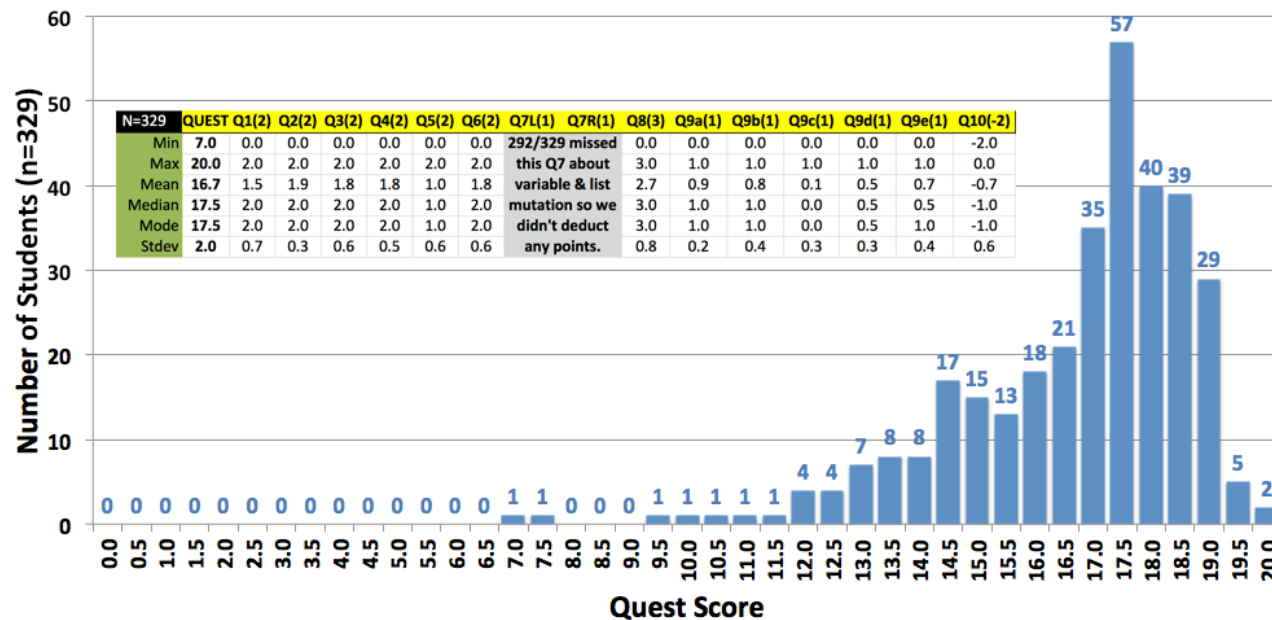# The Beauty and Joy of Computing

## Lecture #8
## Concurrency

UC Berkeley EECS
Lecturer
Gerald Friedland

bjc

### 2013Fa UC Berkeley CS10 Quest Histogram
### Mean = 16.7, Median = 17.5, StDev = 2.0

| N=329 | QUEST | Q1(2) | Q2(2) | Q3(2) | Q4(2) | Q5(2) | Q6(2) | Q7L(1) | Q7R(1) | Q8(3) | Q9a(1) | Q9b(1) | Q9c(1) | Q9d(1) | Q9e(1) | Q10(-2) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Min | 7.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 292/329 missed | | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | -2.0 |
| Max | 20.0 | 2.0 | 2.0 | 2.0 | 2.0 | 2.0 | 2.0 | this Q7 about | | 3.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.0 |
| Mean | 16.7 | 1.5 | 1.9 | 1.8 | 1.8 | 1.0 | 1.8 | variable & list | | 2.7 | 0.9 | 0.8 | 0.1 | 0.5 | 0.7 | -0.7 |
| Median | 17.5 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 2.0 | mutation so we | | 3.0 | 1.0 | 1.0 | 0.0 | 0.5 | 0.5 | -1.0 |
| Mode | 17.5 | 2.0 | 2.0 | 2.0 | 2.0 | 1.0 | 2.0 | didn't deduct | | 3.0 | 1.0 | 1.0 | 0.0 | 0.5 | 1.0 | -1.0 |
| Stdev | 2.0 | 0.7 | 0.3 | 0.6 | 0.5 | 0.6 | 0.6 | any points. | | 0.8 | 0.2 | 0.4 | 0.3 | 0.3 | 0.4 | 0.6 |

# Concurrency: A Definition

Concurrency: A property of computer systems in which several **computations** are **executing** simultaneously, and potentially interacting with each other.

