

Prof. REN-SONG TSAY
NTHU

°

CHAPTER I
COMPUTER
ABSTRACTIONS AND
TECHNOLOGY

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan1

What is Computer Architecture?

Software

Hardware

Application

Compiler

Operating System

Assembler

Instruction Set Architecture

Processor

Memory

I/O System

Datapath & Control

Circuit Design

Transistors

Machine Organization

Computer Architecture =
Instruction Set Architecture
+ Machine Organization

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan2

§1.1 Introduction

The Computer Revolution

- Progress in computer technology
 - Underpinned by Moore’s Law
- Makes novel applications feasible
 - Computers in automobiles
 - Smart phones
 - Human genome project
 - World Wide Web
 - Search Engines
 - Machine learning
- Computers are pervasive

42016/2/15 © Ren-Song Tsay, NTHU, Taiwan3

Evolution of Processor Performance

2X / 1.5-2years

MIPS

1,000

100

10

1

0.1

4 bit 8 bit 16 bit 32 bit

1975 1980 1985 1990 1995 2000

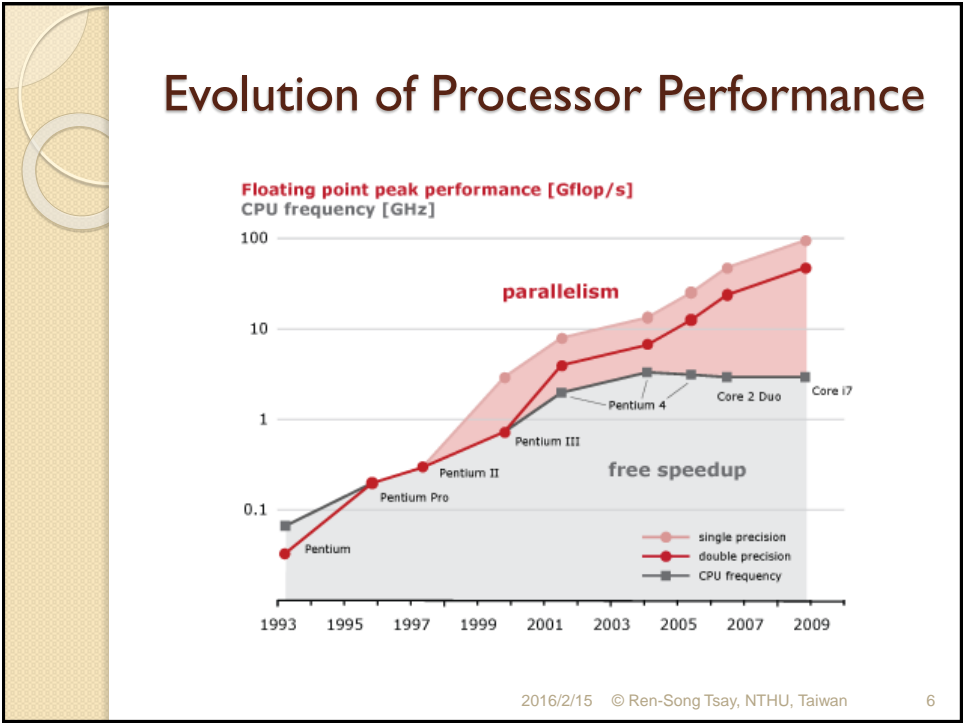
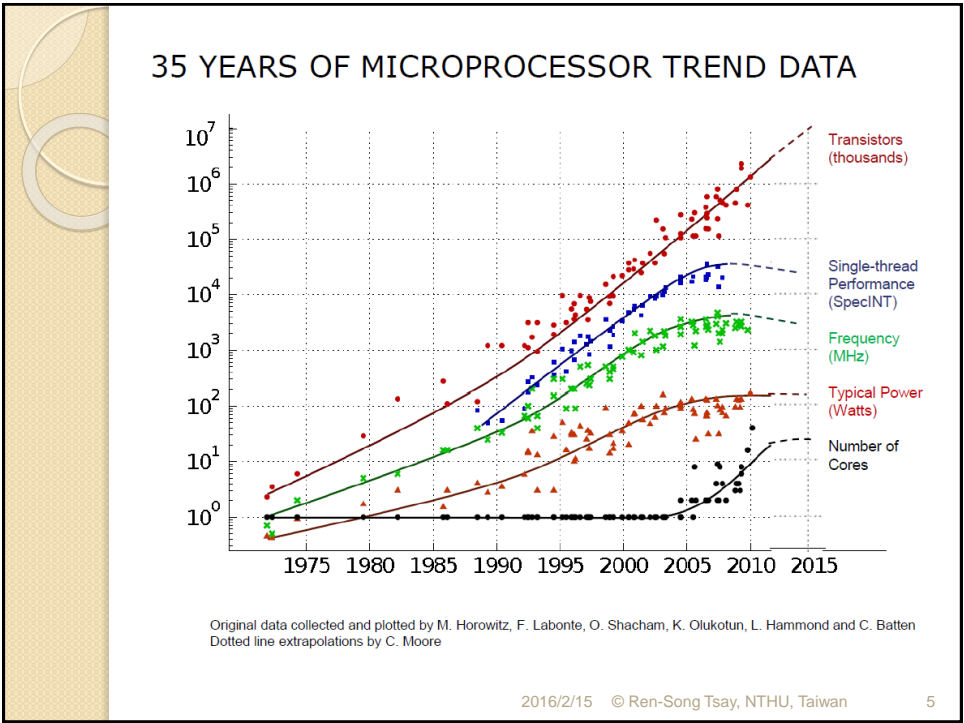
i4004 i8008 i8086 68000 Sparc i486 68040 ARM6 SH1 SH2 SH3 SH4 Pentium Pentium Pro Pentium II StrongARM Enhanced Alpha

CISC

New-generation RISC

(“The Cooler the Better: New Directions in the Nomadic Ages,” *Computer*, April 2001.)

