

# UCR

## CS161 – Design and Architecture of Computer Systems

### Performance Evaluation

UNIVERSITY OF CALIFORNIA, RIVERSIDE

# WHAT IS PERFORMANCE?

# Understanding Performance



- ▶ Algorithm
  - ▶ Determines number of operations executed
- ▶ Programming language, compiler, architecture
  - ▶ Determine number of machine instructions executed per operation
- ▶ Processor and memory system
  - ▶ Determine how fast instructions are executed
- ▶ I/O system (including OS)
  - ▶ Determines how fast I/O operations are executed

# Response Time and Throughput

- ▶ Response time
  - ▶ How long it takes to do a task
- ▶ Throughput
  - ▶ Total work done per unit time
    - ▶ e.g., tasks/transactions/... per hour
- ▶ How are response time and throughput affected by
  - ▶ Replacing the processor with a faster version?
  - ▶ Adding more processors?
- ▶ We'll focus on response time for now...

# Relative Performance

- Define Performance = 1/Execution Time
- “X is  $n$  time faster than Y”

$$\begin{aligned} & \text{Performance}_X / \text{Performance}_Y \\ &= \text{Execution time}_Y / \text{Execution time}_X = n \end{aligned}$$

- Example: time taken to run a program
  - 10s on A, 15s on B
  - $\text{Execution Time}_B / \text{Execution Time}_A$   
 $= 15\text{s} / 10\text{s} = 1.5$
  - So A is 1.5 times faster than B

# Relative Performance

- Define Performance = 1/Execution Time
- “X is  $n$  time faster than Y”

$$\begin{aligned} & \text{Performance}_X / \text{Performance}_Y \\ &= \text{Execution time}_Y / \text{Execution time}_X = n \end{aligned}$$

- Example: time taken to run a program
  - 60s on A, 30s on B
  - $\text{Execution Time}_B / \text{Execution Time}_A = 30\text{s} / 60\text{s} = 0.5$       So A is 0.5 times faster than B
  - or B is 2 times faster than A

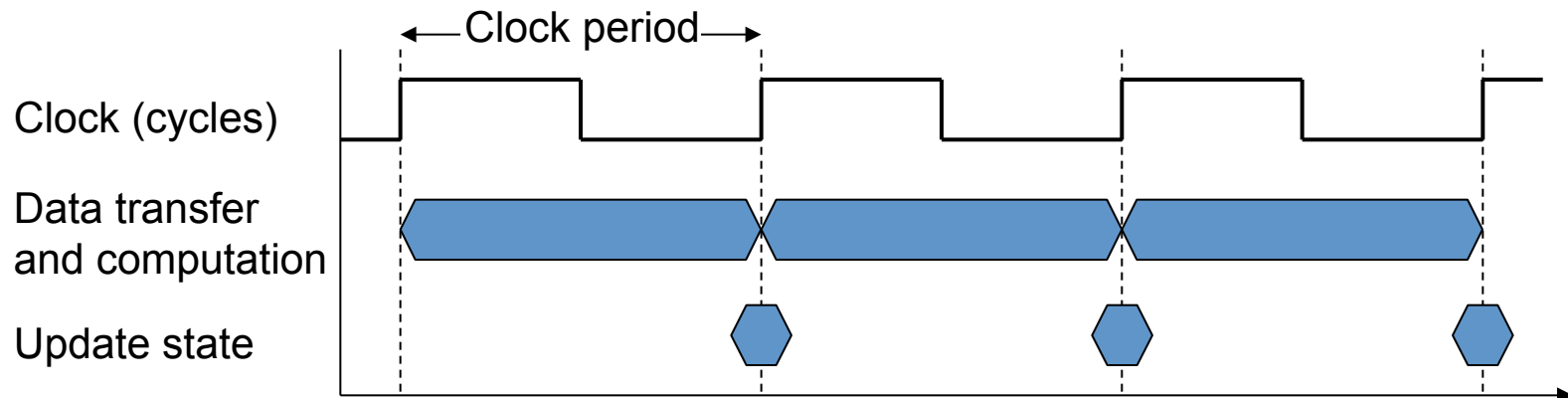
# Measuring Execution Time



- ▶ Elapsed time
  - ▶ Total response time, including all aspects
    - ▶ Processing, I/O, OS overhead, idle time
  - ▶ Determines system performance
- ▶ CPU time
  - ▶ Time spent processing a given job
    - ▶ Discounts I/O time, other jobs' shares
  - ▶ Comprises user CPU time and system CPU time
  - ▶ Different programs are affected differently by CPU and system performance

