



北京大學  
PEKING UNIVERSITY

# 智能硬件体系结构

## 第一讲：智能硬件体系结构简介

主讲：陶耀宇、李萌

2024年秋季

# 课程简介

## • 培养学生初步理解智能时代的硬件芯片的工作原理、设计原理与未来发展方向

课程名称	智能硬件体系结构 Hardware Architecture for Intelligent computing Systems
课程编号	04632042
学分	2 学分
总学时	34 学时
上课时间	1-16周 每周周三 5-6节
上课地点	一教 202
考核方式	- 作业 8 次: 20%, (2 周 1 次, 1 次 3-5 题) - 2 次编程实验: 2*20%, (传统体系结构、简单 AI 加速器) - 期末考试: 30% - 出勤: 10%
授课教师	陶耀宇, 李萌

- 1 智能硬件体系结构简介
- 2 电路基础-1: 晶体管与数字电路设计
- 3 电路基础-2: 芯片的物理设计与验证
- 4 指令集与流水线设计
  - 4.1 指令集基础
  - 4.2 流水线设计原理
  - 4.2 数据/控制冲突及其处理机制

- 5 乱序执行微架构设计
  - 5.1 指令动态发射原理 (MIPS R10K)
  - 5.2 分支预测与超标量设计
- 6 缓存微架构设计
  - 6.1 多级缓存与缓存一致性
  - 6.2 缓存优化与预读取
  - 6.3 虚拟内存技术
- 7 人工智能硬件体系结构
  - 7.1 GPU与FPGA架构
  - 7.2 AI专用加速芯片
  - 7.3 软硬件协同设计
- 8 新型智能计算架构
  - 8.1 感存算一体AI芯片
  - 8.2 未来AI芯片发展趋势



扫描二维码加入【智能硬件体系结构课程】群添加说明: 智能硬件体系结构课程交流-年级-姓名。

课程网站:

<https://aiarchpku.github.io/2024Fall/>

前置知识要求: 无强制要求

编程技能: 简单Python、Verilog

助教: 潘泽伦 (Wx1061758085)

推荐教科书:

Computer Architecture: A Quantitative Approach - John L. Hennessy, David A. Patterson

智能计算系统: 陈云霁, 李玲, 李威, 郭崎, 杜子东

# 目录

CONTENTS



01. 课程简介与体系结构概念
02. 智能芯片历史与发展趋势
03. 智能芯片产业国内外现状
04. 新兴技术与前沿发展趋势

# 什么是硬件体系结构?

- 硬件体系结构这一概念的随着现代计算机的出现而出现，由Amdahl首次提出

“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.”

*Gene Amdahl*, IBM Journal of R&D, April 1964

## 吉恩·阿姆达尔：IBM大型机之父

从1956年的达特茅斯会议开始，人工智能（Artificial Intelligence, AI）作为一个专门的研究领域出现

体系结构作为一个独立研究领域的出现，甚至晚于人工智能



# 为什么要学习硬件体系结构?

- 硬件体系结构是**链接底层微电子器件电路与上层计算任务之间不可或缺的纽带**



**推动现代硬件体系结构的先驱者**

# 为什么要学习硬件体系结构?

- 体系结构是为了解决上世纪60年代出现的实际工程问题 – 如何链接多样算法与单一硬件?

## IBM Compatibility Problem in Early 1960s

By early 1960's, IBM had 4 incompatible lines of computers.

701 ➔ 7094

650 ➔ 7074

702 ➔ 7080

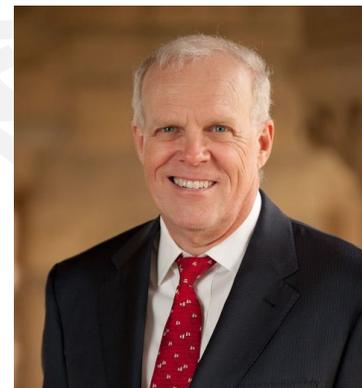
1401 ➔ 7010

Each system had its own:

- Instruction set architecture (ISA)
- I/O system and Secondary Storage: magnetic tapes, drums and disks
- Assemblers, compilers, libraries,...
- Market niche: business, scientific, real time, ...



Stanford



John L. Hennessy

StonyBrook/Stanford

美国科学院/工程院

院士

2017年图灵奖

- 图灵计算理论催生出以图灵机为理论支撑的现代智能芯片体系结构



- 纸带**: 一条无限长的纸带 (TAPE), 被划分为一个接一个小格子, 每个格子上包含一个来自有限字母表的符号
- 笔**: 一个读写头 (HEAD), 可以在纸带上左右移动, 能读出当前所指的格子上的符号, 并能通过写操作改变它
- 运算法则**: 一套规则 (TABLE), 根据当前状态及当前读写头所指格子上的符号来确定读写头下一步的动作
- 状态**: 一个状态寄存器堆栈 (STATE), 保存图灵机当前的状态

模拟人们用纸笔进行数学运算的过程

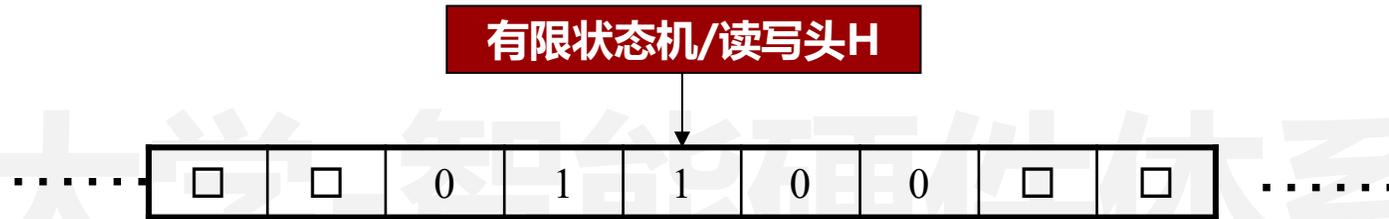
- 图灵机的数学理论框架由一个七元有序组定义

一台图灵机可被定义为  $T = \{Q, \Sigma, \Gamma, q_0, q_{\text{accept}}, q_{\text{reject}}, \delta(q,s)\}$

- $Q$ : 是非空有限状态集合
- $\Sigma$ : 非空有限输入符号表, 其中特殊空白符  $\square \notin \Sigma$
- $\Gamma$ : 非空有限带符号且  $\Sigma \subset \Gamma$ , 空白符  $\square \in \Gamma - \Sigma$ , 也是唯一允许出现无限次的字符
- $q_0 \in Q$  表示图灵机起始状态
- $q_{\text{accept}} \in Q$  表示接受状态
- $q_{\text{reject}} \in Q$  表示拒绝状态, 且  $q_{\text{reject}} \neq q_{\text{accept}}$
- $\delta(q,s): Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$  是转移函数, 根据当前读入符号  $s$  和当前状态  $q$  决定下一个状态、写入的符号、纸带移动方向和距离,  $L, R$  表示读写头是向左移还是向右移,  $-$  表示不移动

- 图灵机的计算方式与工作流程

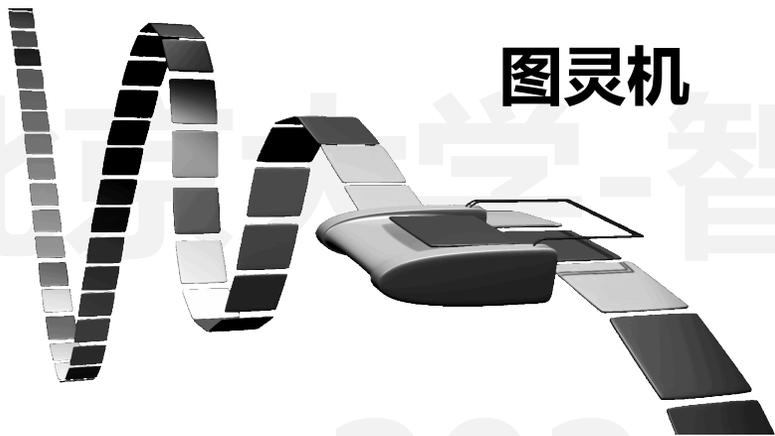
$$T = \{Q, \Sigma, \Gamma, q_0, q_{\text{accept}}, q_{\text{reject}}, \delta(q,s)\}$$



- 初始状态:** 将输入符号串  $\omega = \omega_0\omega_1 \dots \omega_{n-1} \in \Sigma^*$   $\Rightarrow$  纸带第0,1, ..., n-1 号格子
  - 读写头H指向0号格子, T @  $q_0$  状态
- 运行方式:** T按照转移函数所描述的规则进行计算
  - T @  $q$  状态, H =  $x$ , 设  $\delta(q, x) = (q', x', L)$ 
    - T  $\rightarrow q'$ , H  $\rightarrow x'$ , 读写头左移一格
    - 若某时刻H指向0号格子, 但根据  $\delta(q, x)$  将继续左移, 则T原地不动
- 停机情况:**
  - 若某时刻T @  $q_{\text{accept}}$  或  $q_{\text{reject}}$ , T停机, 并接受或拒接 $\omega$ ;
  - $\delta(q,s)$ 对某些 $q$ 和 $s$ 可能无定义, T停机

# 由图灵计算理论衍生出的冯诺依曼体系结构

- 图灵机计算范式中的元素可在冯诺依曼架构中找到对应

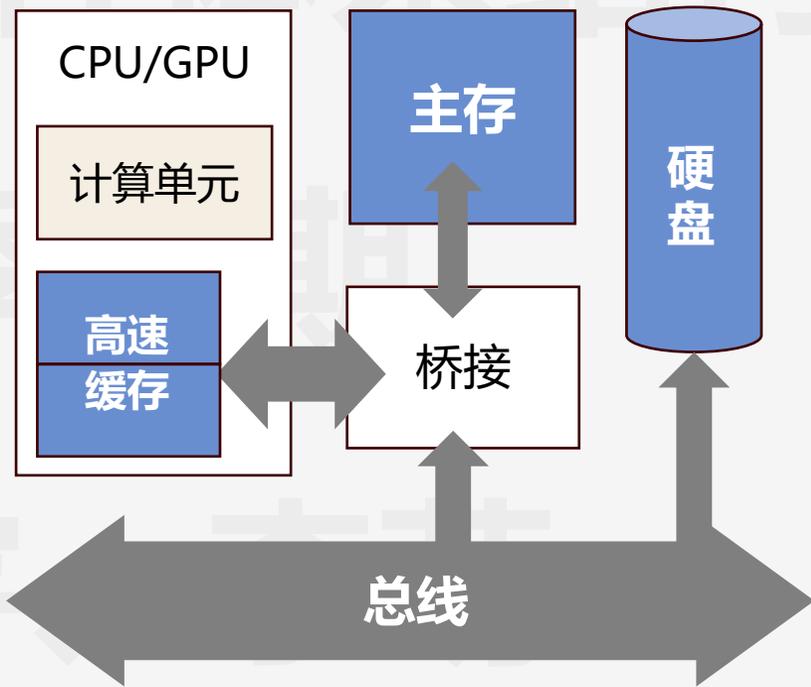


图灵机

- 笔 & 状态 → 数据存储器
- 纸带 → 指令存储器
- 运算法则 → 控制逻辑&计算单元

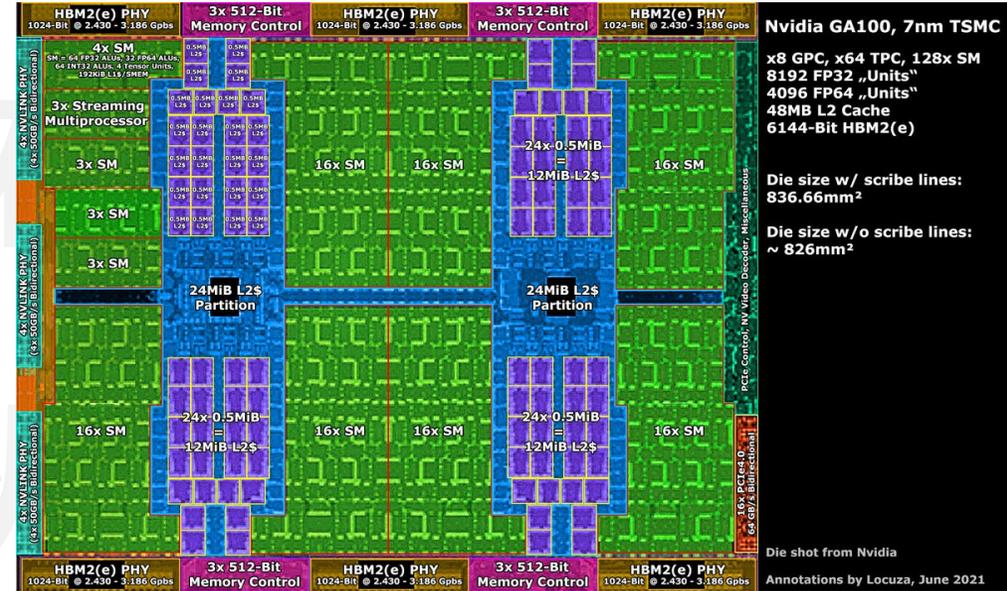
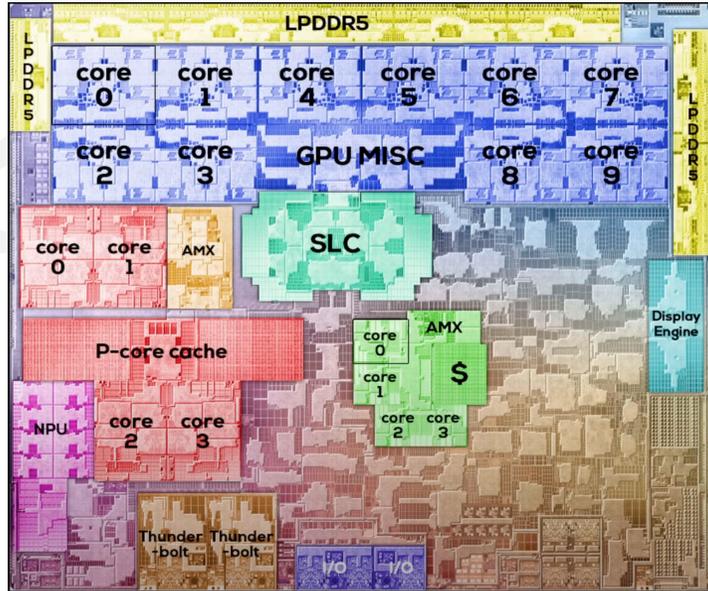