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan4



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

5

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

8

2016.2.1 Microsoft testing underwater data centers



2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

9

Classes of Computers

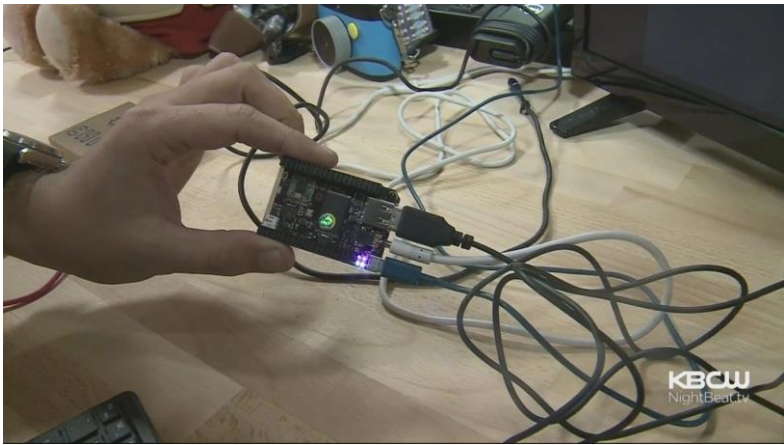
- Supercomputers
 - High-end scientific and engineering calculations
 - Highest capability but represent a small fraction of the overall computer market
- Embedded computers
 - Hidden as components of systems
 - Stringent power/performance/cost constraints

5

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

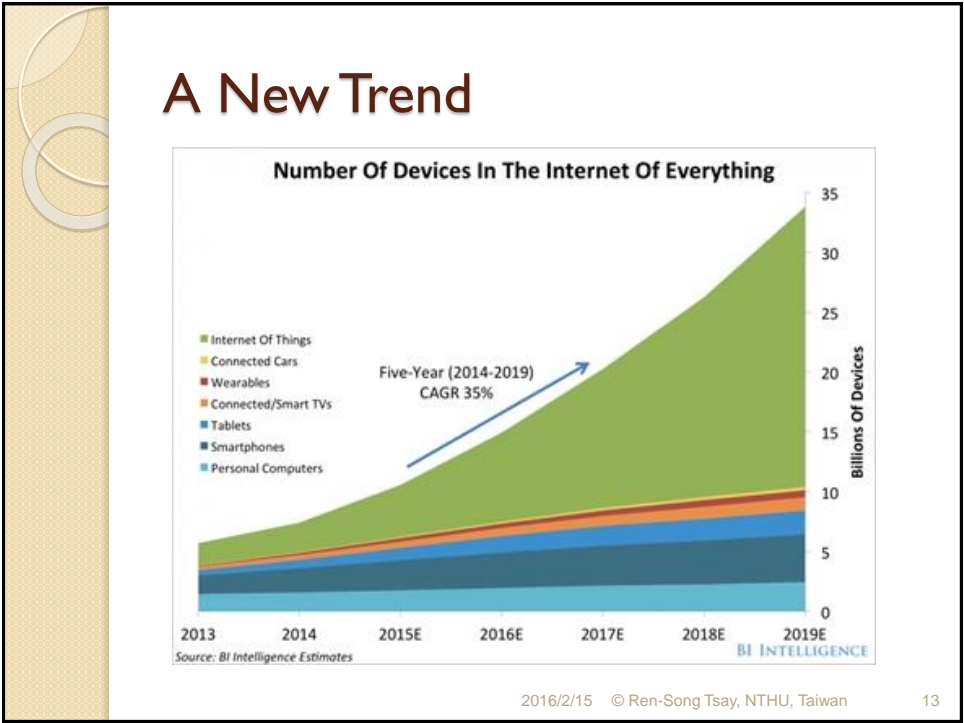
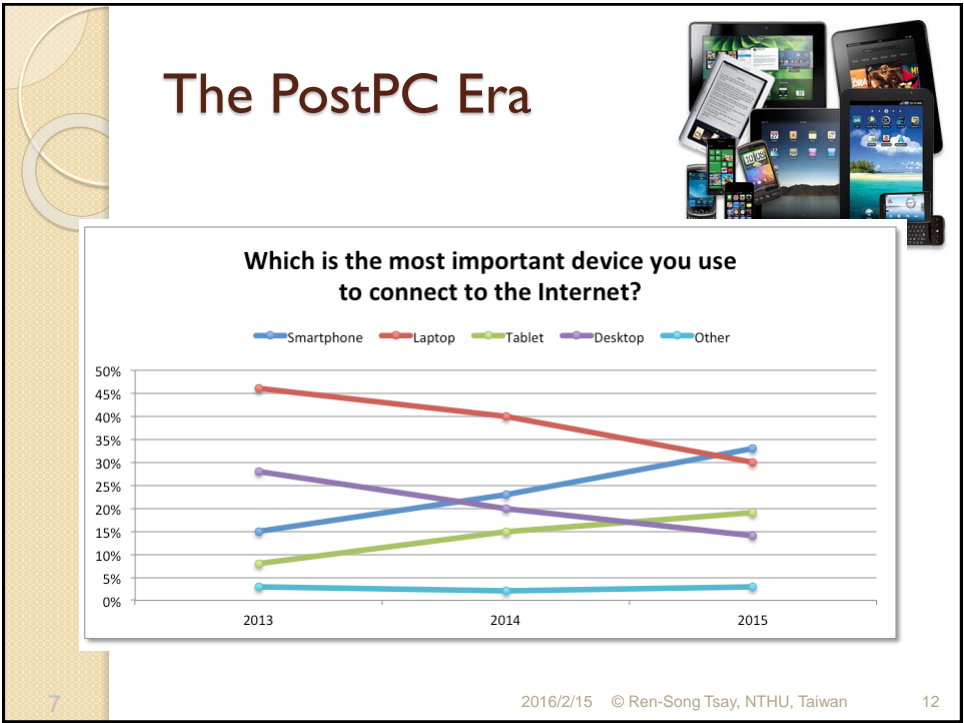
10


