

INTRODUCTION / PERFORMANCE

Computer architecture: Θ the boundary of hardware + software

- parallelism → history
- measuring & evaluation → compiler design
- operating systems → programming languages
- instruction set architecture

DEFINITION: COMPUTER ARCHITECTURE

→ the science and art of selecting and interconnecting hardware components to create computers that meet functional, performance, and cost goals

Performance

- Latency - time to finish a fixed task (execution time)
- Throughput - # tasks per unit time (bandwidth)
 - ↳ exploit parallelism for throughput ← not latency

Speedup

Speedup of A over B

- $X = \frac{\text{latency}(B)}{\text{latency}(A)}$ ← faster in denominator)
- $X = \frac{\text{throughput}(A)}{\text{throughput}(B)}$ ← slower in denom.

→ example:
↳ car capacity: 5 ppl speed: 60 mi/hr
↳ bus capacity: 60 ppl speed: 20 mi/hr

↳ latency: car is 3x (200%) faster than bus
↳ throughput: bus is 4x(300%) faster than car

Comparing performance

→ Given: program A ← 200 cycles
program B ← 350 cycles

→ speedup of A over B?

$$\hookrightarrow \text{speedup} = 350/200 = 1.75 (= 75\%)$$

→ If program C is 50% faster than A, how many cycles does C run for? $200/x = 1.5 \Rightarrow x = \frac{200}{1.5} = 133 \text{ cycles}$

Speedup of A over B

$$x = \text{latency}(B) / \text{latency}(A)$$

$$x = \text{throughput}(A) / \text{throughput}(B)$$

What if $x < 1$?

→ A is slower than B

Percent increase and decrease are not the same.

$$\rightarrow \% \text{ increase of cycles} = \left(\frac{350-200}{200} \right) \cdot 100 = 75\%$$

$$\rightarrow \% \text{ decrease of cycles} = \left(\frac{350-200}{350} \right) \cdot 100 = 42.3\%$$

AVERAGING PERFORMANCE

Arithmetic: $(1/N) \cdot \sum_{P=1..N} \text{latency}(P)$ ← for units proportional to time (e.g. latency)

Harmonic: $\frac{N}{\sum_{P=1..N} \frac{1}{\text{throughput}(P)}}$ ← for units inversely proportional to time (e.g. throughput)

Geometric: $\sqrt[N]{\prod_{P=1..N}}$ ← for unitless quantities (e.g. speedup ratios)

You can add latencies, not throughputs

$$\rightarrow \text{latency}(p_1+p_2, A) = \text{latency}(p_1, A) + \text{latency}(p_2, A)$$

$$\rightarrow \text{throughput}(p_1+p_2, A) \neq \text{throughput}(p_1, A) + \text{throughput}(p_2, A)$$

Example: you drive $\frac{1 \text{ mi}}{30 \text{ mph}}$ @ 30 mph $\frac{1 \text{ mi}}{90 \text{ mph}}$ @ 90 mph average speed?

unit - inversely proportional to time \Rightarrow harmonic mean

$$\text{Mean} = \frac{N}{\sum \frac{1}{30} + \frac{1}{90}} = \frac{2}{\frac{90}{30+90}} = \frac{2 \cdot 90}{4} = 45 \text{ mph} \neq \frac{30+90}{2}$$

CPU Performance Equation

$$\text{latency} = \frac{\text{seconds}}{\text{program}} = \frac{\text{instr}}{\text{program}}$$

\rightarrow instruction count

\rightarrow impacted by program, compiler, ISA

\rightarrow CPI impacted by program, compiler, ISA, micro arch

$$\frac{\text{cycles}}{\text{instr}} \cdot \frac{\text{seconds}}{\text{cycle}}$$

\hookrightarrow clock period (Hz)

\rightarrow impacted by micro-arch, technology

\rightarrow minimize all 3 for low latency (better performance)

$$\text{cycles per instruction} = \text{CPI}$$

$$\rightarrow \text{CPI} = \frac{\text{cycles}}{\text{instr}}$$

$$\hookrightarrow \text{IPC} = \frac{1}{\text{CPI}}$$

\hookrightarrow different instr. have diff. cycle costs

other performance metrics

$$\rightarrow \text{MIPS} \leftarrow \frac{\text{millions instr.}}{\text{sec}}$$

$$\rightarrow \text{FLOPS} \leftarrow \# \text{ floating pt instr/sec}$$

example

a program executes equal # of 3 diff. instr.

