T	10^{-9}	nano	n
10^9	giga	G	10^{-12}	pico	p
10^6	mega	M	10^{-15}	femto	f
10^3	kilo	k	10^{-18}	atto	a
10^2	hecto	h	10^{-21}	zepto	z
10^1	deka	da	10^{-24}	yocto	y

"The 20 SI prefixes used to form decimal multiples and submultiples of SI units."

<https://physics.nist.gov/cuu/Units/prefixes.html>

Quick notes

Getting access to a Linux machine

- download your favorite OS and install
 - standalone or side-by-side with Windows
- built-in Unix subsystem for Windows 10
 - <https://docs.microsoft.com/en-us/windows/wsl/install-win10>
- use a virtual machine

Prefixes for binary multiples

Factor	Name	Symbol	Origin	Derivation
2^{10}	kibi	Ki	kilobinary: $(2^{10})^1$	kilo: $(10^3)^1$
2^{20}	mebi	Mi	megabinary: $(2^{10})^2$	mega: $(10^3)^2$
2^{30}	gibi	Gi	gigabinary: $(2^{10})^3$	giga: $(10^3)^3$
2^{40}	tebi	Ti	terabinary: $(2^{10})^4$	tera: $(10^3)^4$
2^{50}	pebi	Pi	petabinary: $(2^{10})^5$	peta: $(10^3)^5$
2^{60}	exbi	Ei	exabinary: $(2^{10})^6$	exa: $(10^3)^6$

"In December 1998 the International Electrotechnical Commission (IEC), the leading international organization for worldwide standardization in electrotechnology, approved as an IEC International Standard names and symbols for prefixes for binary multiples for use in the fields of data processing and data transmission."

<https://physics.nist.gov/cuu/Units/binary.html>

Examples

Examples and comparisons with SI prefixes

one kibibit $1 \text{ Kibit} = 2^{10} \text{ bit} = 1024 \text{ bit}$

one kilobit $1 \text{ kbit} = 10^3 \text{ bit} = 1000 \text{ bit}$

one byte $1 \text{ B} = 2^3 \text{ bit} = 8 \text{ bit}$

one mebibyte $1 \text{ MiB} = 2^{20} \text{ B} = 1\,048\,576 \text{ B}$

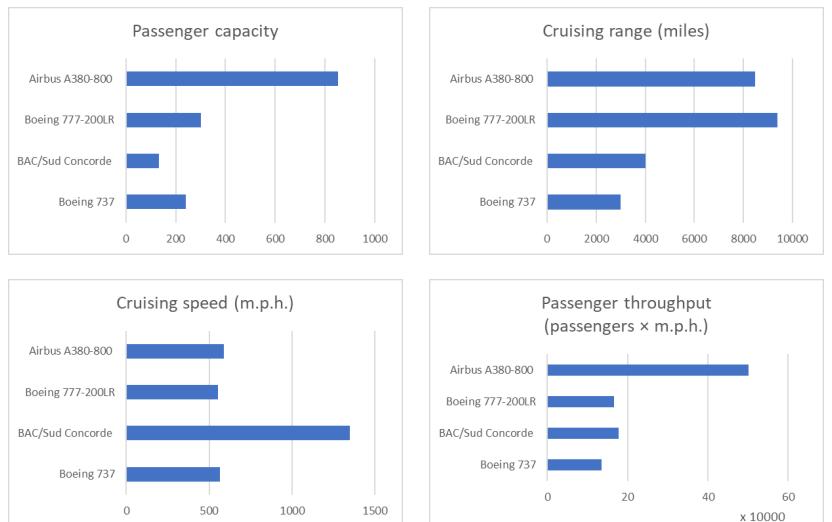
one megabyte $1 \text{ MB} = 10^6 \text{ B} = 1\,000\,000 \text{ B}$

one gibibyte $1 \text{ GiB} = 2^{30} \text{ B} = 1\,073\,741\,824 \text{ B}$

one gigabyte $1 \text{ GB} = 10^9 \text{ B} = 1\,000\,000\,000 \text{ B}$

Performance

Think about performance ...



