

**CS10**  
**The Beauty and Joy of Computing**

**Lecture #8 : Concurrency**

**2010-09-29**

UC Berkeley EECS  
 Lecturer SOE  
 Dan Garcia

**AT A RECENT FIRESIDE CHAT...**

Nvidia CEO Jen-Hsun Huang said:  
 "whatever capability you think you have today, it's nothing compared to what you're going to have in a couple of years ... due to supercomputers in the cloud". Now, cloud computing does what you could do on your PC. Imagine 40,000 times that!




[www.theregister.co.uk/2010/09/24/huang\\_muses\\_at\\_gtc/](http://www.theregister.co.uk/2010/09/24/huang_muses_at_gtc/)

**Happy Confucius Day!**

*Knowledge is recognizing what you know and what you don't.*

*I hear and I forget.  
 I see and I remember.  
 I do and I understand.*




© Garcia, Fall 2010  
 UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency | 3

**Concurrency & Parallelism, 10 mi up...**

**Intra-computer**      Inter-computer


- Today's lecture
- Multiple computing "helpers" are cores within one machine
- Aka "multi-core"
  - Although GPU parallelism is also "intra-computer"

- Week 12's lectures
- Multiple computing "helpers" are different machines
- Aka "distributed computing"
  - Grid & cluster computing



© Garcia, Fall 2010  
 UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency | 3

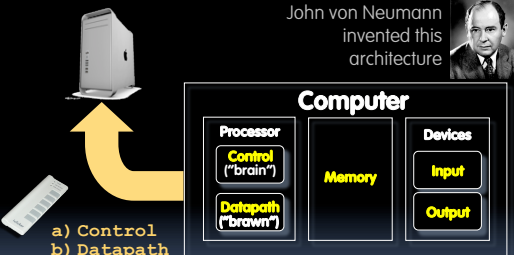
**Anatomy: 5 components of any Computer**



© Garcia, Fall 2010  
 UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency | 4

**Anatomy: 5 components of any Computer**

John von Neumann invented this architecture



**Computer**


- Processor
  - Control ("brain")
  - Datapath ("brawn")
- Memory
- Devices
  - Input
  - Output

a) Control  
 b) Datapath  
 c) Memory  
 d) Input  
 e) Output

**What causes the most headaches for SW and HW designers with multi-core computing?**

© Garcia, Fall 2010  
 UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency | 5

**But what is INSIDE a Processor?**



© Garcia, Fall 2010  
 UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency | 6

## But what is INSIDE a Processor?

- Primarily Crystalline Silicon
- 1 mm – 25 mm on a side
- 2009 "feature size" (aka process) ~ 45 nm =  $45 \times 10^{-9}$  m (then 32, 22, and 16 [by yr 2013])
- 100 - 1000M transistors
- 3 - 10 conductive layers
- "CMOS" (complementary metal oxide semiconductor) - most common
- Package provides:
  - spreading of chip-level signal paths to board-level
  - heat dissipation.
  - Ceramic or plastic with gold wires.

UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (7)

en.wikipedia.org/wiki/Moore's\_Law

## Moore's Law

Predicts: 2X Transistors / chip every 2 years

What is this "curve"?

- Constant
- Linear
- Quadratic
- Cubic
- Exponential

Gordon Moore  
Intel Co-founder  
B.S. Cal 1950!

UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (8)

## Moore's Law and related curves

Transistors (Thousands)

Single-Thread Performance (SpecINT)

Frequency (MHz)

Typical Power (Watts)

Data partially collected by M. Horowitz, F. Labonte, O. Shacham, K. Oukutun, L. Hammond

UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (9)

## Moore's Law and related curves

Transistors (Thousands)

Parallel App Performance

Single-Thread Performance (SpecINT)

Frequency (MHz)

Typical Power (Watts)

Number of Cores

Data partially collected by M. Horowitz, F. Labonte, O. Shacham, K. Oukutun, L. Hammond

UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (10)

## Power Density Prediction circa 2000

Power Density (W/cm<sup>2</sup>)

Year

Source: S. Borkar (Intel)

UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (11)

## Going Multi-core Helps Energy Efficiency

- Power of typical integrated circuit ~  $C V^2 f$ 
  - C = Capacitance, how well it "stores" a charge
  - V = Voltage
  - $f$  = frequency. I.e., how fast clock is (e.g., 3 GHz)

In the same process technology...

