# COMPUTER ORGANIZATION AND DESIGN

## The Hardware/Software Interface

# Chapter 1

## Introduction: Computer Abstractions and Technology

# The Computer Revolution

- Progress in computer technology
  - Underpinned by domain-specific accelerators
- Makes novel applications feasible
  - Computers in automobiles
  - Cell phones
  - Human genome project
  - World Wide Web
  - Search Engines
- Computers are pervasive

# The Computer Revolution

- Progress in computer technology
  - Underpinned by domain-specific accelerators
- Makes novel applications feasible
  - Computers in automobiles
  - Cell phones
  - Human genome project
  - World Wide Web
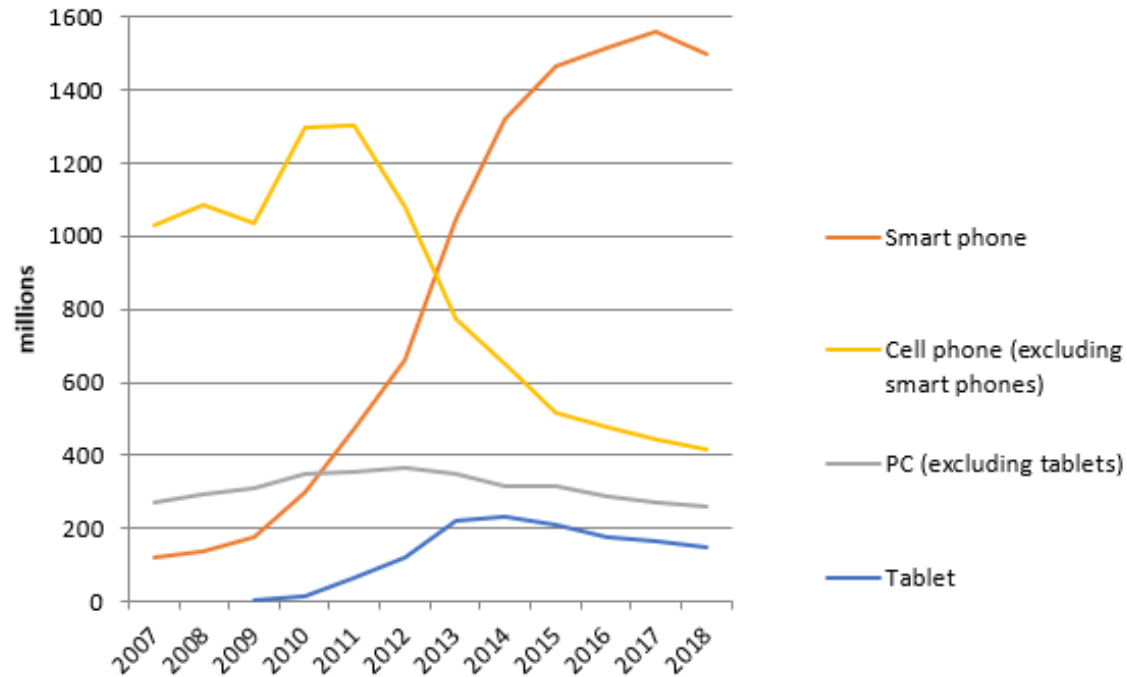  - Search Engines
- Computers are pervasive

# Classes of Computers

- Personal computers
  - General purpose, variety of software
  - Subject to cost/performance tradeoff

- Server computers
  - Network based
  - High capacity, performance, reliability
  - Range from small servers to building sized

# Classes of Computers

- Supercomputers
    - Type of server
    - High-end scientific and engineering calculations
    - Highest capability but represent a small fraction of the overall computer market

- Embedded computers
    - Hidden as components of systems
    - Stringent power/performance/cost constraints

# The PostPC Era



**The number manufactured per year of tablets and smart phones, which reflect the PostPC era, versus personal computers and traditional cell phones.** Smart phones represent the recent growth in the cell phone industry, and they passed PCs in 2011. Tablets are the fastest growing category, nearly doubling between 2011 and 2012. Recent PCs and traditional cell phone categories are relatively flat or declining.

