02-1 Components of CPU Performance and Performance Equation

02-1

02-2

Why is my computer fast (or slow)?

Would it help to improve ____?

CPU performance equation is one way to start answering these questions.

Clock Frequency (φ)

Determined by technology and influenced by organization.

CPU Performance Decomposed into Three Components:

 \bullet Clocks per Instruction (CPI)

Determined by ISA, microarchitecture, compiler, and program.

• Instruction Count (IC)

Determined by program, compiler, and ISA.

These combined to form CPU Performance Equation

$$t_{\rm T} = \frac{1}{\phi} \times {\rm CPI} \times {\rm IC}$$

where $t_{\rm T}$ denotes the execution time.

02-1

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02

02-2

02-4

02-1

02-3

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02.

02-2

02-3

CPU Performance: Simple System

Execution in program order . . .

... one at a time.

Time/cycles: 0 1 2 3 4 5 6 7 8 9 10 11 1,999,996

Time/mms: 0 80 160 39,999,920

| Instr. 1 | Instr. 2 | Instr. 3 | 000 | Instr. 500,000

IC = 500,000;

 $\phi = 50 \, \text{kHz};$

CPI = 4.

Execution time: IC \times CPI. \times clock period.

Here (and only here) CPI is number of cycles for each instruction.

E

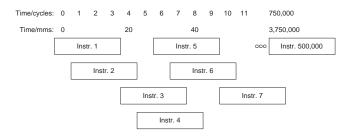
Execution: Pipelined, In Order

02-4

To Run Faster: Overlap Instructions (Pipelined Execution)

Result must be same as one-at-a-time execution . . .

... not too difficult to achieve.



IC = 500,000;

 $\phi = 200 \,\mathrm{kHz};$

 $CPI = \frac{750000}{500000} = 1.5.$

Execution time at best: $IC \times clock period ...$

- ... assuming 1 cycle to start each instruction and ...
- ... instruction can start each cycle. (Slower in illustration.)

02-3

0 Lastum Transmananay Formattad 12:04 24 January 2002 from Isli

02-3

02-4

FE 4790 Lostino Transponency Formatted 12:04 24 January 2002 from Isling

02-5

Execution: Pipelined, Ideal Out of Order

02-6

02-5

Execution: Pipelined, Ideal Out of Order, Superscalar

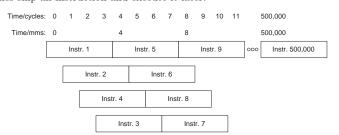
Instr. 500,000

To Run Fastest¹: Overlap, Out-of-Order, Start n per Tick (n-Way Superscalar).

Instr. 10 Instr. 12

To Run Even Faster: Overlap Instructions and Start Out of Order

Sometimes skip an instruction and execute it later.



IC = 500,000;

CPI = 1.

Execution time at best: IC \times clock period \dots

- \dots assuming 1 cycle to start each instruction \dots
- ... instruction can start each cycle.

 $\phi = 200 \,\mathrm{kHz};$

IC = 500,000;

 $\phi = 500 \,\mathrm{MHz};$

Requires about n times as much hardware. (Below, n = 2.)

Time/cycles: 0 1 2 3 4 5 6 7 8 9 10

 $CPI = \frac{1}{2}$.

Instr. 13

Instr 18

Execution time at best: $\frac{1}{n} \times IC \times clock$ period ...

 \dots assuming 1 cycle to start each instruction in struction can start each cycle.

 $^{1}\,$ Using a conventional serial instruction set architecture.

02-5

EE 4720 Lecture Transparency. Formatted $\,$ 13:04, $\,$ 24 January 2003 from lsli02 $\,$

02-6

02-5

02-7

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02.

02-6

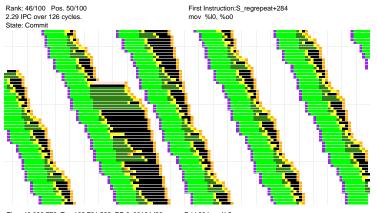
02-6

02-7 Execution: Pipelined, Out of Order, Superscalar

Data from a real program, perl. CPI is 0.44.

Processor can start four instructions per cycle.

Colors show the steps in processing an instruction, yellow is execution.



Time 49,098,778 Tag 102,731,592 PC 0x00164d98 Grid 20 insn X 5 cyc

02-8

Component of CPU Performance: Instruction Count

02-8

Given a program there are two ways instructions could be tallied:

Static Instruction Count:

The number of instructions making up the program.

Dynamic Instruction Count (IC):

The number of instructions executed in a run of the program.

For estimating performance, dynamic instruction count is used.