# CPU Clocking

- Operation of digital hardware governed by a constant-rate clock



- Clock period: duration of a clock cycle
  - e.g.,  $250\text{ps} = 0.25\text{ns} = 250 \times 10^{-12}\text{s}$
- Clock frequency (rate): cycles per second
  - e.g.,  $4.0\text{GHz} = 4000\text{MHz} = 4.0 \times 10^9\text{Hz}$



# CPU Time



$$\begin{aligned}\text{CPU Time} &= \text{CPU Clock Cycles} \times \text{Clock Cycle Time} \\ &= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}\end{aligned}$$

- ▶ Performance improved by
  - ▶ Reducing number of clock cycles
  - ▶ Increasing clock rate
  - ▶ Hardware designer must often trade off clock rate against cycle count

# CPU Time Example

- › Computer A: 2GHz clock, 10s CPU time
- › Designing Computer B
  - › Aim for 6s CPU time
  - › Can do faster clock, but causes  $1.2 \times$  clock cycles
- › How fast must Computer B clock be?

$$\text{Clock Rate}_B = \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6\text{s}}$$

$$\begin{aligned}\text{Clock Cycles}_A &= \text{CPU Time}_A \times \text{Clock Rate}_A \\ &= 10\text{s} \times 2\text{GHz} = 20 \times 10^9\end{aligned}$$

$$\text{Clock Rate}_B = \frac{1.2 \times 20 \times 10^9}{6\text{s}} = \frac{24 \times 10^9}{6\text{s}} = 4\text{GHz}$$

# Instruction Count and CPI

Clock Cycles = Instruction Count  $\times$  Cycles per Instruction

CPU Time = Instruction Count  $\times$  CPI  $\times$  Clock Cycle Time

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

- Instruction Count for a program
  - Determined by program, ISA and compiler
- Average cycles per instruction
  - Determined by CPU hardware
  - If different instructions have different CPI
    - Average CPI affected by instruction mix

# CPI Example

- › Computer A: Cycle Time = 250ps, CPI = 2.0
- › Computer B: Cycle Time = 500ps, CPI = 1.2
- › Same ISA
- › Which is faster, and by how much?

$$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$$

$$= 1 \times 2.0 \times 250\text{ps} = 1 \times 500\text{ps} \leftarrow \text{A is faster...}$$

$$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$$

$$= 1 \times 1.2 \times 500\text{ps} = 1 \times 600\text{ps}$$

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{1 \times 600\text{ps}}{1 \times 500\text{ps}} = 1.2 \leftarrow \text{...by this much}$$

# CPI in More Detail

- If different instruction types take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^n \left( \text{CPI}_i \times \underbrace{\frac{\text{Instruction Count}_i}{\text{Instruction Count}}}_{\text{Relative frequency}} \right)$$

Relative frequency

# CPI Example

- ▶ Alternative compiled code sequences using instructions in type INT, FP, MEM

Type	INT	FP	MEM
CPI for type	1	2	3
IC in Program 1	2	1	2
IC in Program 2	4	1	1

- Program 1: IC = 5
  - Clock Cycles  
 $= 2 \times 1 + 1 \times 2 + 2 \times 3$   
 $= 10$
  - Avg. CPI =  $10/5 = 2.0$

- Program 2: IC = 6
  - Clock Cycles  
 $= 4 \times 1 + 1 \times 2 + 1 \times 3$   
 $= 9$
  - Avg. CPI =  $9/6 = 1.5$