# The PostPC Era

- Personal Mobile Device (PMD)
    - Battery operated
    - Connects to the Internet
    - Hundreds of dollars
    - Smart phones, tablets, electronic glasses
- Cloud computing
    - Warehouse Scale Computers (WSC)
    - Software as a Service (SaaS)
    - Portion of software run on a PMD and a portion run in the Cloud
    - Amazon and Google

# What You Will Learn

- How programs are translated into the machine language
    - And how the hardware executes them
- The hardware/software interface
- What determines program performance
    - And how it can be improved
- How hardware designers improve performance
- What is parallel processing

# 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

# Seven Great Ideas

- Use ***abstraction*** to simplify design

- Make the ***common case fast***

- Performance *via* **parallelism**

- Performance *via* **pipelining**

- Performance *via* **prediction**

- ***Hierarchy*** of memories

- ***Dependability*** *via* redundancy

ABSTRACTION

COMMON CASE FAST

PARALLELISM

PIPELINING

PREDICTION
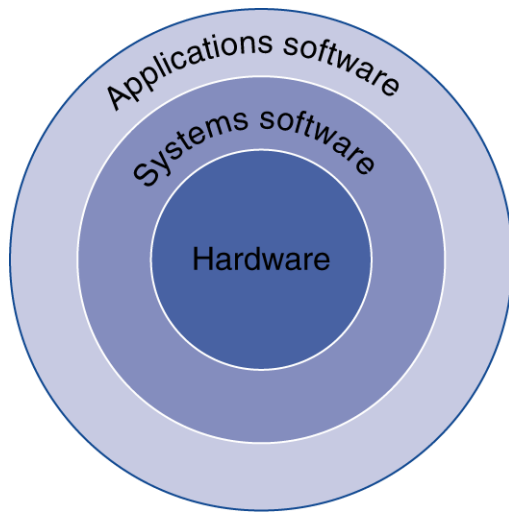
HIERARCHY

DEPENDABILITY

# Below Your Program

- ## Application software
  - ### Written in high-level language
- ## System software
  - ### Compiler: translates HLL code to machine code
  - ### Operating System: service code
    - Handling input/output
    - Managing memory and storage
    - Scheduling tasks & sharing resources
- ## Hardware
  - ### Processor, memory, I/O controllers

Applications software

Systems software

Hardware

# Levels of Program Code

- **High-level language**
  - Level of abstraction closer to problem domain
  - Provides for productivity and portability

- **Assembly language**
  - Textual representation of instructions

- **Hardware representation**
  - Binary digits (bits)
  - Encoded instructions and data

High-level
language
program
(in C)

```
swap(int v[], int k)
{int temp;
    temp = v[k];
    v[k] = v[k+1];
    v[k+1] = temp;
}
```

Compiler

Assembly
language
program
(for MIPS)

```
swap:
    muli $2, $5,4
    add  $2, $4,$2
    lw   $15, 0($2)
    lw   $16, 4($2)
    sw   $16, 0($2)
    sw   $15, 4($2)
    jr   $31
```