## 传统 冯·诺依曼 架构



# 冯诺依曼体系结构

- 目前的成熟商用芯片基本均采用冯诺依曼体系结构



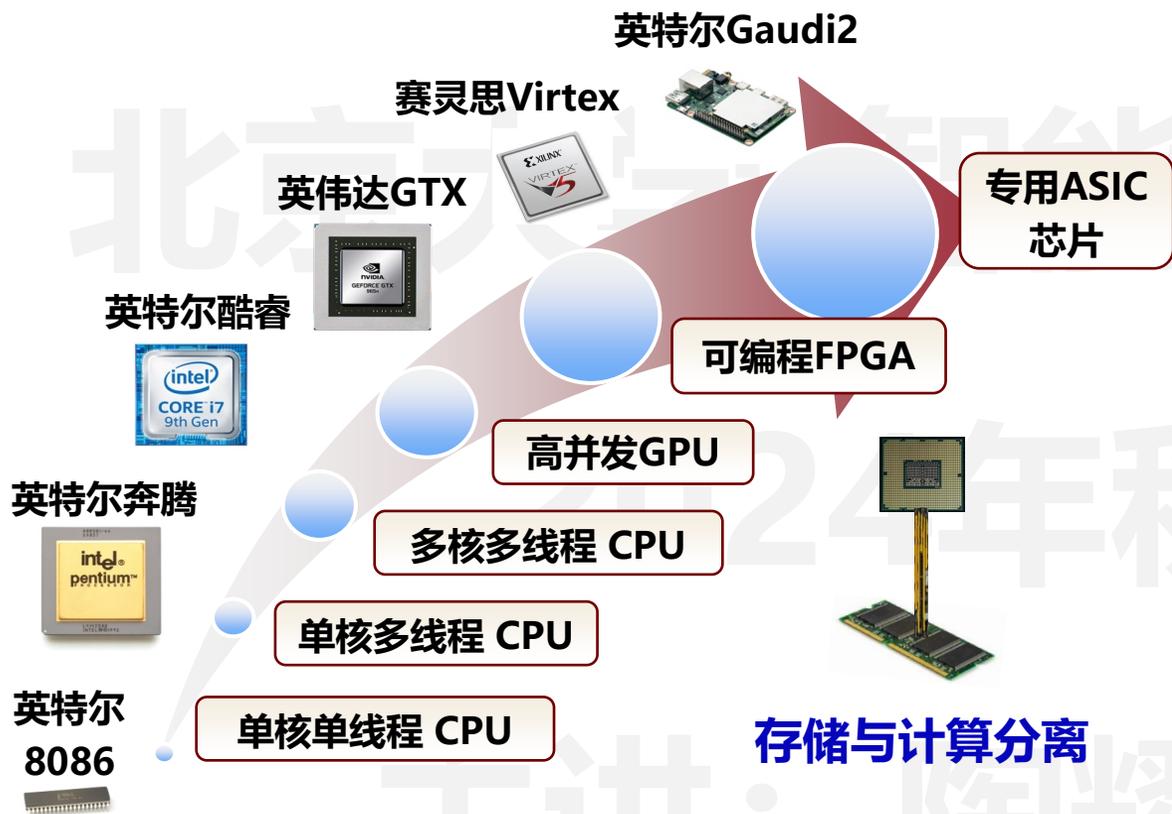
## Apple M3

## NVIDIA H100

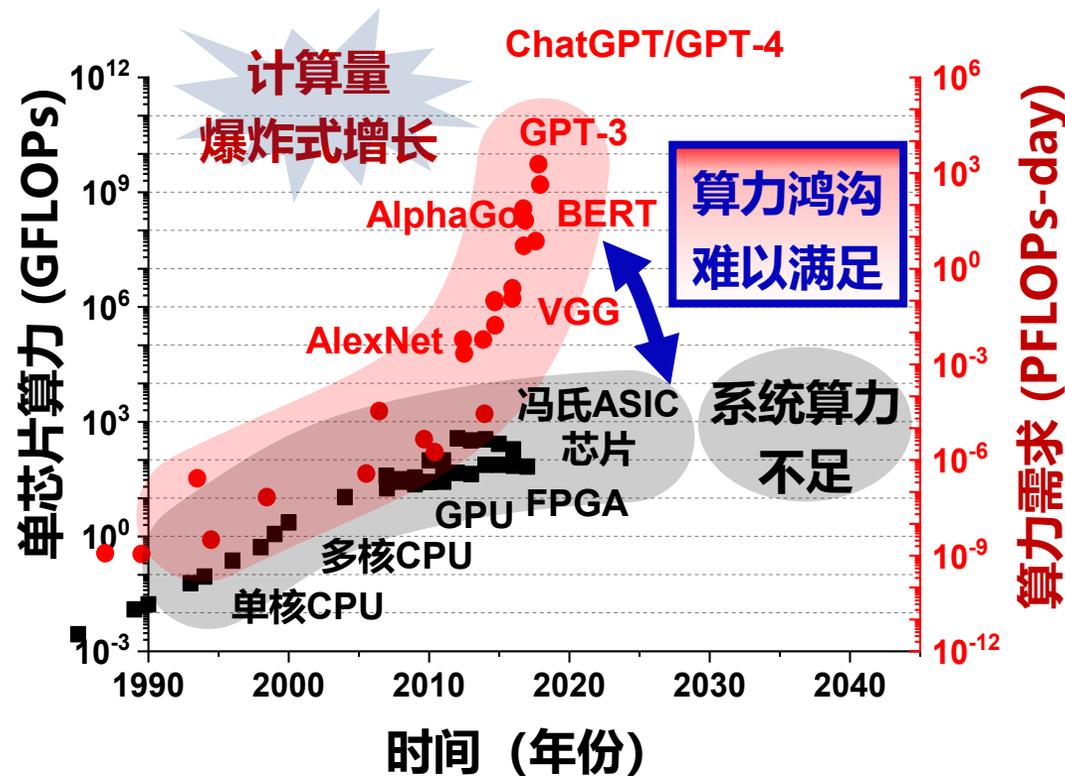
**特点：存储与计算单元分离，依靠总线进行连接，执行程序时需要来回搬运数据（读出→计算→写入）**

