

Chapter 2: Performance

Performance

- **Purchasing perspective**
 - **given a collection of machines, which has the**
 - **best performance ?**
 - **least cost ?**
 - **best performance / cost ?**
- **Design perspective**
 - **faced with design options, which has the**
 - **best performance improvement ?**
 - **least cost ?**
 - **best performance / cost ?**
- **Both require**
 - **basis for comparison**
 - **metric for evaluation**
- **Our goal is to understand cost & performance implications of architectural choices**

Two notions of “performance”

Plane	DC to Paris	Speed	Passengers	Throughput (pmp)
Boeing 747	6.5 hours	610 mph	470	286,700
BAD/Sud Concorde	3 hours	1350 mph	132	178,200

Which has higher performance?

- **Time to do the task (Execution Time)**
 - execution time, response time, **latency**
- **Tasks per day, hour, week, sec, ns. .. (Performance)**
 - **throughput**, bandwidth

Response time and throughput often are in opposition

Definitions

- Performance is in units of things-per-second
 - bigger is better
- If we are primarily concerned with response time
 - $\text{performance}(x) = \frac{1}{\text{execution_time}(x)}$

"X is n times faster than Y" means

$$n = \frac{\text{Performance}(X)}{\text{Performance}(Y)}$$

Example

- Time of Concorde vs. Boeing 747?
 - Concorde is $1350 \text{ mph} / 610 \text{ mph} = 2.2$ **times faster**
 $= 6.5 \text{ hours} / 3 \text{ hours}$
- Throughput of Concorde vs. Boeing 747 ?
 - Concorde is $178,200 \text{ pmph} / 286,700 \text{ pmph} = 0.62$ **“times faster”**
 - Boeing is $286,700 \text{ pmph} / 178,200 \text{ pmph} = 1.6$ **“times faster”**
- Boeing is 1.6 times (“60%”)faster in terms of throughput
- Concorde is 2.2 times (“120%”) faster in terms of flying time

We will focus primarily on execution time for a single job

Basis of Evaluation

Pros

Cons

- representative

Actual Target Workload

- very specific
- non-portable
- difficult to run, or measure
- hard to identify cause

- portable
- widely used
- improvements useful in reality

Full Application Benchmarks

- less representative

- easy to run, early in design cycle

Small “Kernel” Benchmarks

- easy to “fool”