Assembler

Binary machine
language
program
(for MIPS)

```
00000000101000010000000000011000
00000000000110000000011000000100001
10001100011000100000000000000000
10001100111100100000000000000100
10101100111100100000000000000000
10101100011000100000000000000100
00000011111000000000000000001000
```
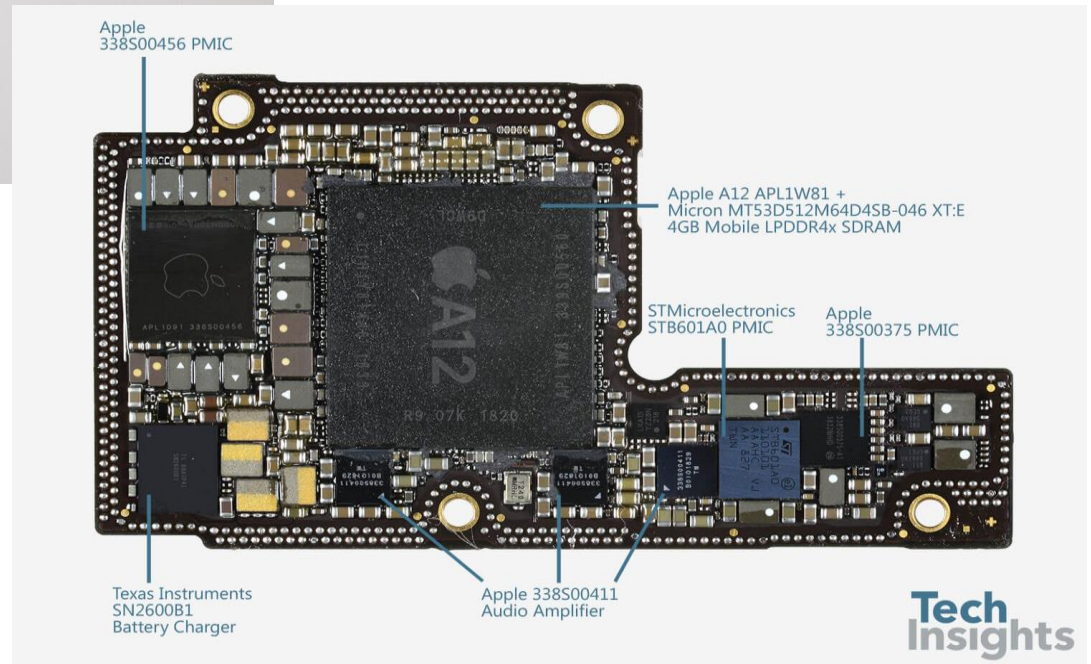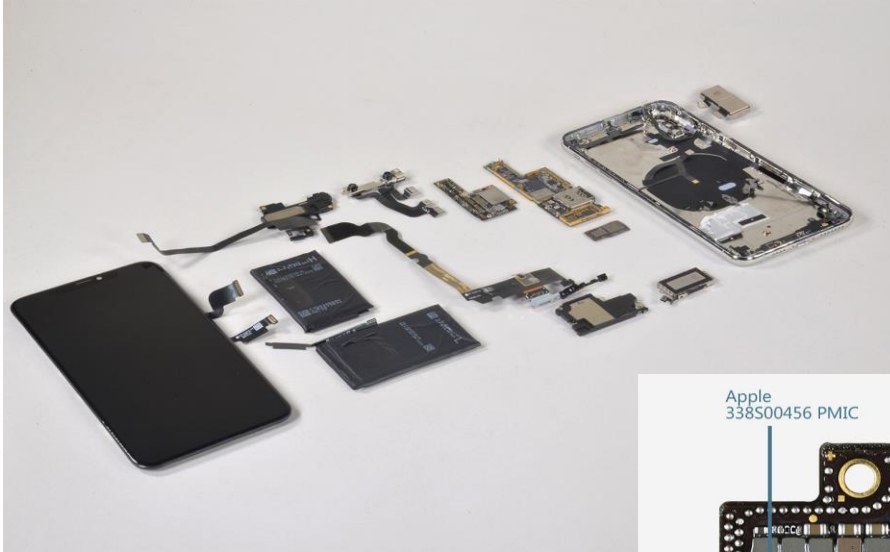
# Components of a Computer

## The BIG Picture



- Same components for all kinds of computer
  - Desktop, server, embedded
- Input/output includes
  - User-interface devices
    - Display, keyboard, mouse
  - Storage devices
    - Hard disk, CD/DVD, flash
  - Network adapters
    - For communicating with other computers

# Opening the Box

# Inside the Processor (CPU)

- Datapath: performs operations on data

- Control: sequences datapath, memory, ...

- Cache memory

  - Small fast SRAM memory for immediate access to data

# Abstractions

- Abstraction helps us deal with complexity
  - Hide lower-level detail
- Instruction set architecture (ISA)
  - The hardware/software interface
- Application binary interface
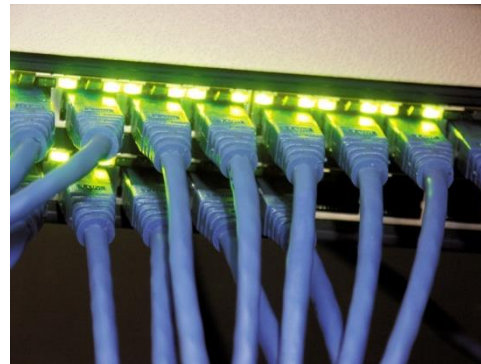  - The ISA plus system software interface
- Implementation
  - The details underlying and interface

