

CS/ECE 5381/7381
Computer Architecture
Spring 2023

Dr. Manikas

Computer Science

Lecture 2: Jan. 19, 2023

Assignments

- Quiz 1 – due Sat., Jan. 21 (11:59 pm)
 - Covers concepts from Module 1 (this week)

Quiz 1 Details

- The quiz is open book and open notes.
- You are allowed 90 minutes to take this quiz.
- You are allowed 2 attempts to take this quiz -
your highest score will be kept.
 - Note that some questions (e.g., fill in the blank) will need to be graded manually
- Quiz answers will be made available 24 hours after the quiz due date.

Fundamentals of Quantitative Design and Analysis

(Chapter 1, Hennessy and Patterson)

Note: some course slides adopted
from publisher-provided material

Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Trends in Technology

- Integrated circuit technology (chip)
 - Transistor density: 35%/year
 - Die size: 10-20%/year
 - Integration overall: 40-55%/year
- DRAM capacity: 25-40%/year (slowing)
 - This is the RAM in your computer
 - 8 GB (2014), 16 GB (2019)

Trends in Technology

- Flash capacity: 50-60%/year
 - 8-10X cheaper/bit than DRAM
Dynamic random access memory.
- Magnetic disk capacity: recently slowed to 5%/year
 - 8-10X cheaper/bit than Flash
 - Eventually to be replaced by Flash (SSD) as Flash costs decrease
 - 200-300X cheaper/bit than DRAM

Bandwidth and Latency

吞吐量

- Bandwidth or throughput
 - Total work done in a given time
 - 32,000-40,000X improvement for processors
 - 300-1200X improvement for memory and disks
- Latency or response time
 - Time between start and completion of an event
 - 50-90X improvement for processors
 - 6-8X improvement for memory and disks

We can process programs faster which means your computer reads faster.

Transistors and Wires

Q(10)

- Feature size
 - Minimum size of transistor or wire in x or y dimension
 - 10 microns in 1971 to .011 microns in 2017
 - 1 micron = 1 micrometer = 10^{-6} meters
 - Transistor performance scales linearly
 - Integration density scales quadratically

속도는 . • Wire delay does not improve with feature size!

<u>Year</u>	<u>Processor</u>	<u># TRANSISTORS</u>
1971	INTEL 4004	2,300
1980	MOTOROLA 68000	68,000
1989	INTEL 80486	~1 MILLION
RECENT	INTEL CORE 2 DUO	~1 BILLION

MORE TRANSISTORS \Rightarrow More Digital Logic
 \Rightarrow More functions on PROCESSORS.

Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Power and Energy

得電力
失電力

- Problem: Get power in, get power out
- Thermal Design Power (TDP)
 - Characterizes sustained power consumption
 - Used as target for power supply and cooling system
 - Lower than peak power, higher than average power consumption
- Clock rate can be reduced dynamically to limit power consumption
- Energy per task is often a better measurement

功率耗能：上端同塊電路中“子純”基底耗能。
也可以“基干次同期每子”來度量。

POWER AND ENERGY

WATT = JOULE / SECONDS

WATT (瓦特) = 焦耳 / 秒

"PHYSICS REVIEW" WATT (W) = JOULE / SECOND

POWER LAPTOP "BATTERY LIFE"

DESKTOP PLUGGED IN

→ POWER ^{耗能} DISSIPATION = "HEAT"

Dynamic Energy and Power

- Dynamic energy
 - Transistor switch from 0 -> 1 or 1 -> 0
 - $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2$
- Dynamic power
 - $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency switched}$
- Reducing clock rate reduces power, not energy

How Does CLOCK RATE AFFECT POWER?

