



دانشکده مهندسی کامپیوتر

## درس معماری کامپیوتر

الهام چشمی خانی

# Introduction

# Introduction

## • الهام چشمی خانی

- ✓ فارغ التحصیل دکترای مهندسی کامپیوتر از دانشگاه شریف
- ✓ زمینه‌های کاری: معماری کامپیوتر، حافظه‌های نوظهور، حافظه‌های نهان، پردازش در حافظه‌ها، سیستم‌های ذخیره‌سازی، قابلیت اطمینان، طراحی سیستم‌های مطمئن، طراحی شتاب‌دهنده‌ها، طراحی توامان سخت‌افزای و نرم‌افزاری
- ✓ مرکز HPC دانشگاه شریف، شرکت Huawei

- ✓ Google Scholar -> Elham Cheshmikhani
- ✓ cheshmikhani\_e@aut.ac.ir

## • تیم دستیاران:

- سردستیار: آقای فرشاد احدی‌نیا
- ✓ farshad.ahn@aut.ac.ir
- دستیاران مسئول: آقایان بردیا اردکانیان و محمدرضا صادقیان
- ✓ b.ardakanian@aut.ac.ir, sadeghian.m79@gmail.com

- سایر دستیاران: خانم‌ها ریحانه آهنی و صبا رمضانی و آقای آراد فیروزکوهی

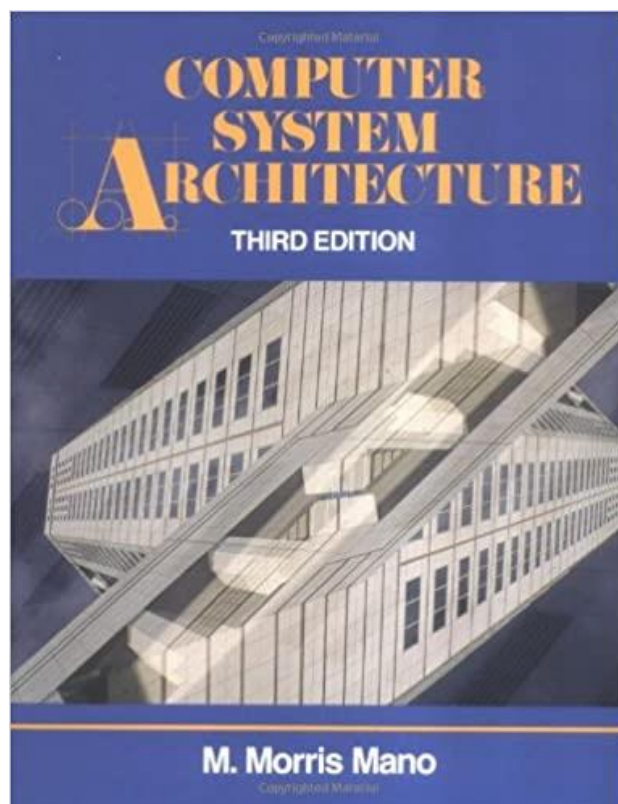
• مراجع اصلی:

- David A Patterson and John L. Hennessy, "Computer Organization and Design RISC-V Edition: The Hardware Software Interface," Morgan Kaufmann Series

- M. Morris Mano, "Computer System Architecture," Pearson Education (US)

• مراجع دیگر:

- Linda Null, "The essentials of computer organization and architecture," Jones & Bartlett Publishers
- David Harris, "Digital Design and Computer Architecture," Elsevier
- M. Mazidi, S. Naimi, "ARM Assembly Language Programming & Architecture," MicroDigitalEd



• تمرین: (شرایط: ۲۴ ساعت - کپی)

➤ ۳۰٪

➤ ۲۰٪

➤ ۲۵٪

➤ ۱۵٪

➤ ۱۰٪

• میان ترم:

• پایان ترم:

• پروژه:

• کوئیز:

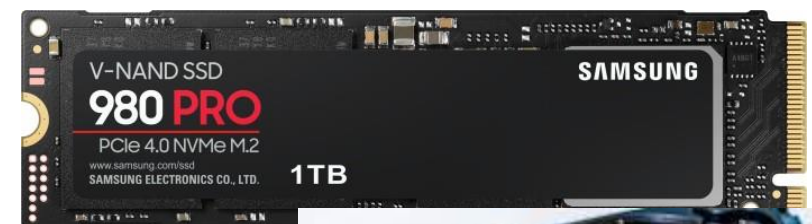


## • طراحی سیستم اساساً بهتر / بهینه

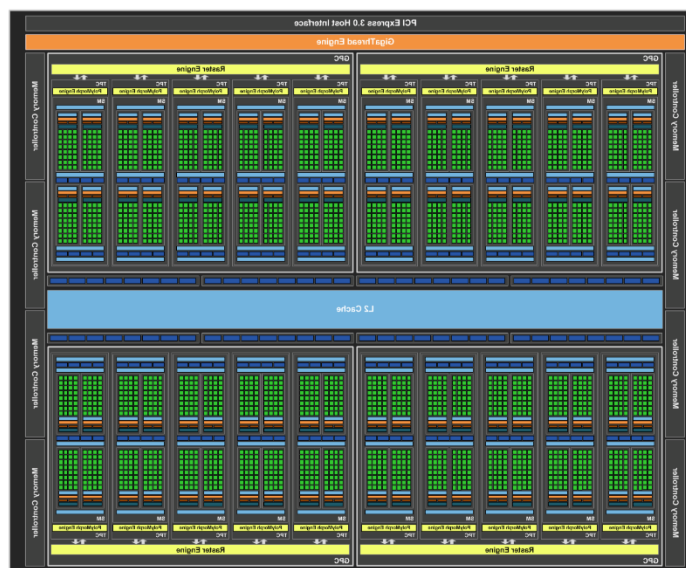


Hierarchical Memory

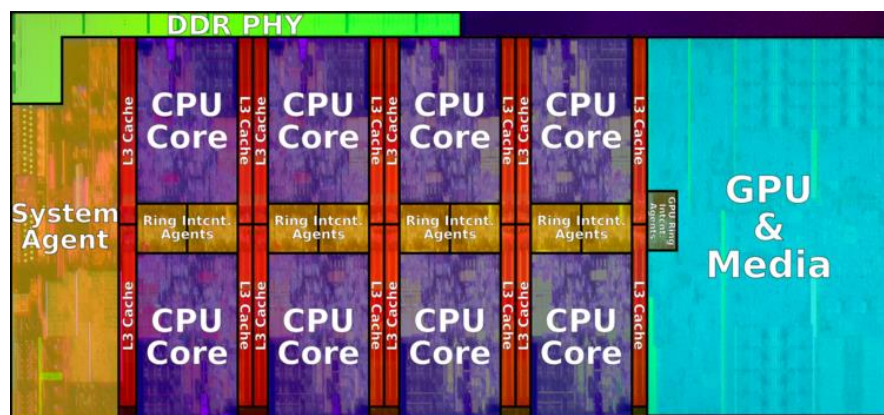
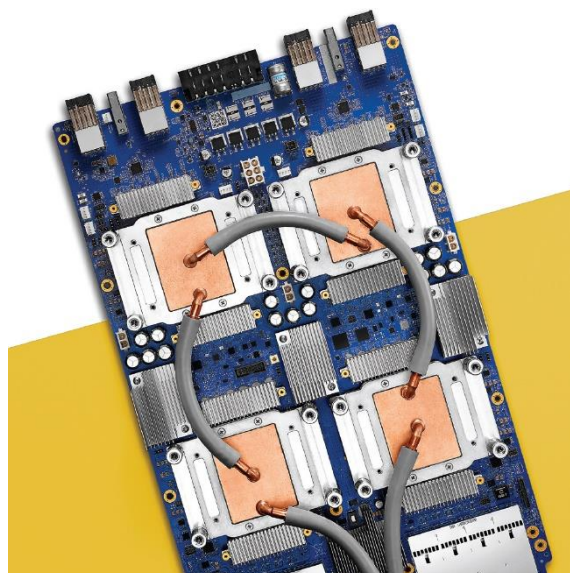
Persistent Storage