# A Safe Place for Data

- Volatile main memory
  - Loses instructions and data when power off
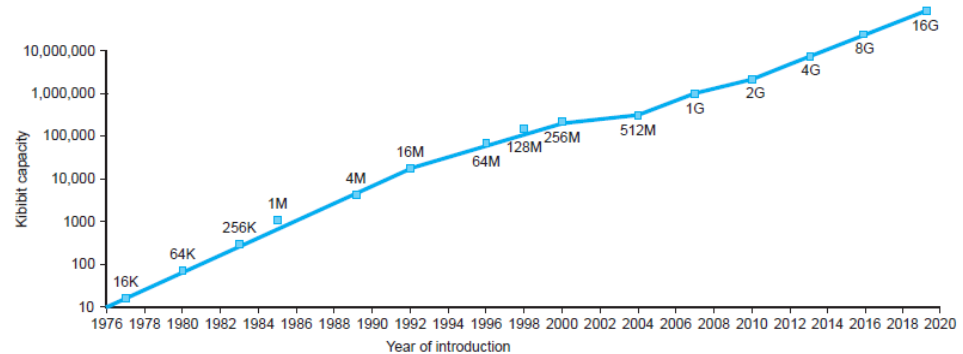- Non-volatile secondary memory
  - Magnetic disk
  - Flash memory
  - Optical disk (CDROM, DVD)

# Networks

- Communication, resource sharing, nonlocal access

- Local area network (LAN): Ethernet

- Wide area network (WAN): the Internet

- Wireless network: WiFi, Bluetooth

# Technology Trends

- **Electronics technology continues to evolve**
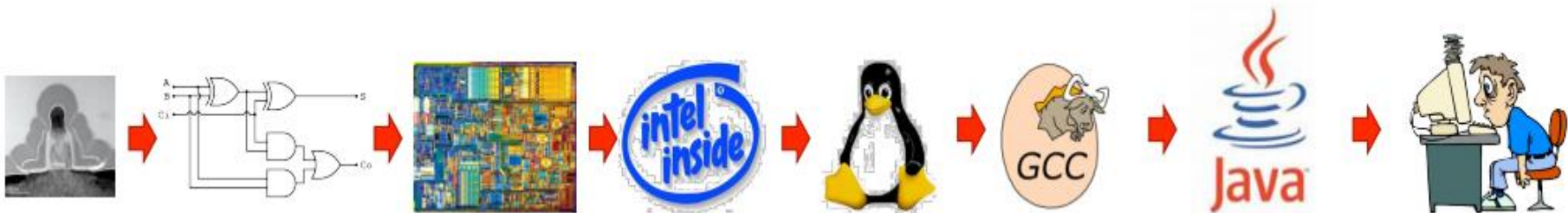  - Increased capacity and performance
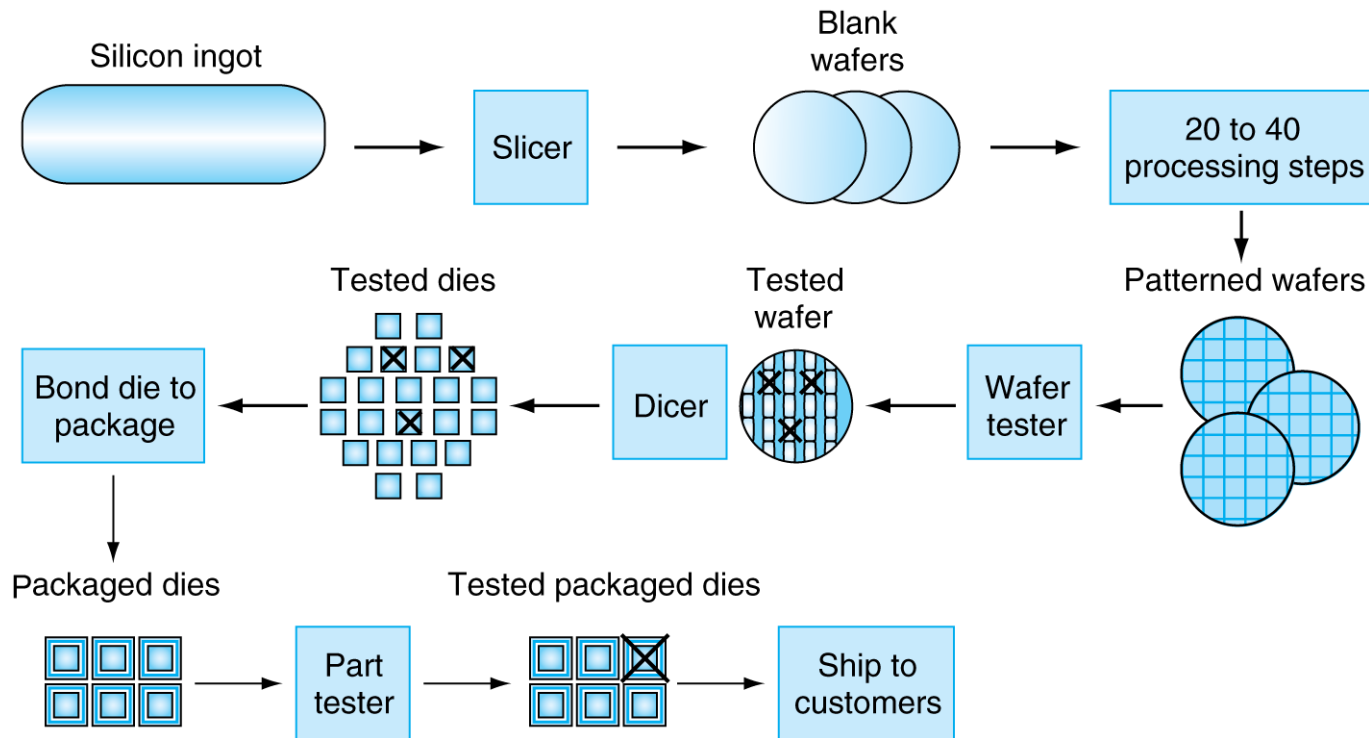  - Reduced cost