# 冯诺依曼体系结构

- 从偏向通用计算任务的CPU、GPU到偏向定制化设计的FPGA、ASIC



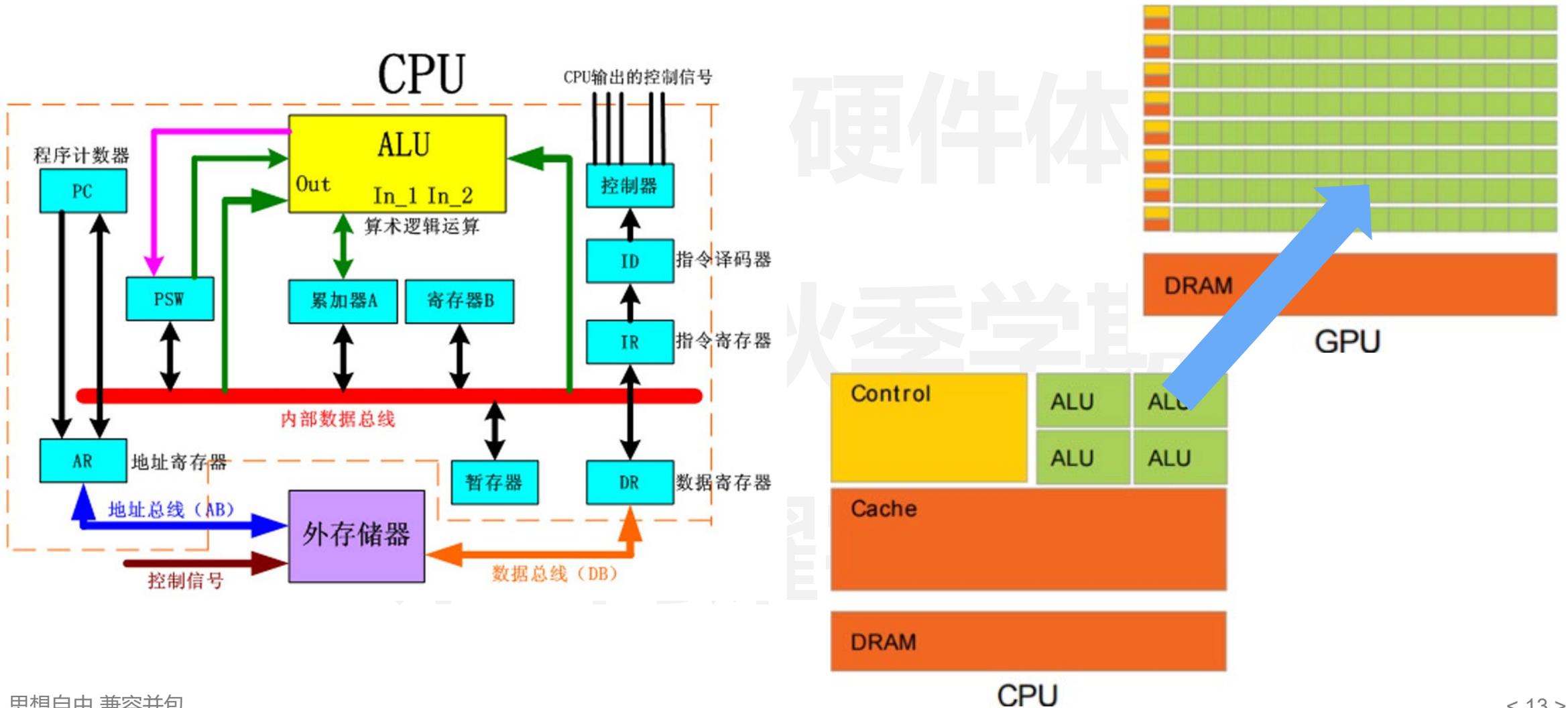
智能芯片体系结构演进图



新兴计算任务所需的计算量呈爆炸式增长

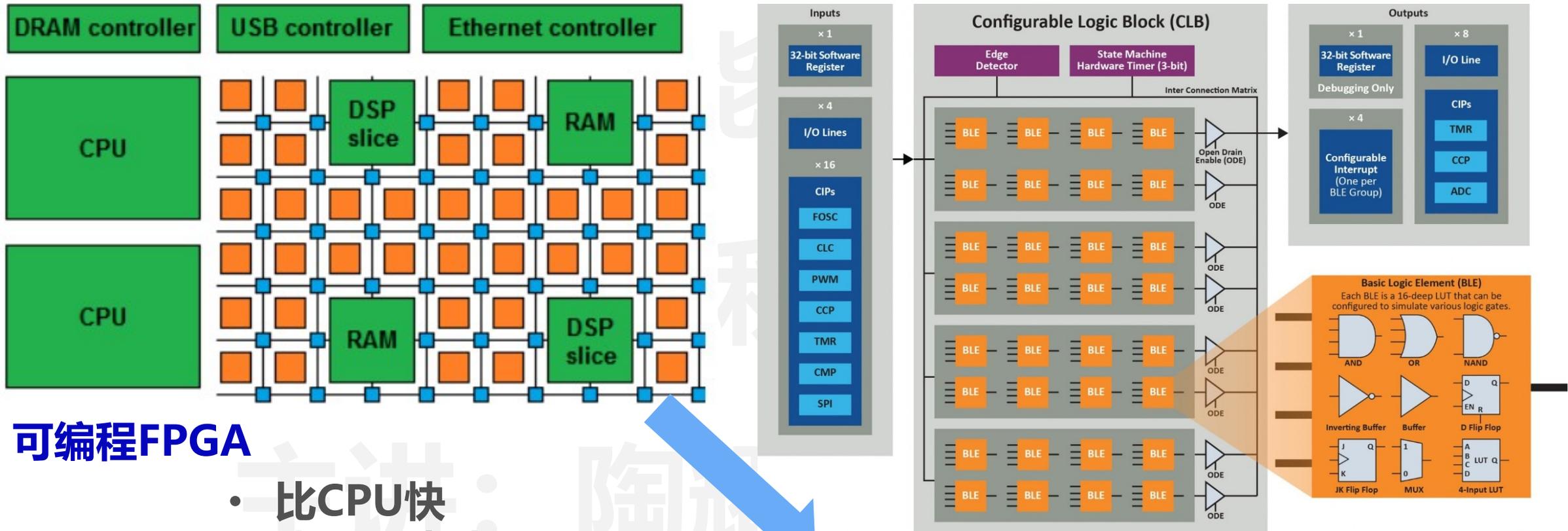
# 典型智能芯片体系结构 – CPU/GPU

- 从偏向通用计算任务的CPU、GPU到偏向定制化设计的FPGA、ASIC



# 典型智能芯片体系结构 – FPGA

- 从偏向通用计算任务的CPU、GPU到偏向定制化设计的FPGA、ASIC



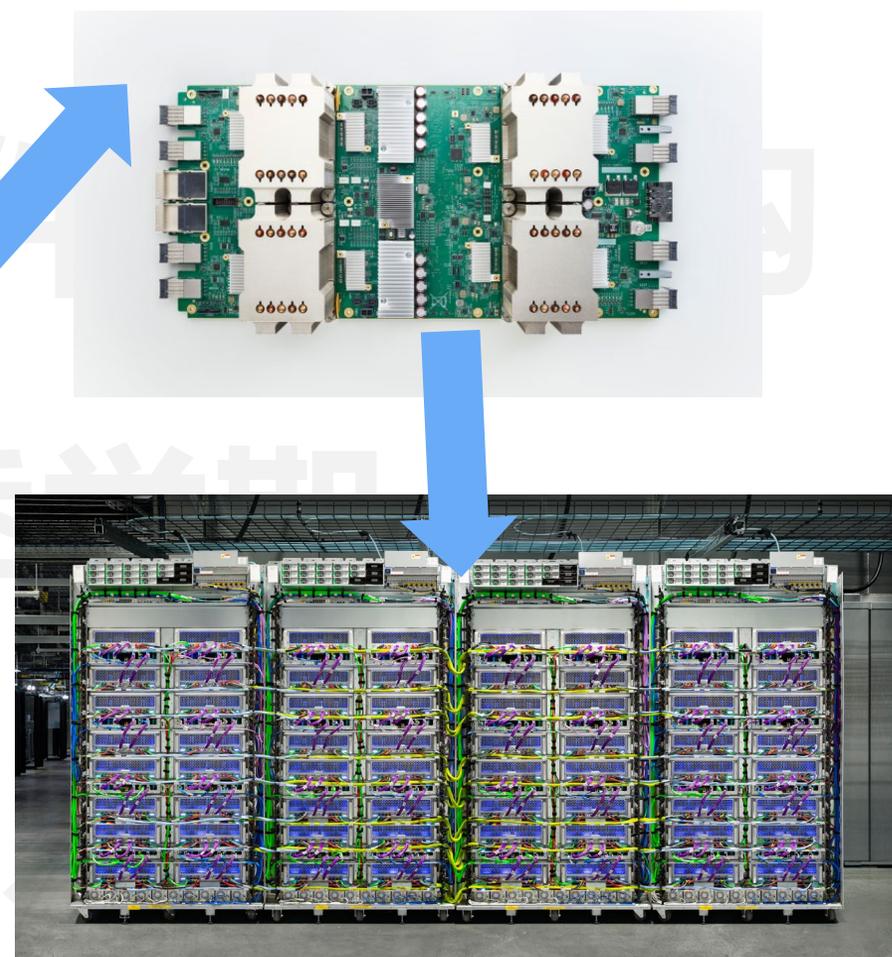
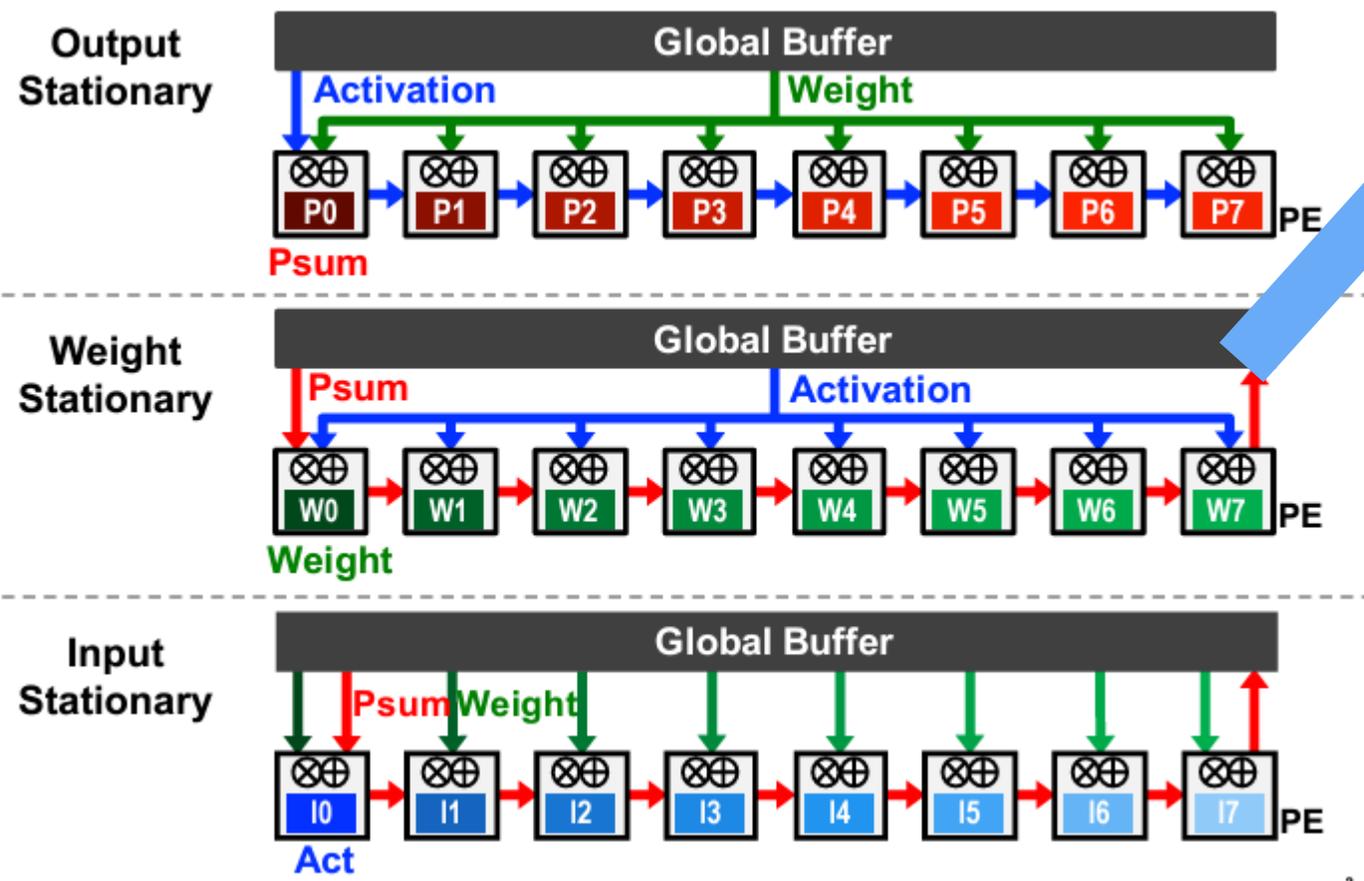
## 可编程FPGA

- 比CPU快
- 比GPU省功耗
- 比ASIC便宜流片周期短

## 可编程逻辑模块CLB

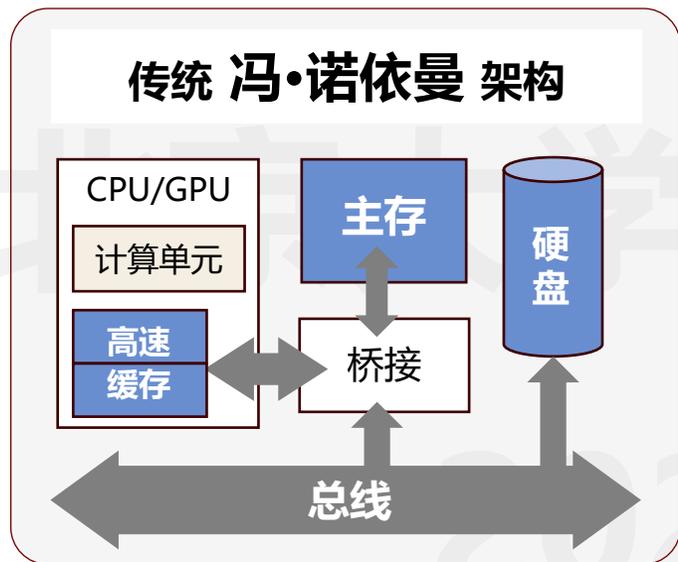
# 典型智能芯片体系结构 – ASIC

- 从偏向通用计算任务的CPU、GPU到偏向定制化设计的FPGA、ASIC

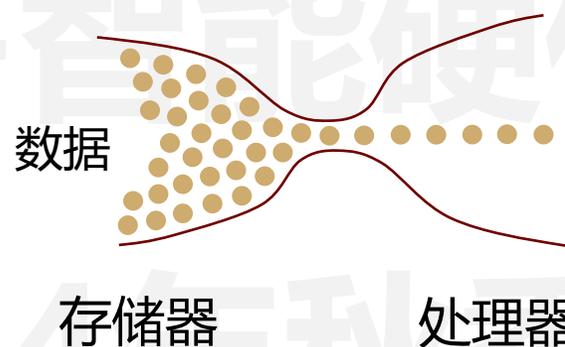


经典的DNN加速器ASIC体系结构

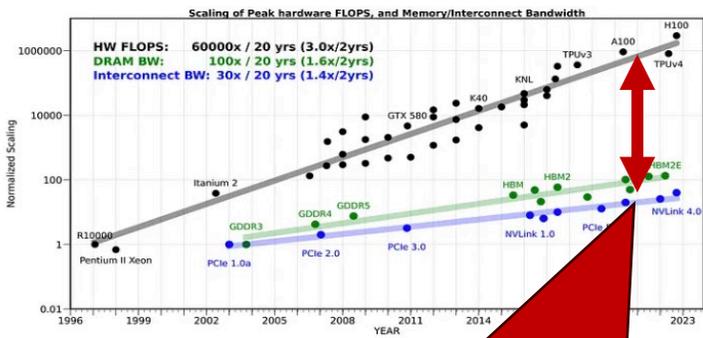
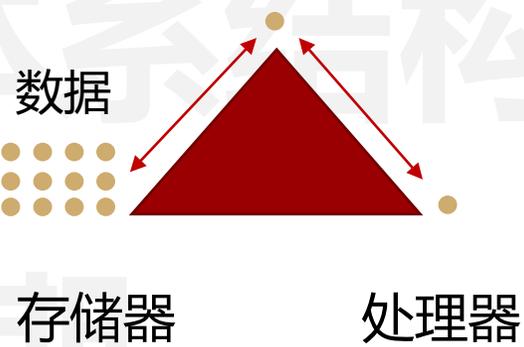
## 当前智能芯片体系结构的瓶颈



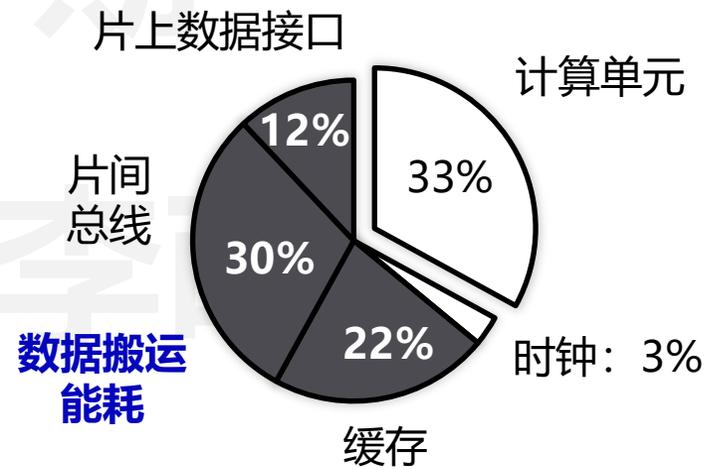
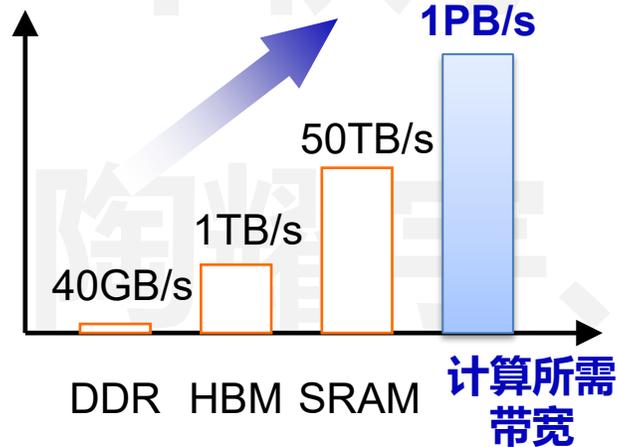
### 存储带宽限制算力



### 数据搬运限制能效



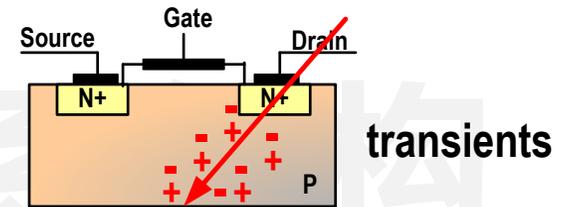
**存储性能无法满足计算需求**



- 当前智能芯片体系结构的瓶颈

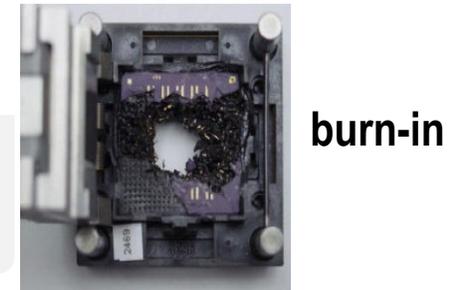
## Transient faults (瞬态故障)

- E.g., high-energy particle strikes



## Manufacturing faults (制造缺陷)

- E.g., broken connections



可靠性问题

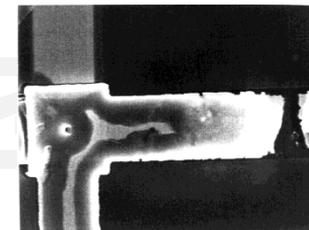
## Wearout faults (老化)

- E.g., Electromigration

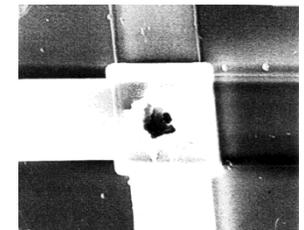
## Device variability (器件涨落)

(not all transistors created equal)

interconnect

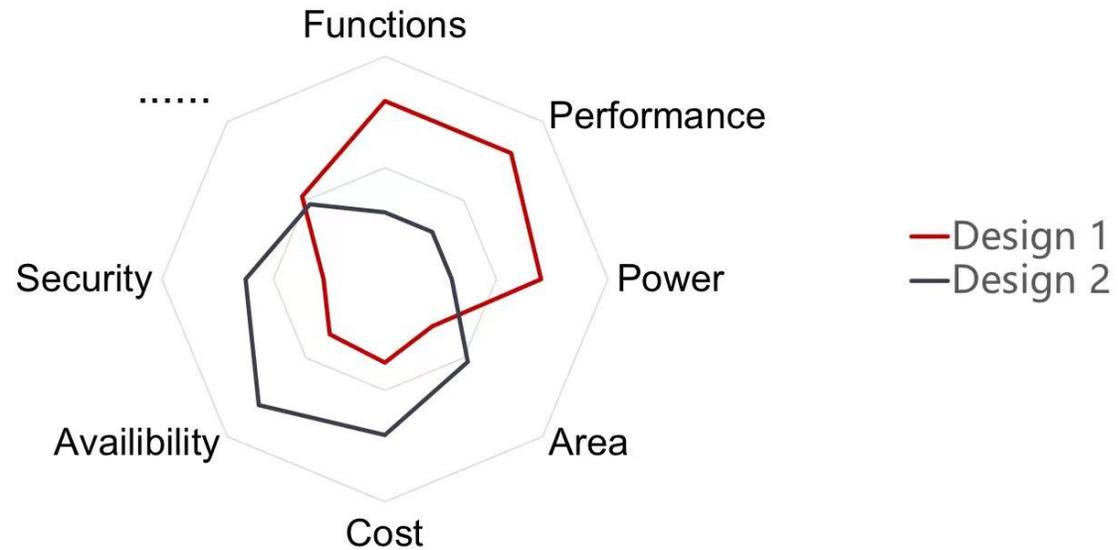


via



- 如何衡量一个芯片的性能

- **Functions**
- **Performance**
- **Power**
- **Area**
- **Cost**
- **Availability**
- **Security**
- ...