COMPUTER HAS DIGITAL LOGIC

\Rightarrow 0 AND 1 SIGNALS

TRANSISTORS IMPLEMENT DIGITAL CIRCUITS

- SWITCH FROM 0 TO 1 CONSUMES POWER
- SWITCH FROM 1 TO 0 DISSIPATES POWER
- SWITCHING SPEED: HOW FAST WE SWITCH

SWITCHING SPEED = CLOCK RATE

EARLY PC (1980's) 5 MHz

$$5 \times 10^6 \frac{\text{cycle}}{\text{sec}}$$

modern PC

MODERN PC
MORE POWER

$$1 \text{ GHz} \quad 1 \times 10^9 \frac{\text{cycle}}{\text{sec}}$$

not making the
clock faster

$$\frac{\text{CYCLES}}{\text{SEC}} \Rightarrow \frac{\text{NUMBER OF SWITCHING}}{\text{sec}}$$

$$0 \rightarrow 1 \\ 1 \rightarrow 0$$

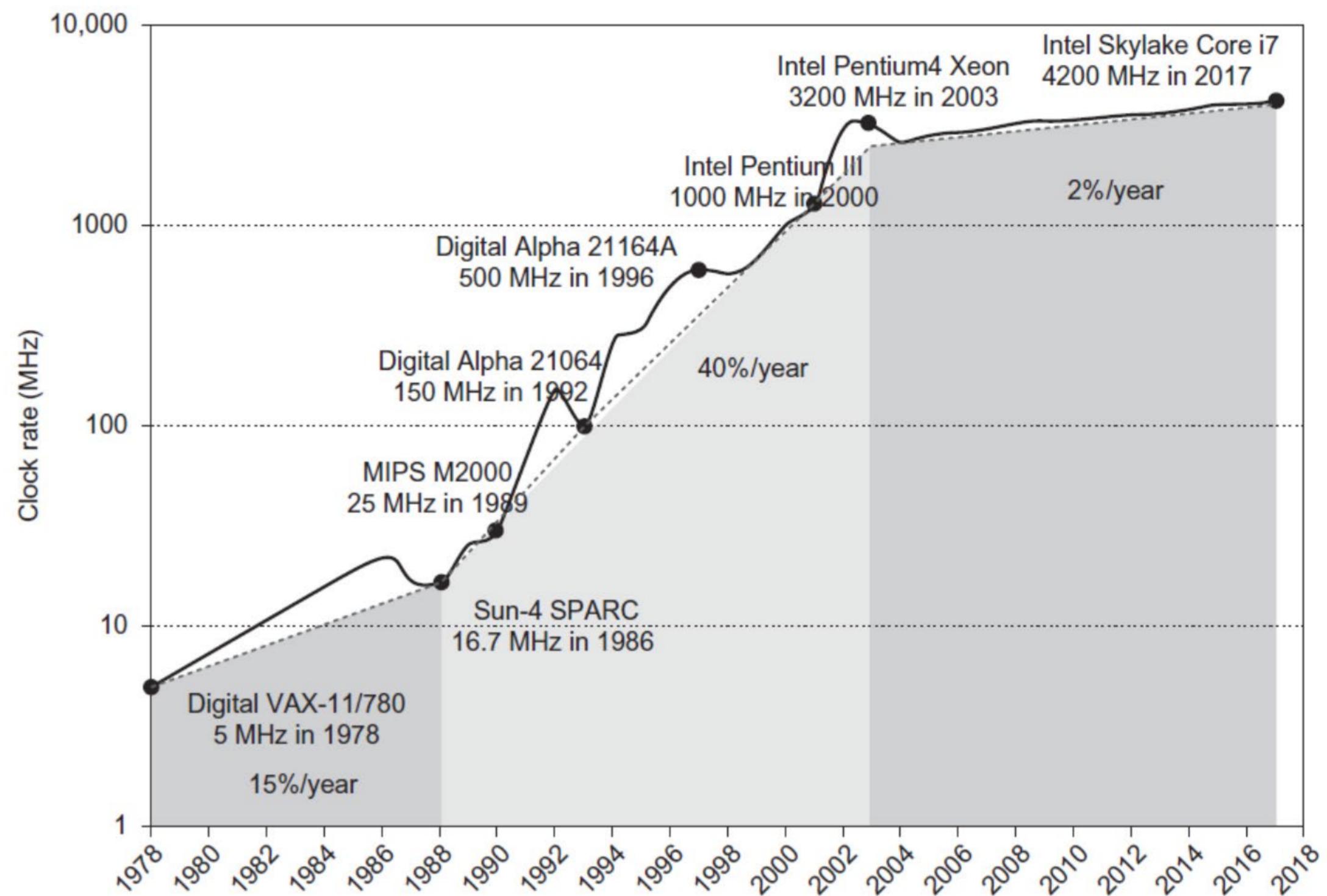
we need to find a way keep instead of increasing the clock rate.

