# CMPE 110: Computer Architecture

## Week 2
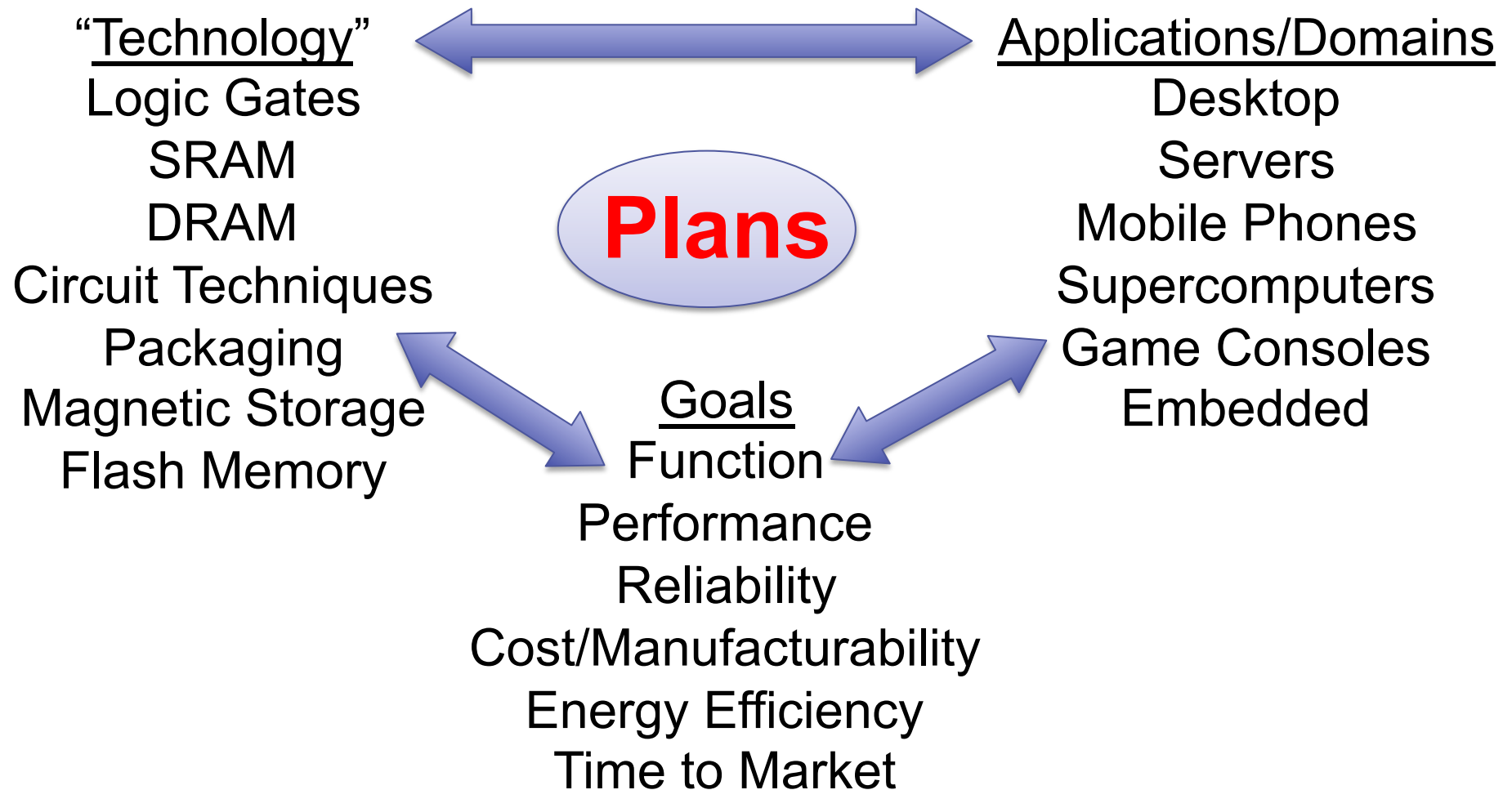## Performance

Jishen Zhao (http://users.soe.ucsc.edu/~jzhao/)