DRAM capacity

| Year | Technology | Relative performance/cost |
|------|------------|:--------------------------:|
| 1951 | Vacuum tube | 1 |
| 1965 | Transistor | 35 |
| 1975 | Integrated circuit (IC) | 900 |
| 1995 | Very large scale IC (VLSI) | 2,400,000 |
| 2013 | Ultra large scale IC | 250,000,000,000 |

# Semiconductor Technology

- Silicon:  semiconductor
- Add materials to transform properties:
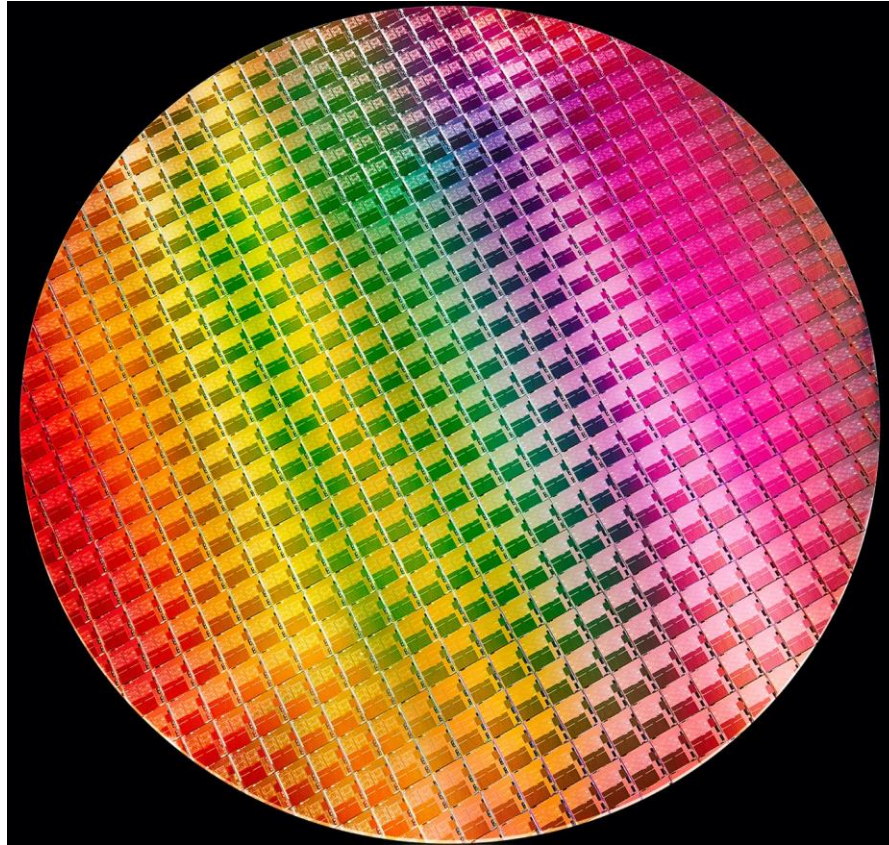  - Conductors
  - Insulators
  - Switch