主讲：陶耀宇、李萌

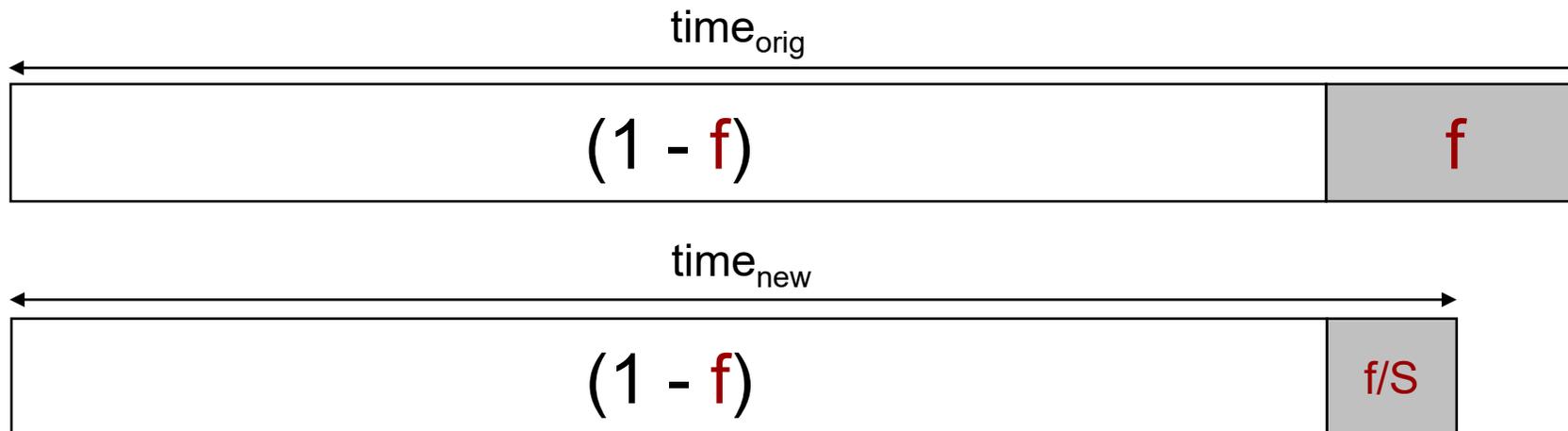
- 速度提升概念 – Amdahl Law

加速比 =  $\text{time}_{\text{without enhancement}} / \text{time}_{\text{with enhancement}}$

Technique speeds up a fraction **f** of a task by a factor of **S**

$$\text{time}_{\text{new}} = \text{time}_{\text{orig}} \cdot ( (1-f) + f/S )$$

$$S_{\text{overall}} = 1 / ( (1-f) + f/S )$$



## 体系结构基本概念 - 2

### • 并行处理的基本概念 - Parallelism law

并行度 - the amount of independent sub-tasks available

Work= $T_S$  - time to complete a computation on a sequential system

Critical Path= $T_\infty$  - time to complete the same computation on an infinitely-parallel system

平均并行度 Average Parallelism

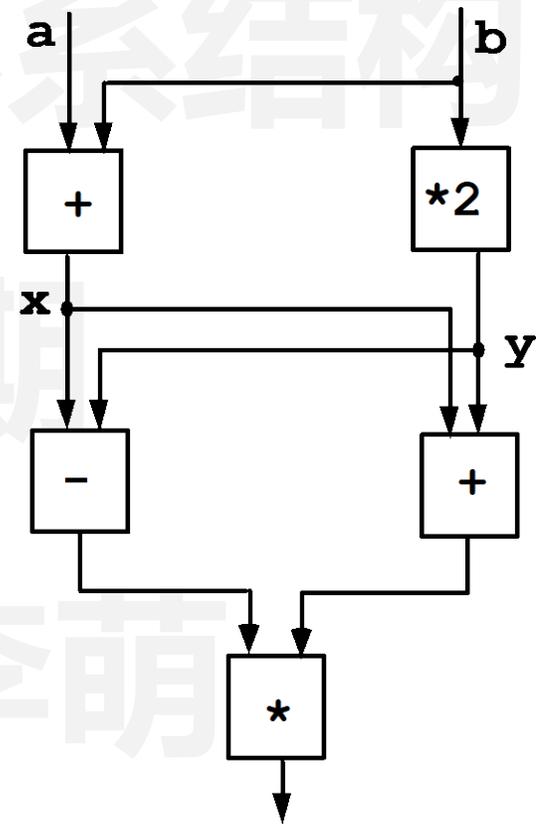
$$P_{avg} = T_S / T_\infty$$

For a  $p$  wide system

$$T_p \geq \max\{ T_S/p, T_\infty \}$$

$$P_{avg} \gg p \Rightarrow T_p \approx T_S/p$$

```
x = a + b;  
y = b * 2  
z = (x-y) * (x+y)
```



- 局部性原理的基本概念 – Locality law

最近发生事件是近期未来发生时间最好的指标.

- **时域局部性 (Temporal Locality)** : If you looked something up, it is very likely that you will look it up again soon
- **空间局部性 (Spatial Locality)** : If you looked something up, it is very likely you will look up something nearby next

Locality == Patterns == Predictability

*Converse:*

*Anti-locality: If you haven't done something for a very long time, it is very likely you won't do it in the near future either*

- 记忆与存储的基本概念 – Memoization law

Dual of temporal locality but for computation

如果某项计算很昂贵（硬件开销、时间、能耗等），那么最好的办法是记住答案一段时间，以备不时之需。

2024年秋季学期

*Why does memoization work??*

Examples

- Trace caches

主讲：陶耀宇、李萌

- 处理开销平摊的基本概念 – Amortization law
  - overhead cost : one-time cost to set something up
  - per-unit cost : cost for per unit of operation

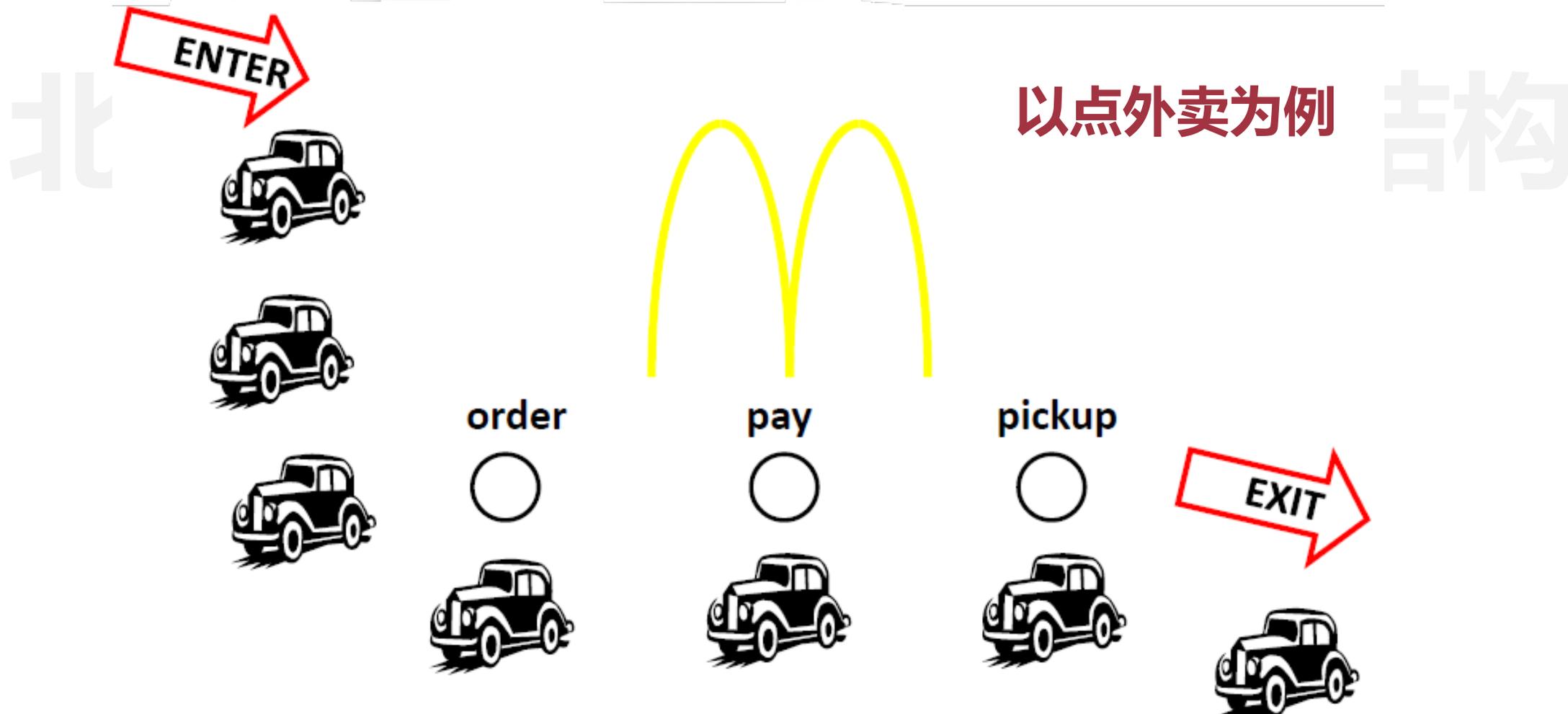
$$\text{total cost} = \text{overhead} + \text{per-unit cost} \times N$$

如果处理可以分摊到大量单位上, 通常高额overhead cost是可以接受

⇒ lower the *average cost*

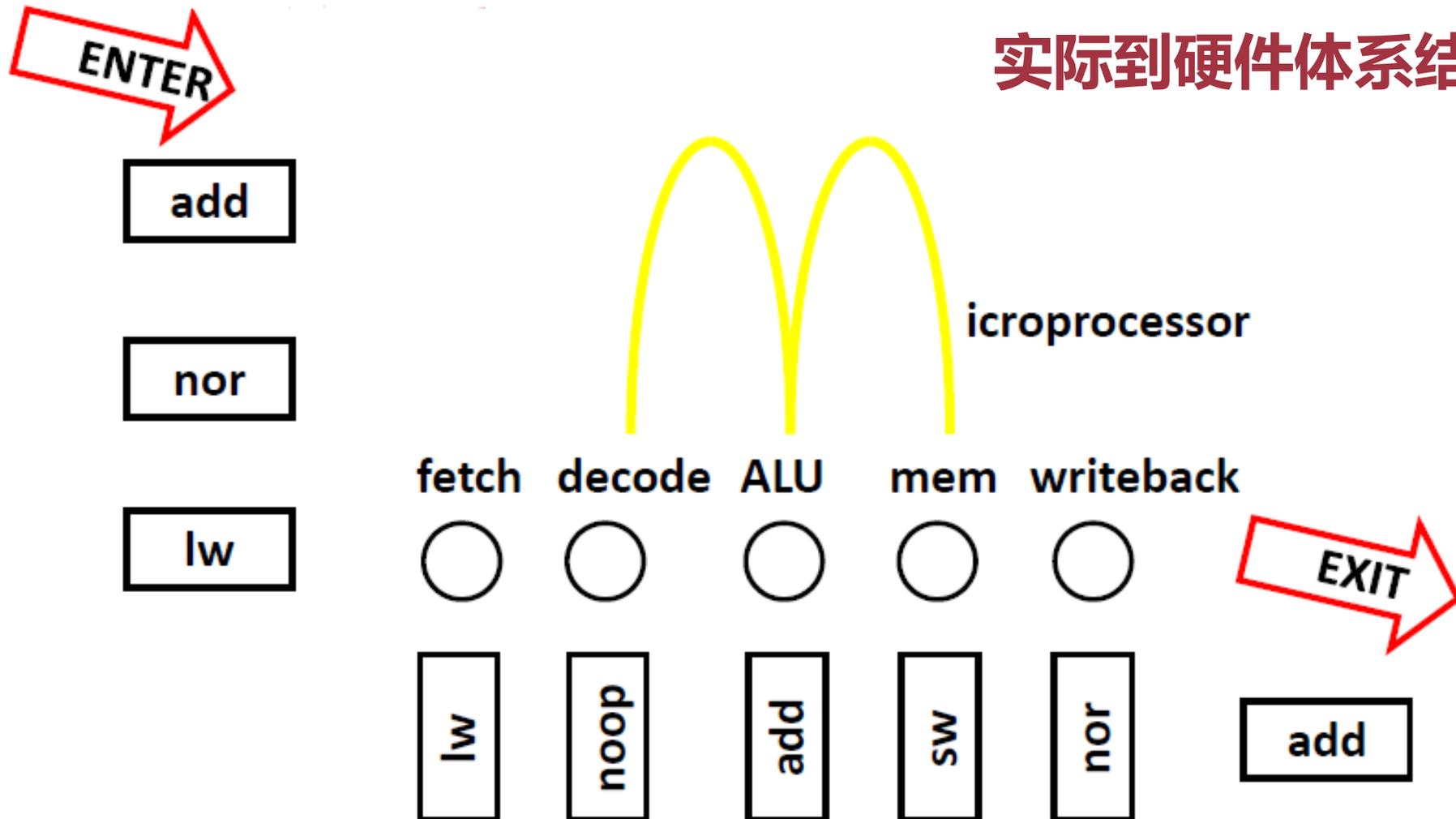
$$\begin{aligned} \text{average cost} &= \text{total cost} / N \\ &= (\text{overhead} / N) + \text{per-unit cost} \end{aligned}$$

- 流水线处理的基本概念 – Pipeline law



- 流水线处理的基本概念 – Pipeline law

实际到硬件体系结构上



## • 指令集的基本概念 – Instruction Set Architecture

### • ISA (instruction set architecture)

#### ▫ 链接软件与硬件的接口

- **Functional definition** of operations, modes, and storage locations supported by hardware
- **Precise description** of how to invoke, and access them

#### ▫ 指令集并不确定内含的指令在硬件上的运行效率，只负责提供功能描述

- Which operations are fast and which are slow and when
- Which operations take more power and which take less

Type	Example Instruction
Arithmetic and logical	and, add
Data transfer	move, load
Control	branch, jump, call, return
System	trap, rett
Floating point	add, mul, div, sqrt
Decimal	addd, convert
String	move, compare

### What operations are necessary?

*What is the minimum complete ISA for a von Neuman machine?*

#### **Too little or too simple → not expressive enough**

difficult to program (by hand)  
programs tend to be bigger

#### **Too much or too complex → most of it won't be used**

too much “baggage” for implementation.  
difficult choices during compiler optimization

## 乱序并行的基本概念 – Out-of-Order Instruction Execution

```
code1:  r1 ← r2 + 1  
        r3 ← r1 / 17  
        r4 ← r0 - r3
```

```
code2:  r1 ← r2 + 1      (1)  
        r3 ← r9 / 17     (2)  
        r4 ← r0 - r10    (3)  
        r11 ← r4 * 16    (4)  
        r28 ← r11 * r3   (5)
```