Graphics Processing

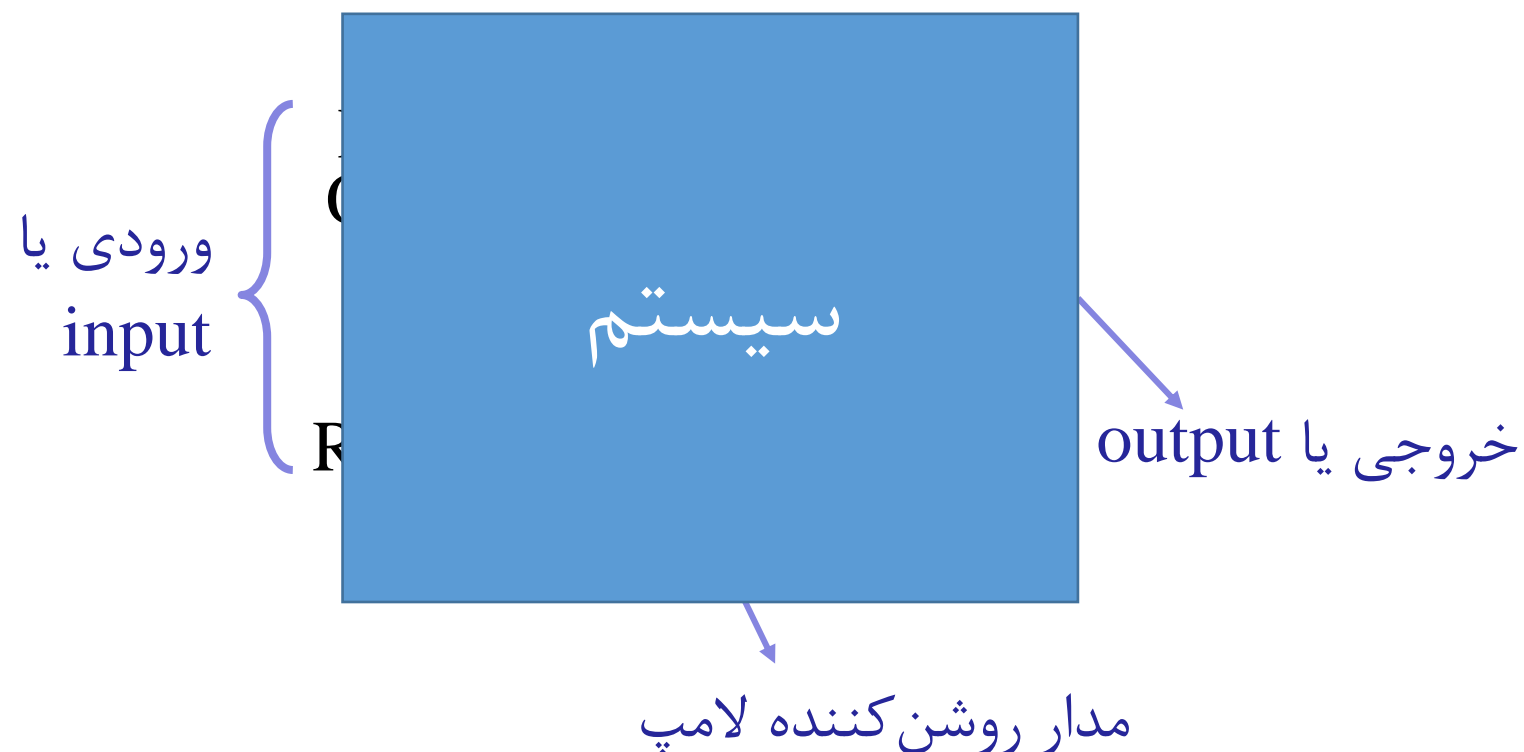


Heterogenous  
processors and  
accelerators

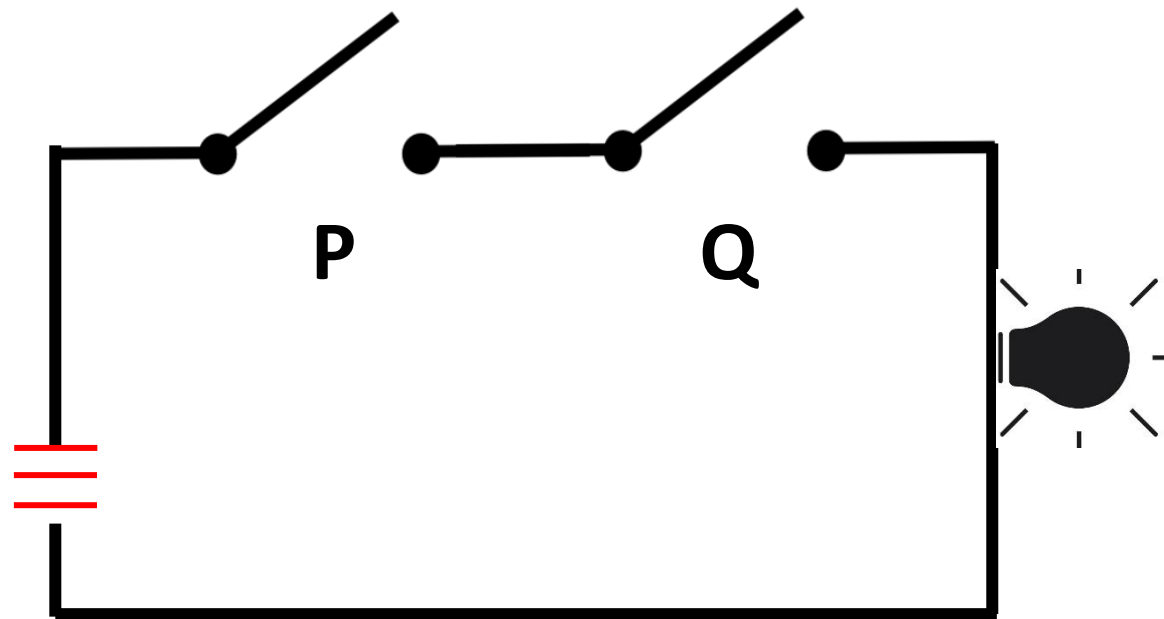


- System is a set of components work as a whole to achieve a goal

➤ Inputs, behavior, and outputs



➤ Behavior is a function that translates input to outputs



لامپ روشن می شود،  
اگر هر دو کلید P و Q بسته باشد

Only if  
اگر و تنها اگر



• سیستم اساساً قابل اطمینان، دارای امنیت، و امن

➤ Reliable, secure, and safe

• سیستم اساساً بهینه از نظر انرژی

➤ Memory-centric

• سیستم اساساً کم تاخیر

➤ Low latency

• سیستم اساساً کارایی بالا

➤ High performance

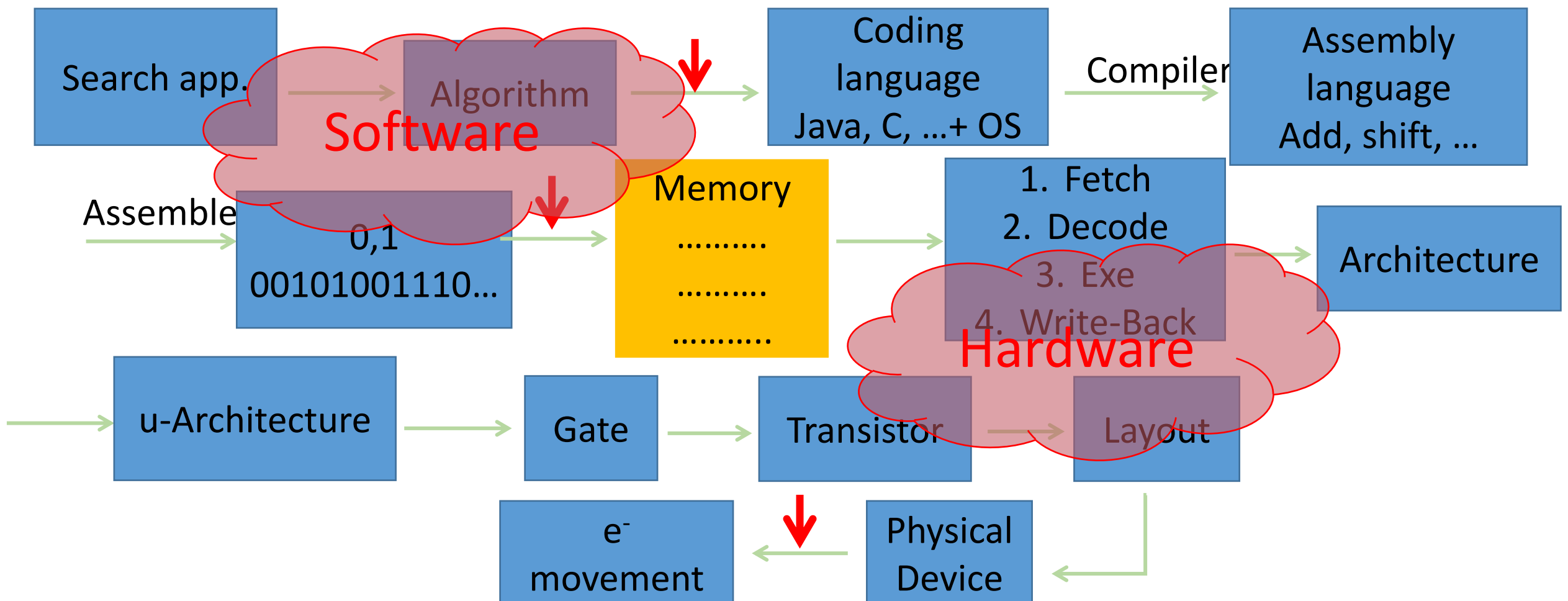
