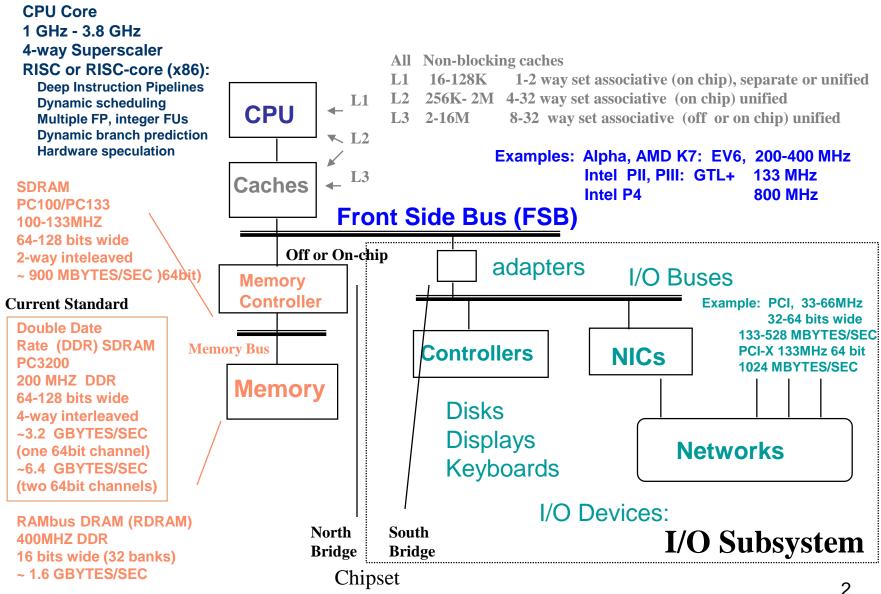
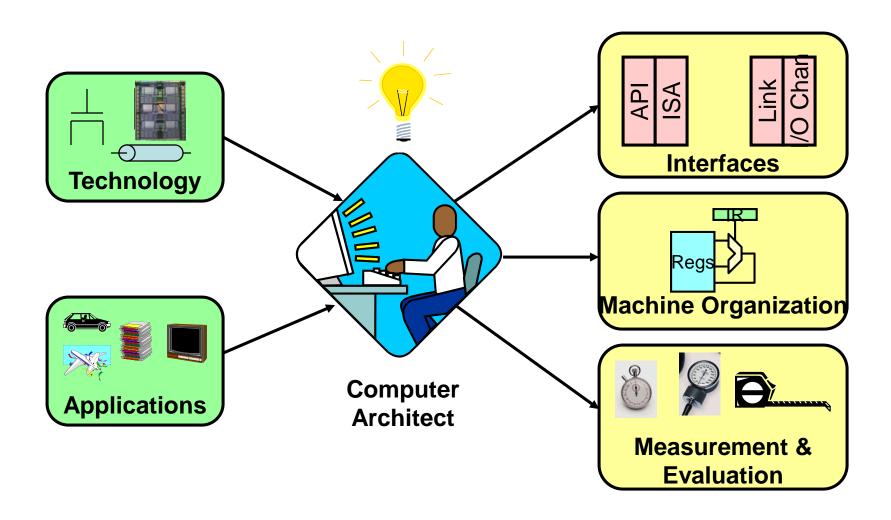
Measuring & Reporting Performance

```
Chapter 1 (4<sup>th</sup> Edition)
OR
Chapter 4 (3<sup>ed</sup> Edition)
```

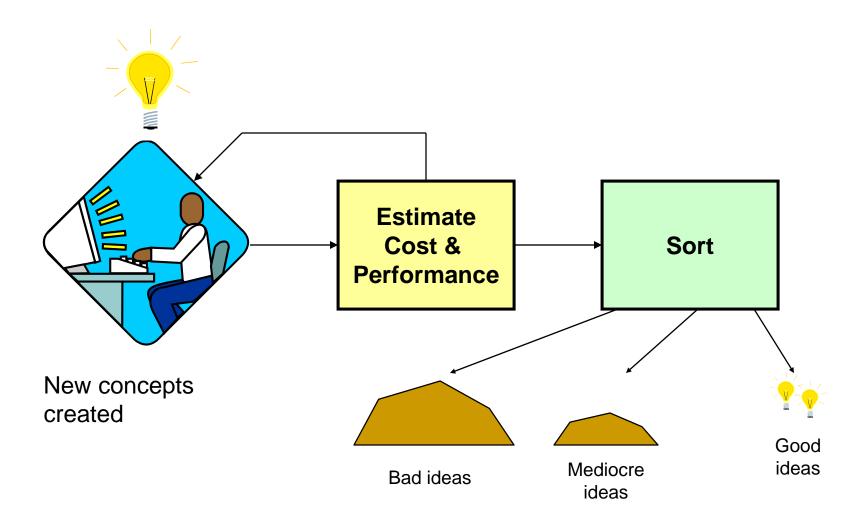
Review: Computer System Components



Review: What is Computer Architecture?



The Architecture Process



Performance

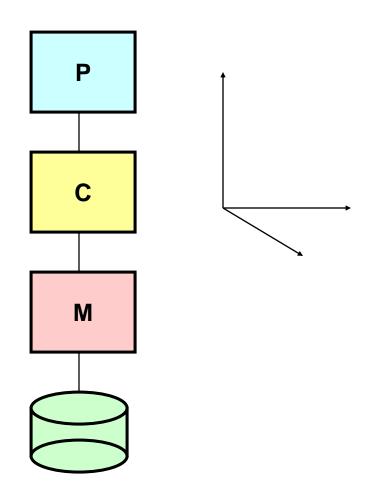
- What is performance?
 - Measuring performance
 - Performance metrics
 - Performance evaluation
 - Why does some hardware perform better
- With different programs?
 - What performance factors are related to hardware?

Measuring performance

- We need measures
 - Comparison of machine properties
 - Comparison of software properties (compilers)
- Purpose
 - Making purchase decisions
 - Development of new architectures
- Is a single measure sufficient?
 - A machine with 600 MHz clock cycle is faster than 500 MHz clock cycle!?
 - Why do we still have mainframes?