ILP: Instruction-Level Parallelism

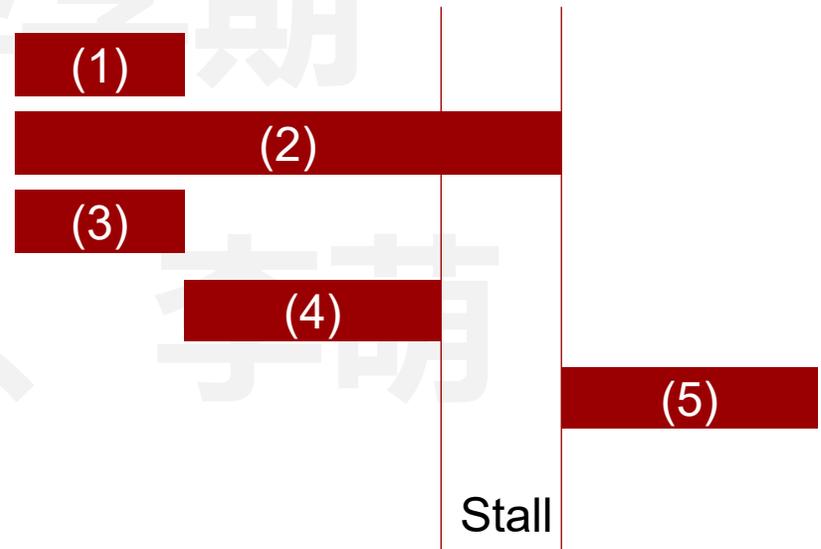
Average ILP = no. instruction / no. cyc required

code1: ILP = 1

*i.e. must execute serially*

code2: ILP = 3

*i.e. can execute at the same time*



- 硬件性能指标的基本概念 – Iron law

Time (latency) 计算延时

- elapsed time vs. processor time

Rate (bandwidth or throughput) 计算速率

- performance = rate = work per time

$$\text{Processor Performance} = \frac{\text{Time}}{\text{Program}}$$

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Time}}{\text{Cycle}}$$

(code size)

算法程序指标

(CPI)

体系结构指标

(cycle time)

体系结构/电路

Architecture --> Implementation --> Realization

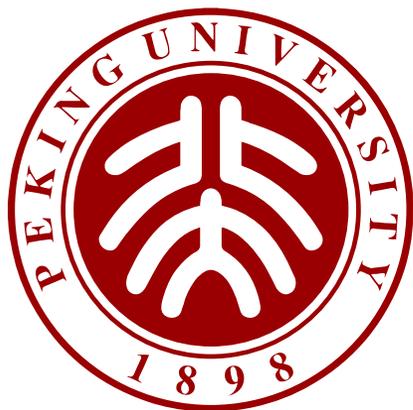
Compiler Designer

Processor Designer

Chip Designer

# 目录

CONTENTS



01. 课程简介与体系结构概念
02. 智能芯片历史与发展趋势
03. 智能芯片产业国内外现状
04. 新兴技术与前沿发展趋势

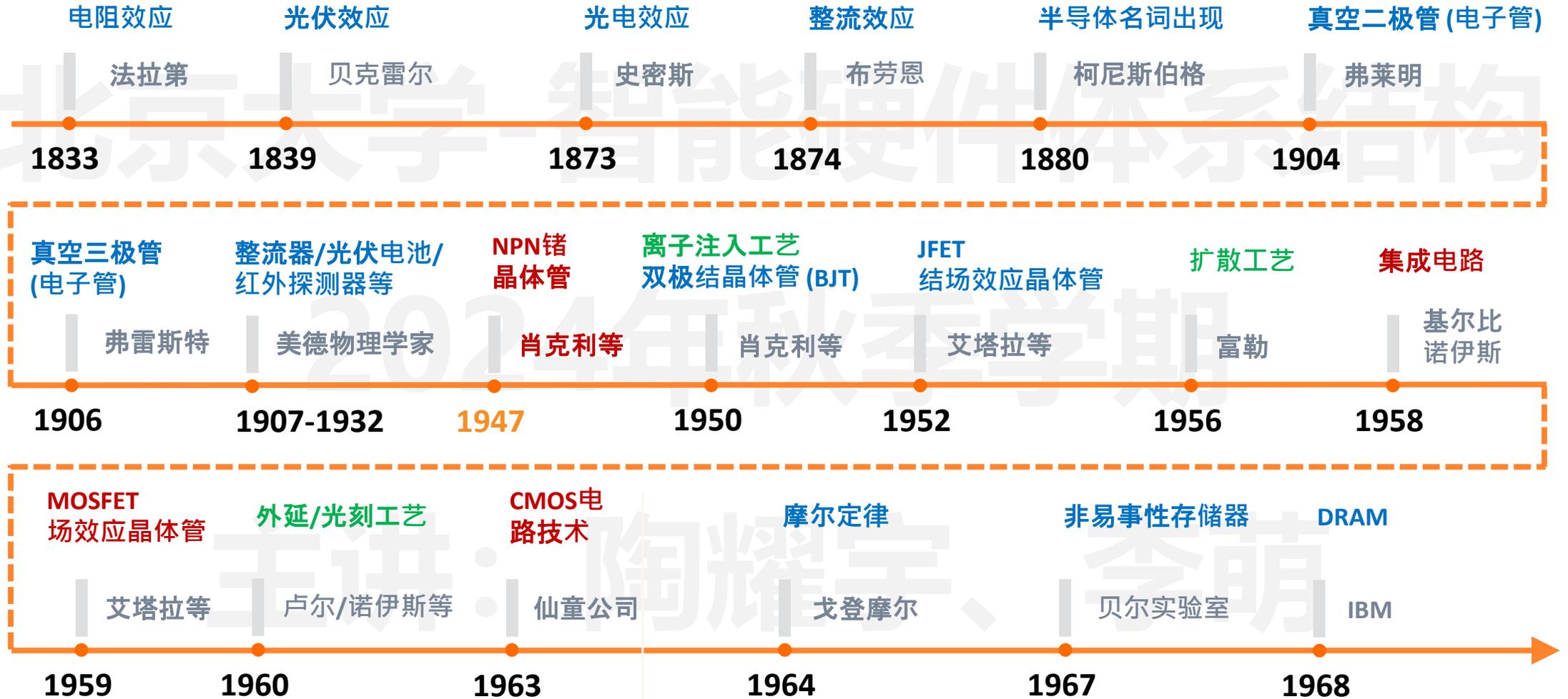
# 智能芯片的计算能力是未来新的生产力

- 数据是新的生产资料，计算能力是新的生产力，是支撑科技发展的源动力



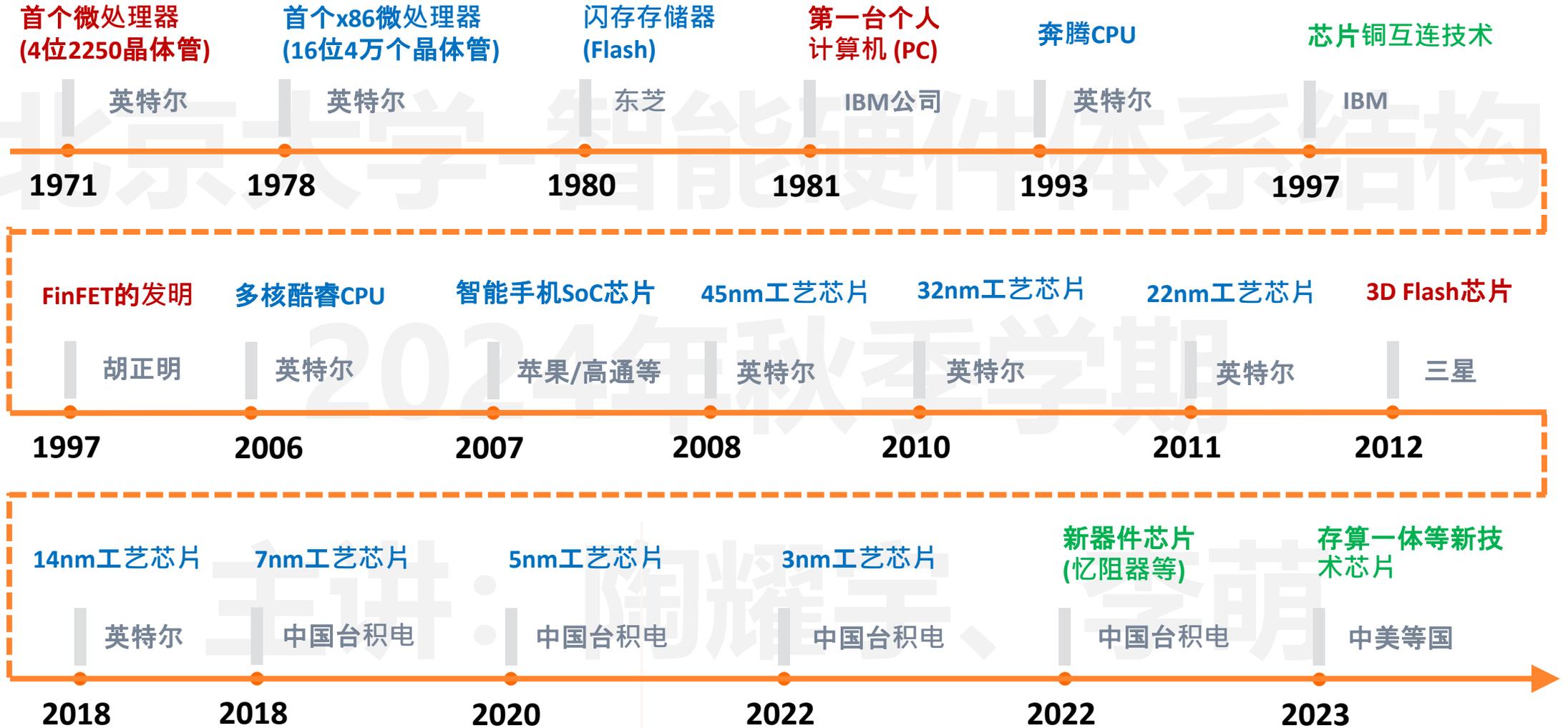
# 波澜壮阔的智能芯片发展史

## 智能芯片的发展历史 (1833 - 1968)



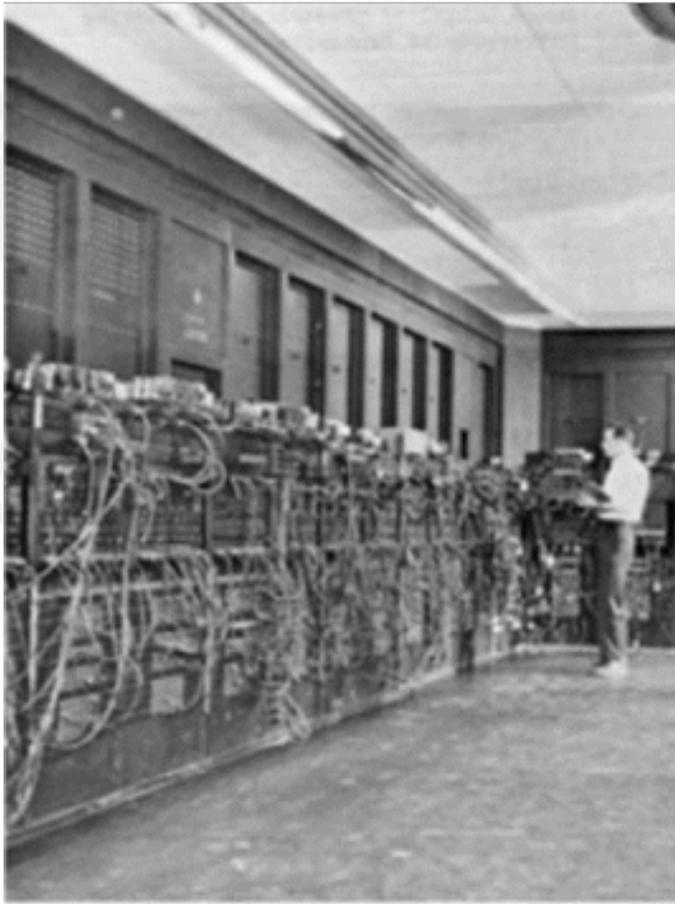
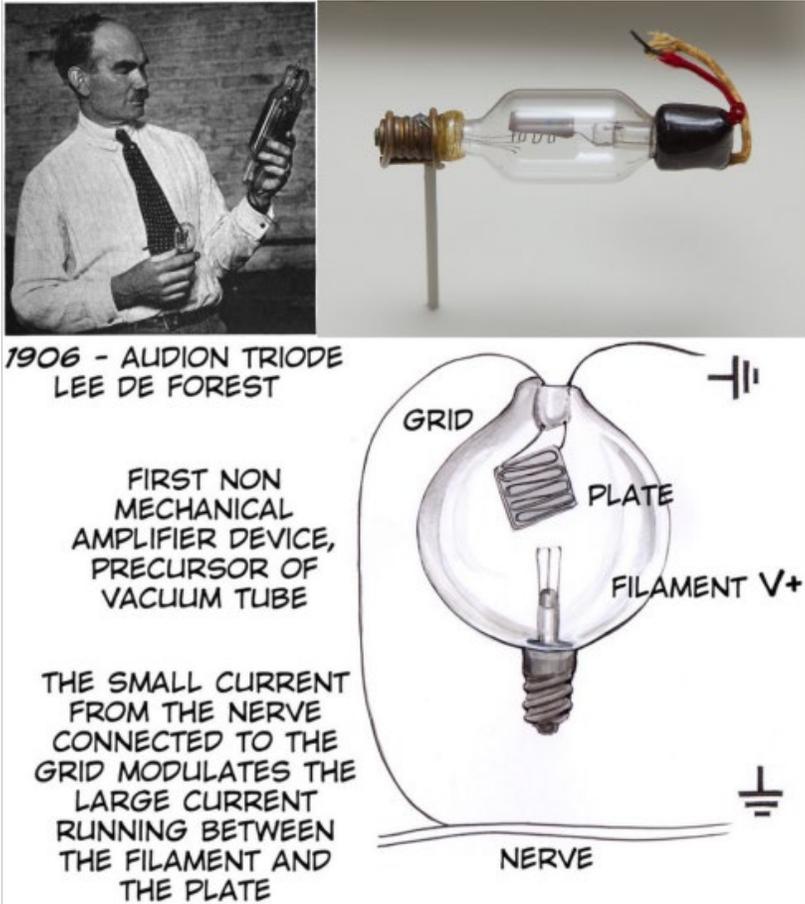
# 波澜壮阔的智能芯片发展史