Cache	Cache
Core	Core
Core	Core


Voltage = 1      Voltage = -15%  
 Freq = 1        Freq = -15%  
 Area = 1        Area = 2  
 Power = 1       Power = 1  
 Perf = 1        Perf = ~1.8

Activity Monitor (on the lab Macs) shows how active your cores are

William Holt, HOT Chips 2005

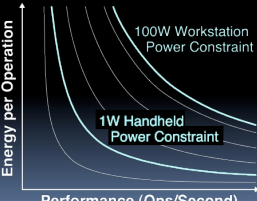
UC Berkeley CS10 "The Beauty and Joy of Computing" | Concurrency (12)

## Energy & Power Considerations



Power =  $\frac{\text{Energy}}{\text{Second}} = \frac{\text{Energy}}{\text{Op}} \times \frac{\text{Ops}}{\text{Second}}$

Power	Energy
Chip Packaging	Battery Life
Chip Cooling	Electricity Bill
System Noise	Mobile Device
Case Temperature	Weight
Data-Center Air Conditioning	



Courtesy: Chris Batten

UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (10)

view.eecs.berkeley.edu

## Parallelism again? What's different this time?


*"This shift toward increasing parallelism is not a triumphant stride forward based on breakthroughs in novel software and architectures for parallelism; instead, this plunge into parallelism is actually a retreat from even greater challenges that thwart efficient silicon implementation of traditional uniprocessor architectures."*

– Berkeley View, December 2006

- HW/SW Industry bet its future that breakthroughs will appear before it's too late

UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (14)

## Background: Threads

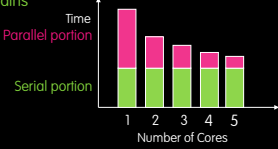
- A **Thread** stands for "thread of execution", is a single stream of instructions
  - A program / process can **split**, or **fork** itself into separate threads, which can (in theory) execute simultaneously.
  - An easy way to describe/think about parallelism
- A single CPU can execute many threads by **Time Division Multiplexing**

- Multithreading** is running multiple threads through the same hardware

UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (18)

en.wikipedia.org/wiki/Amdahl's\_law

## Speedup Issues : Amdahl's Law

- Applications can almost **never** be completely parallelized; some serial code remains



- $s$  is serial fraction of program,  $P$  is # of cores (was processors)
- Amdahl's law:**  

$$\text{Speedup}(P) = \text{Time}(1) / \text{Time}(P)$$

$$\leq 1 / (s + [(1-s) / P]), \text{ and as } P \rightarrow \infty$$

$$\leq 1 / s$$
- Even if the parallel portion of your application speeds up perfectly, your performance may be limited by the sequential portion

UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (14)


## Speedup Issues : Overhead

- Even assuming no sequential portion, there's...
  - Time to think how to divide the problem up
  - Time to hand out small "work units" to workers
  - All workers may not work equally fast
  - Some workers may fail
  - There may be contention for shared resources
  - Workers could overwriting each others' answers
  - You may have to wait until the last worker returns to proceed (the slowest / weakest link problem)
  - There's time to put the data back together in a way that looks as if it were done by one

UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (17)

## Life in a multi-core world...

- This "sea change" to multi-core parallelism means that the computing community has to rethink:
  - Languages
  - Architectures
  - Algorithms
  - Data Structures
  - All of the above



UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (16)

en.wikipedia.org/wiki/Concurrent\_computing

## But parallel programming is hard!

- What if two people were calling withdraw at the same time?
  - E.g., balance=100 and two withdraw 75 each
  - Can anyone see what the problem *could* be?
  - This is a **race condition**
- In most languages, this is a problem.
  - In Scratch, the system doesn't let two of these run at once.

```
withdraw amount
if balance > amount
  set balance to balance - amount
  report true
report false
```

en.wikipedia.org/wiki/Deadlock

## Another concurrency problem ... **deadlock!**

- Two people need to draw a graph but there is only one pencil and one ruler.
  - One grabs the pencil
  - One grabs the ruler
  - Neither release what they hold, waiting for the other to release
- **Livelock** also possible
  - Movement, no progress
  - Dan and Luke demo



## Summary

- "Sea change" of computing because of inability to cool CPUs means we're now in multi-core world
- This brave new world offers lots of potential for innovation by computing professionals, but challenges persist