Performance Measurement and Evaluation

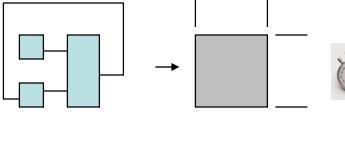
- Many dimensions to computer performance
 - CPU execution time
 - by instruction or sequence
 - floating point
 - integer
 - branch performance
 - Cache bandwidth
 - Main memory bandwidth
 - I/O performance
 - bandwidth
 - seeks
 - pixels or polygons per second
- Relative importance depends on applications



Evaluation Tools

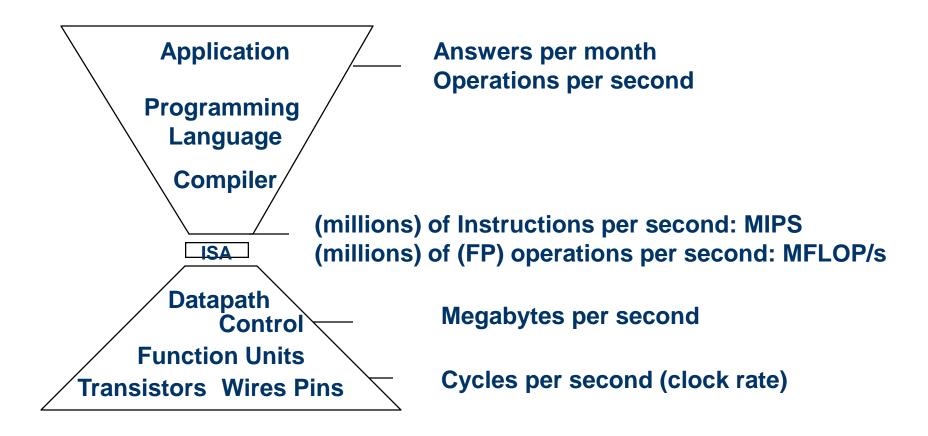
- Benchmarks, traces, & mixes
 - macrobenchmarks & suites
 - application execution time
 - microbenchmarks
 - measure one aspect of performance
 - traces
 - replay recorded accesses
 - cache, branch, register
- Simulation at many levels
 - ISA, cycle accurate, RTL, gate, circuit
 - trade fidelity for simulation rate
- Area and delay estimation
- Analysis
 - e.g., queuing theory
 - Fundamentals Laws

MOVE	39%
BR	20%
LOAD	20%
STORE	10%
ALU	11%





Metrics of Computer Performance



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

Benchmarks and Benchmarking

Some definitions are:

- It is a test that measures the performance of a system or subsystem on a welldefined task or set of task.
- A method of comparing the performance of different computer architecture.
- Or a method of comparing the performance of different software

Some Warnings about Benchmarks

- Benchmarks measure the whole system
 - application
 - compiler
 - operating system
 - architecture
 - implementation
- Popular benchmarks typically reflect yesterday's programs
 - computers need to be designed for tomorrow's programs

- Benchmark timings often very sensitive to
 - alignment in cache
 - location of data on disk
 - values of data
- Benchmarks can lead to inbreeding or positive feedback
 - if you make an operation fast (slow) it will be used more (less) often
 - so you make it faster (slower)
 - and it gets used even more (less)
 - » and so on…

Choosing Programs To Evaluate Performance

Levels of programs or benchmarks that could be used to evaluate performance:

- Actual Target Workload: Full applications that run on the target machine.
- Real Full Program-based Benchmarks:
 - Select a specific mix or suite of programs that are typical of targeted applications or workload (e.g SPEC95, SPEC CPU2000).
- Small "Kernel" Benchmarks:
 - Key computationally-intensive pieces extracted from real programs.
 - Examples: Matrix factorization, FFT, tree search, etc.
 - Best used to test specific aspects of the machine.
- Microbenchmarks:
 - Small, specially written programs to isolate a specific aspect of performance characteristics: Processing: integer, floating point, local memory, input/output, etc.

Types of Benchmarks

Pros

Representative

Actual Target Workload

- Very specific.

Cons

- Non-portable.
- Complex: Difficult to run, or measure.

- Portable.
- Widely used.
- Measurements useful in reality.

Full Application Benchmarks

 Less representative than actual workload.

- Easy to run, early in the design cycle.
- Identify peak performance and potential bottlenecks.

Benchmarks

Small "Kernel"

Microbenchmarks

- Easy to "fool" by designing hardware to run them well.
- Peak performance results may be a long way from real application performance

SPEC: System Performance Evaluation Cooperative

The most popular and industry-standard set of CPU benchmarks.

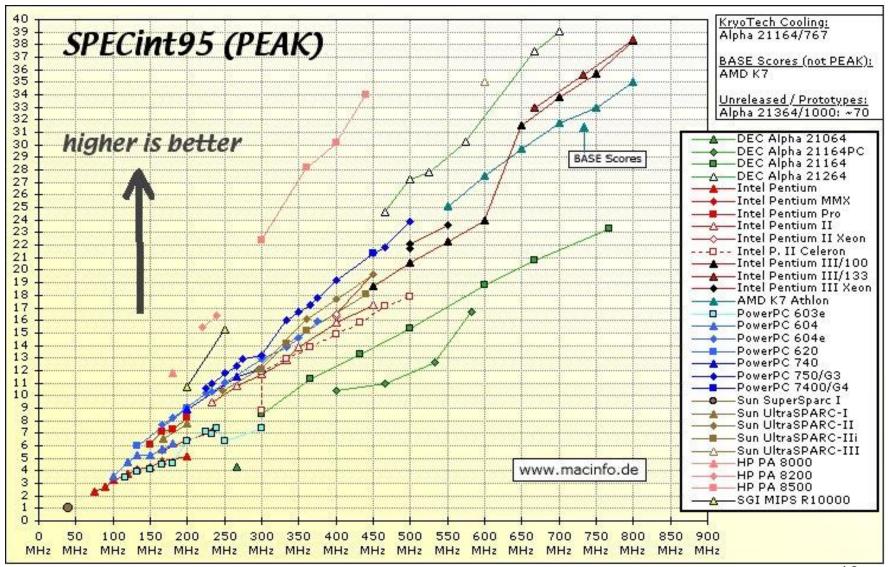
- SPECmarks, 1989:
 - 10 programs yielding a single number ("SPECmarks").
- SPEC92, 1992:
 - SPECInt92 (6 integer programs) and SPECfp92 (14 floating point programs).
- <u>SPEC95, 1995:</u>
 - SPECint95 (8 integer programs):
 - go, m88ksim, gcc, compress, li, ijpeg, perl, vortex
 - SPECfp95 (10 floating-point intensive programs):
 - tomcatv, swim, su2cor, hydro2d, mgrid, applu, turb3d, apsi, fppp, wave5
 - Performance relative to a Sun SuperSpark I (50 MHz) which is given a score of SPECint95 = SPECfp95 = 1
- SPEC CPU2000, 1999:
 - CINT2000 (11 integer programs). CFP2000 (14 floating-point intensive programs)
 - Performance relative to a Sun Ultra5_10 (300 MHz) which is given a score of SPECint2000 = SPECfp2000 = 100

SPEC95 Programs

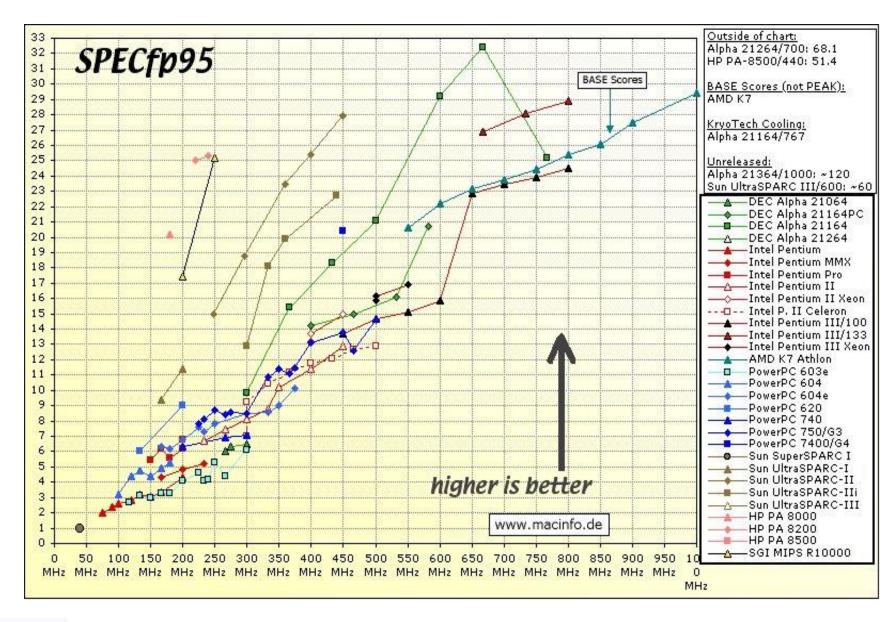
Programs application domain: Engineering and scientific computation