• معماری در حوزه‌های هوش مصنوعی، یادگیری ماشین، درمان و سلامتی

➤ AI, ML, Medicine and health

حوزه تحقیقاتی و تدریس در هم تنیده شده

## • Computer

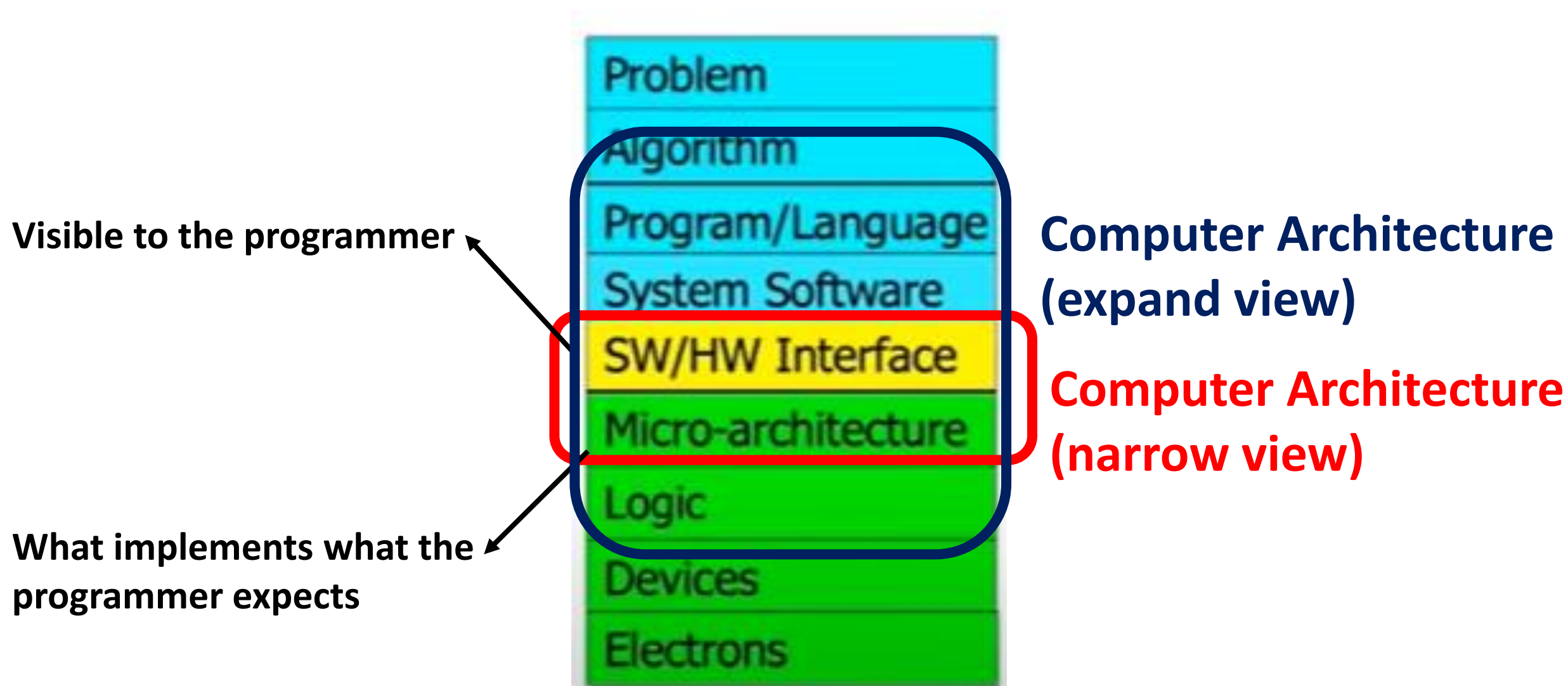
- A **general purpose** device that can be **programmed** to carry out a **set** of arithmetic or logical operations (wiki)
- A system that can execute bunch of instructions



# The Transformation Hierarchy

• ترجمه مشکل به سطوح مختلف تجزید

✓ تبدیل حرکت الکترون‌ها به راه حل مشکل خود



# دسته‌بندی کامپیوترها

- کامپیوترهای رومیزی

- در حال انقراض

- سرورها

- قابلیت بیشتر
- خدمات دادن به کاربران بیشتر

- کامپیوترهای نهفته

- موبایل‌ها
- لوازم خانگی
- در اتومبیل‌ها
- ...