# Manufacturing ICs



- Yield: proportion of working dies per wafer

# Intel® Core 10th Gen



- 300mm wafer, 506 chips, 10nm technology
- Each chip is 11.4 x 10.7 mm

# Integrated Circuit Cost

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

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

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

- 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

# What Is Computer Architecture?

**Computer Architecture:** The science and art of designing, selecting, and interconnecting hardware components and designing the hardware/software interface to create a computing system that meets functional, performance, energy consumption, cost, and other specific goals.

# What Is Computer Architecture?

**Computer Architecture:** The term *architecture* is used here to describe the attributes of a system as seen by the programmer, *i.e.*, the conceptual structure and functional behavior as distinct from the organization of the dataflow and controls, the logic design, and the physical implementation.

# What Is Computer Architecture?

- Computer Architecture
  - Instruction Set Architecture  & Computer Organization

- Instruction Set Architecture (ISA)
  - WHAT the computer does (logical view)

- Computer Organization
  - HOW the ISA is implemented (physical view)

# Current State of Architecture

# Current State of Architecture

Advance of Semiconductors: "Moore's Law"

**Gordon Moore, Founder of Intel**

- 1965: since the integrated circuit was invented, the number of transistors/inch$^2$ in these circuits roughly doubled every year; this trend would continue for the foreseeable future

- 1975: revised - circuit complexity doubles every 18 months



Image credit: http://download.intel.com/research/silicon/Gordon_Moore_ISSCC_021003.pdf