# Performance Summary

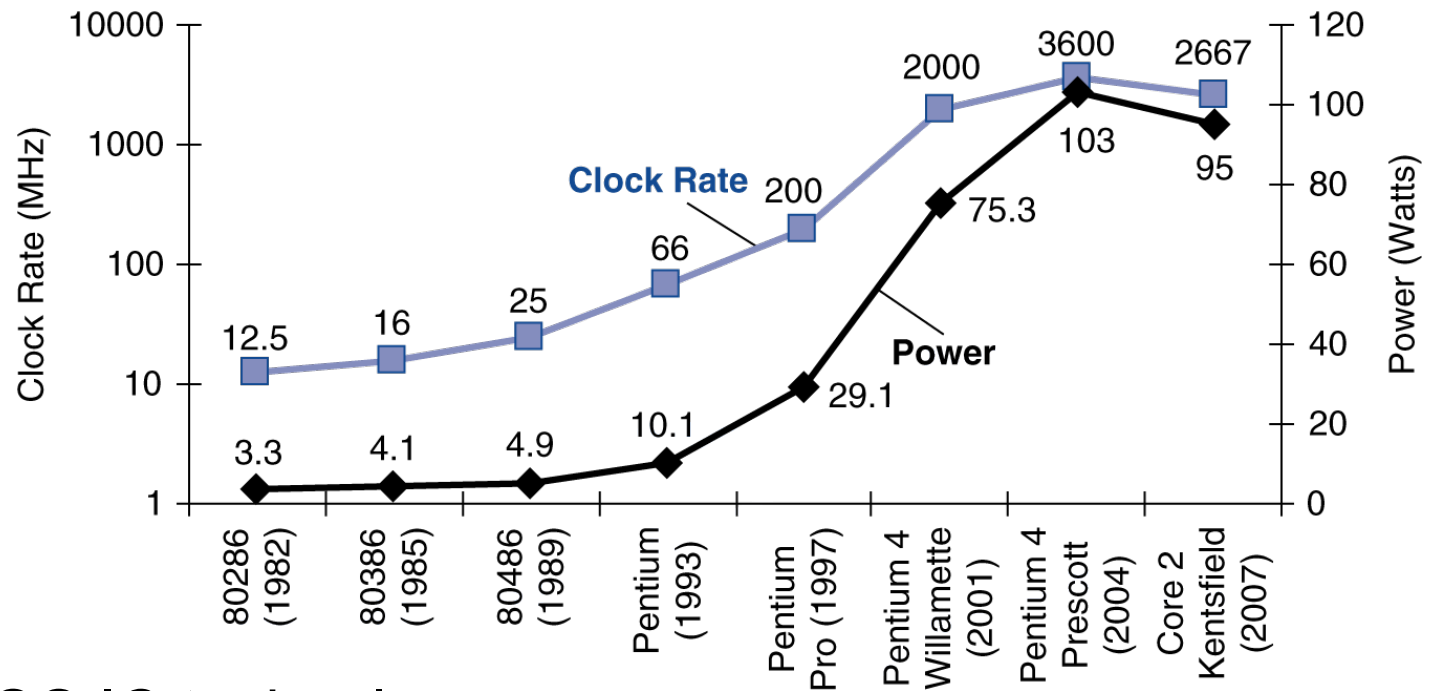


## The BIG Picture

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- Performance depends on
  - Algorithm: affects IC, possibly CPI
  - Programming language: affects IC, CPI
  - Compiler: affects IC, CPI
  - Instruction set architecture: affects IC, CPI,  $T_c$

# Power Trends



- › In CMOS IC technology

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

×30

5V → 1V

×1000



# Reducing Power

- Suppose a new CPU has
  - 85% of capacitive load of old CPU
  - 15% voltage and 15% frequency reduction

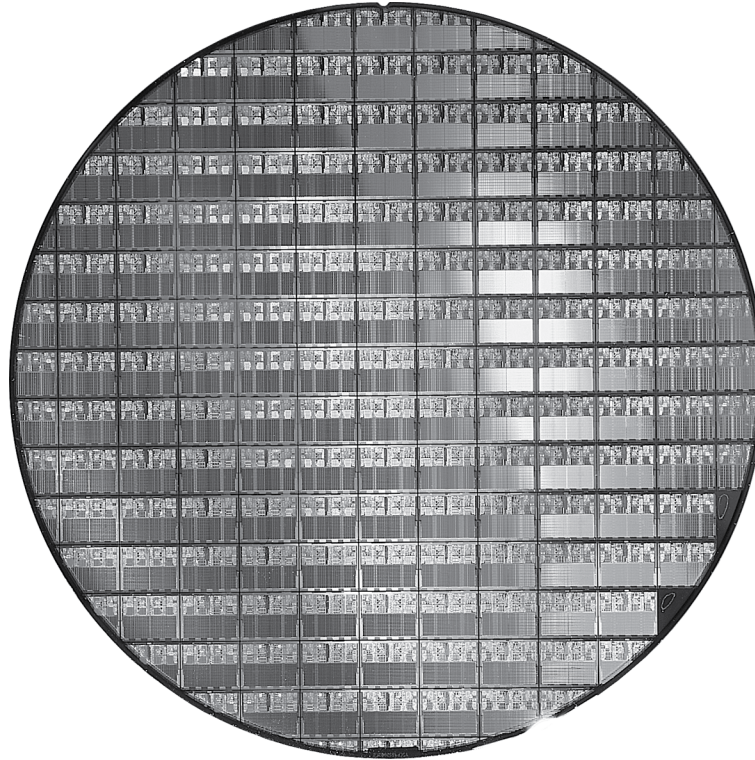
$$\frac{P_{\text{new}}}{P_{\text{old}}} = \frac{C_{\text{old}} \times 0.85 \times (V_{\text{old}} \times 0.85)^2 \times F_{\text{old}} \times 0.85}{C_{\text{old}} \times V_{\text{old}}^2 \times F_{\text{old}}} = 0.85^4 = 0.52$$

- The power wall
  - We can't reduce voltage further
  - We can't remove more heat
- How else can we improve performance?