	Benchmark	Description
	go	Artificial intelligence; plays the game of Go
	m88ksim	Motorola 88k chip simulator; runs test program
	gcc	The Gnu C compiler generating SPARC code
Integer	compress	Compresses and decompresses file in memory
	li	Lisp interpreter
	ijpeg	Graphic compression and decompression
	perl	Manipulates strings and prime numbers in the special-purpose programming language Perl
	vortex	A database program
	tomcatv	A mesh generation program
	swim	Shallow water model with 513 x 513 grid
Floating	su2cor	quantum physics; Monte Carlo simulation
Point	hydro2d	Astrophysics; Hydrodynamic Naiver Stokes equations
	mgrid	Multigrid solver in 3-D potential field
	applu	Parabolic/elliptic partial differential equations
	trub3d	Simulates isotropic, homogeneous turbulence in a cube
	apsi	Solves problems regarding temperature, wind velocity, and distribution of pollutant
	fpppp	Quantum chemistry
	wave5	Plasma physics; electromagnetic particle simulation

Sample SPECint95 (Integer) Results



Sample SPECfp95 (Floating Point) Results

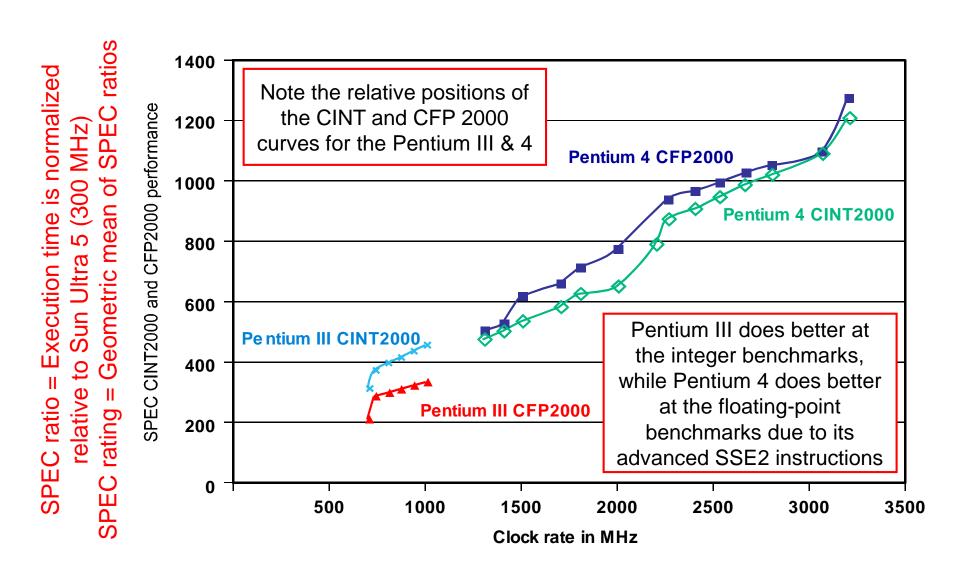


The SPEC CPU2000 Benchmarks

12 Integer benchmarks (C and C++)		14 FP benchmarks (Fortran 77, 90, and C)		
Name	Description	Name	Description	
gzip	Compression	wupwise	Quantum chromodynamics	
vpr	FPGA placement and routing	swim	Shallow water model	
gcc	GNU C compiler	mgrid	Multigrid solver in 3D potential field	
mcf	Combinatorial optimization	applu	Partial differential equation	
crafty	Chess program	mesa	Three-dimensional graphics library	
parser	Word processing program	galgel	Computational fluid dynamics	
eon	Computer visualization	art	Neural networks image recognition	
perlbmk	Perl application	equake	Seismic wave propagation simulation	
gap	Group theory, interpreter	facerec	Image recognition of faces	
vortex	Object-oriented database	ammp	Computational chemistry	
bzip2	Compression	lucas	Primality testing	
twolf	Place and route simulator	fma3d	Crash simulation using finite elements	
		sixtrack	High-energy nuclear physics	
		apsi	Meteorology: pollutant distribution	

- ❖ Wall clock time is used as metric
- ❖ Benchmarks measure CPU time, because of little I/O

SPEC 2000 Ratings (Pentium III & 4)



Common Benchmarking Mistakes

- Only average behavior represented in test workload
- Skewness of device demands ignored
- Loading level controlled inappropriately
- Caching effects ignored
- Buffer sizes not appropriate
- Inaccuracies due to sampling ignored

- Ignoring monitoring overhead
- Not validating measurements
- Not ensuring same initial conditions
- Not measuring transient (cold start) performance
- Using device utilizations for performance comparisons
- Collecting too much data but doing too little analysis

Architectural Performance Laws and Rules of Thumb

Measurement and Evaluation

- Architecture is an iterative process:
 - Searching the space of possible designs
 - Make selections
 - Evaluate the selections made
- Good measurement tools are required to accurately evaluate the selection.

Measurement Tools

- Benchmarks, Traces, Mixes
- Cost, delay, area, power estimation
- Simulation (many levels)
 - ISA, RTL, Gate, Circuit
- Queuing Theory
- Rules of Thumb
- Fundamental Laws

Measuring and Reporting Performance

- What do we mean by one Computer is faster than another?
 - program runs less time
- Response time or execution time
 - time that users see the output
- Elapsed time
 - A latency to complete a task including disk accesses, memory accesses,
 I/O activities, operating system overhead
- Throughput
 - total amount of work done in a given time
- Performance
 - "Increasing and decreasing" ?????
 - We use the term "improve performance" or "improve execution time" When we mean increase performance and decrease execution time.
 - improve performance = increase performance
 - improve execution time = decrease execution time

What is time?