# Current State of Architecture

Leveraging Moore's Law Trends

**From increasing circuit density to performance:**

- More transistors = ↑ opportunities for exploiting parallelism

# The Importance of Architecture

- We design smarter and smarter processors

  - Process technology gives us about 20% performance improvement per year

  - Until 2004, performance grew at about 40% per year
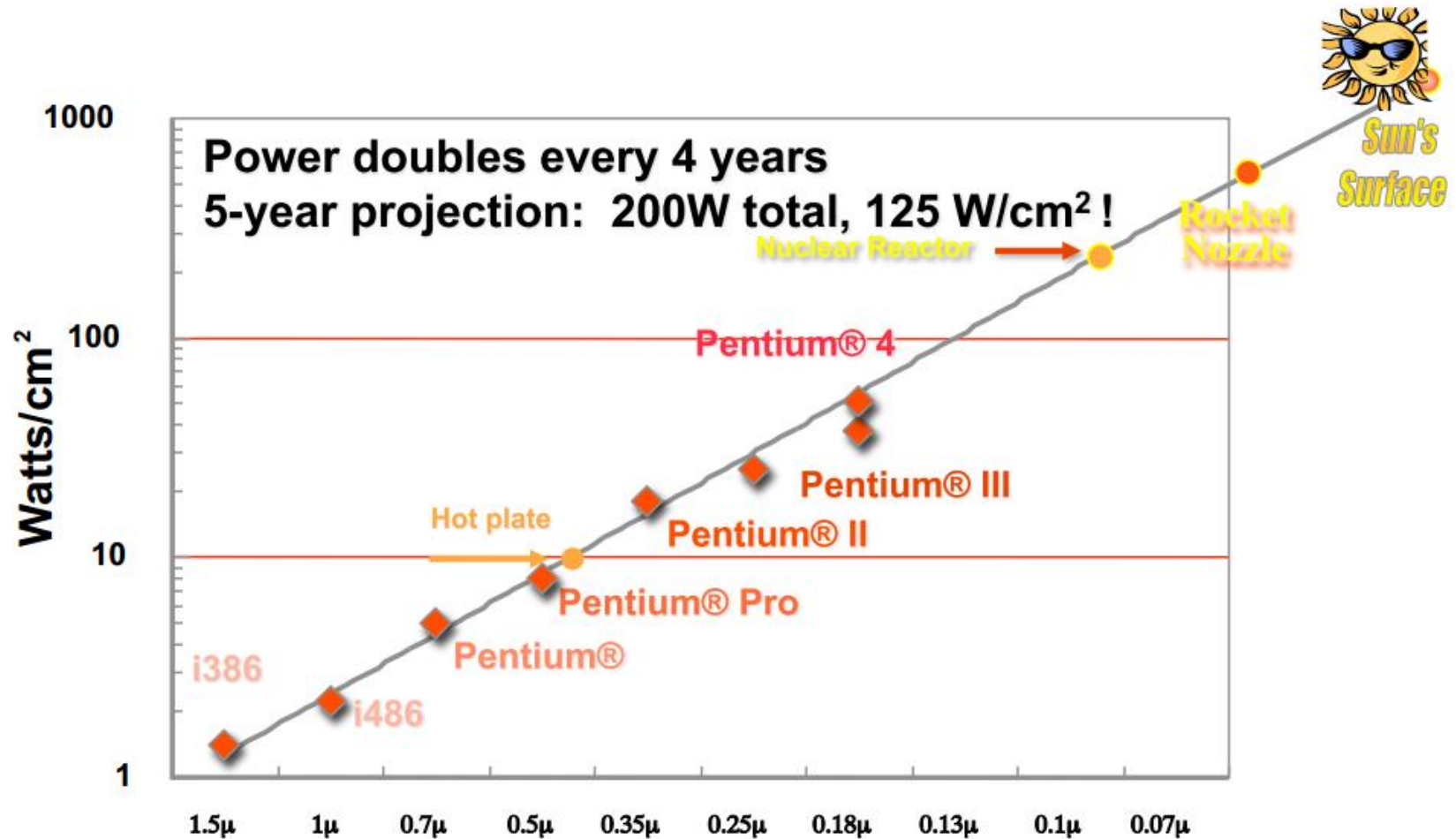
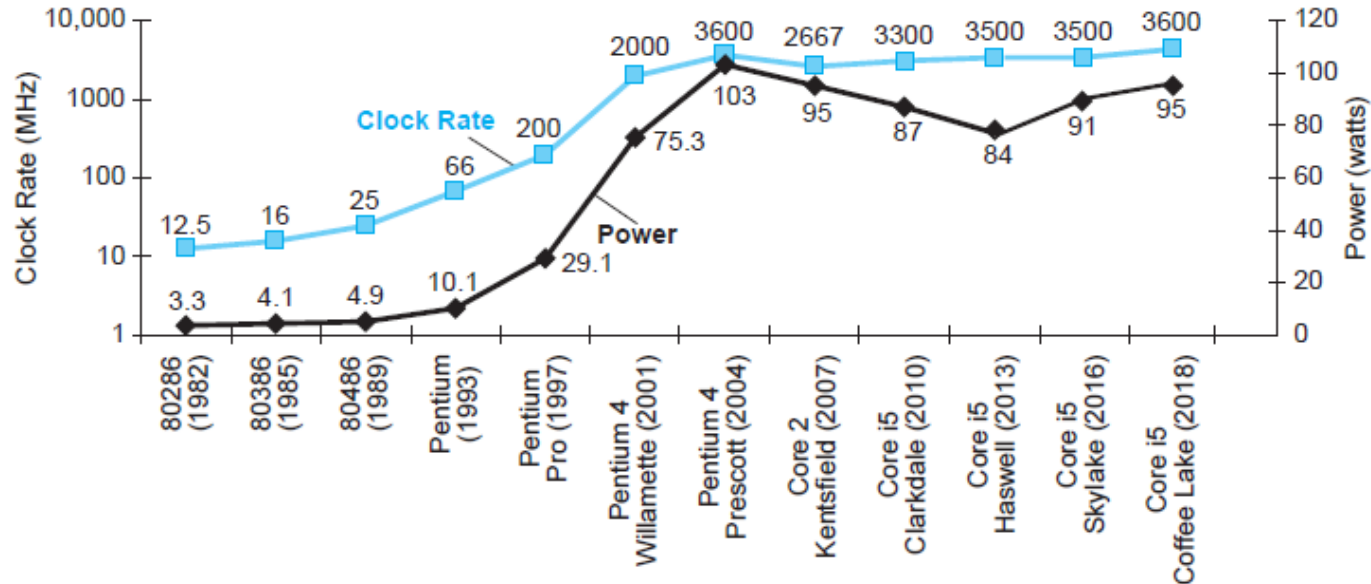- The gap is due to architecture! (and compilers)

# Computer Performance

# Power

- **Clock speed** is the biggest contributor to power

  - Chip manufactures (Intel, esp.) pushed clock speeds very hard in the 90s and early 2000s.

