

2.2 Performance metrics

$$\text{Performance}_x = \frac{1}{\text{Execution time}_x}$$

- Real programs (run)

- Hot
toy programs
synthetic benchmarks
kernels

- Transitivity $A > B \quad \& \quad B > C \Rightarrow A > C$

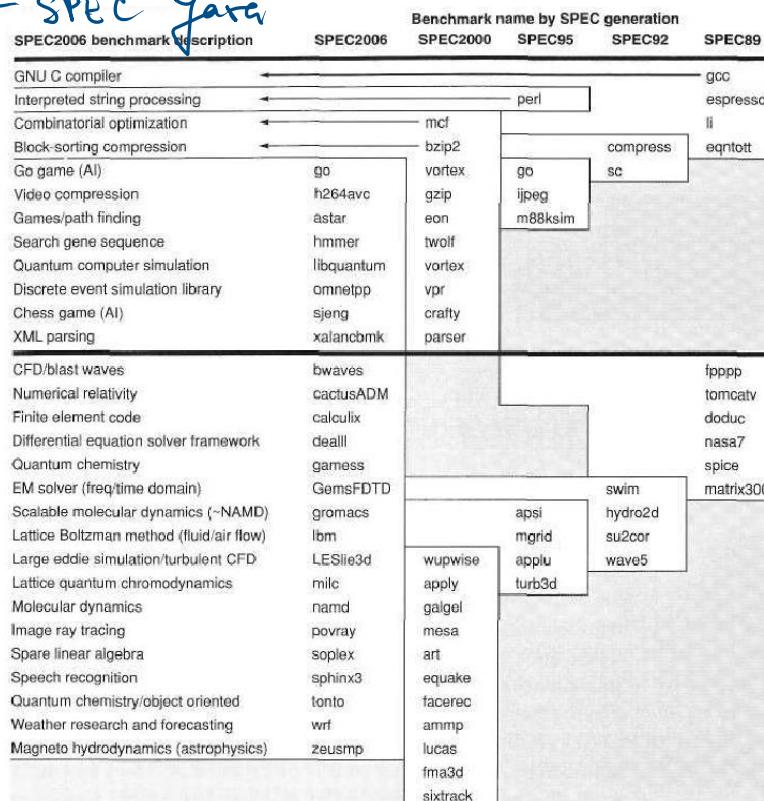
| Dec | Binary | 1 kB | one kilobyte |
|-----------|----------|-------------|--------------|
| 10^3 | 2^{10} | <u>1 kB</u> | |
| 10^6 | 2^{20} | 1 Mi | |
| 10^9 | 2^{30} | 1 GiB | |
| 10^{12} | 2^{40} | 1 TiB | |
| 10^{15} | 2^{50} | 1 PiB | peta |
| 10^{18} | 2^{60} | 1 EiB | exa |
| 10^{21} | 2^{70} | 1 ZiB | zetta |
| 10^{24} | 2^{80} | 1 YiB | yotta |

Benchmark suites

SPEC - Standard Performance Evaluation Corporation www.spec.org

1989, 92, 95, 2000, 2006, 2017

- CPU SPEC
- SPEC graphics
- SPEC Java



Server benchmarks

- Throughput-centric metric
- Correct CPU time ratio aware
- UX CPU SPEC suites

SPEC

2.2.1 Reporting performance results

SPEC Ratio

$$\frac{\text{SPEC Ratio}_A}{\text{SPEC Ratio}_B} = \frac{\frac{\text{Execution time reference}}{\text{Execution time}_A}}{\frac{\text{Execution time reference}}{\text{Execution time}_B}} = \frac{\text{Execution time}_B}{\text{Execution time}_A} \cdot \frac{\text{Performance}_A}{\text{Performance}_B}$$

Geometric mean = $\sqrt[m]{\prod_{i=1}^m \text{SPEC ratio}_i}$

program

- ① The geometric means of ratios = the ratio of geometric means
- ② The ratio of geometric means = the geometric mean of performance ratios

Geometric mean_A

Geometric mean_B

$$\sqrt[m]{\prod_{i=1}^m \text{SPEC ratio}_A i}$$

(m=29)

$$\sqrt[m]{\prod_{i=1}^m \text{SPEC ratio}_B i}$$

$$\frac{\sqrt[m]{\prod_{i=1}^m \text{SPEC ratio}_A i}}{\sqrt[m]{\prod_{i=1}^m \text{SPEC ratio}_B i}}$$

$$= \sqrt[m]{\frac{\frac{\text{Execution time reference}_i}{\text{Execution time}_A i}}{\frac{\text{Execution time reference}_i}{\text{Execution time}_B i}}}$$

$$\sqrt[m]{\frac{\text{Execution time}_B i}{\text{Execution time}_A i}}$$

$$= \sqrt[m]{\frac{\text{Performance}_A i}{\text{Performance}_B i}}$$

$$\text{std dev} = \sqrt{\sum_{i=1}^m (\text{Sample}_i - \text{Mean})^2}$$

↑
SPEC Ratio_i

Geometric standard deviation

$$\text{Geometric mean} = \exp \left(\frac{1}{n} \times \sum_{i=1}^n \ln(\text{Sample}_i) \right)$$

$$\text{gstdev} = \exp \left(\sqrt{\frac{1}{n} \times \sum_{i=1}^n [\ln(\text{Sample}_i) - \ln(\text{Geometric mean})]^2} \right)$$