Understanding performance

Algorithm

- determines number of operations executed

Programming language, compiler, architecture

- determines number of machine instructions executed per operation

Processor and memory system

- determines how fast instructions are executed

I/O system (including OS)

- determines how fast I/O operations are executed

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
- › We'll focus on **response time** for now...

Quizz

- › In the following scenarios, what are the changes:
(A) decreasing response time, (B) increasing throughput, or (C) both?
 - replacing the processor in a computer with a faster version
 - adding additional processors to a system that uses multiple processors for separate tasks

Performance and speedup

- › Performance
 - math display style
$$\text{performance}_X = \frac{1}{\text{execution time}_X}$$
- › Relative performance (speedup):
 - "X is n times faster than Y"
- math display style
$$\frac{\text{performance}_X}{\text{performance}_Y} = \frac{\text{execution time}_Y}{\text{execution time}_X} = n$$
- › Example: time taken to run a program
 - 10s on A, 15s on B
 - $n = 15s / 10s = 1.5$
 - A is **1.5** times faster than B

Measuring execution time

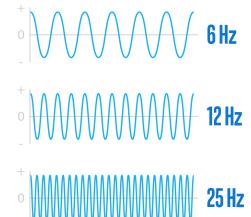
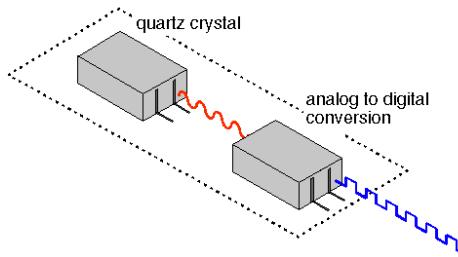
- › **Elapsed time** (wall clock time, response time)
 - total response time, including all aspects: processing, I/O, OS overhead, idle time
 - determines system performance
- › **CPU execution time (CPU time)**
 - time spent processing a job, discounts I/O time, other jobs' shares
 - comprises user **CPU time** and **system CPU time** (difficult to separate accurately)
- › **Different applications are affected by different performance aspects**
 - one must have a clear definition of what performance metric to use to improve system/program performance

```
$ time ls
$ time host uri.edu
$ time (seq 1000 | wc -l)
$ time (seq 10000000 | wc -l)
```

The time command may show different outputs if run on a Linux or Mac systems, in general it shows “elapsed time”, “CPU time”, and “system time”

CPU clocking

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

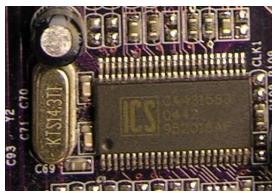


“A clock cycle is a single electronic pulse of a CPU. During each cycle, a CPU can perform a basic operation. Most CPU processes require multiple clock cycles.”

Clock generator

A **clock generator** is an electronic oscillator that produces a **clock signal** for use in synchronizing a circuit's operation. The signal can range from a simple symmetrical **square wave** to more complex arrangements. The basic parts that all clock generators share are a resonant circuit and an amplifier.

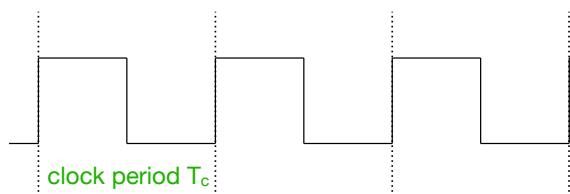
The resonant circuit is usually a **quartz piezo-electric oscillator**, although simpler **tank circuits** and even **RC circuits** may be used.



https://en.wikipedia.org/wiki/Clock_generator

CPU clocking

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



$$f_c = \frac{1}{T_c}$$

CPU execution time

CPU time = CPU clock cycles × 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

Example

- Computer A: 2GHz clock, 10s CPU time
- Designing Computer B aiming for 6s CPU time
 - can do faster clock, but causes 1.2 times the number of clock cycles required by the program
- What clock rate should the designer target?