2016.2.11 \$9 Computer Designed in Oakland



2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

11





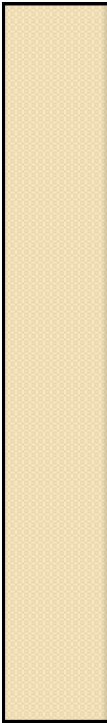
The PostPC Era

- 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

7

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan


14



Global Smartphone Shipments Share Q1 2015

Company	Share (%)
Samsung	24.40%
Apple	17.90%
Lenovo	6.50%
Huawei	5.10%
LG	4.50%
Xiaomi	4.50%
Others	37.20%

Source: Counterpoint Research, April 2015



2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

15

Wearables 1.0

IRONMAN

Wearables 2.0

INVISIBLE MAN

A MUCH More Diversified Market Than Investors Realize

CREDIT SUISSE

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 16

The 4th Industrial Revolution - „Industry 4.0“

Drivers

Quality of life

Engineering Sciences

1st

steam engine

GB

1782

Power generation

Mechanical automation

Mobility

2nd

conveyor belt

US

1913

Industrialization

µelectronics

3rd

Computer, NC, PLC

1954

Electronic Automation

ICT

4th

Cyber Physical Systems

2015

Smart Automation

Simulación

Big DATA

Big Data y análisis

Internet industrial de las cosas

Robots autónomos

Computación en la nube

Realidad aumentada

Ciberseguridad

Sistemas de integración horizontales y verticales

Fabricación aditiva

Industria 4.0

Internet of data

Smart Buildings

Smart Mobility

Smart Grid

Smart Factories

Smart Homes

Social Web

Business Web

Smart Logistics

Internet of people

Internet of things

Internet of services

2016

Chapter 1 — Computer Abstractions and Technology

8

World Largest Companies

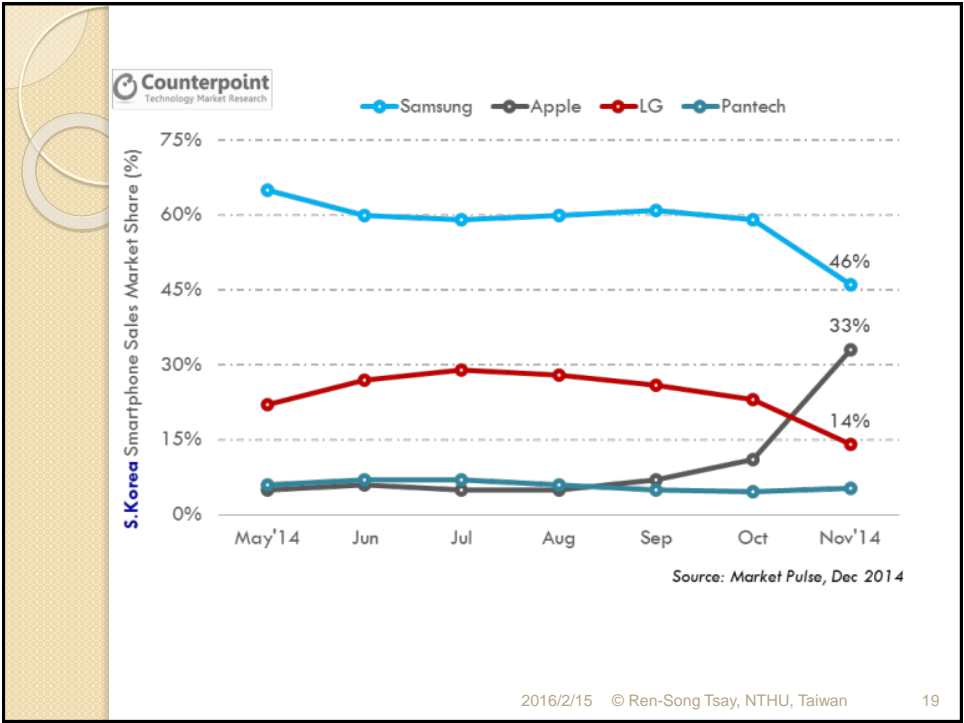


2016.2.15 market value

- Apple \$521B (\$725B in 2015)
- Google \$470B
- Facebook \$290B
 - Taiwan 2015 budget \$60 B
 - Taiwan GDP \$505 B (2014)





2016/2/15 © Ren-Song Tsay, NTHU, Taiwan18







Embedded Computer



Lego Mindstorms



- ◆A computer inside another device used for running one predetermined application
- ◆Uses: control (traffic, printer, disk); consumer electronics (video game, CD player, PDA)

Robotic command explorer:
A “Programmable Brick”,
Hitachi H8 CPU (8-bit), 32KB RAM,
LCD, batteries,
infrared transmitter/receiver,
4 control buttons, 6 connectors




2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

20

Arduino

- 一個開放原始碼的單晶片微控制器
- 使用了Atmel AVR單片機
- 建構於簡易輸出/輸入（ simple I/O ） 介面板，並且具有使用類似Java、C語言的Processing/Wiring開發環境



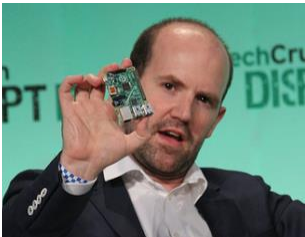
2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

21

Raspberry Pi2

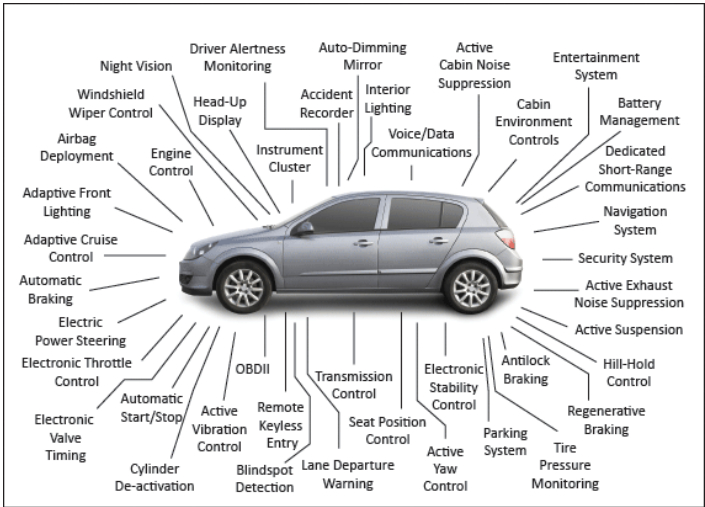


- Since 2012
- 火柴盒大小, \$35
- 可搭載Open Source 的 Linux 系統
- 免費提供Win 10 的開發者套件
- 創始人 Eben Upton