2.3 Quantitative principles of computer design

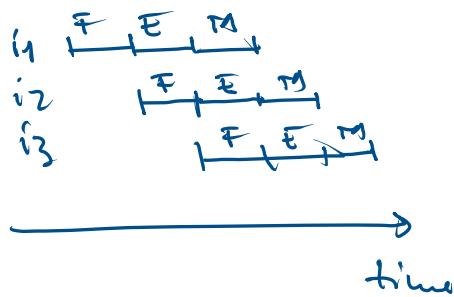
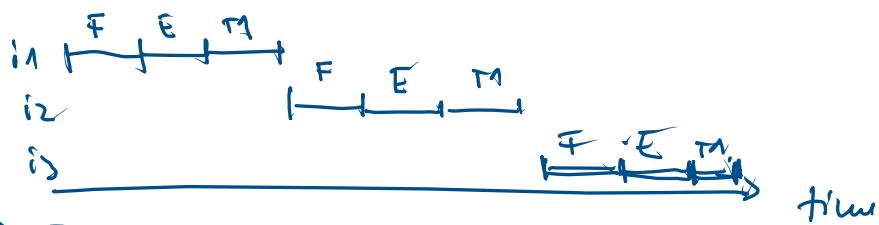
SLD - System-Level Design

① Take advantage of parallelism (when possible)

(CLA)



Instruction-level parallelism - pipelining



② Principle of locality

Obs Programs tend to reuse data and instructions they have used recently

Rule of thumb

A program spends 90% of the time on 10% of the code

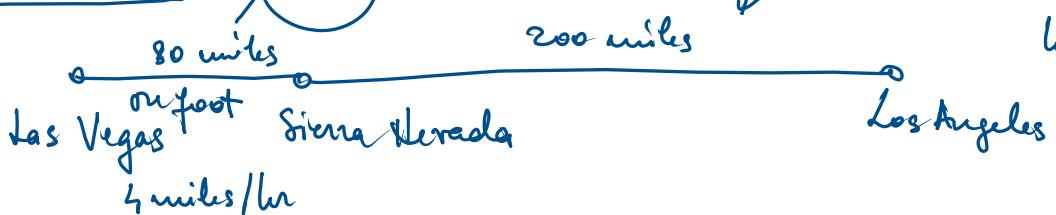
- Mostly applied on instructions

- Also works on data

- 2 types of locality
 temporal locality
 spatial locality

③ Focus on the common case

Amdahl's law



| Vehicle for the 2nd part | Time for the 2nd part | Speedup in the direct | Time for the entire trip | Speedup entire trip |
|--------------------------|-----------------------|-----------------------|--------------------------|---------------------|
| Feet | 50 | 1 | 70 | 1 |
| Bike | 20 | 2.5 | 40 | 1.75 |
| Hyundai | 4 | 12.5 | 24 | 2.9 |
| Ferrari | 1.67 | 30 | 21.67 | 3.023 |
| Rocket car | 1 | 50 | 21 | 3.3 |

$$\text{Speedup entire system} = \frac{\text{Performance of the task using the speedup when possible}}{\text{Performance of the task without the speedup}}$$

$$= \frac{\text{Execution time of task without speedup}}{\text{Execution time of task using speedup when possible}}$$

$$= \frac{\text{Execution time of task without speedup}}{\text{Execution time of task using speedup when possible}}$$

$$= \frac{1}{(1 - \text{Fraction enhanced}) + \frac{\text{Fraction enhanced}}{\text{Speedup}}}$$

$$= \frac{1}{(1 - \text{Fraction enhanced}) + \frac{\frac{200}{280}}{\text{Speedup}}} \rightarrow \frac{200}{280} \rightarrow \frac{\text{Fraction enhanced}}{\text{Speedup}} \rightarrow \text{Speedup in the direct}$$

Limit of system enhancement

$$\lim_{\text{Speedup} \rightarrow \infty} \text{Speedup for the entire system} = \frac{1}{1 - \text{Fraction enhanced}}$$

2.4

Computer Performance equation

$CPU_{time} = \text{Clock cycles for a program} \times \text{Clock cycle time}$

$$\text{Clock cycle time} = \frac{1}{\text{Clock rate}}$$

↑ Hz

Clock cycles for a program = Instruction Count x Clock cycles per instruction

$$\text{CPU Time} = \underbrace{\text{IC}}_{10\%} \times \underbrace{\text{CPI}}_{10\%} \times \text{Clock cycle time} \Rightarrow \text{IC integer CPI average}$$

$$\text{Pulse time} = \frac{\cancel{\text{Instructions}}}{\text{Program}} \times \frac{\cancel{\text{Clock cycles}}}{\cancel{\text{Instruction}}} \times \frac{\cancel{\text{Seconds}}}{\cancel{\text{clock cycle}}} = \frac{\text{Seconds}}{\text{Program}}$$

RISC processors

CISc

$$OPJ = \frac{\sum_{i=1}^n JC_i \times CPJ_i}{JC}$$

Example 1

Measurements for a computer system

frequency of FP operations = 25%

Average CPI of FP operations = 4.0 clock cycles

Average ~~avg~~ of all other instructions = 1.33 clock cycles

Frequency, FPSQRT is 2%

OT of FFSQRT is 20 clock cycles

Assume 2 design alternatives:

A. decrease of CPI of $FPSQR\bar{T}$ to 2.

B. decrease the average CPI of all FP operations to 2.5

Which alternative is best?