# بهینه ← کاری؟

## • کدام سیستم بهتر است؟

- گوشی ساده با پردازنده ARM11 با فرکانس 400MHz
- پردازنده پنتیوم با 66MHz

## • چه آیتم‌هایی وجود دارد که تعیین کننده در کارایی باشد؟

## • برنامه محک

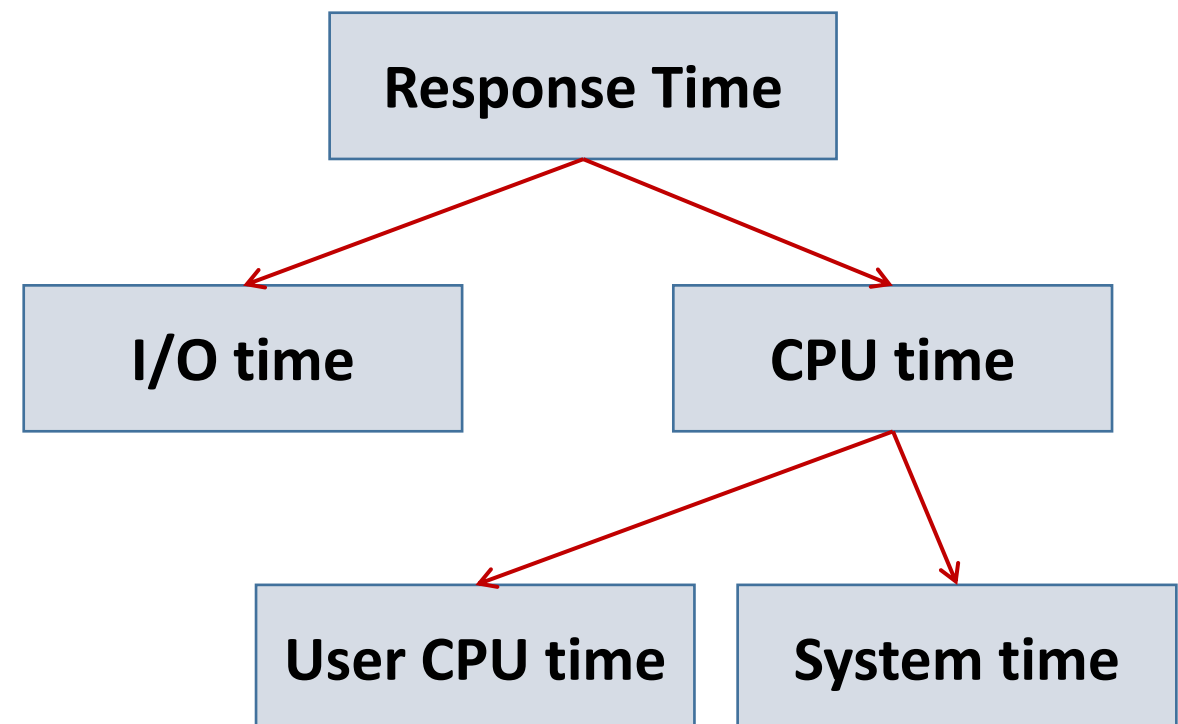
- نحوه نوشتن برنامه، الگوریتم برنامه
- زبان برنامه‌نویسی (مقایسه C و Java)
- نوع پردازنده و سرعت حافظه سیستم
- سیستم عامل
- ارتباطات بین اجزا
- ISA

- Response time
  - Total time to complete a task
  - **OR**: how long it takes to do a task
    - Processing, I/O, OS overhead, idle time
- Response time consists of:
  - CPU time
    - CPU time spent on a program
    - The amount of time used for processing instructions
  - I/O time
    - Time elapsed to wait for I/O operations

- CPU time: Time spent processing a given job
  - Discount I/O time, other jobs' shares
- CPU time comprises
  - User CPU time
    - CPU time spent on a program
  - System time
    - CPU time spent on OS doing tasks on behalf of program
- Performance =  $1/\text{response time}$



- Performance = 1/response time
- System Performance
  - 1/elapsed time
- CPU performance
  - 1/user CPU time



# CPU Time

$$\begin{aligned}\text{CPU Time} &= \# \text{ of Clock Cycles} * \text{Clock Cycle Time} \\ &= \# \text{ of Clock Cycles} * T \\ &= \# \text{ of Clock Cycles} / F\end{aligned}$$

F: CPU Frequency (Hz)

- Example: Which one is faster? A or B?

Computer	Frequency	# Clks
A	3.6 GHz	2000
B	2.6 GHz	1500

# CPU Time

- CPI: # of clocks per instruction
- # of Clock Cycles
  - # of Clks = # of instructions \* CPI
- Example: CPU time = ?

# of instr.	4,000
CPI	3
Freq.	2.6 GHz

# CPU Time

- If different instructions have different CPI
  - CPI: Average CPI (Average cycles per instr.)

$$CPI = \frac{1}{n} \sum_{i=1}^n \# \text{ of Clk for Instr.}_i$$

N: classes of instr.

- CPU time = # of instr. \* CPI \* T
- # of instructions or Instruction Count (IC)
  - code size or lines of code????

# Performance Evaluation

- Consider two CPUs
  - CPU1 runs faster than CPU2 on program A
  - CPU2 runs faster than CPU1 on program B
- Which CPU is faster? or
  - Which program should we choose to compare CPUs?
- **Benchmarks**: a **set** of programs used to compare performance of processors

# Performance Evaluation

## ■ Benchmarks

- Why a **set** of programs?
- Why not just running a simple program?
  - Agreeing on one simple program is very hard
  - Designers can optimize their CPUs towards fast running that simple program

## ■ Standard benchmarks

- **SPEC**: Standard Performance Evaluation Corporation
  - SPCE2000, SPEC2006, SPEC2017

# Performance Metric

- **MIPS**: Million Instruction Per Second
- MIPS Drawbacks
  - Not into account capabilities of instructions
    - Comparing two different ISAs is not fair
  - Not realistic metric even on same CPU
    - $\text{MIPS}(A) > \text{MIPS}(B) \rightarrow \text{Performance}(A) > \text{Performance}(B)??$



# Performance Metric

- **MFLOPS**: Million Floating Point Operation Per Second
  - Popular in scientific computing
  - Drawbacks
    - Ignore other instructions (Load/Store)
    - Not all FP operations have common format
    - Depends on how FP-intensive program is
- So, how determine the best processor?
  - What does the best mean?
    - Performance, cost, power, energy, ISA, .... ?????