Elapsed Time

- counts everything (disk and memory accesses, I/O, etc.)
- a useful number, but often not good for comparison purposes

CPU time

- time the CPU is computing
- doesn't count I/O or time spent running other programs
- User CPU time
 - CPU time spent in the program
- System CPU time
 - CPU time spent in the operating system performing task requested by the program decrease execution time
- CPU time = User CPU time + System CPU time
- Our focus: user CPU time
 - time spent executing the lines of code that are "in" our program

Performance

- System Performance
 - elapsed time on unloaded system
- CPU performance
 - user CPU time on an unloaded system
- Two notions of "performance
 - Response Time (latency)
 - Time to do the task
 - How long does it take for my job to run?
 - How long does it take to execute a job?
 - How long must I wait for the database query?
 - Throughput (bandwidth)
 - Tasks per day, hour, week, sec, ns. ..
 - How many jobs can the machine run at once?
 - What is the average execution rate?
 - How much work is getting done?
- Response time and throughput often are in opposition

CPU Performance Evaluation

• Most computers run synchronously utilizing a CPU clock running at a constant clock rate:

Cycles/sec = Hertz = Hz

Cycles/sec = Hertz = Hz

Cycles/sec = Hertz = Hz

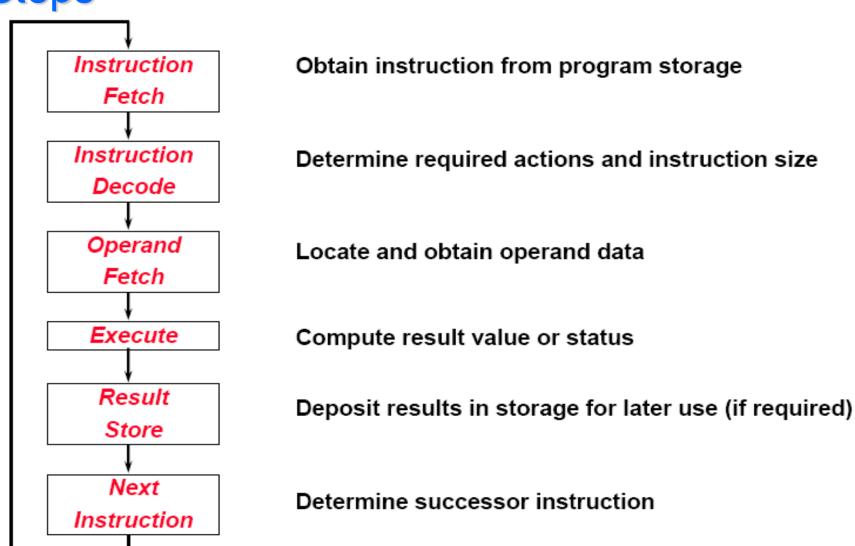
where: Clock rate = 1 / clock cycle

• The CPU clock rate <u>depends</u> on the <u>specific CPU organization</u> (<u>design</u>) and hardware <u>implementation technology (VLSI)</u> used

- A computer machine (ISA) instruction is comprised of a number of elementary or <u>micro operations</u> which vary in number and complexity depending on the <u>instruction</u> and the <u>exact CPU organization</u> (<u>Design</u>)
 - A <u>micro operation</u> is an elementary hardware operation that can be performed during one CPU clock cycle.
 - This corresponds to one micro-instruction in microprogrammed CPUs.
 - Examples: register operations: shift, load, clear, increment, ALU operations: add, subtract, etc.
- Thus a single machine instruction may take one or more CPU cycles to complete termed as the <u>Cycles Per Instruction (CPI)</u>.
 - Average CPI of a program: The average CPI of all instructions executed in the program on a given CPU design.

25

Generic CPU Machine Instruction Execution Steps



Computer Performance Measures: Program Execution Time

- For <u>a specific program</u> compiled to run on <u>a specific</u> machine (CPU) "A", has the following parameters:
 - The total executed instruction count of the program.
 - The average number of cycles per instruction (average CPI).
 - Clock cycle of machine "A"
- How can one measure the performance of this machine (CPU) running this program?
 - Intuitively the machine (or CPU) is said to be faster or has better performance running this program if the total execution time is shorter.
 - Thus the inverse of the total measured program execution time is a possible performance measure or metric:

 $Performance_A = 1 / Execution Time_A$

How to compare performance of different machines? What factors affect performance? How to improve performance?

Comparing Computer Performance Using Execution Time

• To compare the performance of two machines (or CPUs) "A", "B" running <u>a given specific program</u>:

```
Performance_A = 1 / Execution Time_A

Performance_B = 1 / Execution Time_B
```

Machine A is n times faster than machine B means (or slower? if n < 1)

$$Speedup = n = \frac{Performance_A}{Performance_B} = \frac{Execution Time_B}{Execution Time_A}$$

Example:

(i.e Speedup is ratio of performance, no units)

For a given program:

Execution time on machine A: Execution_A = 1 second

Execution time on machine B: Execution_B = 10 seconds

Speedup= Performance_A / Performance_B = Execution Time_B / Execution Time_A =
$$10 / 1 = 10$$

The performance of machine A is 10 times the performance of machine B when running this program, or: Machine A is said to be 10 times faster than machine B when running this program.

CPU Execution Time: The CPU Equation

- A program is comprised of <u>a number of instructions</u> executed, I
 - Measured in: instructions/program
- The average instruction executed takes a number of cycles

per instruction (CPI) to be completed.

Or Instructions Per Cycle (IPC): IPC= 1/CPI

- Measured in: cycles/instruction, CPI
- CPU has a fixed clock cycle time <u>C = 1/clock rate</u>
 - Measured in: seconds/cycle
- CPU execution time is the product of the above three parameters as follows:

```
CPU time = Seconds = Instructions x Cycles x Seconds
Program Program Instruction Cycle
```

T =

IX

CPI

X

C

Execution Time per program in seconds

Number of instructions executed

Average CPI for program

CPU Clock Cycle

CPU Average CPI/Execution Time

For a given program executed on a given machine (CPU):

```
CPI = Total program execution cycles / Instructions count
```

```
    → CPU clock cycles = Instruction count x CPI
    CPU execution time =
    = CPU clock cycles x Clock cycle
    = Instruction count x CPI x Clock cycle
    T = I x CPI x C
```

execution Time per program in seconds

Number of instructions executed

Average CPI for program

CPU Clock Cycle

(This equation is commonly known as the <u>CPU performance equation</u>)

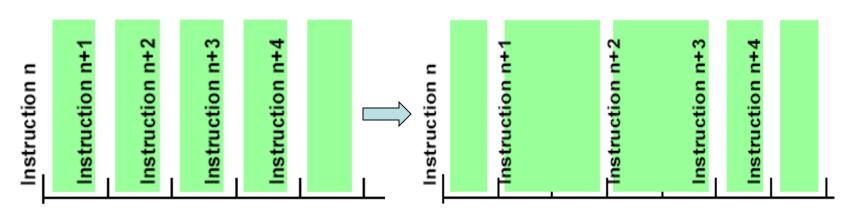
Improving the performance

```
CPU time = Seconds = Instructions x Cycles x Seconds
Program Program Instruction Cycle
```

- Increase the clock frequency =
 - reduce the clock period
- Reduce the number of cycles for the program
- Reduce the number of instructions

Instruction = cycle?