How else we make the programs run faster do things in parallel

Power

- Intel 80386 (1986) consumed ~ 2 W
- 3.3 GHz Intel Core i7 (2017) consumes 130 W
- Heat must be dissipated from 1.5 x 1.5 cm chip
- This is the limit of what can be cooled by air

Power consumption over the years and we've kind of had to flatten that off because we don't want to consume all this power.



Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Trends in Cost

- Impact of Time
- Impact of Volume T_单40
- Cost of an Integrated Circuit
- Cost vs. Price

Impact of Time

- For a given design, manufacturing costs decrease over time
 - “learning curve” – initial design likely to have defects ~~defects~~
 - Later implementations are more reliable

Impact of Volume

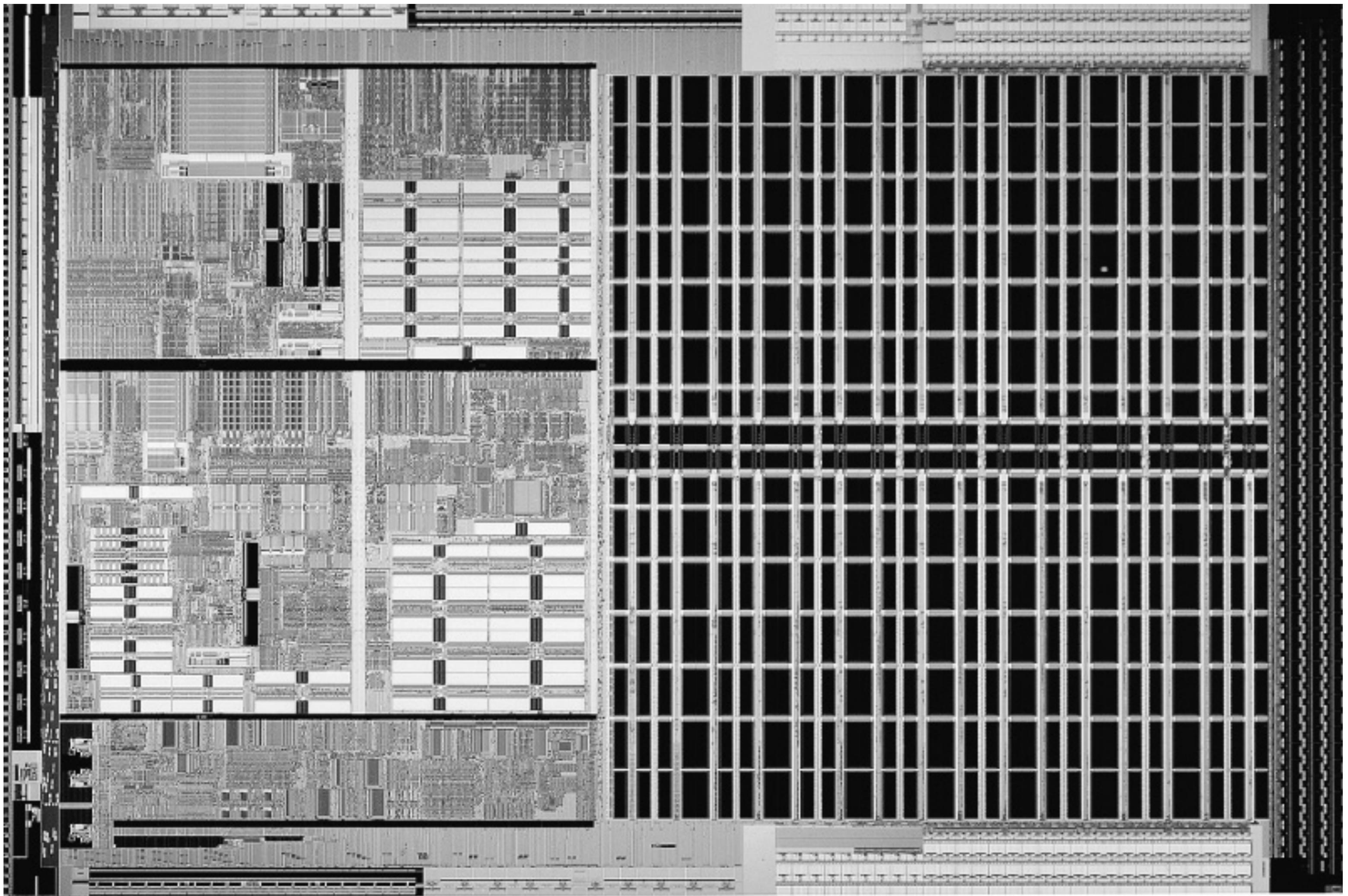
- Chip design and manufacturing cost
 - Cost per part (manufacturing)
 - NRE (Non-Recurring Expense)
 非经常开支
 - One-time cost for design and equipment set-up
 - Usually more expensive than cost per part
- High-volume designs can spread out NRE
 容量
 - Cheaper total cost per part