$$\text{clock cycles}_A = \text{CPU time}_A \times \text{clock rate}_A = 10s \times 2\text{GHz} = 20 \times 10^9$$

$$\text{clock rate}_B = \frac{\text{clock cycles}_B}{\text{CPU time}_B} = \frac{1.2 \times \text{clock cycles}_A}{6s}$$

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

Instruction count and CPI

- Instruction Count for a program
 - previous formulas do not consider number of instructions
 - instructions are determined by program, ISA and compiler
- Average cycles per instruction (CPI)
 - average of clock cycles each instruction takes to execute, over all the instructions executed by the program
 - determined by the CPU hardware

$$\text{CPU clock cycles} = \text{instruction count} \times \text{CPI}$$

$$\begin{aligned}\text{CPU time} &= \text{instruction count} \times \text{CPI} \times \text{clock cycle time} \\ &= \frac{\text{instruction count} \times \text{CPI}}{\text{clock rate}}\end{aligned}$$

Example

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

$$\begin{aligned}\text{CPU time}_A &= \text{ic} \times \text{CPI}_A \times \text{cycle time}_A \\ &= \text{ic} \times 2.0 \times 250\text{ps} = \text{ic} \times 500\text{ps}\end{aligned}$$

$$\begin{aligned}\text{CPU time}_B &= \text{ic} \times \text{CPI}_B \times \text{cycle time}_B \\ &= \text{ic} \times 1.2 \times 500\text{ps} = \text{ic} \times 600\text{ps}\end{aligned}$$

$$\frac{\text{CPU time}_B}{\text{CPU time}_A} = \frac{\text{ic} \times 600\text{ps}}{\text{ic} \times 500\text{ps}} = 1.2$$

CPI in more detail

- If different instruction classes take different numbers of cycles

$$\text{CPU clock cycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{instruction count}_i)$$

- Calculating the (average) CPI using a weighted average

$$\text{CPI} = \frac{\text{CPU clock cycles}}{\text{instruction count}} = \sum_{i=1}^n \left(\text{CPI}_i \times \frac{\text{instruction count}_i}{\text{instruction count}} \right)$$

Quizz

- Consider 3 different processors executing the same instructions
 - P1 has a clock cycle time of 0.33ns and a CPI of 1.5
 - P2 has a clock cycle time of 0.40ns and a CPI of 1
 - P3 has a clock cycle time of 0.25ns and a CPI of 2.2
- Questions
 - which has the highest clock rate? what is it?
 - which is the fastest computer? which is the slowest?

CPI example

- Alternative compiled code sequences using instructions in classes A, B, C. Which code sequence executes the most instructions? Which sequence is faster? What is the CPI for each sequence?

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

- clock cycles = $2 \times 1 + 1 \times 2 + 2 \times 3 = 10$
- average CPI = $10/5 = 2.0$

- Sequence 2

- clock cycles = $4 \times 1 + 1 \times 2 + 1 \times 3 = 9$
- average CPI = $9/6 = 1.5$

Quizz

Quizz

- Case 1:
 - a program takes 4 billion instructions
 - each instruction completes in an average of 1.5 cycles
 - clock speed is 1 GHz
- Case 2:
 - a program takes 2 billion instructions
 - each instruction completes in an average of 6 cycles
 - clock speed is 1.5 GHz
- Which one is better?

Performance summary

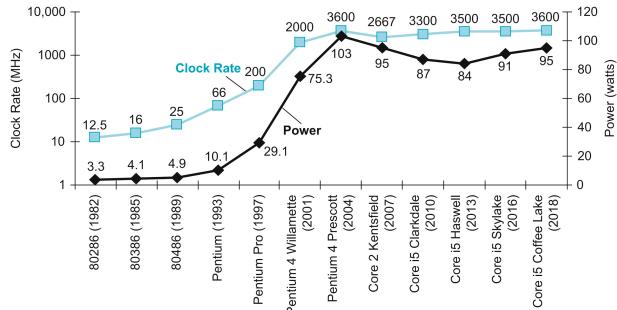
- Performance depends on:
 - algorithm**: affects IC, possibly CPI
 - programming language**: affects IC, CPI
 - compiler**: affects IC, CPI
 - ISA**: affects IC, CPI, T_c