- Is the number of cycles identical with the number of instructions?
 - No!
 - The number of cycles depends on the implementation of the operations in hardware
 - The number differs for each processor
 - Why?
 - Operations take different time
 - Multiplication takes longer than addition
 - Floating point operations take longer than integer operations
 - The access time to a register is much shorter than to memory location



Aspects of CPU Execution Time

CPU Time = Instruction count x **CPI** x **Clock cycle Depends on:** $T = I \times CPI \times C$ **Program Used** Compiler **ISA Instruction Count I** (executed) Depends on: **Depends on: Program Used CPU Organization** Clock **CPI Compiler** Cycle Technology (VLSI) **ISA** (Average **CPU Organization** CPI)

Factors Affecting CPU Performance

 CPU time
 = Seconds
 = Instructions x Cycles
 x Seconds

 Program
 Program
 Instruction
 Cycle

		Instruction Count I	CPI	Clock Cycle C
-	Program	X	X	
	Compiler	X	X	
	ruction Set cture (ISA)	X	X	
	rganization (CPU Design)		X	X
<u>-</u> <u>r</u>	Technology (VLSI)			X

CPU Execution Time: Example

- A Program is running on a specific machine (CPU) with the following parameters:
 - Total executed instruction count: 10,000,000 instructions
 - Average CPI for the program: 2.5 cycles/instruction.
 - CPU clock rate: 200 MHz. (clock cycle = $5x10^{-9}$ seconds)
- What is the execution time for this program:

```
      CPU time
      = Seconds
      = Instructions x
      Cycles
      x
      Seconds

      Program
      Program
      Instruction
      Cycle
```

```
CPU time = Instruction count x CPI x Clock cycle

= 10,000,000 x 2.5 x 1 / clock rate

= 10,000,000 x 2.5 x 5x10^{-9}

= .125 seconds

T = I x CPI x C
```

Performance Comparison: Example

- From the previous example: A Program is running on a specific machine (CPU) with the following parameters:
 - Total executed instruction count, I: 10,000,000 instructions
 - Average CPI for the program: 2.5 cycles/instruction.
 - CPU clock rate: 200 MHz.
- Using the same program with these changes:
 - A new compiler used: New executed instruction count, I: 9,500,000
 New CPI: 3.0
 - Faster CPU implementation: New clock rate = 300 MHz
- What is the speedup with the changes?

```
Speedup = (10,000,000 \times 2.5 \times 5x10^{-9}) / (9,500,000 \times 3 \times 3.33x10^{-9})
= .125 / .095 = 1.32
or 32 % faster after changes.
```

$$T = I x CPI x C$$

Instruction Types & CPI

 Given a program with n types or classes of instructions executed on a given CPU with the following characteristics:

$$i = 1, 2, n$$

 C_i = Count of instructions of type_i executed CPI_i = Cycles per instruction for type_i

Then:

CPI = CPU Clock Cycles / Instruction Count I

Where:

$$CPUclockcycles = \sum_{i=1}^{n} (CPI_i \times C_i)$$

Executed Instruction Count $= \sum_{i} C_{i}$

Instruction Types & CPI: An Example

An instruction set has three instruction classes:

Instruction class	CPI -	
\mathbf{A}	1	For a specific
В	2	CPU design
C	3	

Two code sequences have the following instruction counts:

Instruction counts for instruction of					
Code Sequence	${f A}$	B	\mathbf{C}		
1	2	1	2		
2	4	1	1		

- CPU cycles for sequence 1 = 2 x 1 + 1 x 2 + 2 x 3 = 10 cycles
 CPI for sequence 1 = clock cycles / instruction count
 = 10 /5 = 2
- CPU cycles for sequence 2 = 4 x 1 + 1 x 2 + 1 x 3 = 9 cycles
 CPI for sequence 2 = 9 / 6 = 1.5

$$CPU \ clock \ cycles = \sum_{i=1}^{n} (CPI_i \times C_i)$$
 CPI = CPU Cycles / I

Instruction Frequency & CPI

 Given a program with n types or classes of instructions with the following characteristics:

 C_i = Count of instructions of type_i i = 1, 2, n CPI_i = Average cycles per instruction of type_i F_i = Frequency or fraction of instruction type_i executed = C_i / total executed instruction count = C_i / I

Then:

$$CPI = \sum_{i=1}^{n} (CPI_i \times F_i)$$

Fraction of total execution time for instructions of type $i = \frac{CPI_i \times F_i}{CPI}$

Instruction Type Frequency & CPI: A RISC Example

Program Profile or Executed Instructions Mix Base Machine (Reg / Reg)

CPI; Op Freq, F_i **ALU 50%** Load 20% 5 10% Store **Branch** 20%

CPI_i x F_i .5 1.0

% Time 23% = .5/2.2

45% = **1/2.2** 14% = .3/2.2

CPI_i x F_i

CPI

18% = .4/2.2

$$CPI = \sum_{i=1}^{n} (CPI_i \times F_i)$$

CPI =
$$.5 \times 1 + .2 \times 5 + .1 \times 3 + .2 \times 2 = 2.2$$

= $.5 + 1 + .3 + .4$

Sum = 2.2

Given

Computer Performance Measures: MIPS (Million Instructions Per Second) Rating

 For a specific program running on a specific CPU the MIPS rating is a measure of how many millions of instructions are executed per second:

```
MIPS Rating = Instruction count / (Execution Time x 10<sup>6</sup>)
= Instruction count / (CPU clocks x Cycle time x 10<sup>6</sup>)
= (Instruction count x Clock rate) / (Instruction count x CPI x 10<sup>6</sup>)
```

 Major problem with MIPS rating: As shown above the MIPS rating does not account for the count of instructions executed (I).

= Clock rate / (CPI x 10^6)

- A higher MIPS rating in many cases may not mean higher performance or better execution time. i.e. due to compiler design variations.
- In addition the MIPS rating:
 - Does not account for the instruction set architecture (ISA) used.
 - Thus it cannot be used to compare computers/CPUs with different instruction sets.
 - Easy to abuse: Program used to get the MIPS rating is often omitted.
 - Often the <u>Peak MIPS rating</u> is provided for a given CPU which is obtained using a program comprised entirely of <u>instructions with the lowest CPI</u> for the given CPU design which <u>does not represent real programs.</u>

Computer Performance Measures: MIPS (Million Instructions Per Second) Rating

- Under what conditions can the MIPS rating be used to compare performance of different CPUs?
- The MIPS rating is <u>only valid</u> to compare the performance of different CPUs <u>provided that the following conditions are</u> <u>satisfied:</u>
 - 1 The same program is used (actually this applies to all performance metrics)
 - 2 The same ISA is used
 - 3 The same compiler is used
 - ⇒ (Thus the resulting programs used to run on the CPUs and obtain the MIPS rating are <u>identical</u> at the machine code level including the <u>same instruction count</u>)

Wrong!!!

- 3 significant problems with using MIPS:
 - Problem 1:
 - MIPS is instruction set dependent.
 - (And different computer brands usually have different instruction sets)
 - Problem 2:
 - MIPS varies between programs on the same computer
 - Problem 3:
 - MIPS can vary inversely to performance!
- Let's look at an examples of why MIPS doesn't work...

Compiler Variations, MIPS & Performance: An Example

For a machine (CPU) with instruction classes:

Instruction class	CPI
A	1
В	2
C	3

 For a given high-level language program, two compilers produced the following executed instruction counts:

	Instruction counts (in millions) for each instruction class			
Code from:	\mathbf{A}	В	\mathbf{C}	
Compiler 1	5	1	1	
Compiler 2	10	1	1	

The machine is assumed to run at a clock rate of 100 MHz.