# Multiprocessors

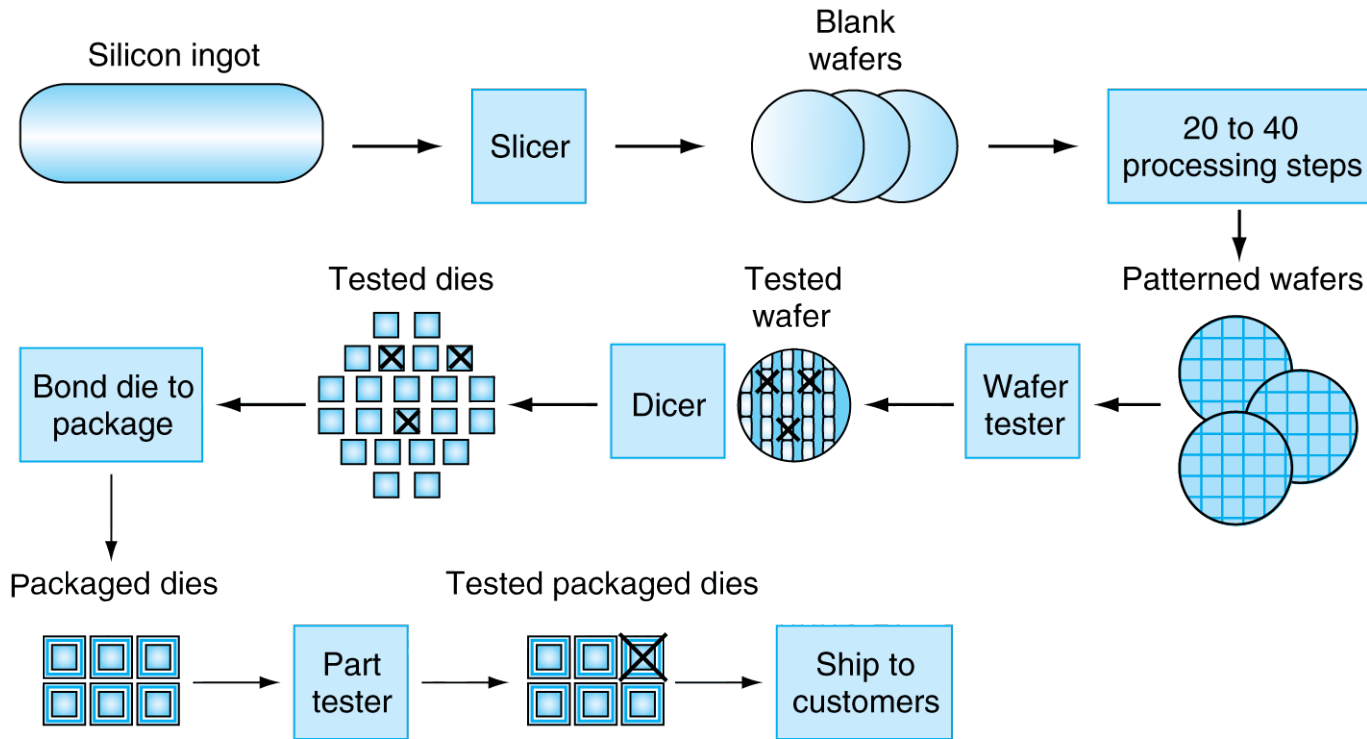
- > Multicore microprocessors
  - > More than one processor per chip
- > Requires explicitly parallel programming
  - > Compare with instruction level parallelism
    - > Hardware executes multiple instructions at once
    - > Hidden from the programmer
  - > Hard to do
    - > Programming for performance
    - > Load balancing
    - > Optimizing communication and synchronization