# Concurrency is Everywhere!

Examples:

- Mouse cursor movement while SNAP! calculates.

- Screen clock advances while typing in a text.

- Busy cursor spins while browser connects to server, waiting for response

Question:  Is this real concurrency?

# Concurrency & Parallelism

## Intra-computer

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

## Inter-computer

- **Future lecture**
- **Multiple computing "helpers" are different machines**
- **Aka "distributed computing"**
  - Grid & cluster computing

# 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

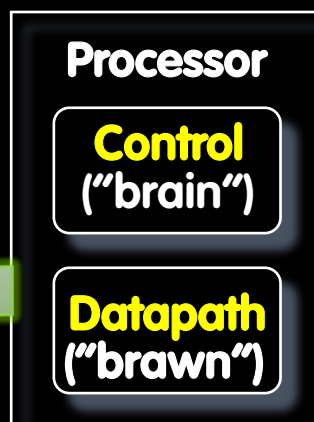**What causes the most headaches for SW and HW designers with multi-core computing?**
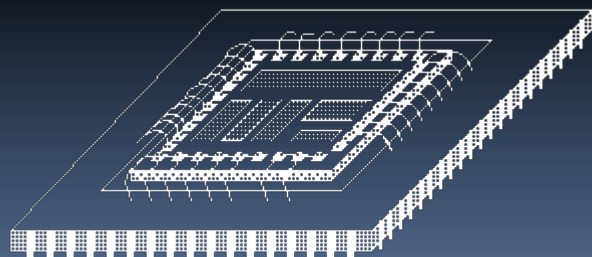
# But what is INSIDE a Processor?