  - Doubling the clock speed increases power by 2-8x

  - Clock speed scaling is essentially finished.

# Power



Power doubles every 4 years
5-year projection: 200W total, 125 W/cm² !

From "New Microarchitecture Challenges in the Coming Generations of CMOS Process Technologies" – Fred Pollack, Intel Corp. Micro32 conference key note - 1999.

# Power Trends

- In CMOS IC technology

$$Power = Capacitive\ load \times Voltage^2 \times Frequency$$

| ×30 | 5V → 1V | ×1000 |

# Reducing Power

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

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

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

# Important Trends

- Historical contributions to performance:
    - Better processes (faster devices) ~20%
    - Better circuits/pipelines ~15%
    - Better organization/architecture ~15%

- In the future, bullet-2 will help little, and bullet-1 will eventually disappear!

|  | Pentium | P-Pro | P-II | P-III | P-4 | Itanium | Montecito |
|---|---|---|---|---|---|---|---|
| Year | 1993 | 95 | 97 | 99 | 2000 | 2002 | 2005 |
| Transistors | 3.1M | 5.5M | 7.5M | 9.5M | 42M | 300M | 1720M |
| Clock Speed | 60M | 200M | 300M | 500M | 1500M | 800M | 1800M |

Moore's Law in action

At this point, adding transistors to a core yields little benefit

37

# What's Next: Parallelism

- You probably own a multi-processor

- They provide some performance, but it's hard to Fully exploit (parallel programming !)
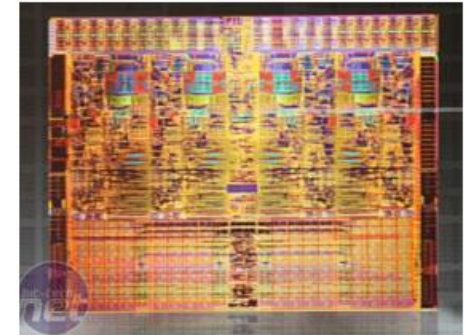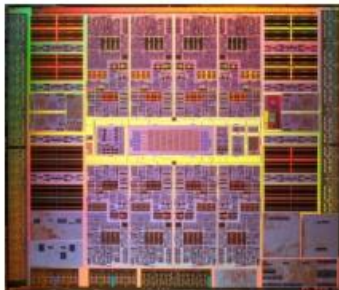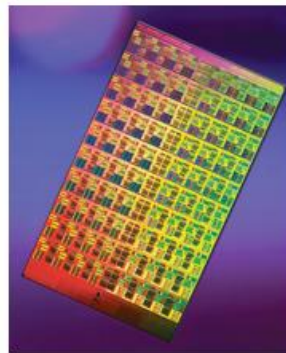
# What's Next: Parallelism
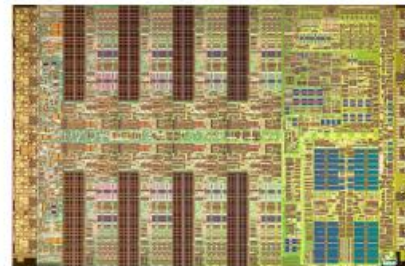
Intel P4
1 core

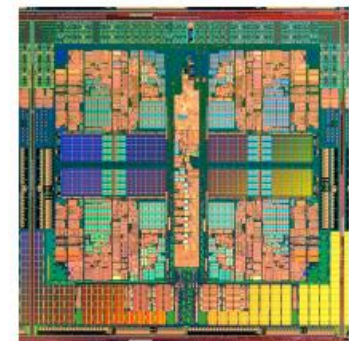Intel Core 2 Duo
2 cores

Intel Nahalem
4 cores

SPARC T1
8 cores

Intel Prototype
80 cores

Cell BE
8 + 1 cores

AMD Barcelona
4 cores