02-9

Instruction Counts

02-10

02-9

Component of CPU Performance: Clock Frequency

02-10

Example, assembler program that computes $a = \sum_{i=0}^{9} i$.

Written in Simplescalar assembler.

```
IC
                           ! r0 is always zero.
1
       move
               r5, r0
               r3. r0
1
       move
      L23:
                           ! Branch label.
10
       addu
               r5, r5, r3 ! Add unsigned.
               r3, r3, 1
10
       addu
10
               r2, r3, 10 ! r2 = r3 < 10
       slt
10
       bne
               r2, r0, L23 ! Branch to L23 if r2 not equal 0.
```

Static count: 6 (number of instructions).

Dynamic count: 42.

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02

Clock frequency determined by critical path.

Critical Path:

Typical Clock Cycle

Logic doing most time consuming work (in a cycle).

• While clock is high work proceeds.

CPUs implemented using synchronous clocked logic.

• When clock switches from low to high work starts.

If clock frequency is too high work will not be completed and so system will not perform properly.

• When clock goes from high to low work should be complete.

For high clock frequencies, keep critical paths short.

02-9

02-11

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from Isli02.

02-10

02-11

02-9

Component of CPU Performance: CPI

02-12

02-10

Review of CPU Performance Equation

02-12

Cycles (clocks) per Instruction (CPI)

Oversimplified definition: CPI:

Average number of cycles needed to execute an instruction.

Better definition: CPI:

Number of cycles to execute some code divided by number of instructions.

Difference:

Interested in rate at which instructions executed in program . . .

... not time time for any one instruction.

 $t_{\rm T} = \frac{1}{\phi} \times {\rm CPI} \times {\rm IC}$

where $t_{\rm T}$ denotes the execution time.

- Clock Frequency (φ)
 Determined by technology and influenced by organization.
- Clocks per Instruction (CPI)

 Determined by organization and instruction mix.
- Instruction Count (IC)

 Determined by program and ISA.

02-13 02-13 02-14 02-14 Interaction of Execution Time Components Example: Trading off Execution Time Components Tradeoffs between Clock Frequency, CPI, and Instruction Count Company X is considering two clock frequencies for its next processor, 500 MHz or 400 MHz. A 500 MHz implementation would execute instructions at 1.7 CPI, the 400 MHz part at 1.1 CPI. Which would be faster?. Increasing Clock Frequency . . . Find time to execute 1 instruction. ... reduces the work that can be done in a clock cycle and possibly limiting instruction overlap. 500 MHz execution time: $\frac{1}{500\times10^6}\times1.7\times1=3.4\,\mu\mathrm{s}$ 400 MHz execution time: $\frac{1}{400 \times 10^6} \times 1.1 \times 1 = 2.75 \,\mu\text{s}$. Reducing IC (by adding "powerful" instructions to ISA) may force implementors to increase CPI or lower clock frequency. The lower clock rate would nevertheless take less time. Perhaps because at 500 MHz too much work had to be split into multiple cycles. Balancing these is an important skill in computer design. Since the ISA is usually fixed, IC is less of a factor. 02-13 02-13 02-14 02-14 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02. 02-15 IC v. CPI Tradeoffs 02-15 02-16 IC v. CPI Tradeoffs, continued. 02-16 Assumption Case 1: Same ISA, different implementation. IC is based on output of a good compiler. Newer implementation may have lower CPI on existing code but even better performance attainable by recompiling ... Compiler is tuned for a particular implementation. ... which may increase CPI. Compiler writer selects instructions based on performance of implementation. Two Cases 1. Same ISA, different implementation. 2. Different ISA, (and of course) different implementation. 02-15 02-15 02-16 02-16 02-17 02-17 02-18 02-18 IC v. CPI Tradeoffs, continued. Consider two implementations: Case 2: Different ISA, (and of course) different implementation. Implementation A: add CPI 1 cycle, mul CPI 5 cycles. Major tradeoffs in complexity and speed. Implementation B: add CPI 1 cycle, mul CPI 2 cycles. Consider two implementations: ! Call original value of r1, x. Code computes 6x. Implementation A: CPI: load, 2; add and store, 1. ! Code For Implementation A Implementation B: CPI: add (doing load and store), 4. add r1, r1, r1 ! r1 = 2xadd r2, r1, r1 ! r2 = 4x! Code for implementation A. add r1, r1, r2 ! r1 = 6x load r1, [r2] ! Load r1 with data at address in r2. add r3, r1, r4 ! r3 = r1 + r4! Code For Implementation B. store [r2], r3 ! Store r3 at address in r2. mul r1, r1, 6 ! r1 = 6x. ! Code for implementation B. add [r2], r4, [r2] Implementation A: IC = 3, CPI = 1 (Computing CPI will be covered later.) Implementation B: IC = 1, CPI = 2. Execution time same. Implementation B is faster despite higher CPI. Implementation A: IC = 3, CPI = $\frac{4}{3}$. Code compiled for B will run slowly on A. Implementation B: IC = 1, CPI = 4. 02-17 02-17 02-18 02-18 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02. EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from Isli02. 02-19 Technological Change 02-19 02-20 02-20 Technological Change and Computer Designer Golden Handcuffs: The need to maintain compatibility in a successful product line. Technology determines "raw materials" for designer. Famously, Intel's IA-32. (Popularly referred to as 80x86.) Raw material: number of gates and their speed. The ISA is the handcuffs... ISA lifetime can be decades. ... and technological change brings the desire to move your arms. Raw materials greatly change over this time. So, design ISA for now and future. 02-19 02-19 02-20 02-20