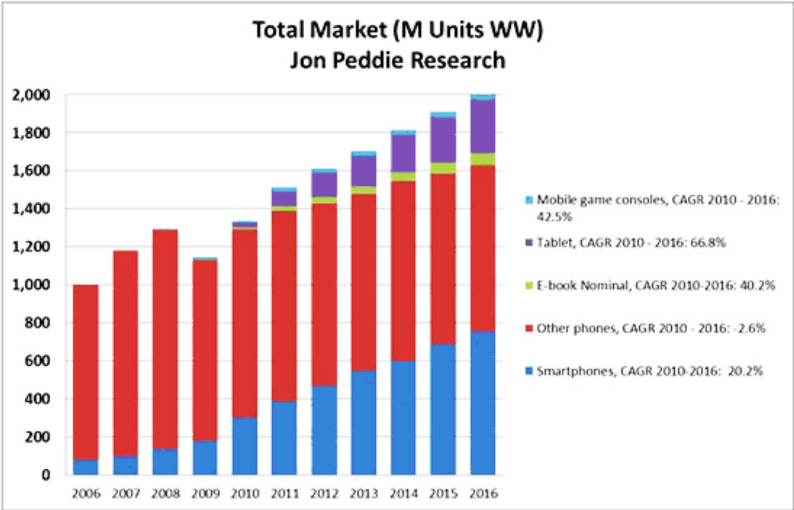
2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 22