**Processor**

**Control** ("brain")

**Datapath** ("brawn")

Friedland

# But what is INSIDE a Processor?

**Processor**

**Control** ("brain")

**Datapath** ("brawn")

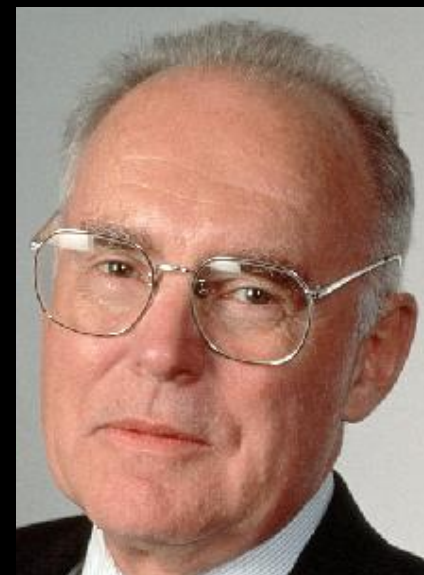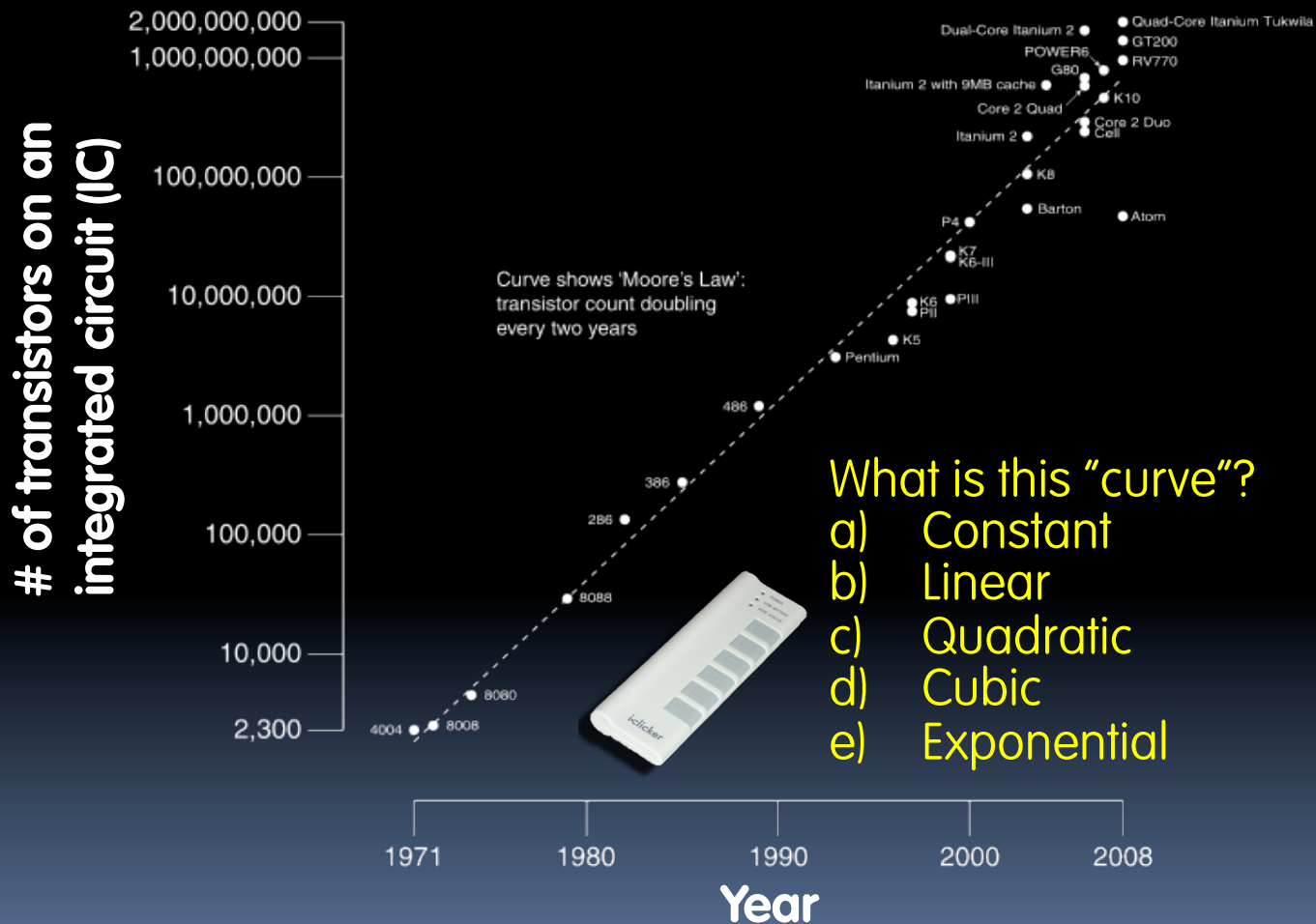**Bare Processor Die**

**Chip in Package**

- Primarily Crystalline Silicon

- 1 mm – 25 mm on a side

- 2009 "feature size" (aka process)
  $\sim 45$ 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.