02-21	02-21	02-22 Benchmarks 02-22
How technological advancement affects processor. Logic Speed, Clock Rate No changes to organization or ISA. Number of Transistors Available for Logic Changes to organization and possible changes to ISA. Memory Size Change ISA to use larger address space. Can use ISA having larger instruction codings. Memory Speed Compared to Processor Speed Include more sophisticated caching in organization.	02-21	Determine effect of design options. Benchmarks O2-22 Benchmark: Program used to evaluate performance. Uses Guide computer design. Guide purchasing decisions. Marketing tool. Guiding Computer Design Measure overall performance. Determine characteristics of programs. E.g., frequency of floating-point operations. Determine effect of design options.
02-21 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from Isli02.	02-21	02-22 EE 4720 Lecture Transparency. Formatted 13.04, 24 January 2003 from lsli02. 02-22
Optimal but unrealistic: The exact set of programs customer will run. Problem: computers used for different applications. Therefore, must model typical users' workload.	02-23	Options: Real Programs: Programs chosen using surveys, for example. + Measured performance improvements apply to customer Large programs hard to run on simulator. (Before system built.) Kernels: Use part of program responsible for most execution time. + Easier to study Not all program have small kernels. Toy Benchmarks: Program performs simplified version of common task. + Easier to study May not be realistic.
02-23 EE 4720 Lecture Transparency: Formatted 13:04, 24 January 2003 from Isli02.	02-23	02-24 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from Isli02.

Synthetic Benchmarks: Program "looks like" typical program, but does nothing useful. + Easier to study. - May not be realistic. Commonly Used Option Overall performance: real programs Test specific features: synthetic benchmarks.	02-25	Benchmark Suite: A named set of programs used to evaluate a system. Typically: Developed and managed by a publication or non-profit organization. E.g., Standard Performance Evaluation Corp., PC Magazine. Tests clearly delineated aspects of system. E.g., CPU, graphics, I/O, application. Specifies a set of programs and inputs for those programs. Specifies reporting requirements for results.
02-25 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from isli02. 02-27 What Suites Might Measure	02-25	02-26 EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lali02. 02-26 02-28
 Application Performance E.g., productivity (office) applications, database programs. Usually tests entire system. CPU and Memory Performance Ignores effect of I/O. Graphics Performance 	02-27	Respected measure of CPU performance. Managed by Standard Performance Evaluation Corporation, a non-profit organization funded by computer companies. Measures CPU and memory performance on integer and FP code. Uses common Unix programs such as perl, gcc, gzip. Requires that results on each program be reported. Programs compiled with publicly available compilers and libraries. Programs compiled with and without expert tuning.

02-29 02-30
SPEC CPU2000 Suites and Measures Other Examples
Suite of integer programs run to determine: (Fall 2001: This list is out of date.)

02-29

• SPECint_base2000, execution time of untuned code.

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from lsli02.

• SPECint_rate2000, throughput of tuned code.

• SPECint2000, execution time of tuned code.

 \bullet SPECint_rate_base2000, throughput of untuned code.

Suite of floating programs run to determine:

• SPECfp2000, execution time of tuned code.

 \bullet SPECfp_base 2000, execution time of untuned code.

 \bullet SPECfp_rate 2000, throughput of tuned code.

02-29

 \bullet SPECfp_rate_rate 2000, throughput of untuned code.

02-30

BAPCO Suites, measure productivity app. performance on Windows 95.

EE 4720 Lecture Transparency. Formatted 13:04, 24 January 2003 from Isli02.

TPC, measure "transaction processing" system performance.

WinMARK, graphics performance.

02-30