Embedded Everywhere

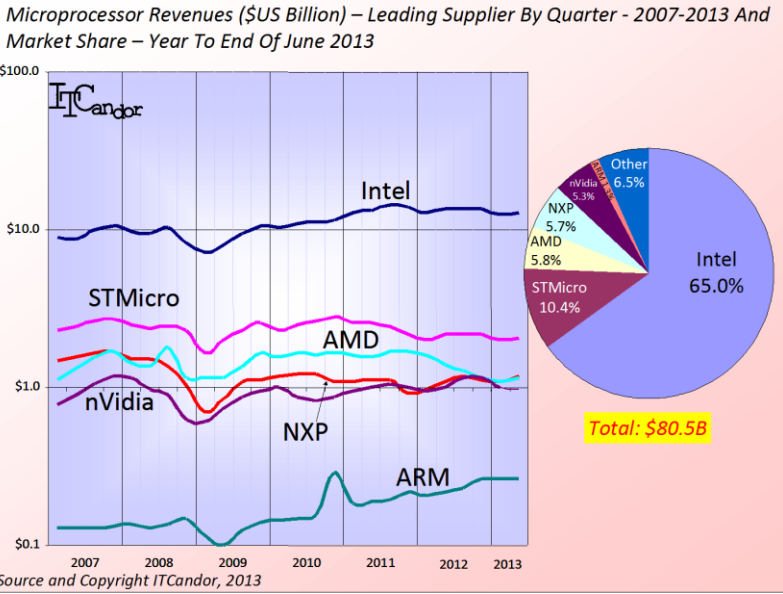


2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 23

The Processor Market

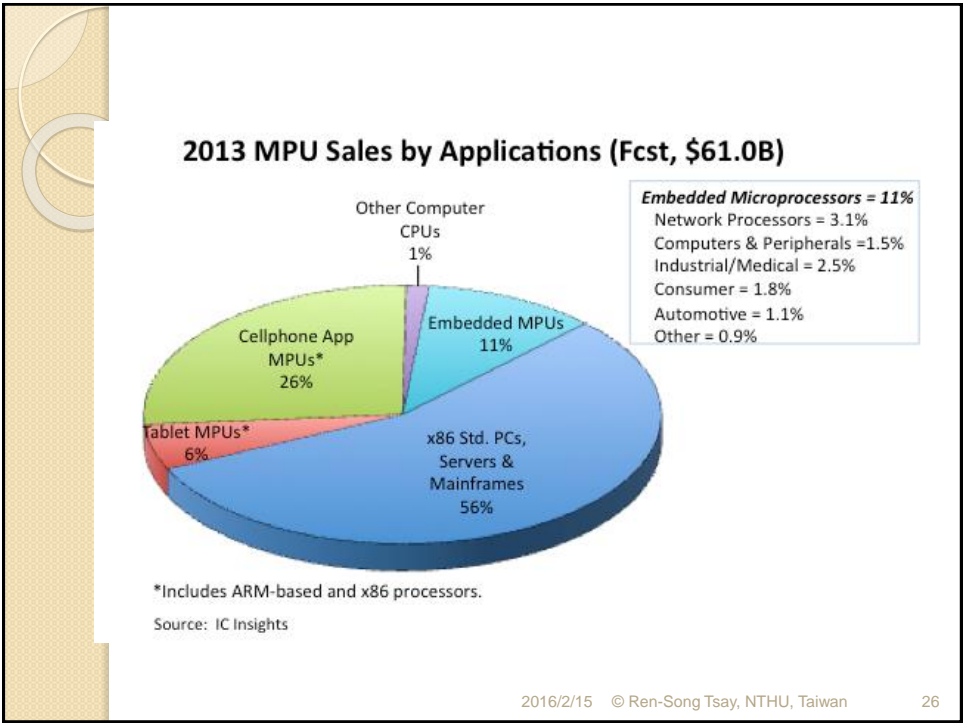


2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 24



Source and Copyright ITCandor, 2013

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 25

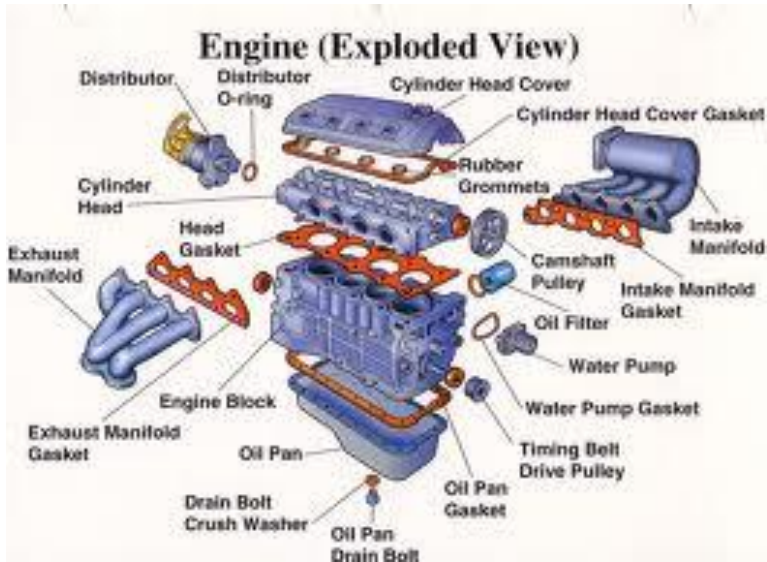


What is coming ...

- Internet-of-things
- Big data
- Intelligent life assistance
- Artificial intelligence
- Industry 4.0

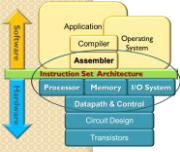
What You Will Learn

- 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



9

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 30

Eight Great Ideas

- Design for **Moore's Law**
- 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



11

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 31

§1.2 Eight Great Ideas in Computer Architecture

§1.3 Below Your Program

Below Your Program

Applications software

Systems software

Hardware

- 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

13

© Ren-Song Tsay, NTHU, Taiwan

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
```

↓
ISA

Assembler

Binary machine language program (for MIPS)