Compiler Variations, MIPS & Performance: An Example (Continued)

MIPS = Clock rate / (CPI x 10^6) = 100 MHz / (CPI x 10^6) CPI = CPU execution cycles / Instructions count

$$CPU \ clock \ cycles = \sum_{i=1}^{n} (CPI_i \times C_i)$$

CPU time = Instruction count x CPI / Clock rate

For compiler 1:

- $CPI_1 = (5 \times 1 + 1 \times 2 + 1 \times 3) / (5 + 1 + 1) = 10 / 7 = 1.43$
- MIPS Rating₁ = $100 / (1.428 \times 10^6) = 70.0 \text{ MIPS}$
- CPU time₁ = $((5 + 1 + 1) \times 10^6 \times 1.43) / (100 \times 10^6) = 0.10$ seconds

For compiler 2:

- CPI₂ = (10 x 1 + 1 x 2 + 1 x 3) / (10 + 1 + 1) = 15 / 12 = 1.25
- MIPS Rating₂ = $100 / (1.25 \times 10^6) = 80.0 \text{ MIPS}$
- CPU time₂ = $((10 + 1 + 1) \times 10^6 \times 1.25) / (100 \times 10^6) = 0.15$ seconds

MIPS rating indicates that compiler 2 is better while in reality the code produced by compiler 1 is faster

MIPS (The ISA not the metric) Loop

Performance Example

For the loop:

```
for (i=0; i<1000; i=i+1){
 x[i] = x[i] + s; }
```

MIPS assembly code is given by:

```
$3, 8($1)
       lw
                       ; load s in $3
                                                                compute
                                                         Low Memory
           $6, $2, 4000; $6 = address of last element + 4
             $4, 0($2)
                           ; load x[i] in $4
      lw
loop:
       add $5, $4, $3
                        ; $5 has x[i] + s
       sw \$5, 0(\$2) ; store computed x[i]
       addi $2, $2, 4 ; increment $2 to point to next x[] element
             $6, $2, loop ; last loop iteration reached?
      bne
```

The MIPS code is executed on a specific CPU that runs at 500 MHz (clock cycle = $2ns = 2x10^{-9}$ seconds) with following instruction type CPIs:

For this MIPS code running on this CPU find:

Instruction type	CPI	1- Fraction of total instructions executed for each instruction type
ALU	4	2- Total number of CPU cycles
Load	5	3- Average CPI
Store	7	4- Fraction of total execution time for each instructions type
Branch	3	5- Execution time
		6- MIPS rating, peak MIPS rating for this CPU

X[] array of words in memory, base address in \$2, s a constant word value in memory, address in \$1

High Memory

← Last element to compute

←First element to

X[999]

X[998]

\$6 points here →

\$2 initially

points here $\rightarrow X[0]$

MIPS (The ISA) Loop Performance Example (continued)

- The code has 2 instructions before the loop and 5 instructions in the body of the loop which iterates 1000 times,
- Thus: Total instructions executed, I = 5x1000 + 2 = 5002 instructions
- 1 Number of instructions executed/fraction F_i for each instruction type:
- 2 $CPU \ clock \ cycles = \sum_{i=1}^{n} (CPI_i \times C_i)$

$$= 2001x4 + 1001x5 + 1000x7 + 1000x3 = 23009$$
 cycles

- 3 Average CPI = CPU clock cycles / I = 23009/5002 = 4.6
- 4 Fraction of execution time for each instruction type:
 - Fraction of time for ALU instructions = $CPI_{ALL} \times F_{ALL} / CPI = 4x0.4/4.6 = 0.348 = 34.8\%$
 - Fraction of time for load instructions = CPI_{load} x F_{load} / CPI= 5x0.2/4.6 = 0.217 = 21.7%
 - Fraction of time for store instructions = $CPI_{store} \times F_{store} / CPI = 7x0.2/4.6 = 0.304 = 30.4\%$
 - Fraction of time for branch instructions = CPI_{branch} x F_{branch} / CPI = 3x0.2/4.6 = 0.13 = 13%
- 5 Execution time = I x CPI x C = CPU cycles x C = $23009 \times 2 \times 10^{-9} =$

$$= 4.6x \cdot 10^{-5} \text{ seconds} = 0.046 \text{ msec} = 46 \text{ usec}$$

- 6 MIPS rating = Clock rate / (CPI x 10^6) = 500 / 4.6 = 108.7 MIPS
 - The CPU achieves its peak MIPS rating when executing a program that only has instructions of the type with the lowest CPI. In this case branches with CPI_{Branch} = 3
 - Peak MIPS rating = Clock rate / $(CPI_{Branch} \times 10^6) = 500/3 = 166.67 MIPS$

CPI

Instruction type

ALU Load

Store

Branch

Computer Performance Measures : MFLOPS

- A floating-point operation is an addition, subtraction, multiplication, or division operation applied to numbers represented by a single or a double precision floating-point representation.
- MFLOPS, for a specific program running on a specific computer, is a measure of millions of floating point-operation (megaflops) per second:

MFLOPS = Number of floating-point operations / (Execution time x 10⁶)

- MFLOPS rating is a better comparison measure between different machines (applies even if ISAs are different) than the MIPS rating.
 - Applicable even if ISAs are different

Computer Performance Measures : MFLOPS

- Program-dependent: Different programs have different percentages of floating-point operations present. i.e compilers have no floating-point operations and yield a MFLOPS rating of zero.
- Dependent on the type of floating-point operations present in the program.
 - Peak MFLOPS rating for a CPU: Obtained using a program comprised entirely of the <u>simplest floating point</u> <u>instructions</u> (with the lowest CPI) for the given CPU design which <u>does not represent real floating point programs.</u>

Quantitative Principles of Computer Design

- Amdahl's Law:
 - The performance gain from improving some portion of a computer is calculated by:

```
Speedup = Performance for entire task using the enhancement
Performance for the entire task without using the enhancement
```

```
or Speedup = Execution time without the enhancement Execution time for entire task using the enhancement
```

Performance Enhancement Calculations: Amdahl's Law

- The performance enhancement possible due to a given design improvement is limited by the amount that the improved feature is used
- Amdahl's Law:
 - Performance improvement or speedup due to enhancement E:

 Suppose that enhancement E accelerates a fraction F of the execution time by a factor S and the remainder of the time is unaffected then:

Execution Time with $E = ((1-F) + F/S) \times Execution Time without E$ Hence speedup is given by:

F (Fraction of execution time enhanced) refers to original execution time before the enhancement is applied

Pictorial Depiction of Amdahl's Law

Enhancement E accelerates fraction F of original execution time by a factor of S

Before:

Execution Time without enhancement E: (Before enhancement is applied)

• shown normalized to 1 = (1-F) + F = 1



After:

Execution Time with enhancement E:

Speedup(E) =
$$\begin{array}{c} & \text{Execution Time without enhancement E} & 1 \\ & \text{Speedup(E)} & \text{Execution Time with enhancement E} & (1 - F) + F/S \\ \end{array}$$

Example of Amdahl's Law

Floating point instructions improved to run 2X;
 but only 10% of actual instructions are FP

ExTime_{new} = ExTime_{old} x
$$(0.9 + .1/2) = 0.95$$
 x ExTime_{old}

Speedup_{overall} = $\frac{1}{0.95}$ = 1.053

Performance Enhancement Example

For the RISC machine with the following instruction mix given earlier:

Op	Freq	Cycles	CPI(i)	% Time	
ALU	50%	1	.5	23%	$\mathbf{CPI} = 2.2$
Load	20%	5	1.0	45%	
Store	10%	3	.3	14%	
Branch	20%	2	.4	18%	

• If a CPU design enhancement improves the CPI of load instructions from 5 to 2, what is the resulting performance improvement from this enhancement:

Fraction enhanced =
$$F = 45\%$$
 or .45
Unaffected fraction = 1- $F = 100\%$ - 45% = 55% or .55
Factor of enhancement = $S = 5/2 = 2.5$
Using Amdahl's Law:

Speedup_(E) = ----- = ----- = 1.37

$$(1 - F) + F/S$$
 $.55 + .45/2.5$

An Alternative Solution Using CPU Equation

Op	Freq	Cycles	CPI(i)	% Time	
ALU	50%	1	.5	23%	
Load	20%	5	1.0	45%	$\mathbf{CPI} = 2.2$
Store	10%	3	.3	14%	
Branch	20%	2	.4	18%	

• If a CPU design enhancement improves the CPI of load instructions from 5 to 2, what is the resulting performance improvement from this enhancement:

Old CPI =
$$2.2$$

New CPI = $.5 \times 1 + .2 \times 2 + .1 \times 3 + .2 \times 2 = 1.6$

Which is the same speedup obtained from Amdahl's Law in the first solution.

Performance Enhancement Example

 A program runs in 100 seconds on a machine with multiply operations responsible for 80 seconds of this time. By how much must the speed of multiplication be improved to make the program four times faster?

 \rightarrow Execution time with enhancement = 100/4 = 25 seconds 25 seconds = (100 - 80 seconds) + 80 seconds / S 25 seconds = 20 seconds + 80 seconds / S \rightarrow 5 = 80 seconds / S \rightarrow S = 80/5 = 16

Alternatively, it can also be solved by finding enhanced fraction of execution time:

and then solving Amdahl's speedup equation for desired enhancement factor S Hence multiplication should be 16 times faster to get an overall speedup of 4.

Performance Enhancement Example

 For the previous example with a program running in 100 seconds on a machine with multiply operations responsible for 80 seconds of this time. By how much must the speed of multiplication be improved to make the program five times faster?

 \rightarrow Execution time with enhancement = 100/5 = 20 seconds

20 seconds =
$$(100 - 80 \text{ seconds}) + 80 \text{ seconds} / s$$

20 seconds = 20 seconds + $80 \text{ seconds} / s$
 $\Rightarrow 0 = 80 \text{ seconds} / s$

No amount of multiplication speed improvement can achieve this.

Extending Amdahl's Law To Multiple Enhancements

Suppose that enhancement E_i accelerates a fraction F_i of the original execution time by a factor S_i and the remainder of the time is unaffected then:

$$Speedup = \frac{Original Execution Time}{\left((1 - \sum_{i} \mathbf{F}_{i}) + \sum_{i} \frac{\mathbf{F}_{i}}{\mathbf{S}_{i}}\right) X Original Execution Time}$$

$$Unaffected fraction$$

$$Speedup = \frac{1}{\left((1 - \sum_{i} \mathbf{F}_{i}) + \sum_{i} \frac{\mathbf{F}_{i}}{\mathbf{S}_{i}}\right)}$$

Note: All fractions $\mathbf{F}_{\mathbf{i}}$ refer to original execution time before the enhancements are applied.

Amdahl's Law With Multiple Enhancements: Example

 Three CPU performance enhancements are proposed with the following speedups and percentage of the code execution time affected:

Speedup₁ =
$$S_1$$
 = 10 Percentage₁ = F_1 = 20%
Speedup₂ = S_2 = 15 Percentage₁ = F_2 = 15%
Speedup₃ = S_3 = 30 Percentage₁ = F_3 = 10%

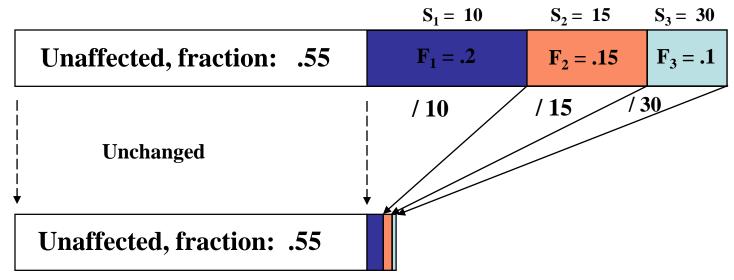
- While all three enhancements are in place in the new design, each enhancement affects a different portion of the code and only one enhancement can be used at a time.
- What is the resulting overall speedup?

$$Speedup = \frac{1}{\left(\left(1 - \sum_{i} \mathbf{F}_{i}\right) + \sum_{i} \frac{\mathbf{F}_{i}}{\mathbf{S}_{i}}\right)}$$

Pictorial Depiction of Example

Before:

Execution Time with no enhancements: 1



After:

Execution Time with enhancements: .55 + .02 + .01 + .00333 = .5833

Speedup = 1 / .5833 = 1.71

Note: All fractions refer to original execution time.

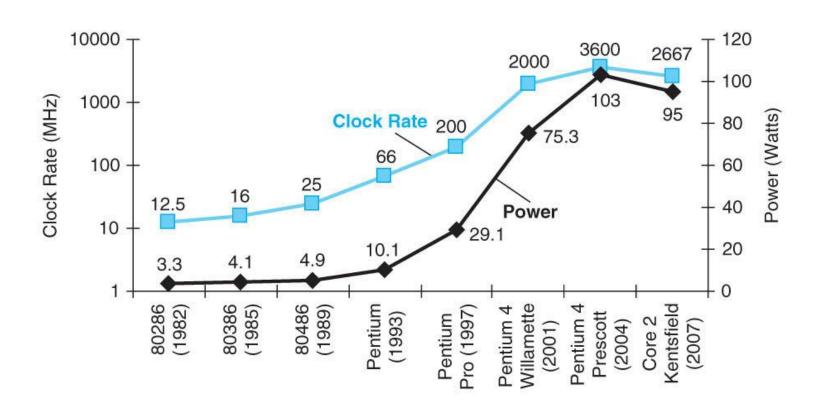
Performance and Power

- Power is a key limitation
 - Battery capacity has improved only slightly over time
- Need to design power-efficient processors
- Reduce power by
 - Reducing frequency
 - Reducing voltage
 - Putting components to sleep
- Energy efficiency
 - Important metric for power-limited applications
 - Defined as performance divided by power consumption