- **علم و هنر طراحی سیستم‌هایی برای محاسبات**

- سخت‌افزار

- رابط

- SystemSW

- مدل برنامه‌نویسی

- **برای دستیابی به مجموعه‌ای از اهداف طراحی**

- مثال: دستیابی به بالاترین کارایی برای برنامه‌های X, Y, Z

- مثال: طراحی سیستمی با بالاترین عمر باتری

- ...

- **نوع ماشین و سیستم مطرح نیست بلکه هدف طراحی مطرح است**

- طراحی یک سوپرکامپیوتر متفاوت از طراحی یک گوشی هوشمند ← بسیاری از نکات طراحی اصلی یکسان است

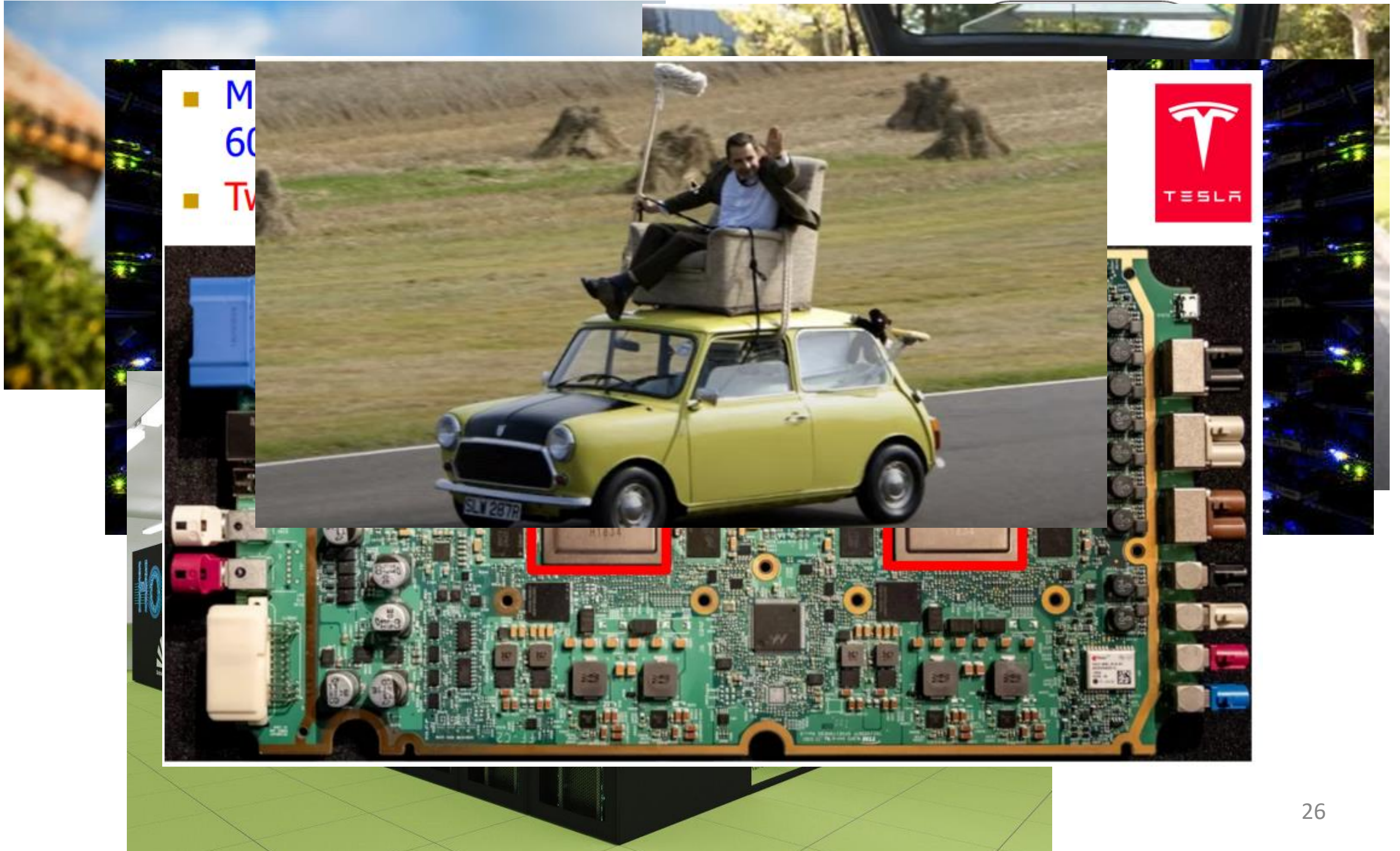
- استفاده از مفاهیم برای دسترسی به طراحی موردنیاز