Cost of an Integrated Circuit (IC)

集成電路

- Computers contain chips (integrated circuits)
 - MPU, DRAM, etc.
- Thus, IC cost affects computer cost
- A given IC is manufactured as a *die*
はり
- Several dies (IC copies) are manufactured on a waffer
はり

晶片



Die for an AMD Opteron microprocessor (MPU)

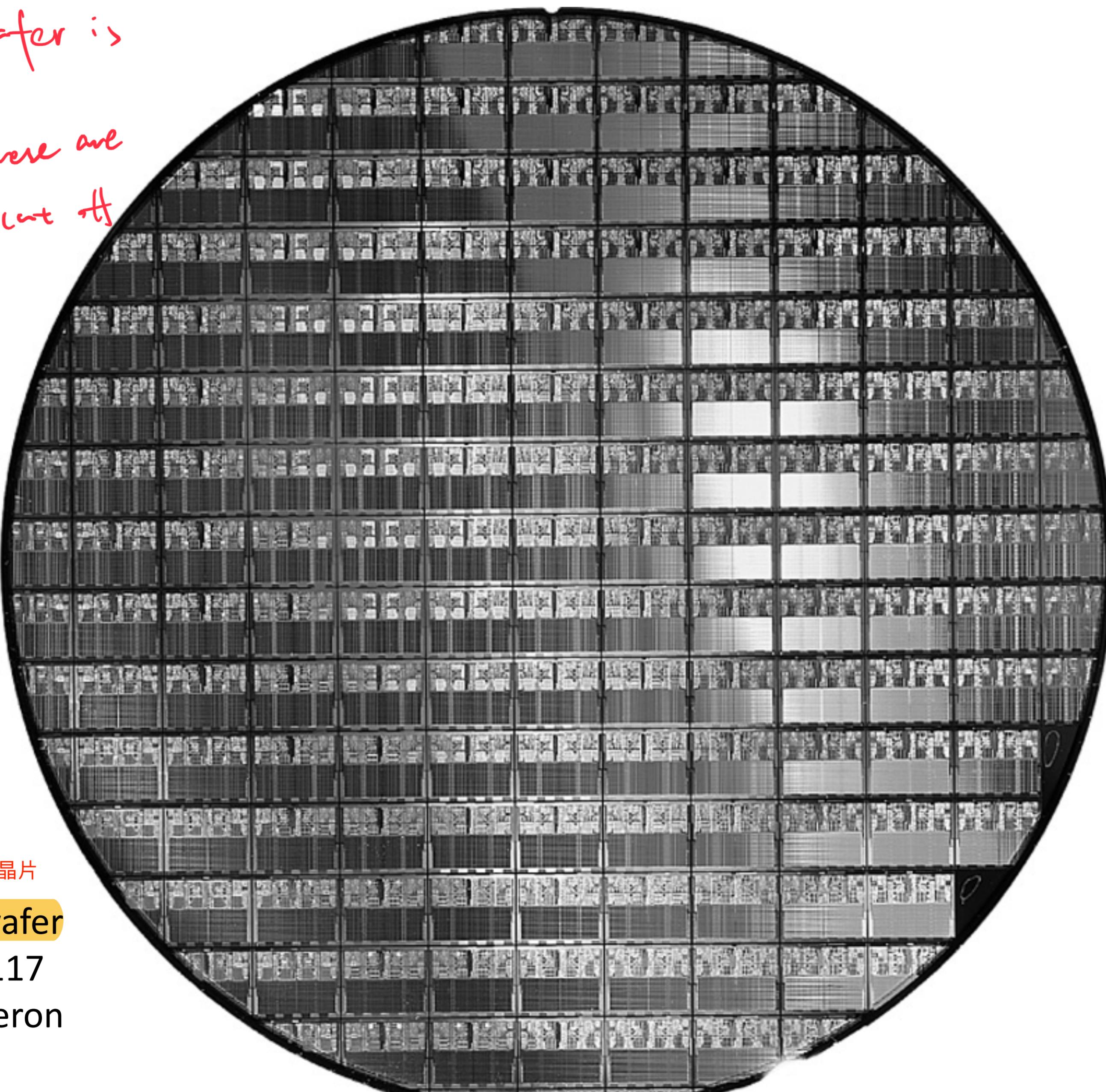
—
looks like a rectangle

die is the virtual processor with
all the transistors.



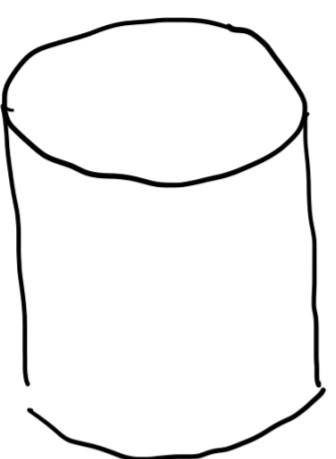
but this wafer is
around.
some of these are
going to be cut off
here right.

why?
why they
do that way?



Why Is A Wafer Round?

What happens is how they prepare it.



(99.9% pure)
Silicon is extruded ^{to} _{in}

So it actually looks like a cylinder

⇒ CUT INTO wafers



But!

chips are
rectangle.

designing it's a lot easier to
design on a rectangle than
a circle.

What they're doing is the extruding process if they recast it to a rectangle they can lose some of the risk as they lose some of that purity and you name that silicon as pure as possible what happens every time you have an impurity you have a defect. If you have a defect that affects your chip that means your chip doesn't work.

晶圆(Wafer): 晶圆圆是半导体集成电路的核心材料,是一种圆形的板

晶粒(Die): 很多四边形都聚集在圆形晶圆上。这些四边形都是集成电子电路的IC芯片

Dies per Wafer

→ rectangle.