[Adapted in part from Jose Renau, Mary Jane Irwin, Joe Devietti, Onur Mutlu, and others]

# Reminders

- ## TA office hours
  - Rebecca Rashkin: Mon 4-5pm, E2-480
  - Xin Li: TBD
  - Narendra Kumar Govinda Raju: Mon 5:30-6:30pm, BE-312B
  - Aziz (Abdulazaz) Albalawi: Thu 5-6pm, BE-312B

- ## Discussion sessions
  - Section 1: Tue 08:30 - 09:35am in Kresge Clrm 327, TA: Aziz
  - Section 2: Wed 04:00 - 05:05pm in Kresge Clrm 327, TA: Rebecca
  - Section 3: Wed 05:20 - 06:25pm in Kresge Clrm 327, TA: Narendra
  - Section 4: TBA, TA: Xin

# Review: What is computer architecture?

**"Technology"**
Logic Gates
SRAM
DRAM
Circuit Techniques
Packaging
Magnetic Storage
Flash Memory

**Plans**

**Applications/Domains**
Desktop
Servers
Mobile Phones
Supercomputers
Game Consoles
Embedded

Goals
Function
Performance
Reliability
Cost/Manufacturability
Energy Efficiency
Time to Market

# Today: Performance

- Metrics
  - Latency and throughput
  - Speedup
  - Averaging

- CPU Performance

# Performance Metrics

# Performance: Latency vs. Throughput

- **Latency (execution time)**: time to finish a fixed task
- **Throughput (bandwidth)**: number of tasks per unit time
  - Different: exploit parallelism for throughput, not latency



- Choose definition of performance that matches your goals
  - Scientific program? latency.
  - Web server? throughput.

# Examples

- How to measure the performance of moving people for 10 miles round trip?

  - Car: capacity = 5, speed = 60 miles/hour
  - Bus: capacity = 60, speed = 20 miles/hour
  - Latency: how long does ~ take to move one person for 20 miles?
  - Throughput: how many people can ~ move per hour?

  Answer:

  - Latency: **car = 20 min**, bus = 60 min
  - Throughput: car = 15 PPH, **bus = 60 PPH**

# Examples

- Fastest way to send 10TB of data from US to UK?

  FTP, SMB, Rsync / Robocopy, other?

  Used FedEx overnight to deliver the drive

  Even 1 Gbps data transfer takes days!

# Amazon Does This...

| Available Internet Connection | Theoretical Min. Number of Days to Transfer 1TB at 80% Network Utilization | When to Consider AWS Import/Export? |
|---|---|---|
| T1 (1.544Mbps) | 82 days | 100GB or more |
| 10Mbps | 13 days | 600GB or more |
| T3 (44.736Mbps) | 3 days | 2TB or more |
| 100Mbps | 1 to 2 days | 5TB or more |
| 1000Mbps | Less than 1 day | 60TB or more |

**amazon** web services™ **AWS IMPORT/EXPORT CALCULATOR**

Amazon Web Services  »  AWS Import/Export  »  AWS Import/Export Calculator

| | | |
|---|---|---|
| Operation Type | | Import to S3 ▾ |
| Location | AWS Region | US Standard Region ▾ |
| AWS Import/Export Data Load | Total Terabytes to Load | 1 TB |
| | Number of Devices | 1 |
| | Wipe Device After Import | No ▾ |
| Estimated Transfer Speed | Average File Size* | 1 MB |
| | Interface Type | eSATA ▾ |
| | Transfer Speed** | 22.51 MB/sec |

# What we learned

Measuring performance

# **Latency & throughput**

# Comparing Performance - Speedup

- Speedup of A over B
  - X = Latency(B)/Latency(A) (divide by the faster)
  - X = Throughput(A)/Throughput(B) (divide by the slower)
- A is X% faster than B if
  - X = ((Latency(B)/Latency(A)) – 1) * 100
  - X = ((Throughput(A)/Throughput(B)) – 1) * 100
  - Latency(A) = Latency(B) / (1+(X/100))
  - Throughput(A) = Throughput(B) * (1+(X/100))

- Car/bus example
  - Latency?
  - Throughput?
  - See next slide…

# Car/bus example

- Latency: **car = 20 min**, bus = 60 min
- Throughput: car = 15 PPH, **bus = 60 PPH**

Speedup?