# هدف معماری کامپیوتر؟

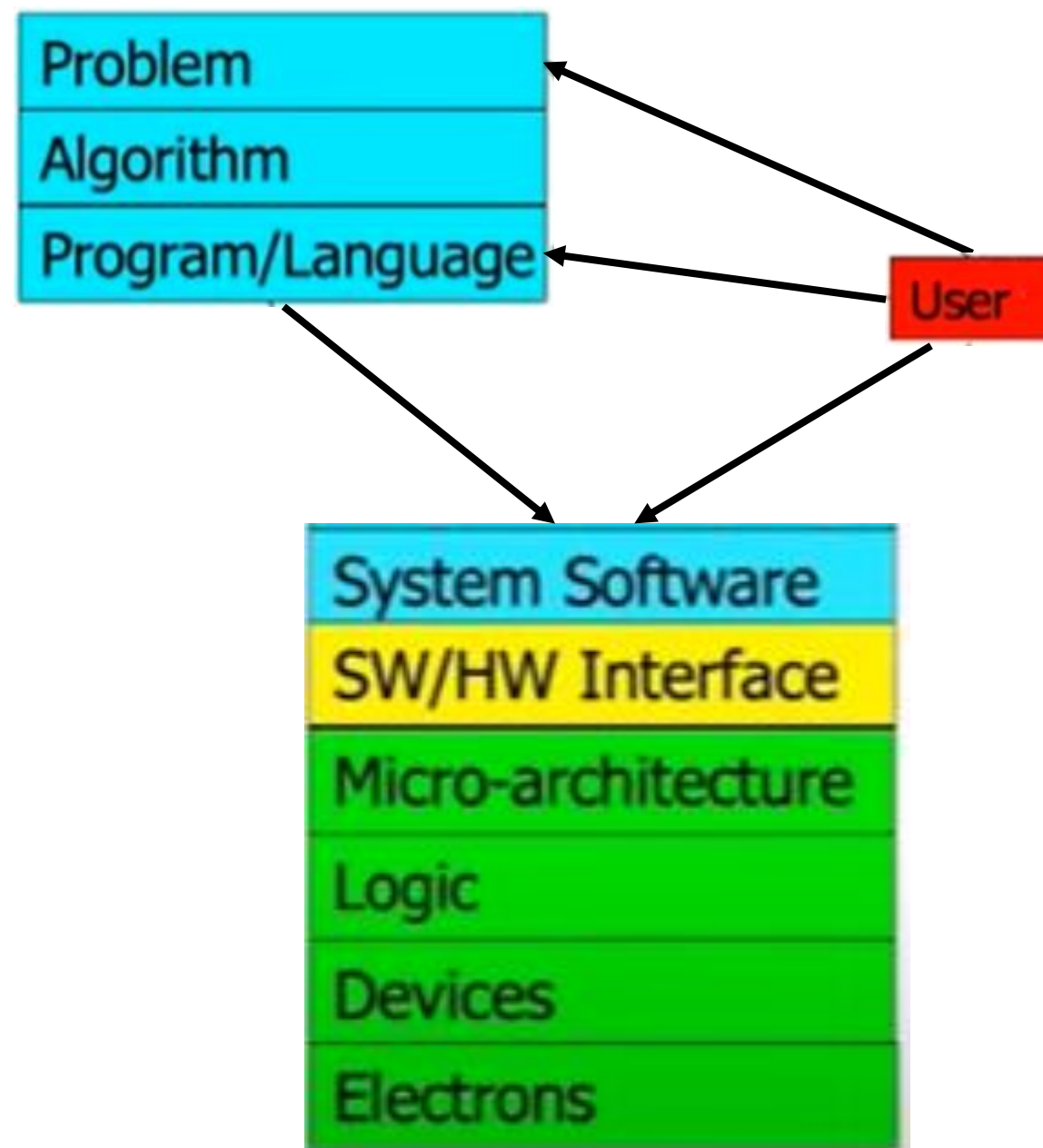
- **علم و هنر طراحی، انتخاب، و اتصال اجزای سخت‌افزاری، و طراحی رابط سخت‌افزاری/نرم‌افزاری برای ساخت سیستم محاسباتی برای رسیدن به اهداف خاص مانند انرژی، هزینه، کارایی، مساحت و ...**
- **طراحی سیستم بهتر**
  - کامپیوتر سریع، ارزان، کوچک، کم مصرف، قابل اطمینان و ...
  - **فعال کردن اپلیکیشن‌ها/نیازهای جدید**
  - خودروهای خودران، واقعیت مجازی، سلامتی، و ...
  - **درک اینکه چرا کامپیوترها به این شکل امروزی کار می‌کند**



# Different platforms and goals

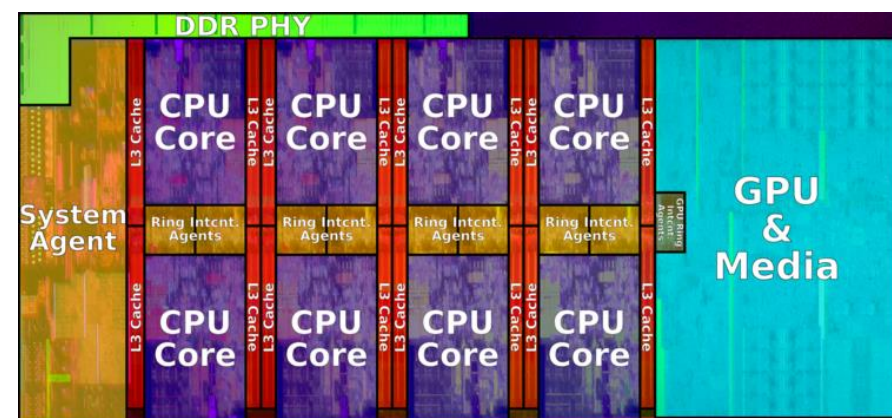


# هدف معماری کامپیوتر؟





# Different platforms and goals



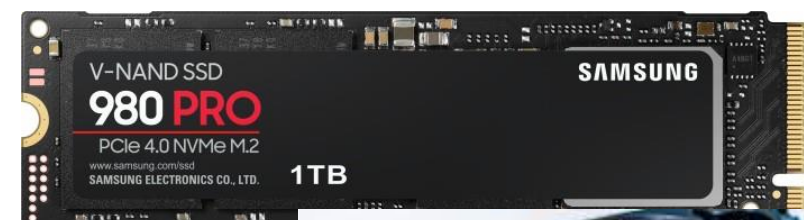
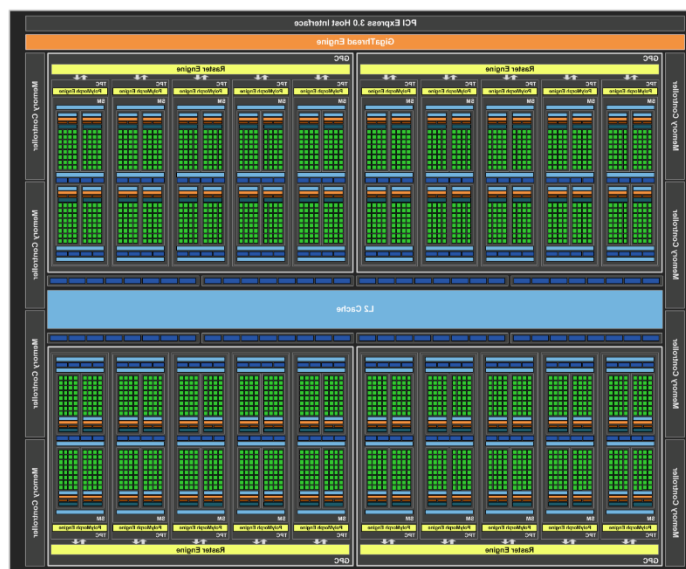
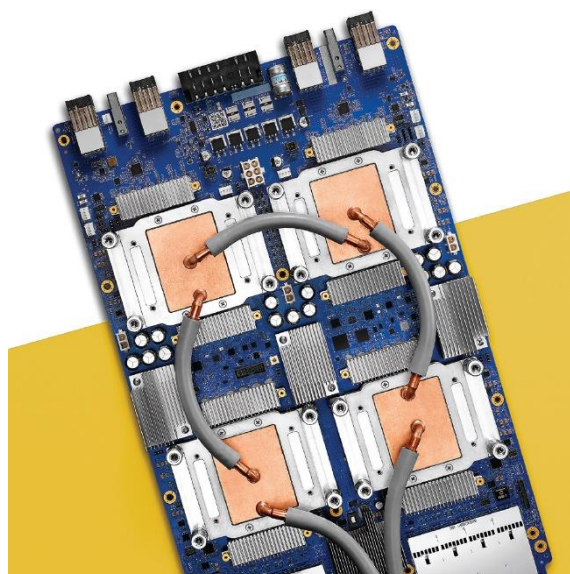
**Heterogenous  
processors and  
accelerators**