# Moore's Law

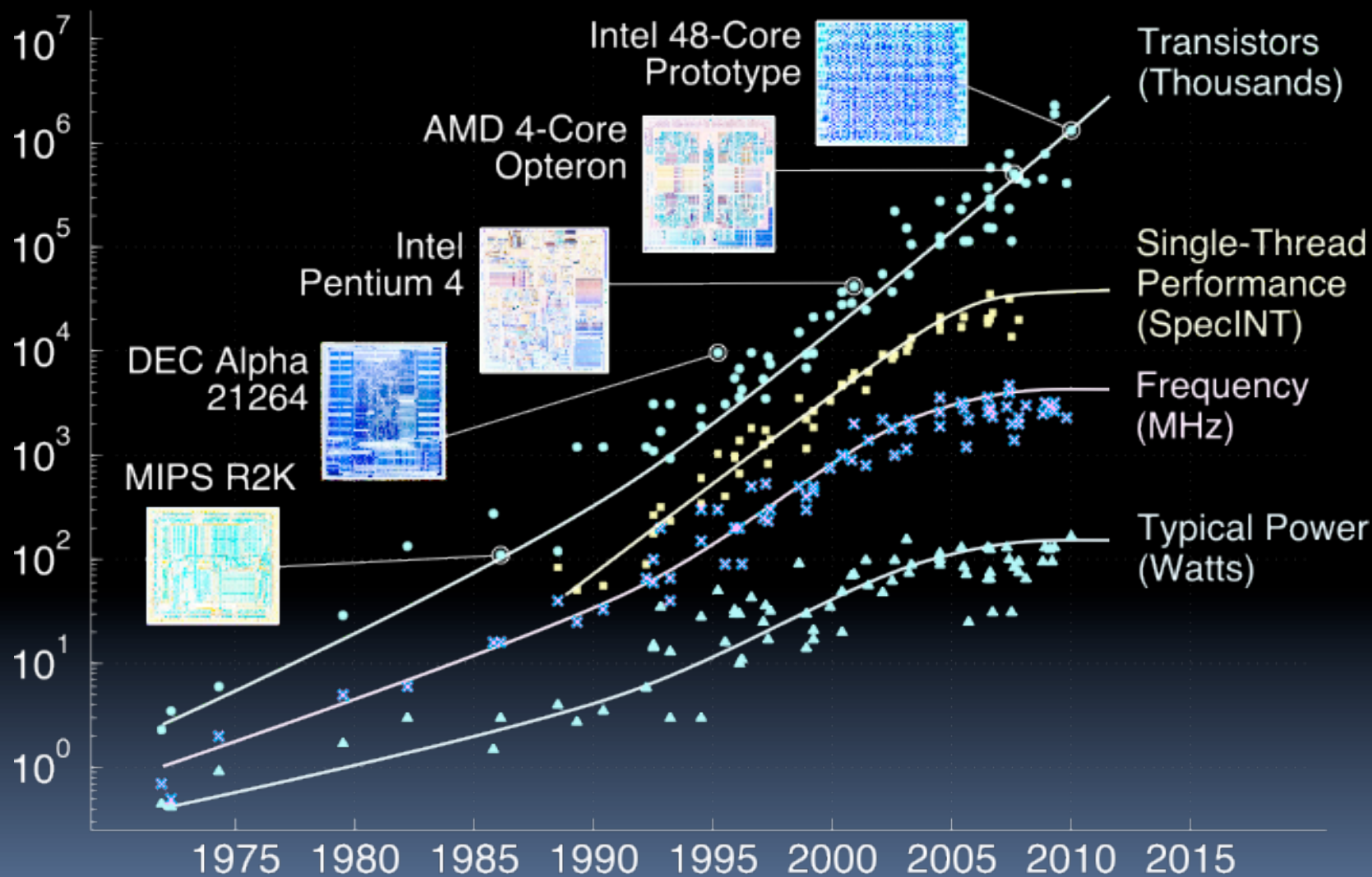## Predicts: 2X Transistors / chip every 2 years