- Area of wafer / area of die
晶圆面积 晶体面积
- Wafer is round, but die is square
- Estimate number of dies as

$$d = \pi r^2 = \pi (\frac{R}{2})^2$$

$$\frac{Dies}{Wafer} = \frac{\pi (Diameter_{wafer} / 2)^2}{Area_{die}} - \frac{\pi (Diameter_{wafer})}{\sqrt{2 (Area_{die})}}$$

Die Yield

- silicon was 99.9% pure not 100% plus
also you know they
have these clean
rooms but not
100% clean.
little speck
of dust could
get down there.
- Fraction of “good” dies on wafer
 - Dies without manufacturing defects

$$Yield_{die} = \frac{Yield_{wafer}}{[1 + (defects / unitArea)(area_{die})]^N}$$

where $N = \text{process - complexity factor}$

Typical contemporary values (for 40 nm process):

$$N = 12$$

$$\text{defects/unit_area} = 0.04 \text{ defects/cm}^2$$

Effects on IC cost

What's the cost per die

$$\text{cost}_{\text{die}} = \frac{\text{wafer}_{\text{cost}}}{(\text{dies/wafer})\text{yield}_{\text{die}}}$$

$$\text{cost}_{\text{IC}} = \frac{\text{cost}_{\text{die}} + \text{cost}_{\text{testing_die}} + \text{cost}_{\text{packaging_and_final_test}}}{\text{yield}_{\text{final_test}}}$$

Cost vs. Price

say:::

- Margin = sales price – manufacturing cost
- Margins cover overhead costs
 - Salaries, benefits, utilities, equipment, maintenance
 - R & D, sales, manufacturing

Example 1.6.1

- New chip with code name “Peruna”
 - Die area is 250 mm^2
 - To be fabricated on wafer with diameter of 300 mm
12寸
- Fabrication parameters: *製造参数*
 - Estimated defect rate = 0.03 per cm^2
缺陥率
 - Wafer yield = 100%
 - Process-complexity factor $N = 12$

Example 1.6.1

- a. How many Peruna dies can we fabricate on a wafer?
- b. What is the die yield? 没有缺陷
- c. If we can make a \$20 profit per defect-free chip, how much profit can we make for a wafer of Pernuna dies?

(a) DIES PER WAFER?

GIVEN: DIE AREA = 250 mm^2

WAFER DIAMETER = 300 mm

DIES PER WAFER

$$= \frac{\pi [(WAFER \text{ DIAM})/2]^2}{DIE \text{ AREA}} - \frac{\pi (WAFER \text{ DIAM})}{\sqrt{2} (\text{DIE AREA})}$$

$$= \frac{\pi [(300 \text{ mm})/2]^2}{250 \text{ mm}^2} - \frac{\pi (300 \text{ mm})}{\sqrt{2} (250 \text{ mm}^2)}$$

$$= \frac{\pi (150)^2}{250} - \frac{\pi (300)}{\sqrt{500}}$$

$$= \pi \left[\frac{150^2}{250} - \frac{300}{\sqrt{500}} \right] = \pi \left[\frac{3 \times 50 \times 30 \times 5}{250} - \frac{300}{\sqrt{5} \times 10} \right]$$

$$= \pi \left[90 - \frac{30}{\sqrt{5}} \right] = \pi [90 - 13.4164] \approx 240.594$$

$$\frac{\text{DIES}}{\text{WAFER}} = \text{COMPLETE DIES} \Rightarrow \text{INTEGER} \Rightarrow \boxed{240}$$

(b) DIE YIELD

$$= \frac{WAFER}{[1 + (\text{DEFECTS/UNIT AREA}) (\text{DIE AREA})]^N}$$

GIVEN: WAFER YIELD = 100% = 1

• DEFECTS PER UNIT AREA

$$= \text{DEFECT RATE} = 0.03 \text{ per/cm}^2$$

$$\bullet N = 12$$

$$\bullet \text{DIE AREA} = 250 \text{ mm}^2$$

$$\text{DIE YIELD} = \frac{1}{[1 + (0.03 \text{ cm}^{-2}) \cdot (250 \text{ mm}^2)]^{12}}$$

CONVERT DEFECTS PER UNIT AREA

$$\frac{0.03}{\text{cm}^2} \left[\frac{\text{cm}^2}{(100 \text{ mm})^2} \right] = \frac{0.03}{100 \text{ mm}^2}$$

$$\text{DIE YIELD} = \frac{1}{[1 + (\frac{0.03}{100}) (250)]^{12}} = \frac{1}{(1 + 0.075)^{12}} \approx 0.420$$