**Hierarchical Memory**

**Persistent Storage**

**Graphics Processing**



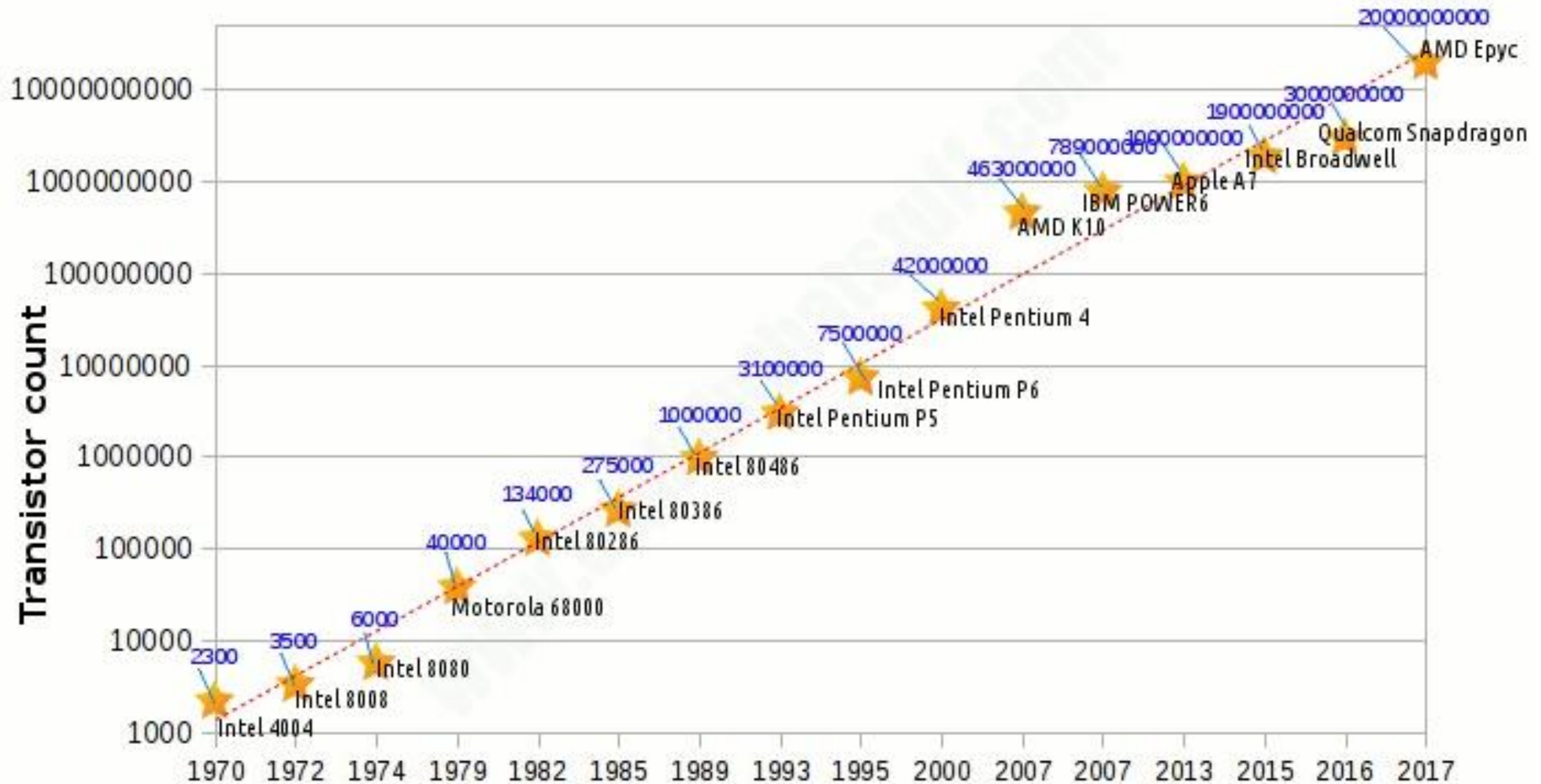
# There's Plenty of Room at the Bottom

From Wikipedia, the free encyclopedia

**"There's Plenty of Room at the Bottom: An Invitation to Enter a New Field of Physics"** was a lecture given by [physicist Richard Feynman](#) at the annual [American Physical Society](#) meeting at [Caltech](#) on December 29, 1959.<sup>[1]</sup> Feynman considered the possibility of direct manipulation of individual atoms as a more powerful form of synthetic chemistry than those used at the time. Although versions of the talk were reprinted in a few popular magazines, it went largely unnoticed and did not inspire the conceptual beginnings of the field. Beginning in the 1980s, nanotechnology advocates cited it to establish the scientific credibility of their work.



## 50 Years of Moore's law



Feynman considered some ramifications of a general ability to manipulate matter on an atomic scale. He was particularly interested in the possibilities of denser computer circuitry and microscopes that could see things much smaller than is possible with scanning electron microscopes. These ideas were later realized by the use of the scanning tunneling microscope, the atomic force microscope and other examples of scanning probe microscopy and storage systems such as Millipede, created by researchers at IBM.

Feynman also suggested that it should be possible, in principle, to make nanoscale machines that "arrange the atoms the way we want" and do chemical synthesis by mechanical manipulation.

He also presented the possibility of "swallowing the doctor", an idea that he credited in the essay to his friend and graduate student Albert Hibbs. This concept involved building a tiny, swallowable surgical robot.