Inst	Cycles/Instr
int	1
Mem	2
FP	3

FIND CPI

use arithmetic mean

$$\frac{1+2+3}{3} = 2 \text{ cycles/instr}$$

CPI example assume a processor w/instruction frequencies + cost

instr	freq	cost (cycles)
ALU	50%	1
load	20%	5
store	10%	1
branch	20%	2

which change would improve performance more?

- (A) branch prediction which reduce branch cost to 1 cycle
- (B) faster data memory to reduce load cost to 3 cycles

base: $1 \cdot 0.5 + 5 \cdot 0.2 + 1 \cdot 0.1 + 2 \cdot 0.2 = 2 \text{ CPI}$

B
3 · 0.2
 $A 1 \cdot 0.2 = 1.8 \text{ CPI}$

$= 1.6 \text{ CPI} \leftarrow B \text{ is fastest}$

Example

instr	freq	CPI-M1	CPI-M2	average MIPS for M1 + M2? clock: 1GHz
ALU	50%	1	2	
branch	15%	2	1	
load	20%	2	1	
store	15%	1	1	

① find IPC = $\frac{\text{instr}}{\text{cycle}}$

② MIPS = $\frac{\text{instr}}{\text{cycle}} \cdot \frac{\text{cycle}}{\text{sec}} \cdot 10^6$

$0.5 + 0.15(2) + 0.2(2) + 0.15(1) = 1.35 \text{ M1 } \frac{\text{cycles}}{\text{instr}}$

$0.5(2) + 0.15 + 0.2 + 0.15 = 1.5 \text{ M2 } \frac{\text{cycles}}{\text{instr}}$

MIPS M1 = $\frac{1}{1.35} \frac{\text{instr}}{\text{cycle}} \cdot \frac{\text{cycle}}{10^{-9} \text{ sec}} = 741 \text{ MIPS}$

MIPS M2 = $\frac{1}{1.5} \frac{\text{instr}}{\text{cycle}} \cdot \frac{\text{cycle}}{10^{-9} \text{ sec}} = 667 \text{ MIPS}$

Peak MIPS, $CPI = 1 = 1PC$

$$1 \cdot \frac{1}{1E-9} \cdot \frac{1}{E6} = \frac{1}{1E-3} = 1E3 = 1000 \text{ MIPS}$$

Amdahl's Law

$$\frac{1}{(1-P) + \frac{P}{S}}$$

P=proportion of running time affected by optimization

S=speedup

example: speedup 25% of a program's execution by 2x?

$$P=0.25, S=2 \quad \frac{1}{(0.75) + \frac{0.25}{2}} = 1.14X$$

example: $P=25\%$, $S=\infty$

$$\frac{1}{0.75 + \frac{0.25}{\infty}} = 1.33X$$

Amdahl's Law for Parallelization

$$\frac{1}{(1-P) + \frac{P}{N}}$$

How much will parallelization improve performance?

P=portion of parallel code

N=threads

example: What's the Max speedup for a program that's 10% serial?

$$\frac{1}{(1-0.9) + \frac{0.9}{\infty}} = \frac{1}{0.1} = 10X \text{ speedup}$$

$\Rightarrow P=0.9$

example: max speedup, 1% serial?

$$\frac{1}{1-0.99} = \frac{1}{0.01} = 100X \text{ speedup}$$

another example

program has:

instr	latency	freq
mem	4	40%
add	2	50%
mult	16	10%

If you could pick one type of instr to make twice as fast, which instr type would you pick?

$$\begin{array}{l} \text{Mem : } 4 \cdot 0.4 = 1.6 \\ \text{add : } 2 \cdot 0.5 = 1 \\ \text{Mult : } 16 \cdot 0.1 = 1.6 \end{array}$$

choose Mem or Mult

Performance must be associated w/a workload ← set of tasks

benchmarks: standard workloads, representative of real programs

micro-benchmarks: non-standard workloads, tiny programs used
to evaluate specific aspects of performance

(e.g. binary tree
search, towers of
hanoi, 8-queens, etc)