- Latency:
    - Speedup of car over bus is 3
    - Car is 200% faster than bus
- Throughput:
    - Speedup of bus over car is 4
    - Bus is 300% faster than car

# Comparing Performance - Speedup

- Program A runs for 200 cycles
- Program B runs for 350 cycles

*What is "cycle"?*
*Execution time * clock frequency*
*i.e., second * Hz*

- Speedup of A over B?
  - Speedup = 350/200 = 1.75
  - As a percentage: (1.75 − 1) * 100 = 75% (Program A runs 75% faster than program B)

- If program C is 50% faster than A, how many cycles does C run for?
  - 133 cycles

# Note

- Speedup of A over B
  - X = Latency(B)/Latency(A)
  - X = Throughput(A)/Throughput(B)

# What if X < 1?

-- means A is slower than B

# Speedup and % Increase and Decrease

- Program A runs for 200 cycles
- Program B runs for 350 cycles
- Percent increase and decrease are <span style="color:red">not the same</span>.
  - % increase of cycles: ((350 – 200)/200) * 100 = 75%
  - % decrease of cycles: ((350 - 200)/350) * 100 = 42.3%

# What we learned

Comparing performance

## **Speedup**

Performance metrics

## **Latency, throughput, speedup**

# Averaging performance

# Mean (Average) Performance Numbers

- **Arithmetic**: $(1/N) * \Sigma_{P=1..N}$ Latency(P)
  - For units that are proportional to time (e.g., latency)

- **Harmonic**: $N / \Sigma_{P=1..N}$ 1/Throughput(P)
  - For units that are inversely proportional to time (e.g., throughput)

- You can add latencies, but not throughputs
  - Latency(P1+P2, A) = Latency(P1, A) + Latency(P2, A)
  - Throughput(P1+P2, A) != Throughput(P1, A) + Throughput(P2, A)

- **Geometric**: $\sqrt[N]{\Pi_{P=1..N}}$ Speedup(P)
  - For unitless quantities (e.g., speedup ratios)

# For Example...

1 mile @ 30 miles/hour + 1 mile @ 90 miles/hour

- You drive two miles
  - 30 miles per hour for the first mile
  - 90 miles per hour for the second mile

- Question: what was your average speed?
  - Hint: the answer is not 60 miles per hour
  - Why?

# Answer: 45 miles/hour

- You drive two miles
  - 30 miles per hour for the first mile
  - 90 miles per hour for the second mile

- Question: what was your average speed?
  - Hint: the answer is not 60 miles per hour
  - 0.03333 hours per mile for 1 mile
  - 0.01111 hours per mile for 1 mile
  - 0.04444 hours for 2 miles
  - = 45 miles per hour
  - != (30 + 90) / 2

# What we learned

Averaging performance

# Arithmetic mean for latency
# Harmonic mean for throughput
# Geometric mean for speedup

# CPU Performance

How to evaluate

Latency, throughput, and speedup

# CPU Performance Equation

- Latency = seconds / program =
  - (insns / program) * (cycles / insn) * (seconds / cycle)

  - **Insns / program**: insn count
    - Impacted by program, compiler, ISA
  - **Cycles / insn**: **CPI**
    - Impacted by program, compiler, ISA, **micro-arch**
  - **Seconds / cycle**: **clock period (Hz)**
    - Impacted by micro-arch, technology

- For low latency (better performance) minimize all three
  - Difficult: often pull against one another
  - Example we have seen: RISC vs. CISC ISAs
    - ± RISC: low CPI/clock period, high insn count
    - ± CISC: low insn count, high CPI/clock period