## RESEARCH

---

### REVIEW SUMMARY

#### COMPUTER SCIENCE

## There's plenty of room at the Top: What will drive computer performance after Moore's law?

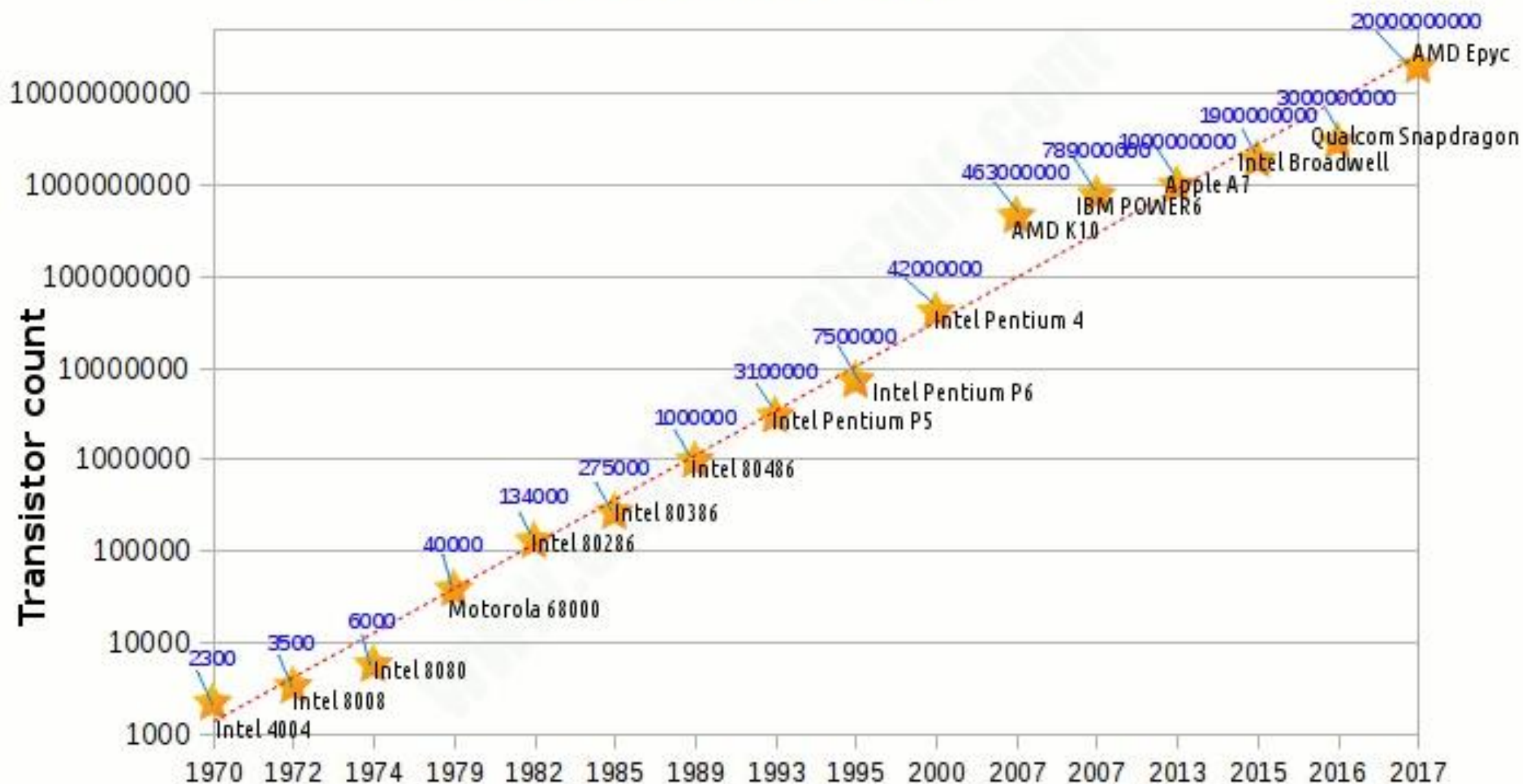
Charles E. Leiserson, Neil C. Thompson\*, Joel S. Emer, Bradley C. Kuszmaul, Butler W. Lampson, Daniel Sanchez, Tao B. Schardl

**BACKGROUND:** Improvements in computing power can claim a large share of the credit for many of the things that we take for granted in our modern lives: cellphones that are more powerful than room-sized computers from 25 years ago, internet access for nearly half the world, and drug discoveries enabled by powerful supercomputers. Society has come to rely on computers whose performance increases exponentially over time.

in computing power stalls, practically all industries will face challenges to their productivity. Nevertheless, opportunities for growth in computing performance will still be available, especially at the “Top” of the computing-technology stack: software, algorithms, and hardware architecture.

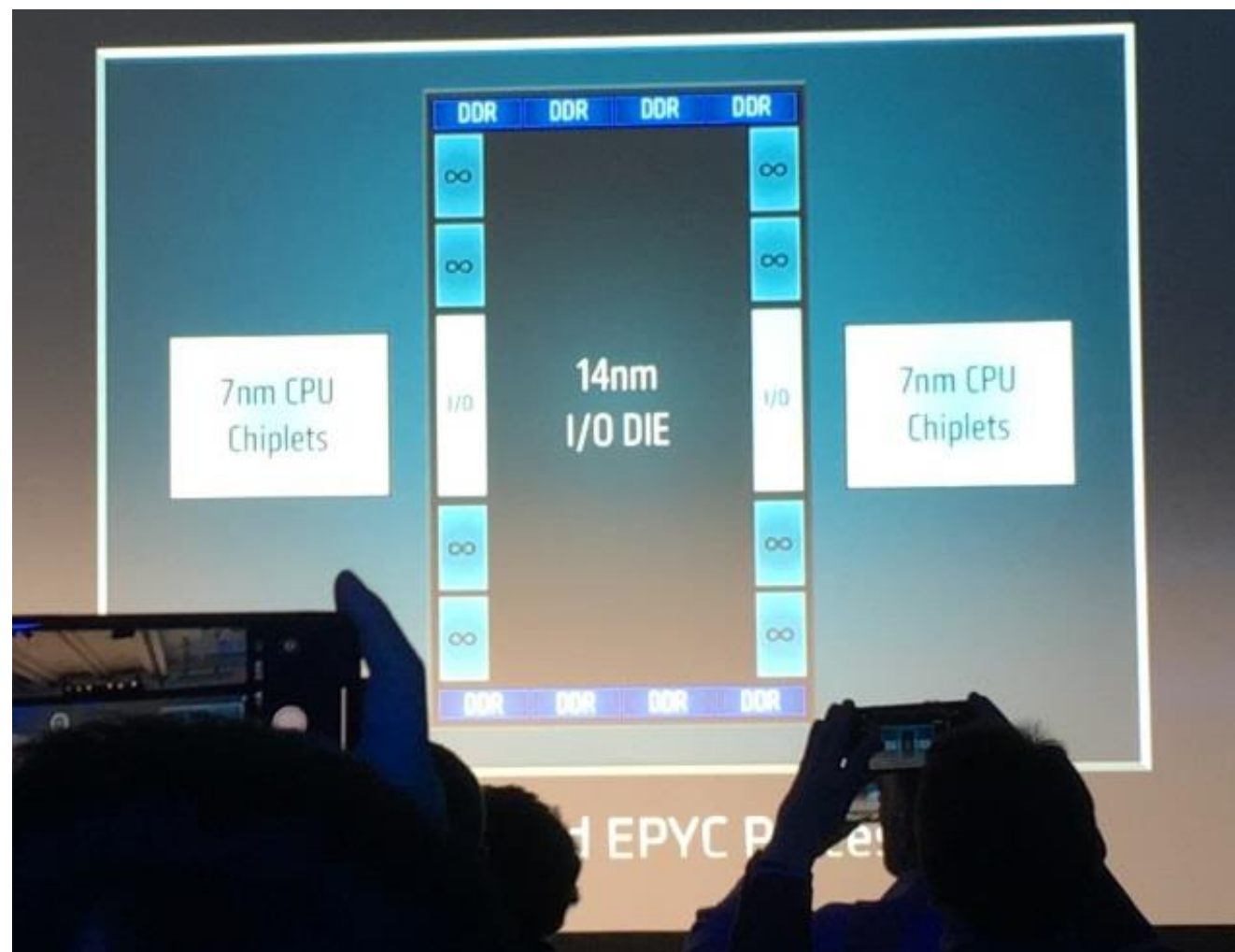
**ADVANCES:** Software can be made more efficient by performance engineering: restructur-

## 50 Years of Moore's law





- **Lisa Su:** “Forget about whether Moore’s Law is alive or dead. That’s an interesting conversation. But the more interesting thing is applications need for us to be above Moores’s Law. We need to do more than the industry has done in the past because the applications and data are such that there’s a lot more computation that’s necessary.”



There is plenty of room both at the top and at the bottom  
but much more so  
when you communicate well between and optimize across  
the top and the bottom.