Curve shows 'Moore's Law': transistor count doubling every two years

**What is this "curve"?**
a) Constant
b) Linear
c) Quadratic
d) Cubic
e) Exponential

Gordon Moore
Intel Cofounder
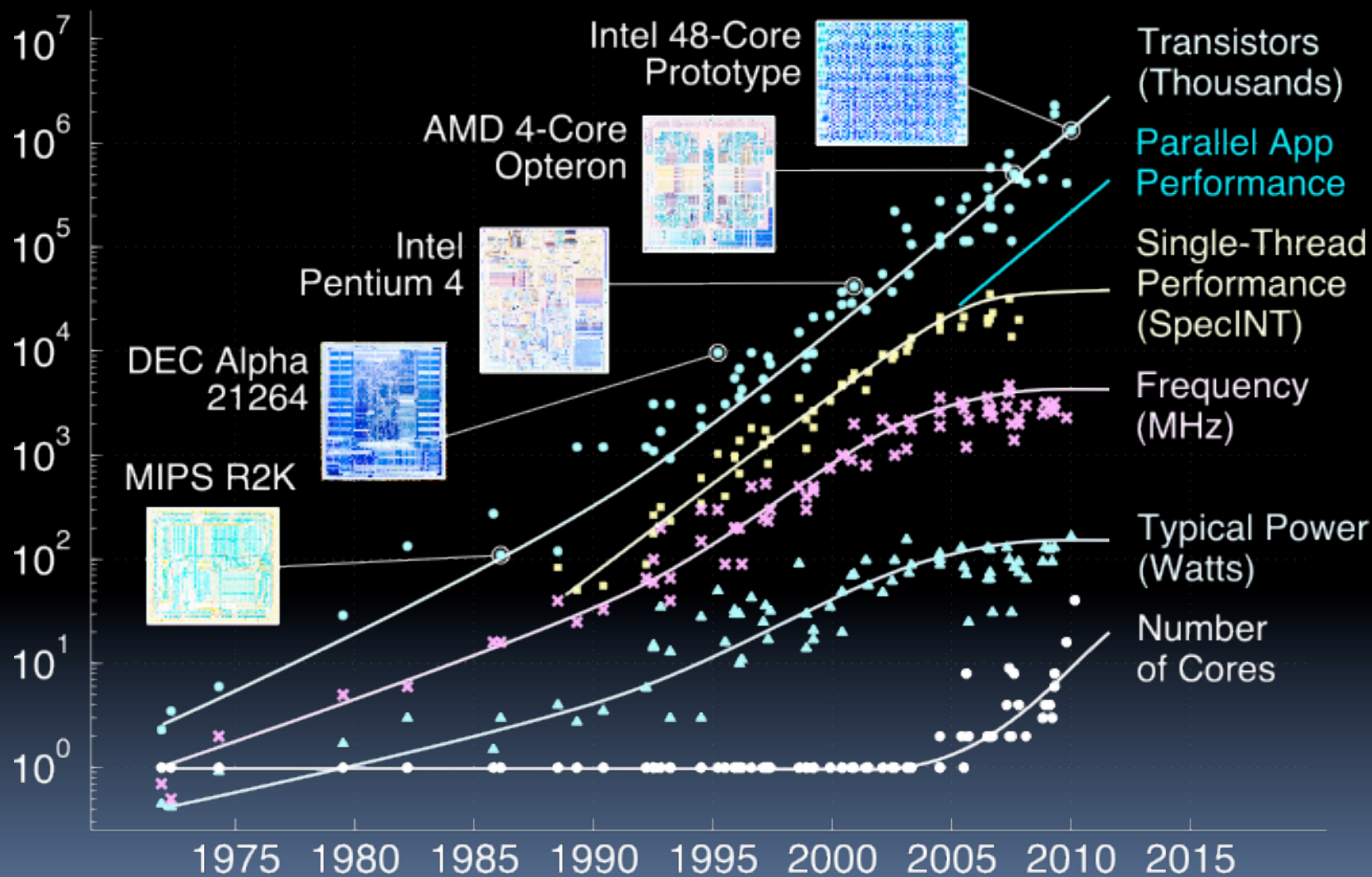B.S. Cal 1950!