- identify peak capability and potential bottlenecks

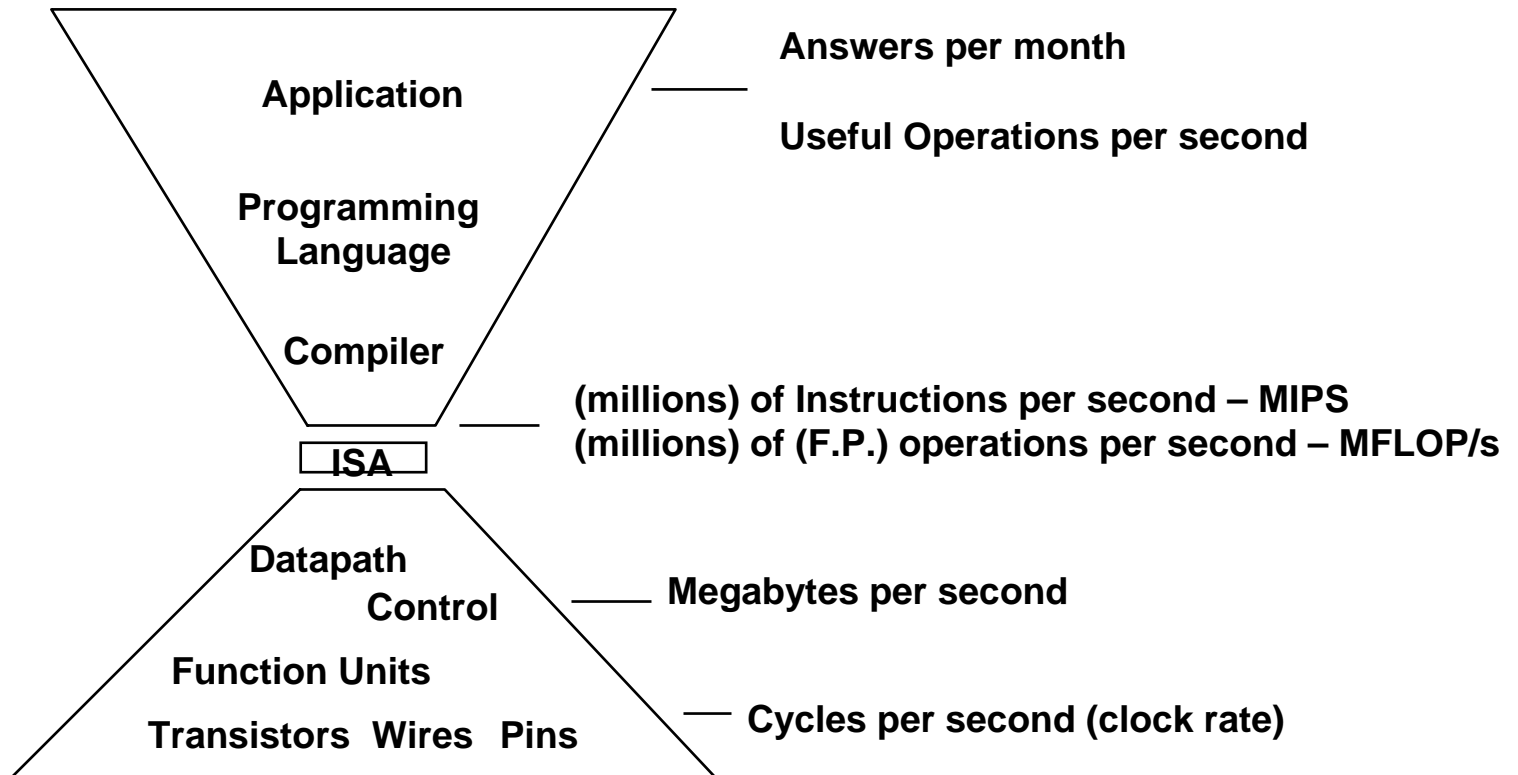
Microbenchmarks

- “peak” may be a long way from application performance

SPEC95

- **Eighteen application benchmarks (with inputs) reflecting a technical computing workload**
- **Eight integer**
 - **go, m88ksim, gcc, compress, li, jpeg, perl, vortex**
- **Ten floating-point intensive**
 - **tomcatv, swim, su2cor, hydro2d, mgrid, applu, turb3d, apsi, fppp, wave5**
- **Must run with standard compiler flags**
 - **eliminate special undocumented incantations that may not even generate working code for real programs**

Metrics of performance



Each metric has a place and a purpose, and each can be misused

Aspects of CPU Performance

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

	instr. count	CPI	clock rate
Program			
Compiler			
Instr. Set Arch.			
Organization			
Technology			

Aspects of CPU Performance

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

	instr count	CPI	clock rate
Program	X		
Compiler	X	X	
Instr. Set	X	X	
Organization		X	X
Technology			X

CPI

“Average cycles per instruction”

$$\begin{aligned} \text{CPI} &= (\text{CPU Time} * \text{Clock Rate}) / \text{Instruction Count} \\ &= \text{Clock Cycles} / \text{Instruction Count} \end{aligned}$$

n

$$\text{CPU time} = \text{ClockCycleTime} * \sum_{i=1}^n \text{CPI}_i * I_i$$

n

$$\text{CPI} = \sum_{i=1}^n \text{CPI}_i * F_i \quad \text{where } F_i = \frac{I_i}{\text{Instruction Count}}$$

"instruction frequency"

Invest Resources where time is Spent!

Example (RISC processor)

Base Machine (Reg / Reg)

Op	Freq	Cycles	CPI(i)	% Time
ALU	50%	1	.5	23%
Load	20%	5	1.0	45%
Store	10%	3	.3	14%
Branch	20%	2	.4	18%
			<hr/>	
			2.2	

Typical Mix

How much faster would the machine be if a better data cache reduced the average load time to 2 cycles?