CPU time = instruction count \times CPI \times clock cycle time

The Power Wall

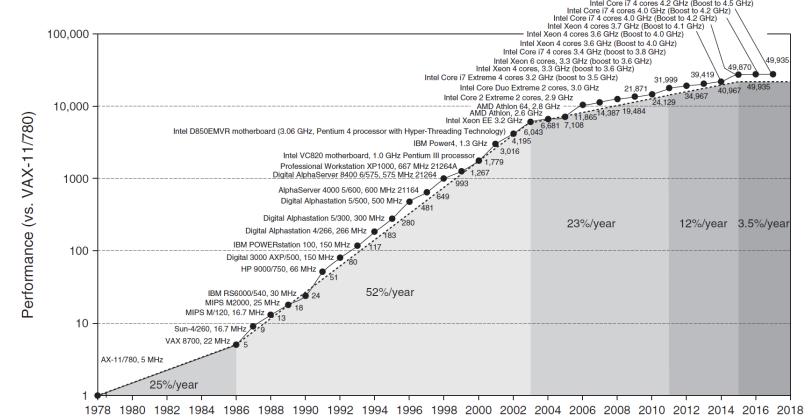
Power trends

- Clock rate and power are correlated
 - $\text{Power} \propto 1/2 \times \text{capacitive load} \times \text{voltage}^2 \times \text{frequency}$
- We have run into a practical limit for cooling



Uniprocessor performance (response time)

- Constrained by power, instruction-level parallelism, and 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

The multicore effect

processors	exec. time/ processor	time w/overhead	speedup	actual speedup/ideal speedup
1	100			
2	50	54	$100/54 = 1.85$	$1.85/2 = .93$
4	25	29	$100/29 = 3.44$	$3.44/4 = 0.86$
8	12.5	16.5	$100/16.5 = 6.06$	$6.06/8 = 0.75$
16	6.25	10.25	$100/10.25 = 9.76$	$9.76/16 = 0.61$

Benchmarking Processors

<https://www.spec.org/>



Standard Performance Evaluation Corporation

Home Benchmarks Tools Results Contact Blog Site Map Search Help

Benchmarks

- Cloud
- CPU
- Graphics/Workstations
- High Performance Computing
- Java Client/Server
- Machine Learning
- Storage
- Power
- Virtualization
- Results Search

■ Submitting Results

- Cloud/CPUs/Java/Power
- Storage/Virtualization
- ACCEL/HPC/CPU/OMP
- SPECap/SPECviewperf
- SPECworkstation

The Standard Performance Evaluation Corporation (SPEC) is a non-profit corporation formed to establish, maintain and endorse standardized benchmarks and tools to evaluate performance and energy efficiency for the newest generation of computing systems. SPEC develops benchmark suites and also reviews and publishes submitted results from our member organizations and other benchmark licensees.

What's New:

December 10, 2021 - The IEEE Micro Special Issue on the Microprocessor at 50 includes an article by SPEC's John Henning, "How Many VAXes Fit in the Palms of Your Hands?" The article, and a [companion video](#), use the original SPECmark benchmark to show relative performance on a modern system vs. a system from the 1970s.

December 9, 2021: SPECgpc has released the [SPECviewperf 2020 v3.0 benchmark](#). The worldwide standard for measuring graphics performance based on professional applications, the SPECviewperf 2020 benchmark measures the 3D graphics performance of systems running under the OpenGL and Direct X application programming interfaces.

Summarizing performance

- Summaries as **geometric mean** of performance ratios (speedup)
 - not affected by the baseline performance

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

- Summaries as **arithmetic mean** of execution times

$$\frac{1}{n} \sum_{i=1}^n \text{execution time}_i$$

SPECpower

Target Load %	Performance (ssj_ops)	Average Power (watts)
100%	4,864,136	347
90%	4,389,196	312
80%	3,905,724	278
70%	3,418,737	241
60%	2,925,811	212
50%	2,439,017	183
40%	1,951,394	160
30%	1,461,411	141
20%	974,045	128
10%	485,973	115
0%	0	48
Overall Sum	26,815,444	2,165
$\sum \text{ssj_ops} / \sum \text{power} =$		12,385

SPECpower_ssj2008 running on a dual socket 2.2 GHz Intel Xeon Platinum 8276L with 192 GiB of DRAM and one 80 GB SSD disk.

