# Lecture 3: Performance

Sunday, January 14, 2018        7:21 AM

## Outline

- Latency vs. bandwidth
- Benchmarks
- Iron law
- Speedup
- Amdahl's law

## Latency vs bandwidth: Moving data around

Grace Hopper - Nanoseconds
SeHouMusic



https://youtu.be/JEpsKnWZrJ8

Nanosecond => latency →

Goodness of a system?

Transfer rate → bandwidth

Inst. per second (IPS) → throughput →

Accuracy            Security          Usability
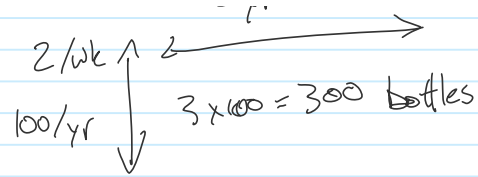
Storage space   ☆ Price          Stability

Energy              Size

Little's Law

2/wk ↑   3yr →

~ 300 bottles

## Little's Law

Resources = bandwidth * latency

2/wk ↑  ⌐ → 

100/yr ↓  3 × 100 = 300 bottles

## Measure a system

Benchmarking
→ tests → fast    Geekbench
→ standard        SPEC →
→ "score"         Kraken
→ realistic

  ↳ predict real-world performance/other metric

☆ take results with a grain of salt

### Problems

Rarely representative of real world
Not full applications

## Iron Law of performance

Simple model

$$Time = Inst \cdot \frac{cycles}{inst.} \cdot \frac{s}{cycle}$$

Inst · CPI · Cycletime
                        ↳ $\frac{1}{freq.}$

$1 \cdot 10^6 \cdot 5 \cdot \frac{1}{500MHz}$          Hz = ¹/s

$1 \cdot 10^6 \cdot 5 \cdot 2 \cdot 10^{-9} s = 10 \cdot 10^{-3} = 10ms$

Application: 1 million inst

System: 500 MHz

5 cycles per instruction (CPI)

Separate out different design considerations

$$Time = Inst * \boxed{CPI} * Cycle\ time \longrightarrow \begin{array}{l} reduce\ cycle\ time \Rightarrow technology \\ incr.\ freq. \end{array}$$

  ↳ Change your alg.          → Reduced Inst set (RISC)
  ↳ change ISA (x86, ARM, MIPS)
                ↳ Complex Inst set (CISC)          ld $rI
                                                    add $rI, $rZ
                add ($rax), $rbx → Micro-ops ↗
  ↳ change compiler

      to reduce cycles per inst → improve microarchitecture
  → pipelining  or  diff. CPU design

      Today's systems  CPI → ¹/₆   6 inst/cycle
                → ¹/₂

Latency → how to compare system A      system B
                              AMD        Intel

Latency → how to compare system A    system B
                              AMD         Intel
    gcc
                              2s           1s
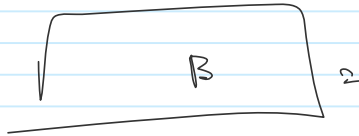
Intel is twice as fast for gcc compiling program X
    takes 1 less second
    half as long
    AMD is 50% slower

Speedup → it's the ratio that gives a big number
         → $\frac{\text{time on old system}}{\text{time on new system}}$      $\frac{2s}{1s}$ = Intel has a 2x speedup

| Program | A | B | Baseline | |
|---|---|---|---|---|
| 30% loads | 1x | 1x | 1x | |
| 20% stores | 1x | 1x | 1x | Majority is ALU |
| → 10% divides | 5x | 1x | 1x | |
| 40% ALU | 1x | 1.5x | 1x | |

**Make the common case fast**
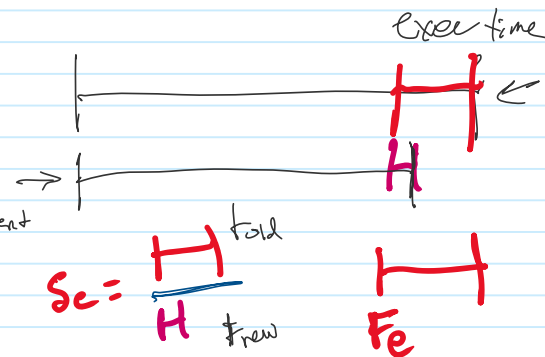
System A speedup = $\frac{T_{old}}{T_{new}}$ = $\frac{1}{.3 + .2 + .1 \cdot \frac{1}{5} + .4}$ = $\frac{1}{.92}$ = 1.09x

                                                = $\frac{1}{.3 + .2 + .1 + \frac{.4}{1.5}}$ = $\frac{1}{.6 + \frac{.8}{3}}$ = 1.16x


B          2

$T_{new} = T_{old} * \left[ (1 - F_e) + \frac{F_e}{S_e} \right]$
                      ↓                    ↳ speedup of enhancement
              fraction enhanced

exec time

$S_e = \frac{T_{old}}{T_{new}}$

Amdahl's Law

Speedup of app. on new system

Speedup = $\frac{T_{old}}{T_{new}}$ = $\frac{1}{(1 - F_e) + \frac{F_e}{S_E}}$       $F_e = 0.5$

                                                    Max speedup?

$\frac{1}{1 - .5 + \frac{.5}{\infty}}$ = $\frac{1}{.5}$ = 2