*not good
lower than half.*

(c) PROFIT PER WAFER?

a) 240 DIES

b) DIE YIELD = 0.420

GIVEN: PROFIT per DEFECT-FREE CHIP = \$20.

$$\begin{aligned}\text{PROFIT PER WAFER} &= (\text{DIE YIELD}) (\# \text{DIES}) (\text{PROFIT}) \\ &= (0.42) (240) (\$20) \\ &= \boxed{\$2016}\end{aligned}$$

Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Dependability

- Module reliability
 - Mean time to failure (MTTF)
 - Mean time to repair (MTTR)
 - Mean time between failures (MTBF) = MTTF + MTTR
 - Availability = MTTF / MTBF
- Failures in Time (FIT)
 - Rate of failures per *billion* hours
 - $\text{MTTF} = 10^9/\text{FIT}$

Example 1.7.1

- Our Peruna chip has the following parameters
 - FIT = 150
 - MTTR = 2 days
- a. What is the MTTF of our chip?
 - b. What is the availability of our chip?

Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Measuring, Reporting, and Summarizing Performance

- Performance = “speed” of computer
- User view – how fast does my program run?
 - Execution (response) time: time between start and end of an event
- Web administrator view – how many transactions/hour?
 - Throughput: total amount of work done in a given time

Measuring Execution Time

- User view: “Wall-clock” time
 - E.g., process started at 2:00 pm, ended at 2:10 pm
- Actual execution time: CPU time
 - Time that processor is actually computing
 - Ignores waiting time for I/O, other programs

Benchmarks

- Common set of programs to run on different computers
- SPEC (Standard Performance Evaluation Corporation)
 - Started in 1988 by group of workstation vendors.
 - Non-profit org - used as indep. testing source.
 - Goal: produce benchmarks that measure "real" performance.
 - Becoming the standard for performance measurement.

SPEC Benchmarks

- Benchmark types:
 - Open Systems: benchmarks for PCs, servers
 - High Performance: supercomputers
 - Graphics: high-end graphical workstations (CAD, simulators, games)

Reporting Performance Results

- Experiments should be *reproducible*
 - Report sufficient information such that another researcher can get the same results (assuming he/she follows your exact approach)
- SPEC benchmark reports require detailed info on computer, compiler, program parameters, etc.

Summarizing Performance Results

- We want to compare two computers (A and B)
- Run several different SPEC benchmarks on both computer A and B
- For computer j, SPECRatio for a given benchmark is

$$SPECRatio_j = \frac{Exec_Time_{ref}}{Exec_Time_j}$$

where ref = a reference computer (baseline)

shorter exec time \Rightarrow larger SPECRatio

Outline

- 1.1 Introduction
- 1.2 Classes of Computers
- 1.3 Defining Computer Architecture
- 1.4 Trends in Technology
- 1.5 Trends in Power and Energy in Integrated Circuits
- 1.6 Trends in Cost
- 1.7 Dependability
- 1.8 Measuring, Reporting, and Summarizing Performance
- 1.9 Quantitative Principles of Computer Design

