



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/

CS10 The Beauty and Joy of Computing

Lecture #8 : Concurrency

2010-09-29



Concurrency & Parallelism, 10 mi up...

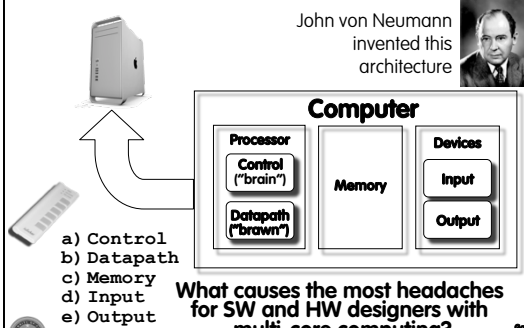
- | | |
|---|--|
| <p>Intra-computer</p> <ul style="list-style-type: none"> Today's lecture Multiple computing "helpers" are cores <u>within one machine</u> Aka "multi-core" <ul style="list-style-type: none"> Although GPU parallism is also "intra-computer" | <p>Inter-computer</p> <ul style="list-style-type: none"> Week 12's lectures Multiple computing "helpers" are <u>different machines</u> Aka "distributed computing" <ul style="list-style-type: none"> Grid & cluster computing |
|---|--|



UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (2)



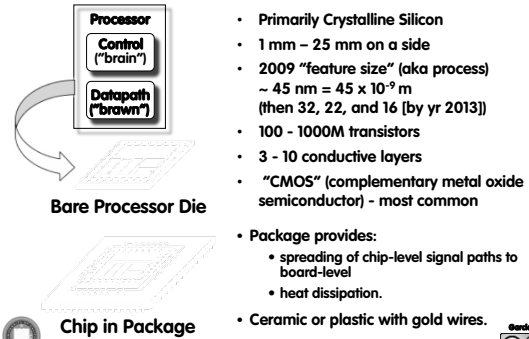
Anatomy: 5 components of any Computer



UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (4)



But what is INSIDE a Processor?



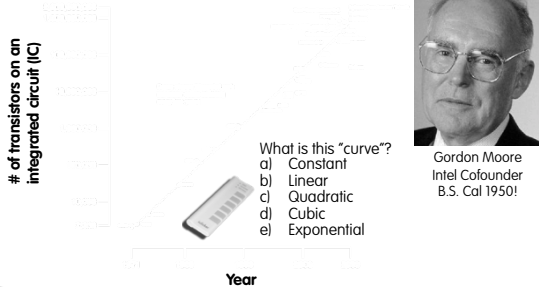
UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (6)



Moore's Law

en.wikipedia.org/wiki/Moore's_law

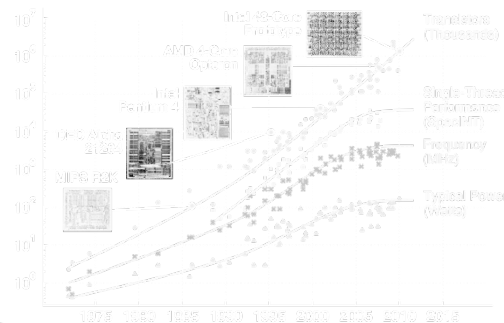
Predicts: 2X Transistors / chip every 2 years



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



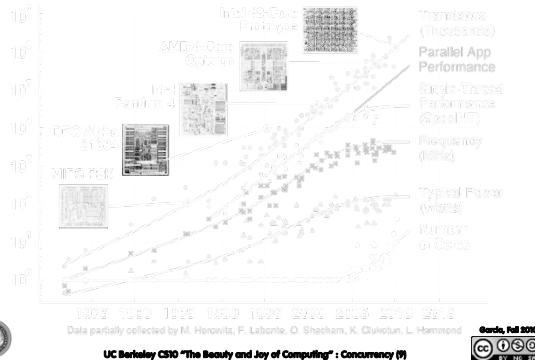
Moore's Law and related curves



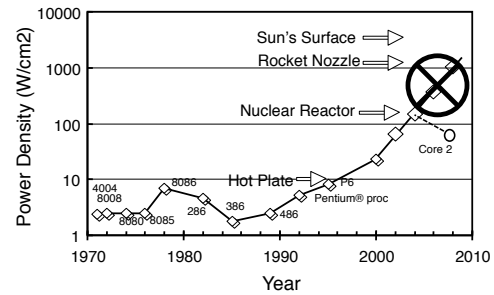
UC Berkeley CS10 "The Beauty and Joy of Computing" : Concurrency (8)



Moore's Law and related curves

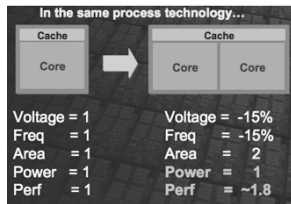


Power Density Prediction circa 2000



Going Multi-core Helps Energy Efficiency

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

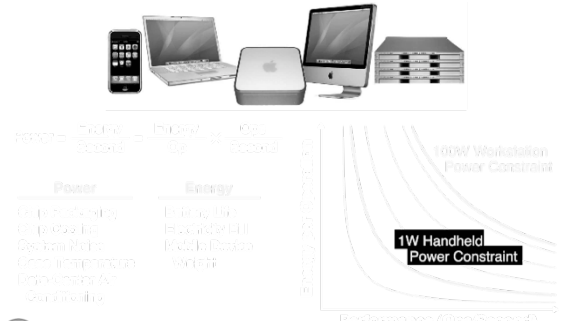


Activity Monitor (on the lab Macs) shows how active your cores are

William Holt, HOT Chips 2005

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

Energy & Power Considerations



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

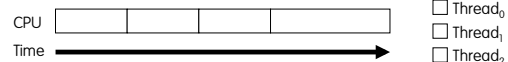
view.eecs.berkeley.edu

UC Berkeley CS10 "The Beauty and Joy of Computing": Concurrency (13)

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**



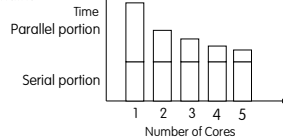
- Multithreading** is running multiple threads through the same hardware

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

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

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



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

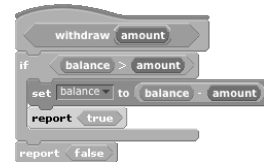


en.wikipedia.org/wiki/Concurrent_computing

But Concurrent 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.

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 also challenges