```
0000000001010000100000000000011000
0000000000001100000001100000100001
1000110001100010000000000000000000
1000110011110010000000000000000100
1010110011110010000000000000000000
1010110001100010000000000000000100
000000111110000000000000000001000
```

15

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

33

Chapter 1 — Computer Abstractions and Technology

16

§1.4 Under the Covers

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

17

© Ren-Song Tsay, NTHU, Taiwan

Even thinner: only 0.29-inch

Retina Display

Stereo Speakers

New Antennas

FaceTime HD Camera

5-megapixel camera Full HD 1080p

LED Flash

New Curved-Glass Back

iPad 3

Concept by: Guilherme Martins Schasiepen

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

35

Opening the Box

Capacitive multitouch LCD screen

3.8 V, 25 Watt-hour battery

Computer board

20

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan 36

iPhone components

- A4 processor (1GHz ARM Cortex A8)
- 512 MB RAM
- Gorilla Glass
- 1420 mAh Li-Ion
- 5M Pixel CCD
- SAMSUNG flash
- TI Touch Panel Controller
- Cirrus Logic Codec




37

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

19

Touchscreen

- PostPC device
 - Supersedes keyboard and mouse
- Resistive and Capacitive types
 - Most tablets, smart phones use capacitive
 - Capacitive allows multiple touches simultaneously
- Human Interface
 - sensors

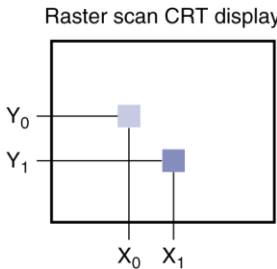
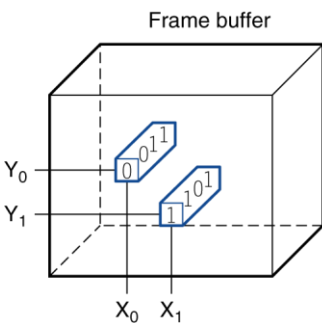


© Ren-Song Tsay, NTHU, Taiwan

18

Through the Looking Glass

- LCD screen: picture elements (pixels)
 - Mirrors content of frame buffer memory



2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

39

Inside the Processor (CPU)

- Datapath: performs operations on data
- Control: sequences datapath, memory, ...
- Cache memory
 - Small fast SRAM memory for immediate access to data

Inside the Processor Apple A5



22

Abstractions

The BIG Picture

- 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 underlying details and interface

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

42

22

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)





© Ren-Song Tsay, NTHU, Taiwan

23

Networks

- Communication and resource sharing
- Local area network (LAN): Ethernet
 - Within a building
- Wide area network (WAN): the Internet
- Wireless network: WiFi, Bluetooth



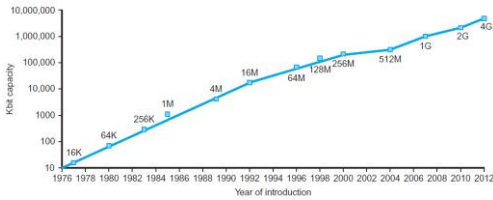
2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

44

25

Technology Trends

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



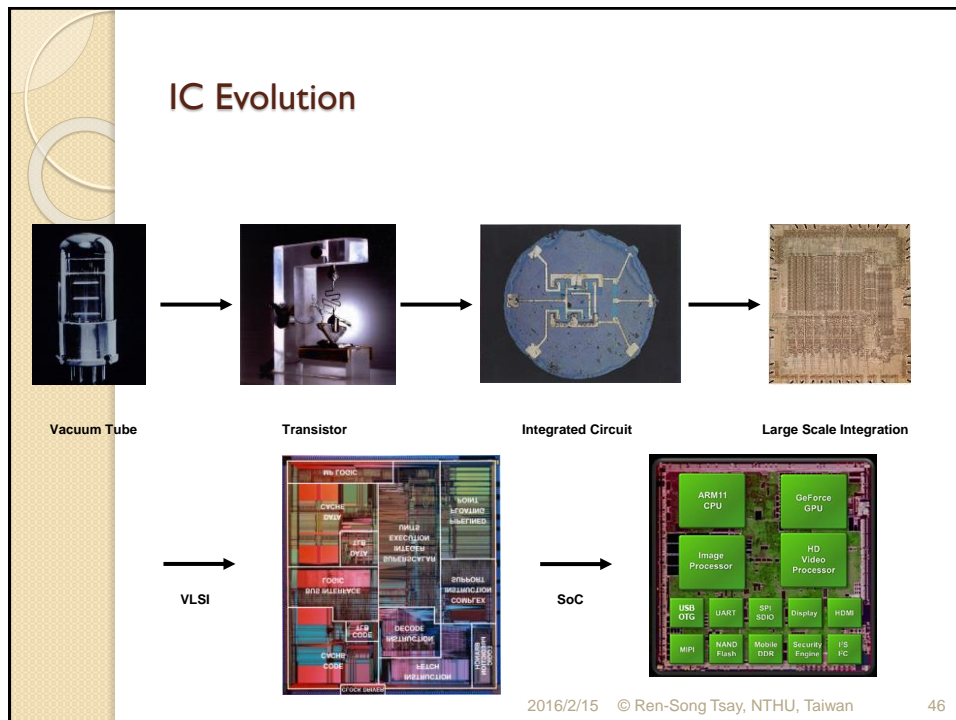
DRAM capacity

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

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

45

§1.5 Technologies for Building Processors and Memory



Semiconductor Technology

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

Manufacturing ICs

The flowchart illustrates the IC manufacturing process. It begins with a 'Silicon ingot' (cylinder) which is processed by a 'Slicer' to produce 'Blank wafers' (stack of circles). These wafers undergo '20 to 40 processing steps' to become 'Patterned wafers' (stack of circles with grid patterns). A 'Wafer tester' then examines a 'Tested wafer' (circle with grid and 'X' marks). The wafer is then processed by a 'Dicer' to produce 'Tested dies' (grid of squares, some with 'X' marks). These dies are then 'Bond die to package' to create 'Packaged dies' (grid of squares). A 'Part tester' tests these to produce 'Tested packaged dies' (grid of squares, some with 'X' marks), which are finally 'Ship to customers'.

- Yield: proportion of working dies per wafer

26

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

48

Intel Core i7 Wafer

A circular image showing a silicon wafer with a dense grid of square dies.

- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

27

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

49

28

Integrated Circuit Cost

$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$
$$\text{Dies per wafer} \approx \text{Wafer area} / \text{Die area}$$
$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area} / 2))^2}$$

- 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


2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

50

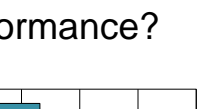
29

Defining Performance

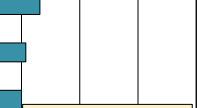
- Which airplane has the best performance?



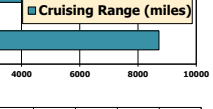
\$350M



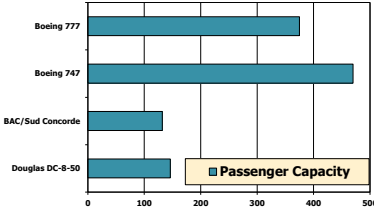
\$250M



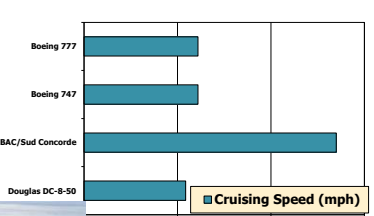
\$150M



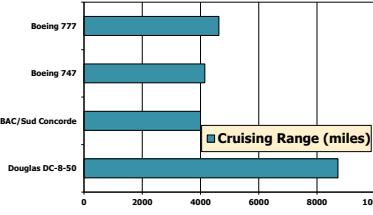
\$55M



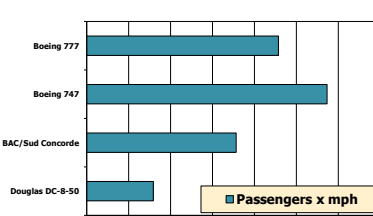
Passenger Capacity



Cruising Speed (mph)



Cruising Range (miles)



Passengers x mph

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan

51

30

Response Time and Throughput

- Response time
 - How long it takes to do a task
- Throughput
 - Total work done per unit time
 - e.g., tasks/transactions/... per hour
- How are response time and throughput affected by