## 智能芯片的发展历史 (1968 - 2023)



# 前半导体时代的霸主：真空三极管（电子管） - 1906年

- 佛雷斯特进一步再真空二极管中加入了栅极，提供额外电场调控阴极热电子向阳极运动的行为

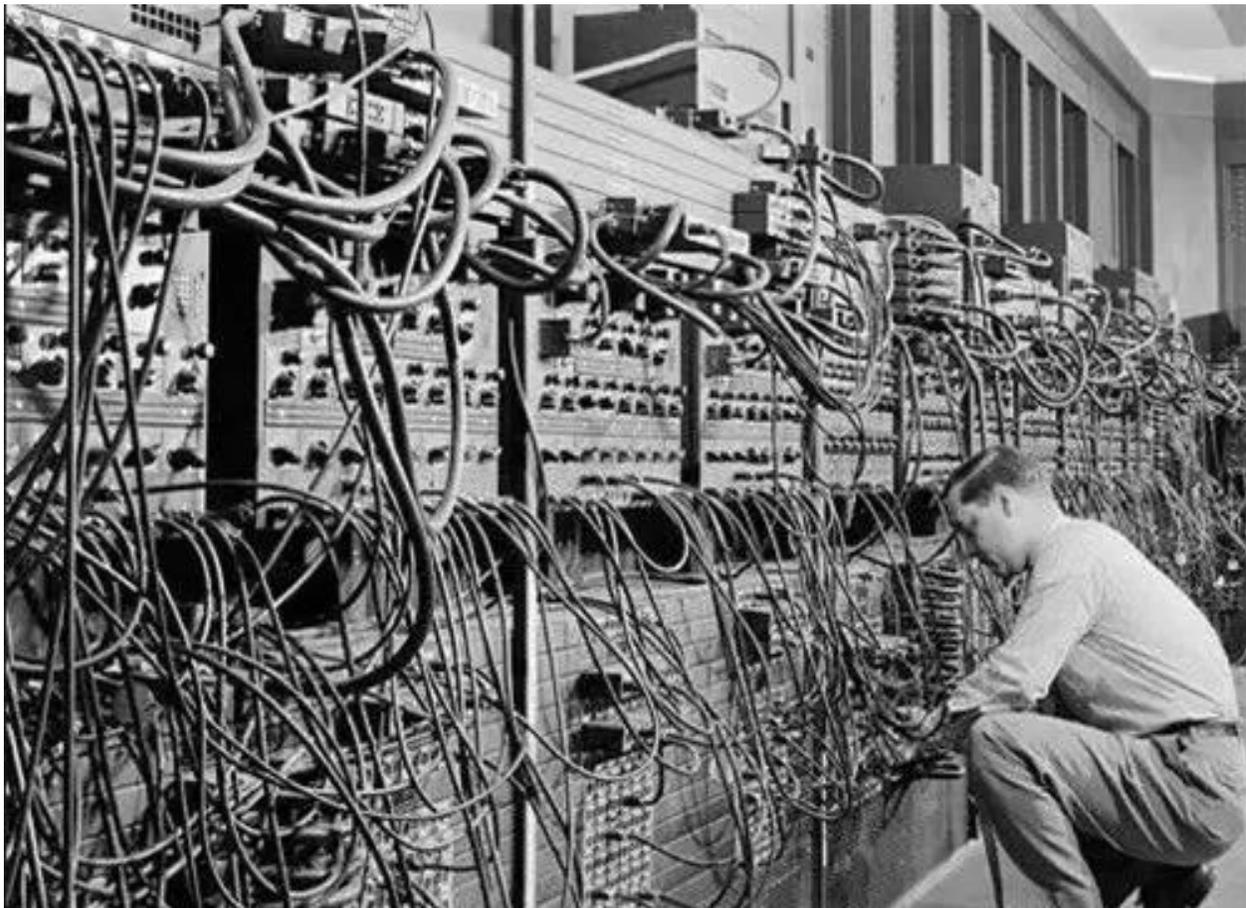


1906年，佛雷斯特进一步再真空二极管中加入了栅极，提供额外电场调控阴极热电子向阳极运动的行为，栅极电压就可以调控阴极的发射电流。这种新型真空管被称为三极管。三极管具备了检波、放大和振荡的功能，其应用场景被大大扩展，并促使了第一台现代意义的电子计算机埃尼阿克的诞生。

美国电子管计算机ENIAC

# 电子管的发展瓶颈 – 二十世纪中叶开始

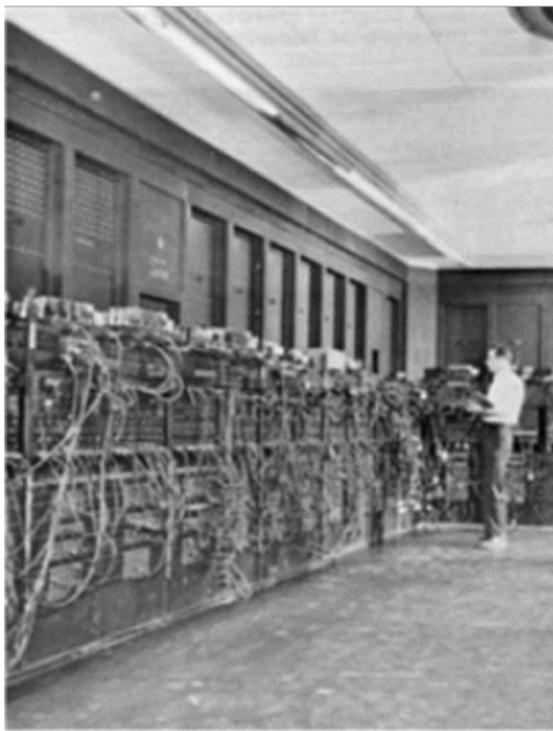
- 消失的电子管 – 根本原因是体积过大无法大规模集成、寿命较低难以长时间工作



基于热电子发射的真空管寿命较短、功耗高、体积大、成本高。埃尼阿克有一半的机时都浪费在检修损坏的真空管上，这导致它难以长时间地处理复杂的计算任务。

# 电子管的发展瓶颈 – 前苏联/俄罗斯半导体芯片产业发展的教训

- 前苏联的计算机起步与美国几乎同时代，但在**电子管与晶体管的路线选择**上出现重大失误



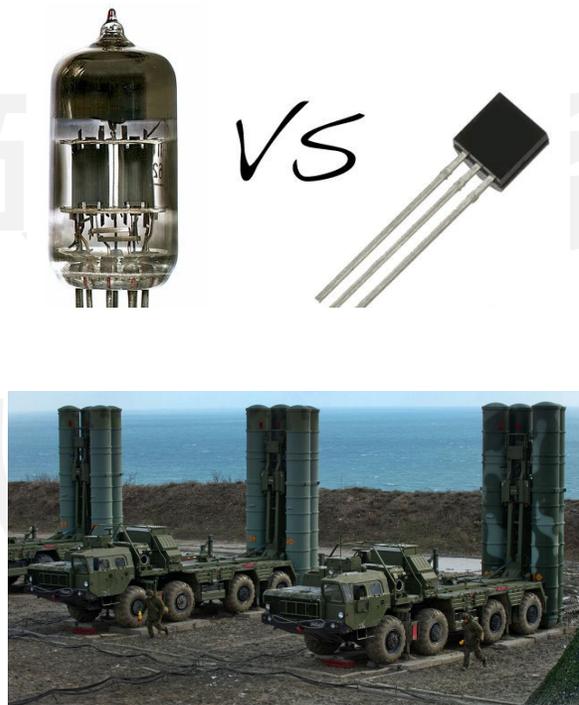
**美国电子管计算机ENIAC**

重达 30 吨，占地 170 平米  
每分钟能执行 5000 次运算



**前苏联电子管计算机MESM**

6000 个电子管每分钟3000  
次运算，算力稍弱，但耐用和  
省电上有一些优势



前苏联选择把主要精力放在了电子管的小型化上，在半导体晶体管时代逐渐落后

电子管在特定的**军事应用领域与国防工业**中仍具有一定的作用

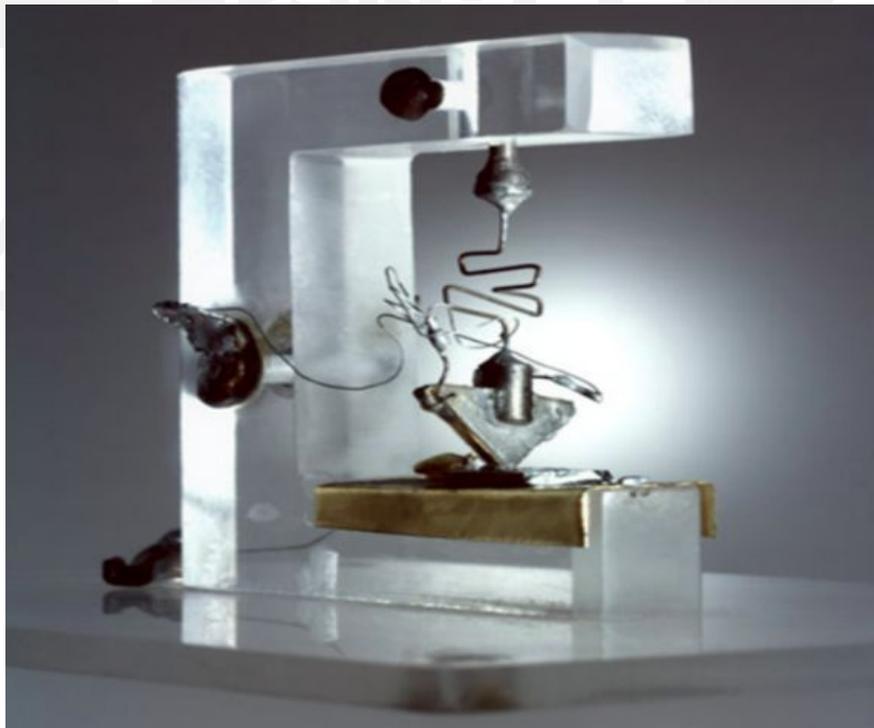
在电子管小型化方面，俄罗斯的实力在目前世界是最强的。俄罗斯S300/400等防空导弹系统极强的抗干扰能力，其实就来源于前苏联/俄罗斯的电子管小型化技术

# 芯片发展的重要历史节点：半导体锗晶体管的发明 – 1947年

- 半导体晶体管被誉为“21世纪最伟大的发明”，深刻的改变了人类历史发展进程



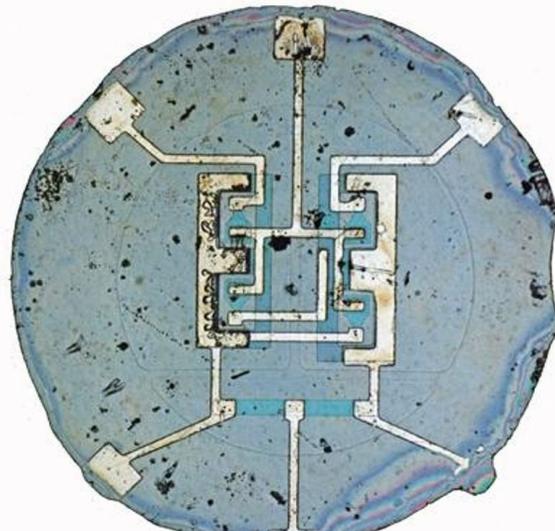
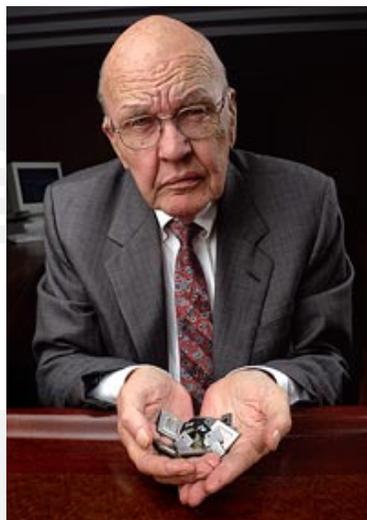
肖克利（前）、巴丁（后一）、布拉顿（后二），因为晶体管的发明，共同获得了1956年的诺贝尔物理学奖



点接触式晶体管：把间距为 $50\ \mu\text{m}$ 的两个金电极压在锗半导体上，微小的电信号由一个金电极（发射极）进入锗半导体（基极）并被显著放大，然后通过另一个金电极（集电极）输出，这个器件在 $1\text{kHz}$ 的增益为4.5

# 重要历史节点：集成电路的发明 – 1958年/1959年

- 德州仪器公司的工程师基尔比 (Jack Kilby) 发明了第一块集成电路



1958年8月28日世界第一块集成电路  
尺寸7/16×1/16英寸

基尔比获**2000年**  
**诺贝尔奖**

罗伯特-诺伊斯于1959年8月发明第一块  
**硅集成电路**

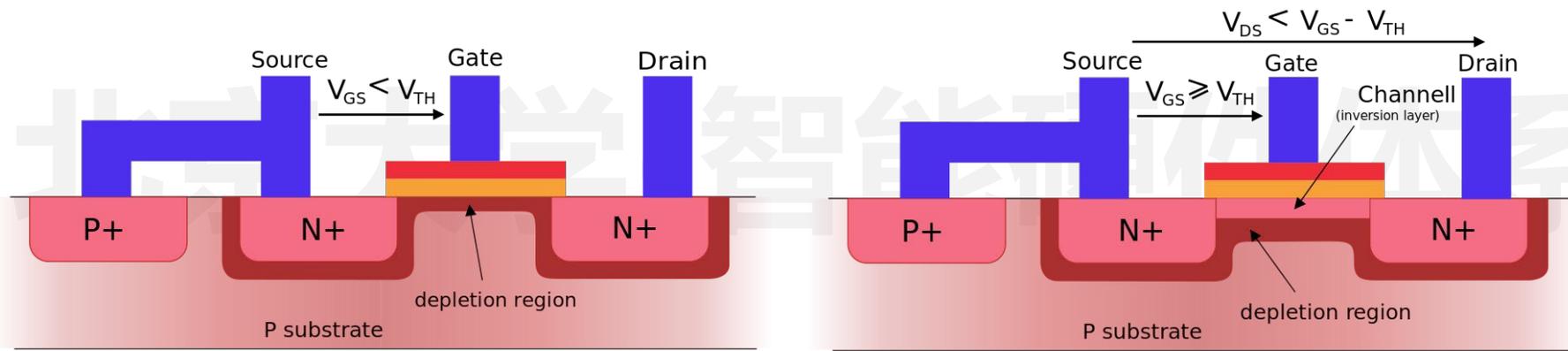
将包括锗晶体管在内的**五个元器件**集成在一起，基于锗材料制作了一个叫做**相移振荡器的简易集成电路**