Friedland

# Moore's Law and related curves



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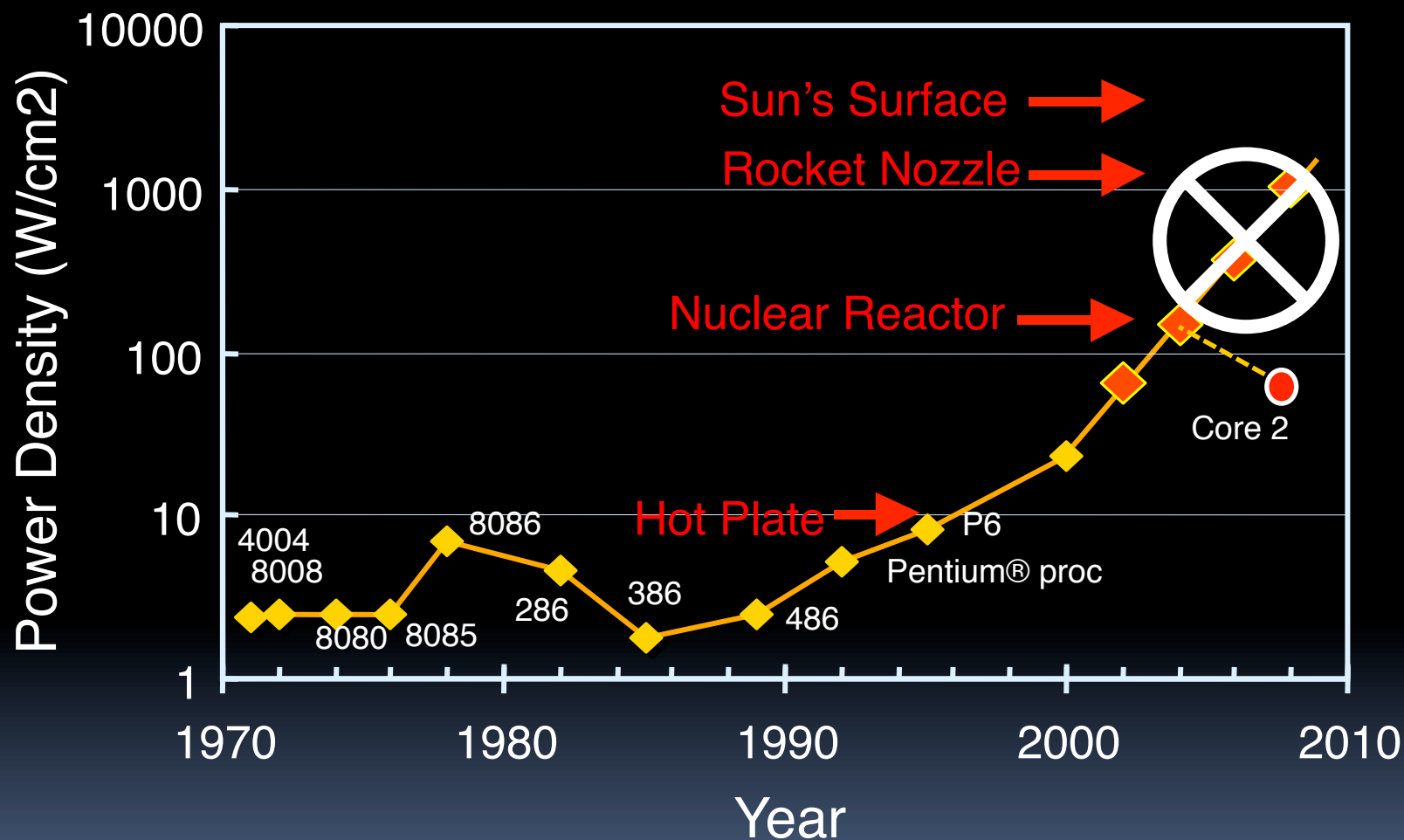
Friedland

# Moore's Law and related curves

Friedland

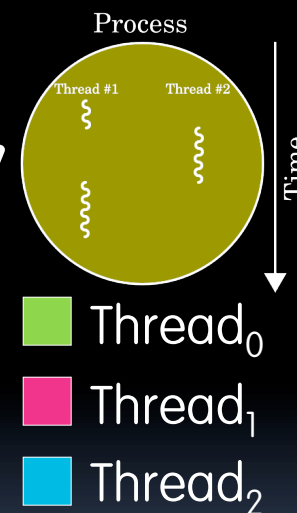# Power Density Prediction circa 2000



Source: S. Borkar (Intel)