 - Replacing the processor with a faster version?
 - Adding more processors?
- We'll focus on response time for now...

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan52

31

Relative Performance

- Define Performance = 1/Execution Time
- “X is n time faster than Y”

$$\frac{\text{Performance}_X}{\text{Performance}_Y}$$
$$= \frac{\text{Execution time}_Y}{\text{Execution time}_X} = n$$
- Example: time taken to run a program
 - 10s on A, 15s on B
 - $\text{Execution Time}_B / \text{Execution Time}_A$
 $= 15\text{s} / 10\text{s} = 1.5$
 - So A is 1.5 times faster than B

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan53

32

Idle

IO

System

Measuring Execution Time

- Elapsed time
 - Total response time, including all aspects
 - Processing, I/O, OS overhead, idle time
 - Determines system performance
- CPU time
 - Time spent processing a given job
 - Discounts I/O time, other jobs' shares
 - Comprises user CPU time and system CPU time
 - Different programs are affected differently by CPU and system performance

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

54

33

CPU Clocking

- Operation of digital hardware governed by a constant-rate clock

The diagram illustrates the timing of CPU operations relative to a clock. It features three horizontal tracks: 'Clock (cycles)' at the top, 'Data transfer and computation' in the middle, and 'Update state' at the bottom. The clock track shows a square wave with a double-headed arrow labeled 'Clock period' spanning one full cycle. The data transfer and computation track shows a continuous blue bar with 'X' marks at the boundaries of each clock cycle. The update state track shows blue hexagons positioned at the midpoint of each clock cycle.

- Clock period: duration of a clock cycle
 - e.g., 250ps = 0.25ns = 250×10⁻¹²s
- Clock frequency (rate): cycles per second
 - e.g., 4.0GHz = 4000MHz = 4.0×10⁹Hz

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

55

34

CPU Time

CPU Time = CPU Clock Cycles × Clock Cycle Time
$$= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}$$

- Performance improved by
 - Reducing number of clock cycles
 - Increasing clock rate
 - Hardware designer must often trade off clock rate against cycle count

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

56

35

CPU Time Example

- Computer A: 2GHz clock, 10s CPU time
- Designing Computer B
 - Aim for 6s CPU time
 - Can do faster clock, but requires 1.2 × clock cycles
- How fast must Computer B clock be?

$$\text{Clock Rate}_B = \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6s}$$
$$\text{Clock Cycles}_A = \text{CPU Time}_A \times \text{Clock Rate}_A$$
$$= 10s \times 2\text{GHz} = 20 \times 10^9$$
$$\text{Clock Rate}_B = \frac{1.2 \times 20 \times 10^9}{6s} = \frac{24 \times 10^9}{6s} = 4\text{GHz}$$

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

57

Chapter 1 — Computer Abstractions and Technology

28

Instruction Count and CPI

Clock Cycles = Instruction Count × Cycles per Instruction

CPU Time = Instruction Count × CPI × Clock Cycle Time

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

- Instruction Count for a program
 - Determined by program, ISA and compiler
- Average cycles per instruction
 - Determined by CPU hardware
 - If different instructions have different CPI
 - Average CPI affected by instruction mix

36

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan58

CPI Example

- Computer A: Cycle Time = 250ps, CPI = 2.0
- Computer B: Cycle Time = 500ps, CPI = 1.2
- Same ISA
- Which is faster, and by how much?

$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$
 $= 1 \times 2.0 \times 250\text{ps} = 1 \times 500\text{ps}$ ← A is faster...

$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$
 $= 1 \times 1.2 \times 500\text{ps} = 1 \times 600\text{ps}$

$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{1 \times 600\text{ps}}{1 \times 500\text{ps}} = 1.2$ ← ...by this much

36

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan59

CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^n \left(\text{CPI}_i \times \underbrace{\frac{\text{Instruction Count}_i}{\text{Instruction Count}}}_{\text{Relative frequency}} \right)$$

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan60

CPI Example

- Alternative compiled code sequences using instructions in classes A, B, C

Class	A	B	C
CPI for class	1	2	3
IC in sequence 1	2	1	2
IC in sequence 2	4	1	1

- Sequence 1: IC = 5
 - Clock Cycles = $2 \times 1 + 1 \times 2 + 2 \times 3 = 10$
 - Avg. CPI = $10/5 = 2.0$
- Sequence 2: IC = 6
 - Clock Cycles = $4 \times 1 + 1 \times 2 + 1 \times 3 = 9$
 - Avg. CPI = $9/6 = 1.5$

372016/2/15 © Ren-Song Tsay, NTHU, Taiwan61

Chapter 1 — Computer Abstractions and Technology

30

Performance Summary

The BIG Picture

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- Performance depends on
 - Algorithm: affects IC, possibly CPI
 - Programming language: affects IC, CPI
 - Compiler: affects IC, CPI
 - Instruction set architecture: affects IC, CPI, T_c