参与创立**仙童半导体 (Fairchild)** 和**英特尔 (Intel)** 公司，奠定了硅谷的基石

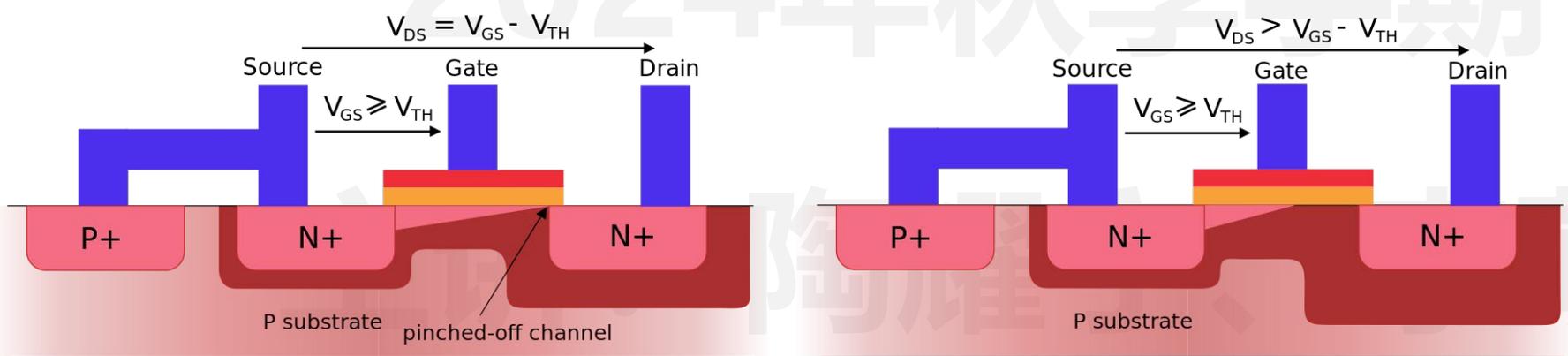


# 重要历史节点：MOSFET晶体管工作原理 – 1959年/1960年

- MOSFET有三个工作区间：断开、线性（欧姆区间）、饱和（电压不随电流线性增加）

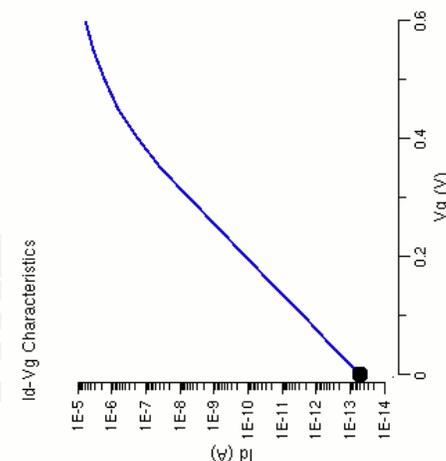
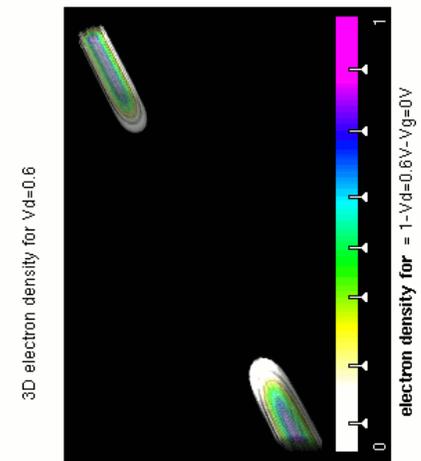


Linear operating region (ohmic mode)



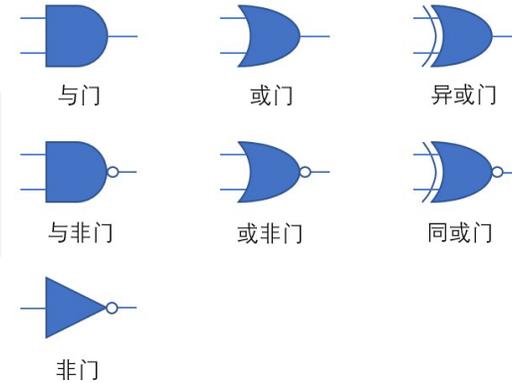
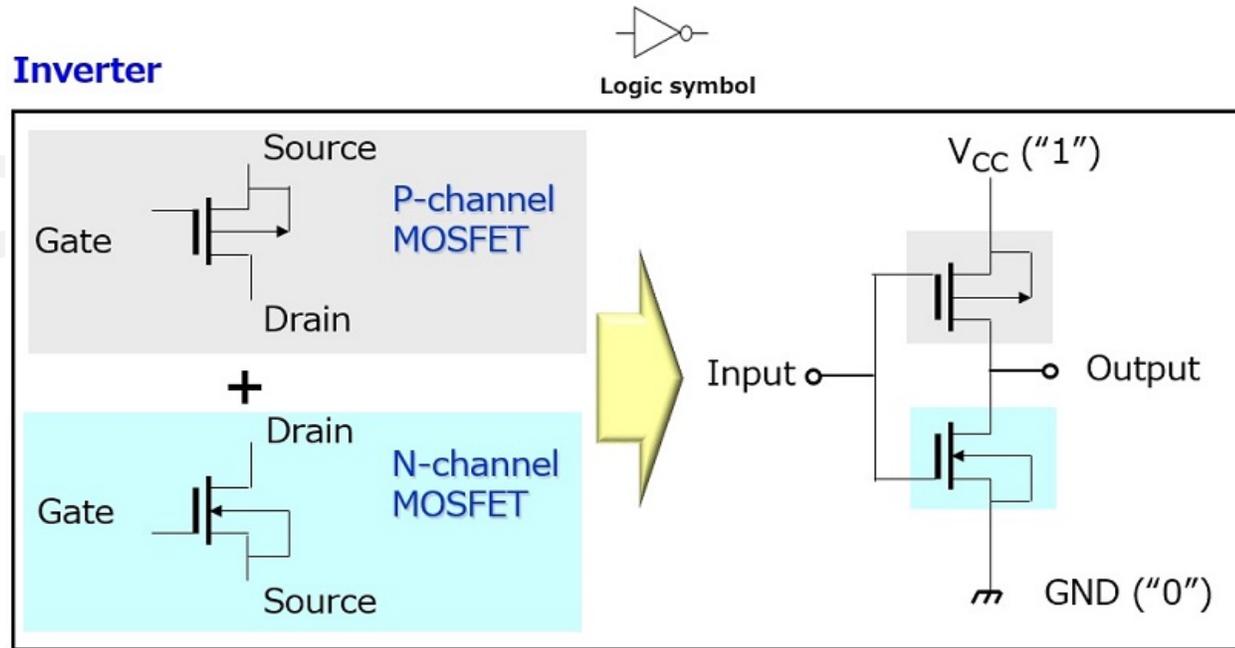
Saturation mode at point of pinch-off

Saturation mode



# 重要历史节点：CMOS电路的发明 – 1963年

- 仙童半导体于1963年首次发明互补金属氧化物半导体 (Complementary Metal Oxide Sem.)



CMOS非门电路图

互补式金属氧化物半导体具有**只有在晶体管需要切换启动与关闭时才需消耗能量的优点，因此非常节省电力且发热量少，且工艺上也是最基础而最常用的半导体器件**

硅质晶圆模板上制出NMOS (n-type MOSFET) 和PMOS (p-type MOSFET) 的基本器件，由于NMOS与PMOS在物理特性上为互补性，因此被称为CMOS

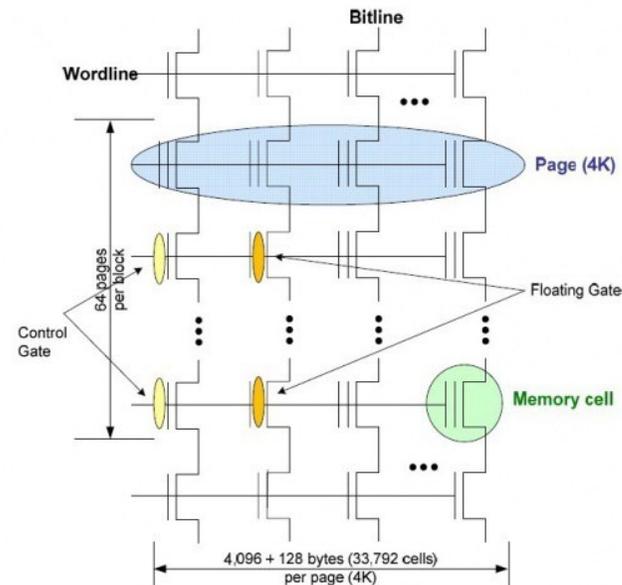


# 重要历史节点：非易失性存储器Flash的发明 – 1967年

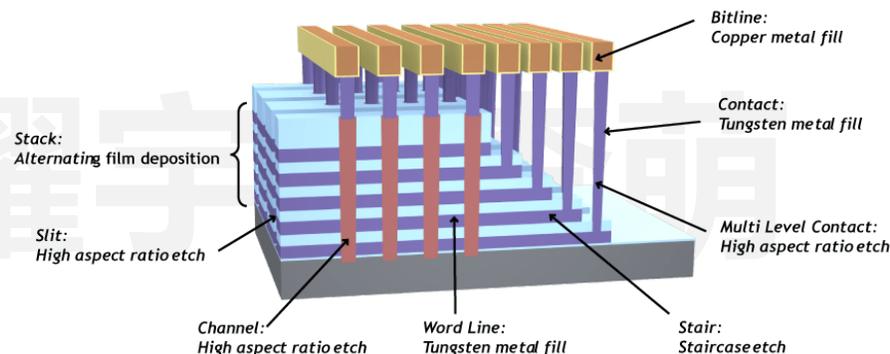
- 除了计算场景之外，存储也是占据智能芯片重要份额的典型应用场景



Dawon Kahng (韩) 和 Simon Sze (华裔) 在贝尔实验室发明了非易失性存储器浮动门 (Floating Gate) 本文发表为 “A Floating Gate and Its Application to Memory Devices” (贝尔系统技术期刊)



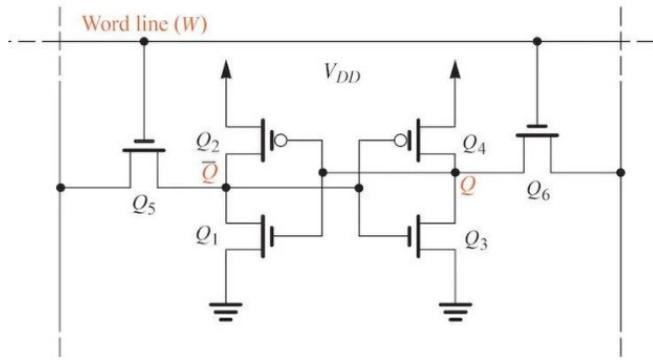
吉构  
传统平面型NAND  
Flash非易失性存储器



现代三维堆叠  
NAND Flash

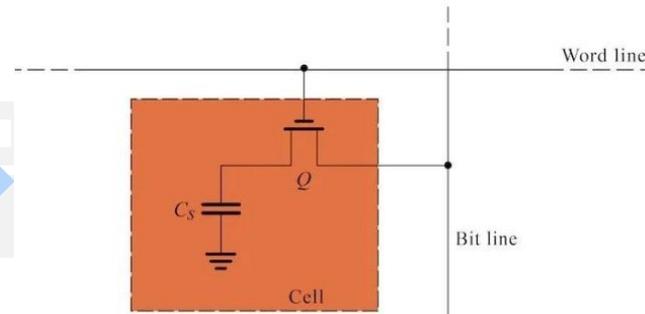
# 重要历史节点：易失性存储器DRAM的发明 – 1968年

- SRAM/DRAM是两种最常用的易失性存储器件，广泛应用于现代智能芯片中



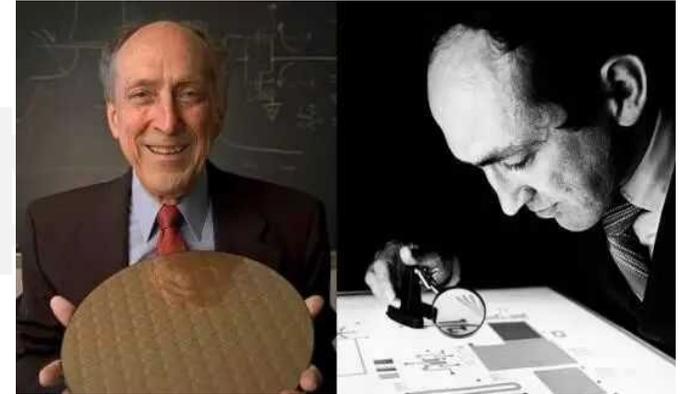
SRAM需要6个CMOS  
晶体管来存储数据

SRAM（**静态随机存取存储器**）的优点是它的速度快，它的存取速度比DRAM（**动态随机存取存储器**）快得多，因为它不需要每次访问数据都要重新刷新电容。



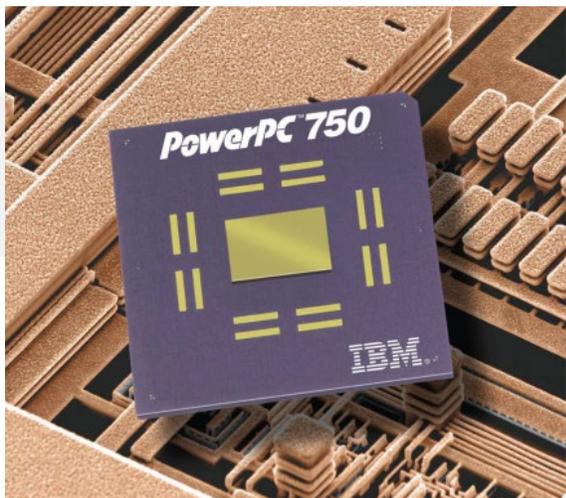
罗伯特·丹纳德发明了DRAM（**动态随机存取存储器**）存储器

与SRAM相比，DRAM的优势在于结构简单——**每比特都只需一个电容跟一个晶体管来处理**，相比之下在SRAM上一个比特通常需要六个晶体管。正因这缘故，**DRAM拥有非常高的密度，单位体积的容量较高因此成本较低**。但相反的，DRAM也有访问速度较慢，耗电量较大的缺点。