# AMD Opteron X2 Wafer



- X2: 300mm wafer, 117 chips, 90nm technology
- X4: 45nm technology

# Manufacturing ICs



- › Yield: proportion of working dies per wafer

# Integrated Circuit Cost



$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$

$$\text{Dies per wafer} \approx \text{Wafer area} / \text{Die area}$$

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

- Nonlinear relation to area and defect rate
  - Wafer cost and area are fixed
  - Defect rate determined by manufacturing process
  - Die area determined by architecture and circuit design

# SPEC CPU Benchmark

- ▶ Programs used to measure performance
  - ▶ Supposedly typical of actual workload
- ▶ Standard Performance Evaluation Corp (SPEC)
  - ▶ Develops benchmarks for CPU, I/O, Web, ...
- ▶ SPEC CPU2006
  - ▶ Elapsed time to execute a selection of programs
    - ▶ Negligible I/O, so focuses on CPU performance
  - ▶ Normalize relative to reference machine
  - ▶ Summarize as geometric mean of performance ratios
    - ▶ CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$

# CINT2006 for Opteron X4 2356

Name	Description	IC×10 <sup>9</sup>	CPI	Tc (ns)	Exec time	Ref time	SPECratio
perl	Interpreted string processing	2,118	0.75	0.40	637	9,777	15.3
bzip2	Block-sorting compression	2,389	0.85	0.40	817	9,650	11.8
gcc	GNU C Compiler	1,050	1.72	0.47	24	8,050	11.1
mcf	Combinatorial optimization	336	10.00	0.40	1,345	9,120	6.8
go	Go game (AI)	1,658	1.09	0.40	721	10,490	14.6
hmmer	Search gene sequence	2,783	0.80	0.40	890	9,330	10.5
sjeng	Chess game (AI)	2,176	0.96	0.48	37	12,100	14.5
libquantum	Quantum computer simulation	1,623	1.61	0.40	1,047	20,720	19.8
h264avc	Video compression	3,102	0.80	0.40	993	22,130	22.3
omnetpp	Discrete event simulation	587	2.94	0.40	690	6,250	9.1
astar	Games/path finding	1,082	1.79	0.40	773	7,020	9.1
xalancbmk	XML parsing	1,058	2.70	0.40	1,143	6,900	6.0
Geometric mean							11.7

High cache miss rates



# SPEC Power Benchmark

- ▶ Power consumption of server at different workload levels
  - ▶ Performance: ssj\_ops/sec
  - ▶ Power: Watts (Joules/sec)

$$\text{Overall ssj\_ops per Watt} = \left( \sum_{i=0}^{10} \text{ssj\_ops}_i \right) / \left( \sum_{i=0}^{10} \text{power}_i \right)$$



# SPECpower\_ssj2008 for X4

Target Load %	Performance (ssj_ops/sec)	Average Power (Watts)
100%	231,867	295
90%	211,282	286
80%	185,803	275
70%	163,427	265
60%	140,160	256
50%	118,324	246
40%	92,035	233
30%	70,500	222
20%	47,126	206
10%	23,066	180
0%	0	141
Overall sum	1,283,590	2,605
$\Sigma$ ssj_ops/ $\Sigma$ power		493

# Fallacy: Low Power at Idle

- > Look back at X4 power benchmark
  - > At 100% load: 295W
  - > At 50% load: 246W (83%)
  - > At 10% load: 180W (61%)
- > Google data center
  - > Mostly operates at 10% – 50% load
  - > At 100% load less than 1% of the time
- > Consider designing processors to make power proportional to load

# Pitfall: Amdahl's Law

- ▶ Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{\text{improvement factor}} + T_{\text{unaffected}}$$

- Example: multiply accounts for 80s/100s
  - How much improvement in multiply performance to get 5× overall?

$$20 = \frac{80}{n} + 20 \quad \text{■ Can't be done!}$$

- Corollary: make the common case fast

# Pitfall: MIPS as a Performance Metric

- MIPS: Millions of Instructions Per Second
  - Doesn't account for
    - Differences in ISAs between computers
    - Differences in complexity between instructions

$$\begin{aligned} \text{MIPS} &= \frac{\text{Instruction count}}{\text{Execution time} \times 10^6} \\ &= \frac{\text{Instruction count}}{\frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}} \times 10^6} = \frac{\text{Clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

- CPI varies between programs on a given CPU

# Concluding Remarks

- › Cost/performance is improving
  - › Due to underlying technology development
- › Hierarchical layers of abstraction
  - › In both hardware and software
- › Instruction set architecture
  - › The hardware/software interface
- › Execution time: the best performance measure
- › Power is a limiting factor
  - › Use parallelism to improve performance