38

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan62

Power Trends

Processor	Year	Clock Rate (MHz)	Power (watts)
80286	1982	12.5	3.3
80386	1985	16	4.1
80486	1989	25	4.9
Pentium	1993	66	10.1
Pentium Pro	1997	200	29.1
Pentium 4	2001	2000	75.3
Pentium 4 Prescott	2004	3600	103
Core 2	2007	2667	95
Core i5	2010	3300	87
Core i5 Ivy Bridge	2012	3400	77

- In CMOS IC technology

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

x30

5V → 1V

x1000

40

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan63

Reducing Power

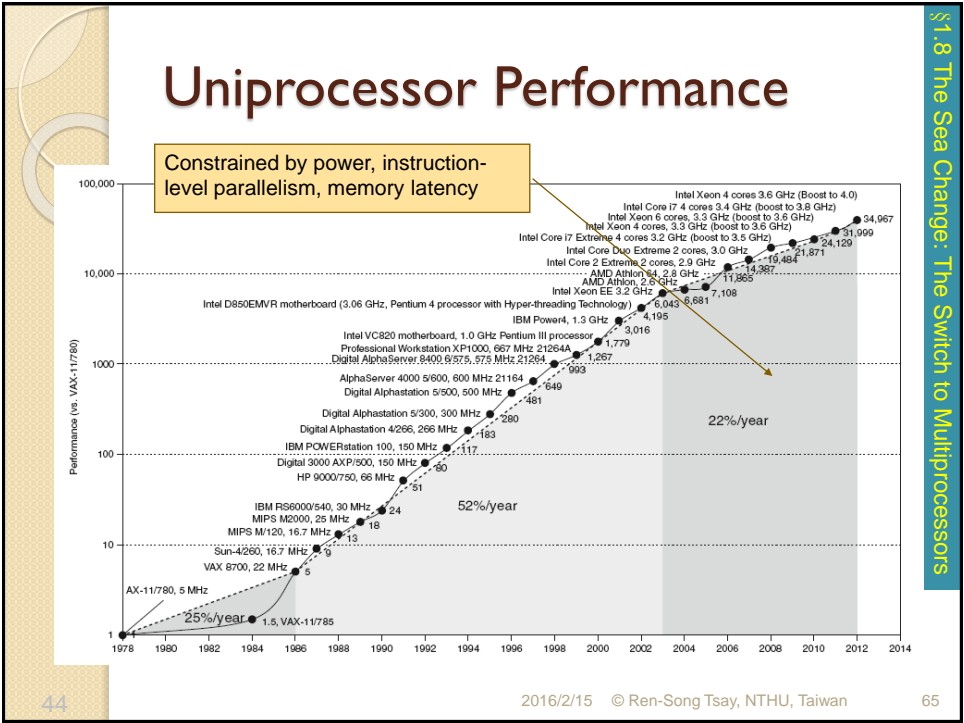
- Suppose a new CPU has
 - 85% of capacitive load of old CPU
 - 15% voltage and 15% frequency reduction

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

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

41

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan64



45

Multiprocessors

- Multicore microprocessors
 - More than one processor per chip
- Requires explicitly **parallel** programming
 - Compare with instruction level parallelism
 - Hardware executes multiple instructions at once
 - Hidden from the programmer
 - Hard to do
 - Programming for performance
 - Load balancing
 - Optimizing communication and synchronization

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan66

46

SPEC CPU Benchmark

- Programs used to measure performance
 - Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
 - Develops benchmarks for CPU, I/O, Web, ...
- SPEC CPU2006
 - Elapsed time to execute a selection of programs
 - Negligible I/O, so focuses on CPU performance
 - Normalize relative to reference machine
 - Summarize as geometric mean of performance ratios
 - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan67

§1.9 Real Stuff Benchmarking the Intel Core i7

CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 ⁹	CPI	Clock cycle time (seconds x 10 ⁻⁹)	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	bzip2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	221	2.66	0.376	221	9120	41.2
Go game (AI)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	hmmer	2616	0.60	0.376	590	9330	15.8
Chess game (AI)	sjeng	1948	0.80	0.376	586	12100	20.7
Quantum computer simulation	libquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264avc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	omnetpp	367	2.10	0.376	290	6250	21.5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalancbmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7

High cache miss rates

47

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan68

SPEC Power Benchmark

- Power consumption of server at different workload levels
 - Performance: ssj_ops/sec
 - Power: Watts (Joules/sec)

$$\text{Overall ssj_ops per Watt} = \left(\sum_{i=0}^{10} \text{ssj_ops}_i \right) / \left(\sum_{i=0}^{10} \text{power}_i \right)$$

48

2016/2/15 © Ren-Song Tsay, NTHU, Taiwan69

SPECpower_ssj2008 for Xeon X5650

Target Load %	Performance (ssj_ops)	Average Power (Watts)
100%	865,618	258
90%	786,688	242
80%	698,051	224
70%	607,826	204
60%	521,391	185
50%	436,757	170
40%	345,919	157
30%	262,071	146
20%	176,061	135
10%	86,784	121
0%	0	80
Overall Sum	4,787,166	1,922
Σssj_ops/Σpower =		2,490

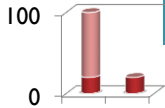
482016/2/15 © Ren-Song Tsay, NTHU, Taiwan70

1.10 Fallacies and Pitfalls

Pitfall (陷阱): Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{\text{improvement factor}} + T_{\text{unaffected}}$$



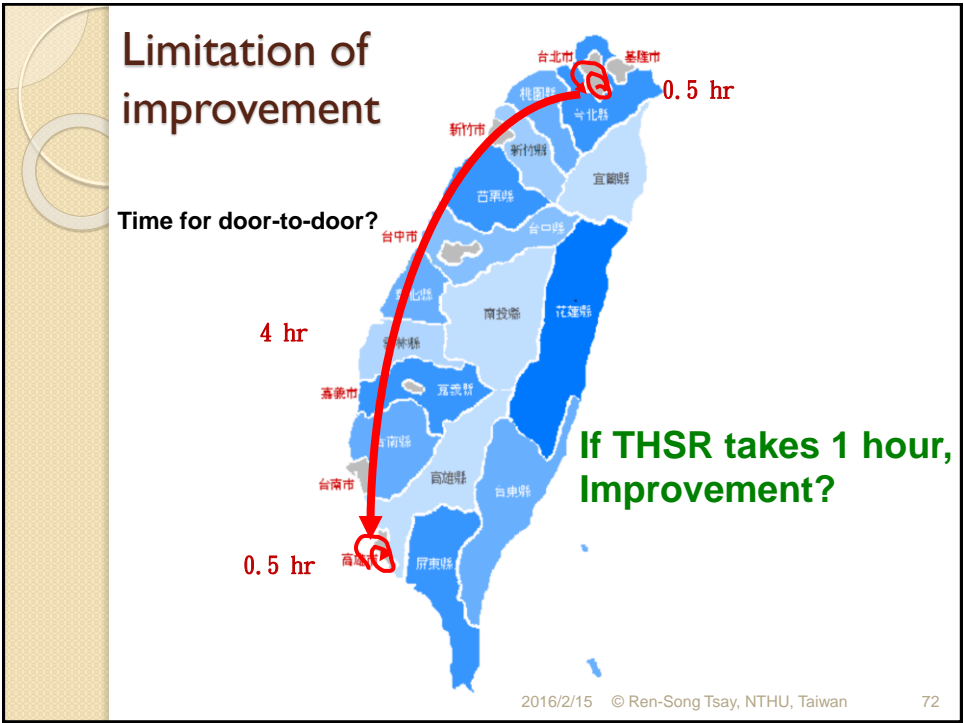
- Example: multiply accounts for 80s/100s
 - How much improvement in multiply performance to get 5x overall?

$$20 = \frac{80}{n} + 20$$

- Can't be done!

- Corollary: make the common case fast

492016/2/15 © Ren-Song Tsay, NTHU, Taiwan71



Fallacy (迷思): Low Power at Idle

- Look back at i7 power benchmark
 - At 100% load: 258W
 - At 50% load: 170W (66%)
 - At 10% load: 121W (47%)
- Google data center
 - Mostly operates at 10% – 50% load
 - At 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

51

Pitfall: MIPS as a Performance Metric

- MIPS: Millions of Instructions Per Second
 - Doesn't account for
 - Differences in ISAs between computers
 - Differences in complexity between instructions

$$\begin{aligned} \text{MIPS} &= \frac{\text{Instruction count}}{\text{Execution time} \times 10^6} \\ &= \frac{\text{Instruction count}}{\frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}} \times 10^6} = \frac{\text{Clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

- CPI varies between programs on a given CPU

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

74

52

Concluding Remarks

- Cost/performance is improving
 - Due to underlying technology development
- Hierarchical layers of abstraction
 - In both hardware and software
- Instruction set architecture
 - The hardware/software interface
- Execution time: the best performance measure
- Power is a limiting factor
 - Use parallelism to improve performance

§1.11 Concluding Remarks

75

2016/2/15

© Ren-Song Tsay, NTHU, Taiwan

75