# 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 Multipexing***
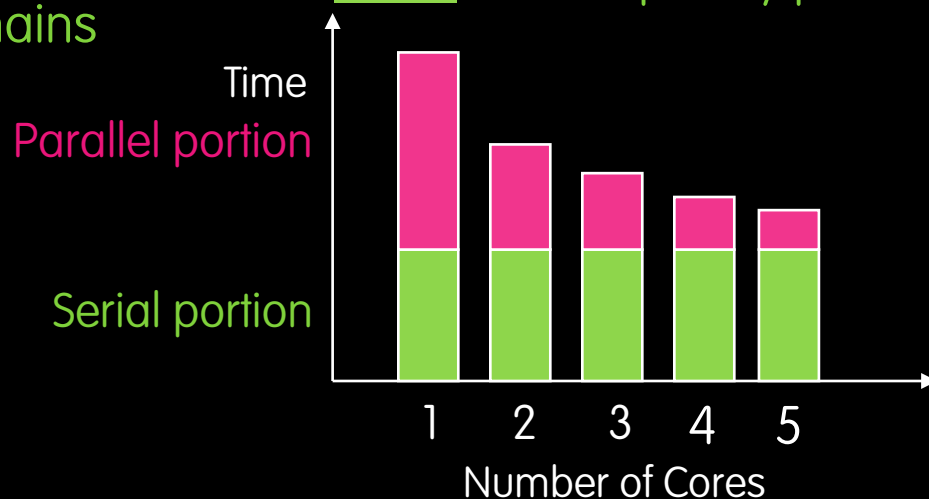
Process

Thread #1    Thread #2

Time

CPU

Time

$Thread_0$

$Thread_1$

$Thread_2$

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

# Speedup Issues : Amdahl's Law

- Applications can almost <u>never</u> be completely parallelized; some serial code remains

Time
Parallel portion

Serial portion

1   2   3   4   5
Number of Cores

- s is serial fraction of program, P is # of cores (was processors)

- **Amdahl's law:**

Speedup(P) = Time(1) / Time(P)

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

$$\leq 1 / s$$

- Even if the parallel portion of your application speeds up perfectly, your performance may be limited by the sequential portion

Friedland

# 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
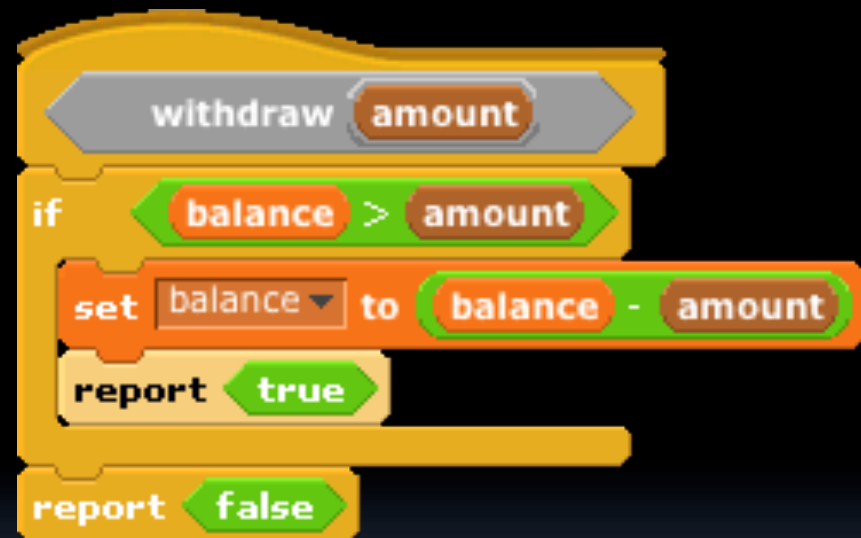
Friedland

# Life in a multi-core world…

- **This "sea change" to multi-core parallelism means that the computing community has to rethink:**
  a) Languages
  b) Architectures
  c) Algorithms
  d) Data Structures
  e) All of the above

# 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.

# 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

# 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