SPEC CPU2017 Results

These results have been submitted to SPEC; see [the disclaimer](#) before studying any results.

Available Results

Browse

The following are sets of available results since the announcement of the benchmark in June 2017.

- [All CPU2017 Results](#)

Results from all publication quarters, broken out by reported metric:

- Speed:
[SPECspeed 2017 Integer, SPECspeed 2017 Floating Point]
- Throughput:
[SPECrate 2017 Integer, SPECrate 2017 Floating Point]

<https://www.spec.org/cpu2017/results/>

Pitfalls and Myths

Amdahl's law (pitfall)

› Amdahl's Law

- "the **overall performance** improvement gained by **optimizing a single part** of a system is **limited by the fraction of time** that the improved part is actually used"
- improving an aspect of a computer does not translate automatically into a proportional improvement in overall performance

› Architecture design is bottleneck-driven

- **make the common case fast**
- do not waste resources on a component that has little impact on overall performance/power

Which enhancement is better?

Two independent parts A B

Original process



Make B 5x faster



Make A 2x faster



https://en.wikipedia.org/wiki/Amdahl%27s_law

Amdahl's law

› Assume:

- S: overall system speedup, F: fraction enhanced, I: factor of improvement

$$S = \frac{1}{(1 - F) + \frac{F}{I}}$$

non-enhanced enhanced

› Example

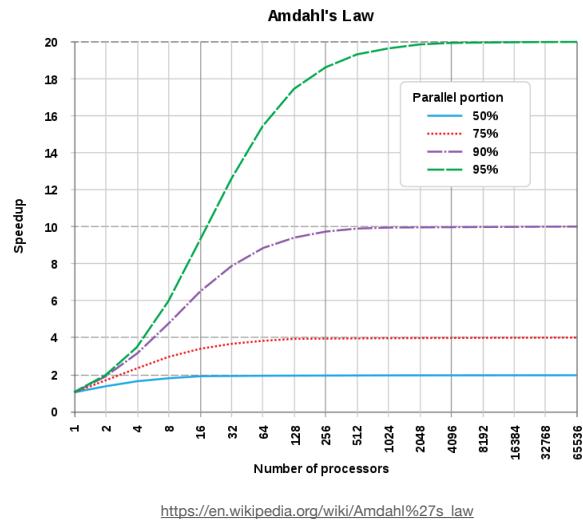
- 40% of instructions in program A are multiplications, and a new processor runs multiplications faster by a factor of 3.4. What is the overall system speedup?

Amdahl's law (example)

- A program spends 40% of time in the CPU and 60% doing I/O
 - a new processor **10x faster** results in a reduction in total execution time (what is the speedup?)

- Amdahl's Law states that the **maximum speedup** will be 1.66)

Amdahl's law and parallel computing



https://en.wikipedia.org/wiki/Amdahl%27s_law

Low power at idle (fallacy)

- Look back at Xeon's SPECpower benchmark
 - at 100% load: 347W
 - at 50% load: 183W (52.7%)
 - at 10% load: 115W (33.1%)
- Google data center
 - mostly operates at 10% – 50% load
 - at 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

Datacenter tour



https://www.youtube.com/watch?v=80aK2_iwMOs

MIPS as a Performance Metric (pitfall)

- **MIPS:** Millions of Instructions Per Second
 - doesn't account for differences in ISAs between computers, nor differences in complexity between instructions

$$\begin{aligned} \text{MIPS} &= \frac{\text{instruction count}}{\text{execution time} \times 10^6} \\ &= \frac{\text{instruction count}}{\frac{\text{instruction count} \times \text{CPI}}{\text{clock rate}} \times 10^6} \\ &= \frac{\text{clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

CPI varies between programs on a given CPU