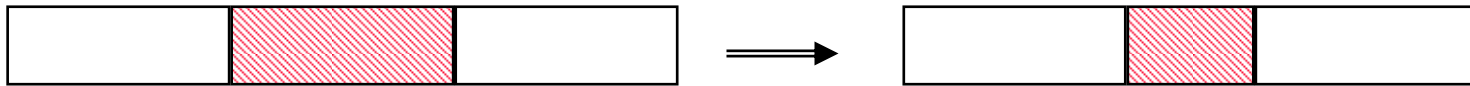
How does this compare with using branch prediction to shave a cycle off the branch time?

What if two ALU instructions could be executed at once?

Amdahl's Law

Speedup due to enhancement E:

$$\text{Speedup}(E) = \frac{\text{ExTime w/o E}}{\text{ExTime w/ E}} = \frac{\text{Performance w/ E}}{\text{Performance w/o E}}$$



Suppose that enhancement E accelerates a fraction F of the task by a factor S and the remainder of the task is unaffected then,

$$\text{ExTime}(\text{with E}) = ((1-F) + F/S) \times \text{ExTime}(\text{without E})$$

$$\text{Speedup}(\text{with E}) = \frac{1}{(1-F) + F/S}$$

Summary: Salient features of MIPS I

- **32-bit fixed format inst** (3 formats)
- **32 32-bit GPR** (R0 contains zero) and 32 FP registers (and HI LO)
 - partitioned by software convention
- **3-address, reg-reg arithmetic instr.**
- **Single address mode for load/store:** base+displacement
 - no indirection, scaled
- 16-bit immediate plus LUI**
- **Simple branch conditions**
 - compare against zero or two registers for =, °
 - no integer condition codes
- **Delayed branch**
 - execute instruction after the branch (or jump) even if the branch is taken (Compiler can fill a delayed branch with useful work about 50% of the time)

Summary: Instruction set design (**MIPS**)

- Use general purpose registers with a load-store architecture: **YES**
- Provide at least 16 general purpose registers plus separate floating-point registers: **31 GPR & 32 FPR**
- Support basic addressing modes: displacement (with an address offset size of 12 to 16 bits), immediate (size 8 to 16 bits), and register deferred; : **YES: 16 bits for immediate, displacement (disp=0 => register deferred)**
- All addressing modes apply to all data transfer instructions : **YES**
- Use fixed instruction encoding if interested in performance and use variable instruction encoding if interested in code size : **Fixed**
- Support these data sizes and types: 8-bit, 16-bit, 32-bit integers and 32-bit and 64-bit IEEE 754 floating point numbers: **YES**
- Support these simple instructions, since they will dominate the number of instructions executed: load, store, add, subtract, move register-register, and, shift, compare equal, compare not equal, branch (with a PC-relative address at least 8-bits long), jump, call, and return: **YES, 16b**
- Aim for a minimalist instruction set: **YES**

Summary: Evaluating Instruction Sets?

Design-time metrics:

- Can it be implemented, in how long, at what cost?
- Can it be programmed? Ease of compilation?

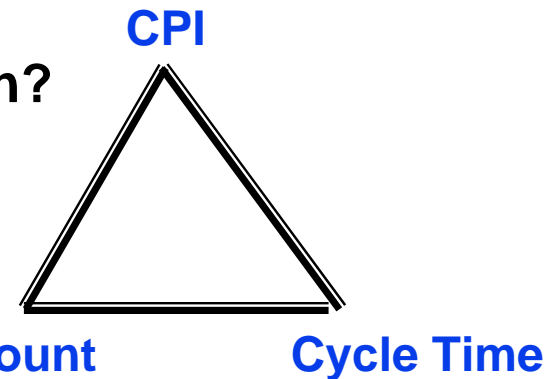
Static Metrics:

- How many bytes does the program occupy in memory?

Dynamic Metrics:

- How many instructions are executed?
- How many bytes does the processor fetch to execute the program?
- How many clocks are required per instruction?
- How "lean" a clock is practical?

Best Metric: Time to execute the program!



NOTE: this depends on instructions set, processor organization, and compilation techniques.