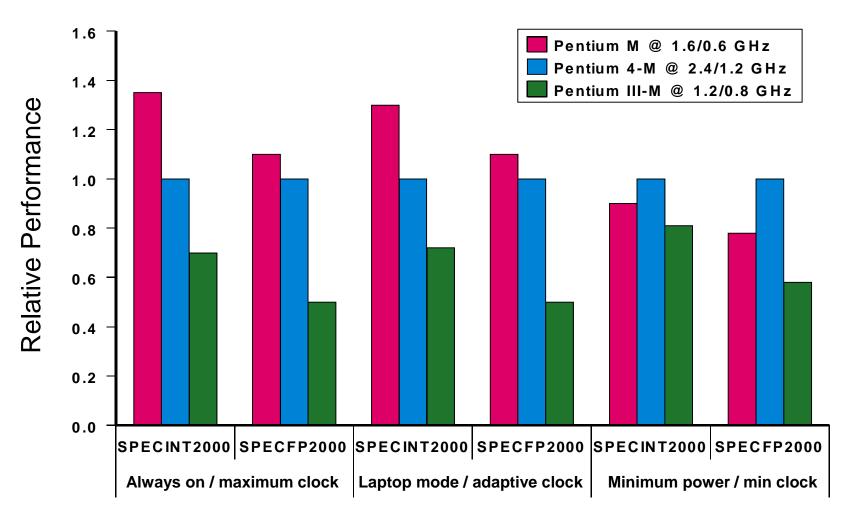
Dynamic Energy and Power

- Dynamic energy
 - Transistor switch from 0 -> 1 or 1 -> 0
 - ½ x Capacitive load x Voltage²
- Dynamic power
 - ½ x Capacitive load x Voltage² x Frequency switched
- Reducing clock rate reduces power, not energy

Power & Clock Rate

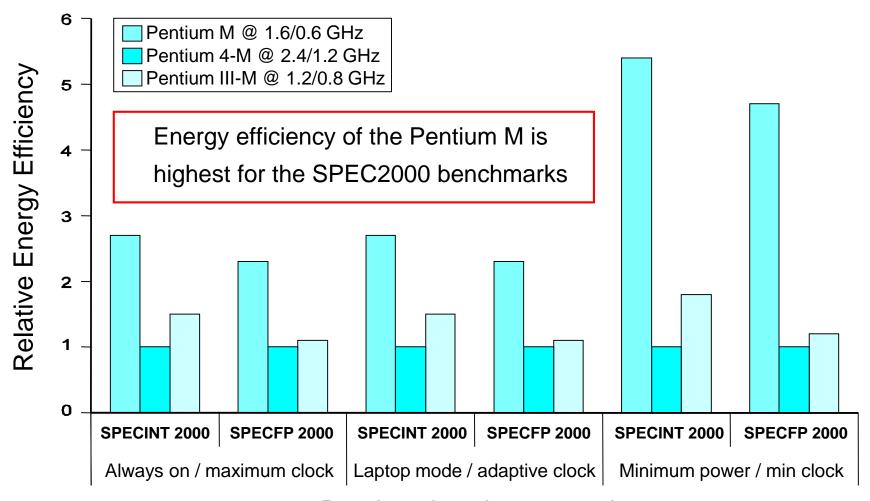


Performance and Power



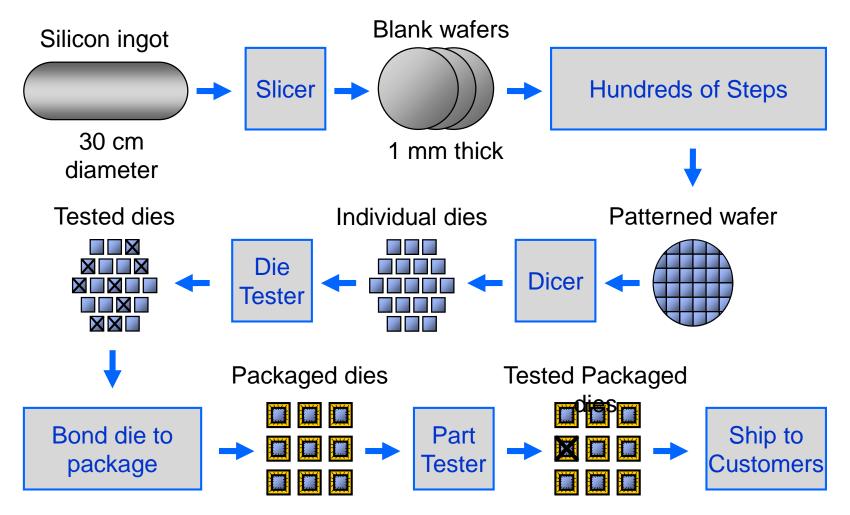
Benchmark and Power Mode

Energy Efficiency

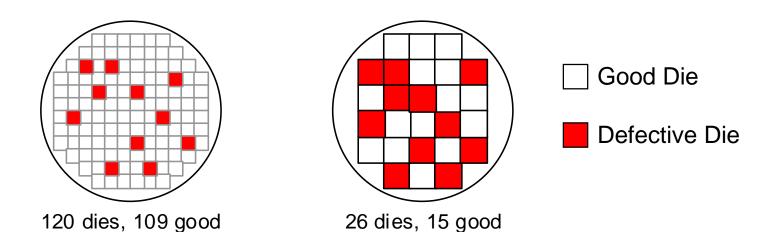


Benchmark and power mode

Chip Manufacturing Process



Effect of Die Size on Yield



Dramatic decrease in yield with larger dies

Yield = (Number of Good Dies) / (Total Number of Dies)

Yield = $\frac{1}{(1 + (\text{Defect per area} \times \text{Die area} / 2))^2}$ Die Cost = (Wafer Cost) / (Dies per Wafer × Yield)

STUDENTS-HUB.com

Integrated Circuit Cost

Integrated circuit

$$Cost of integrated circuit = \frac{Cost of die + Cost of testing die + Cost of packaging and final test}{Final test yield}$$

Cost of die =
$$\frac{\text{Cost of wafer}}{\text{Dies per wafer} \times \text{Die yield}}$$

Dies per wafer =
$$\frac{\pi \times (\text{Wafer diameter/2})^2}{\text{Die area}} - \frac{\pi \times \text{Wafer diameter}}{\sqrt{2 \times \text{Die area}}}$$

Yield = (Number of Good Dies) / (Total Number of Dies)

Yield =
$$\frac{1}{(1 + (\text{Defect per area} \times \text{Die area} / 2))^2}$$

Things to Remember

- Performance is specific to a particular program
 - Any measure of performance should reflect execution time
 - Total execution time is a consistent summary of performance
- For a given ISA, performance improvements come from
 - Increases in clock rate (without increasing the CPI)
 - Improvements in processor organization that lower CPI
 - Compiler enhancements that lower CPI and/or instruction count
 - Algorithm/Language choices that affect instruction count
- Pitfalls (things you should avoid)
 - Using a subset of the performance equation as a metric
 - Expecting improvement of one aspect of a computer to increase performance proportional to the size of improvement

You are going to enhance a machine and there are two types of possible improvements: either (i) make multiply instructions run 4 times faster, or (ii) make memory access instructions run two times faster than before. You repeatedly run a program that takes 100 seconds to execute (on the original machine) and find that of this time 25% is used for multiplication, 50% for memory access instructions, and 25% for other tasks.

- 1. What will the speedup be if you improve both multiplication and memory access?
- 2. Assume the program you run has 10 billions instructions and runs on the machine that has a clock rate of 1GHz. Calculate the CPI for this machine. Assume further that the CPI for multiplication instructions is 20 cycles and the CPI for memory access instructions is 6 cycles. Compute the CPI for all other instructions.
- 3. What is the CPI for the improved machine when improvements on both multiplication and memory access instructions are made?