Quantitative Principles of Computer Design

1. Parallelism
2. Principle of locality
3. Focus on common case
4. Amdahl's Law
5. Processor Performance Equation

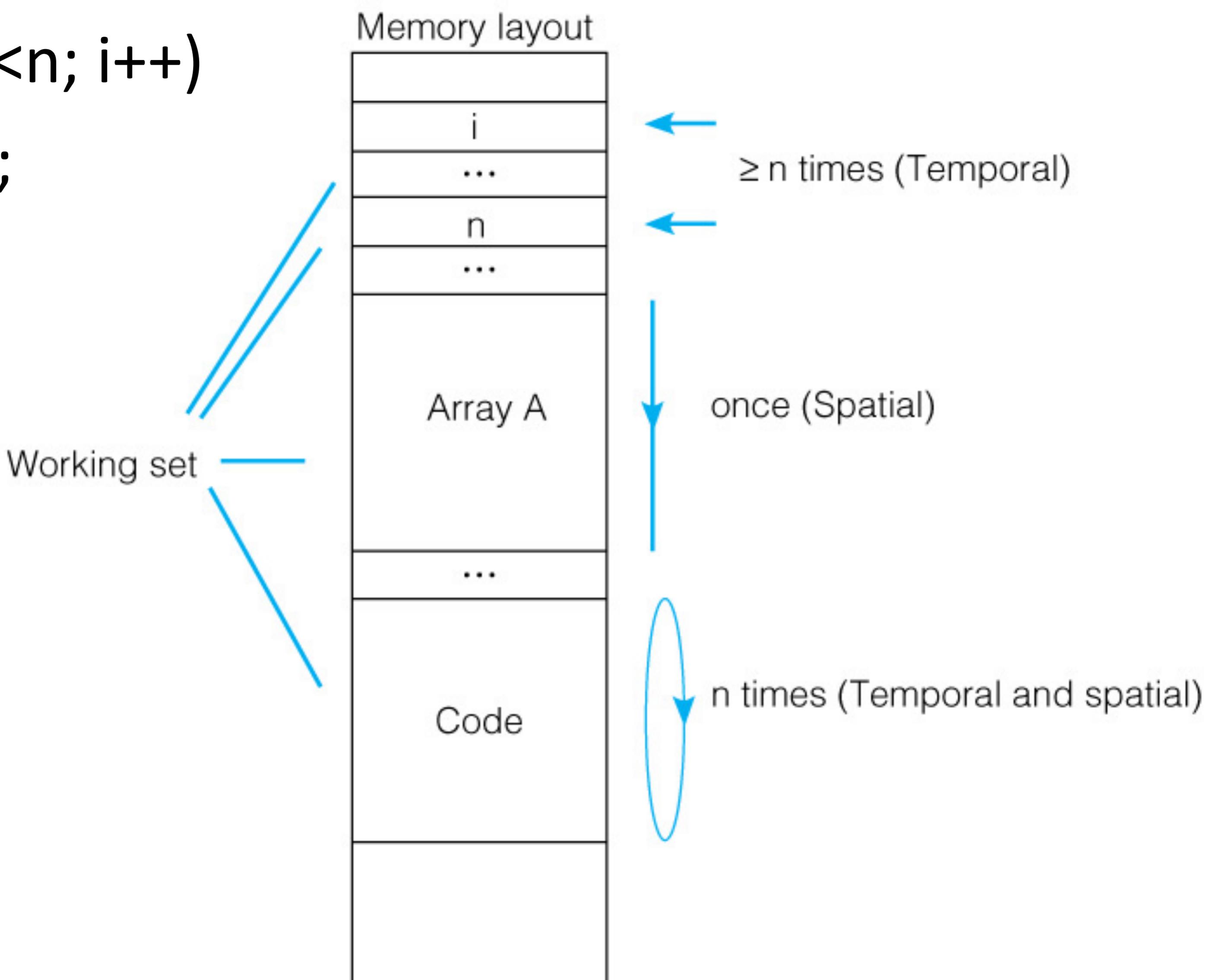
Take Advantage of Parallelism

- Parallel operations = faster execution time
- System level – multiple processors
- Processor level – can we do two (or more) tasks at the same time? (pipelining)

Principle of Locality

- Programs tend to reuse recent data and instructions
- **Temporal Locality:** items that have been recently accessed will likely be accessed again soon
- **Spatial Locality:** items that are near each other (memory address) will likely be accessed close together in time

```
for (i=0; i<n; i++)  
    A[i] = 0;
```



Copyright © 2004 Pearson Prentice Hall, Inc.

Focus on the Common Case

- Most designs have trade-offs: optimizing for one objective will degrade another objective
- For a given design, optimize for main purpose
 - Graphics workstation – optimize mathematical operation speed, at expense of power
 - Laptop – optimize battery life, at expense of operation speed