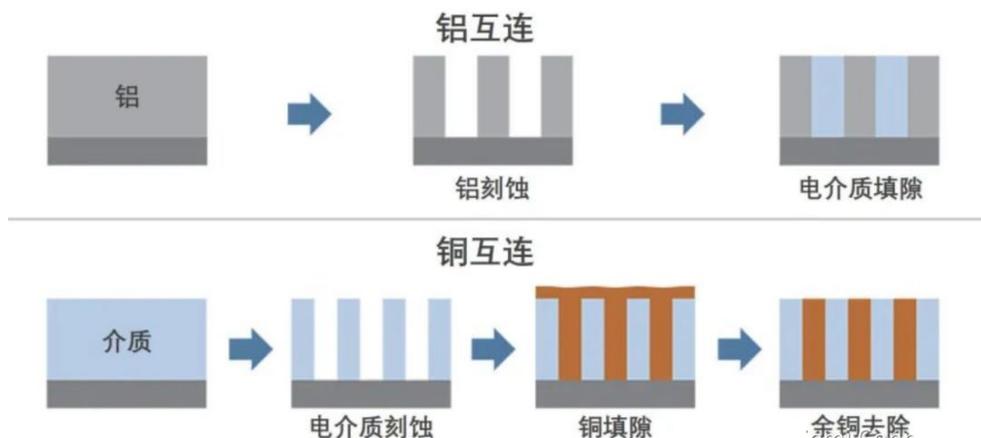
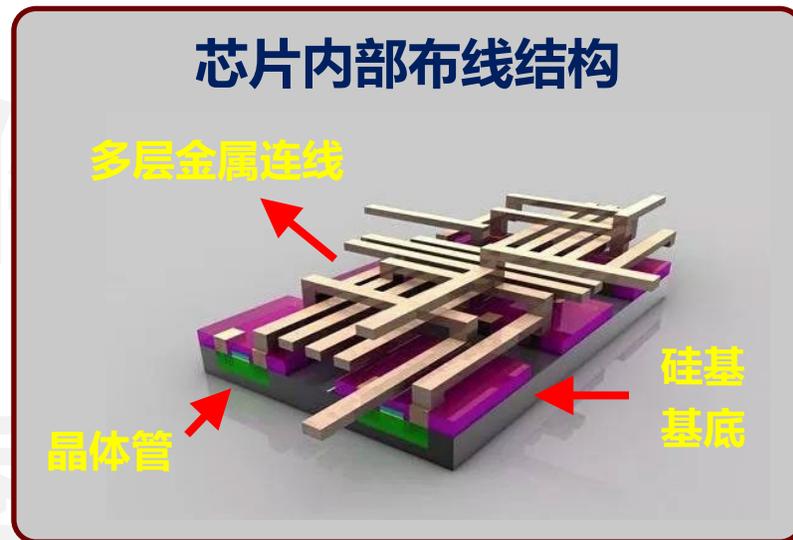


# 重要历史节点：智能芯片的铜互连技术 – 1997年

- IBM率先从铝互连转向铜互连，并推出了第一个铜基微处理器 IBM PowerPC 750



IBM PowerPC 750 最初是采用铝设计的，其工作频率高达300 MHz，采用铜互连之后，同一芯片的速度至少能达到400MHz，提高了33%



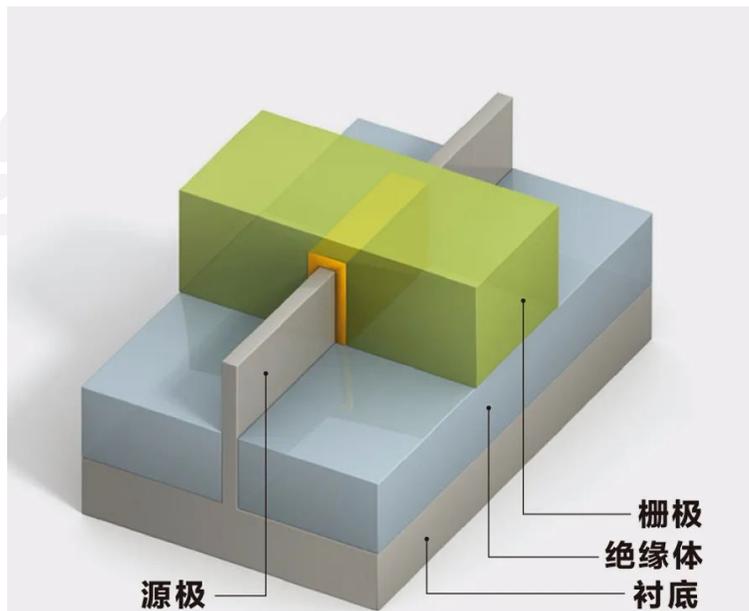
集成电路金属互连线制造工艺达到纳米级后，因为超高纯铜具有更佳的电阻率和抗电迁徙能力，很快高纯铜就替代超高纯铝合金成为金属互连线的主要材料

# 拯救摩尔定律的发明：鳍式三维晶体管FinFET – 1999年

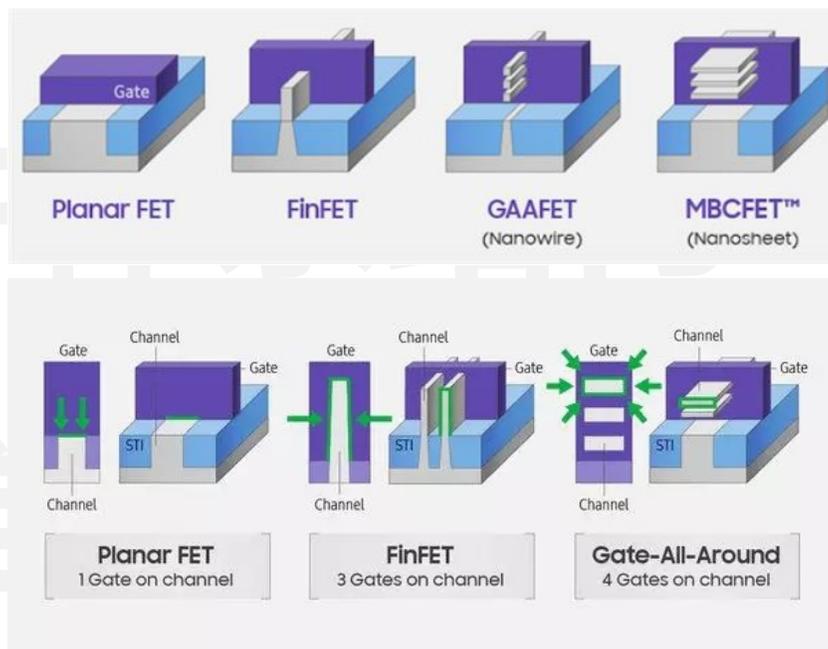
- 原本预计2010年后，传统CMOS工艺技术在20nm走到尽头，胡正明的发明拯救了摩尔定律



加州大学伯克利分校的**胡正明教授**  
(IEEE Fellow, 美国工程院院士,  
中国科学院外籍院士)



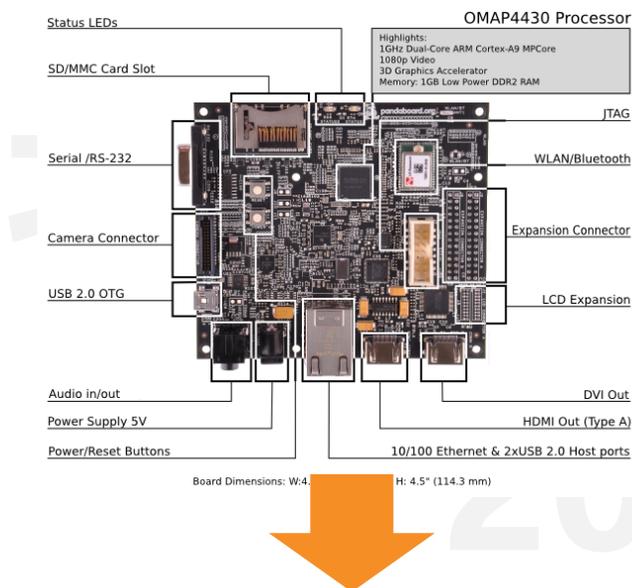
**FinFET的特点：**每个晶体管包含一个源极、一个漏极、一个连接两者的导电通道和一个控制电流沿通道流动的栅极。在FinFET中，将通道抬高，如同鲨鱼鳍高出芯片表面，使得栅极三面环绕，从而给予栅极更大的控制力。



由FinFET演化出多种**三维晶体**  
**管构型**，推动制程向  
**3nm/1nm**演进

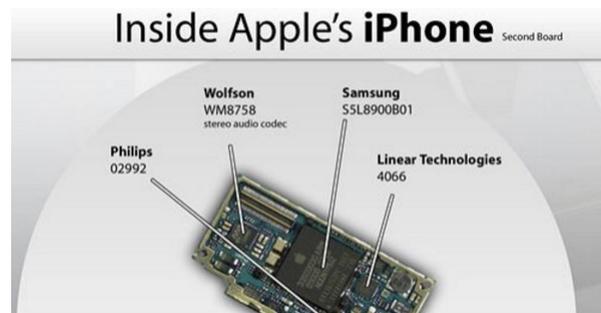
# 推动移动互联网飞速发展：移动SOC芯片 – 2007年至今

- 移动电话SOC芯片成为推动移动互联网飞速发展的算力基石，引领过去十几年的技术革命



## 德州仪器OMAP手机芯片

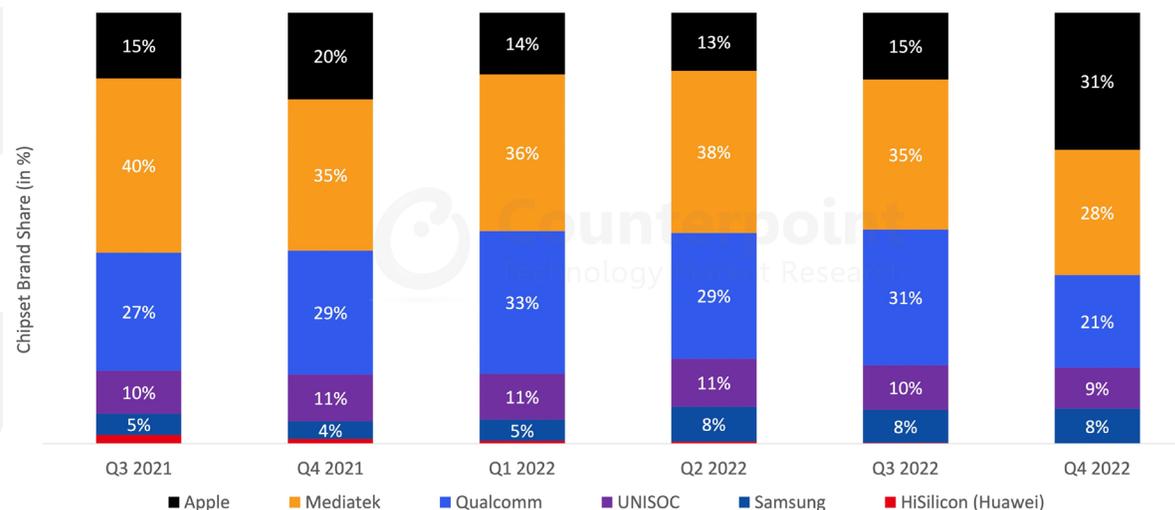
诺基亚 6630、6680、6681、E50、E60、E61、E62、E65、E70、N70、N71、N72、N73、N80、N90、N91和N92 等



## 三星S5L8900 SOC芯片

2007年乔布斯发布了第一代iPhone采用90nm制程三星SOC芯片

Global Smartphone Chipsets Market Share (Q3 2021 – Q4 2022)



This data is based on the smartphone AP/Soc Shipments | Note: Totals may not add up due to rounding.

苹果、高通、联发科、三星、紫光展锐、华为海思占据移动SOC市场的前列

# 推动制程不断向前发展：中国台湾台积电/英特尔/三星 – 2008年至今

- 过去十几年，中国台湾积体电路公司、英特尔公司、三星公司是推动芯片制程发展的主要力量

### 全球主要晶圆厂制程节点技术路线图

晶圆代工厂	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024
台积电	28nm			20nm	16nm		10nm	7nm	7nm +	5nm 6nm		3nm		2nm
三星		28nm		22nm	14nm		10nm	8nm	7nm EUV 6nm	5nm	3nm			
英特尔	22nm			14nm		14nm +	14nm ++		10nm	10nm +	7nm 10nm ++	7nm +	7nm ++	
格罗方德			28nm		14nm		12nm							
联电				28nm			14nm							
中芯国际	40nm				28nm				14nm					

备注：以上信息整理自网络，如有错漏欢迎指正。

中国台湾积体电路公司后来追上，超越英特尔与三星

# 推动智能时代飞速发展：AI芯片 – 2014/2015年至今

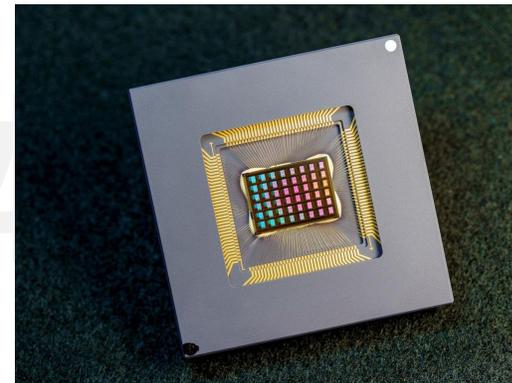
- 高性能AI芯片成为推动智能时代发展的算力基石，将引领未来十几年的技术革命



Google  
TPU



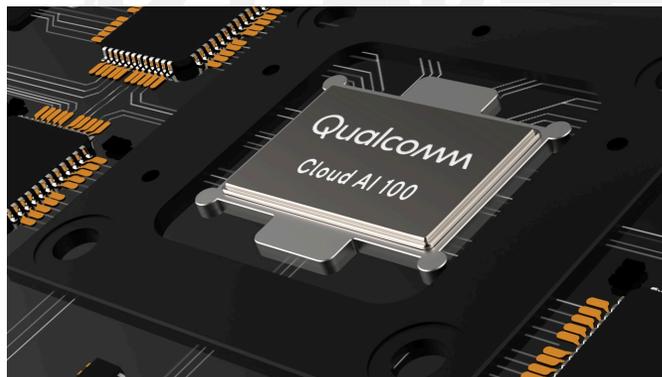
Tesla Dojo



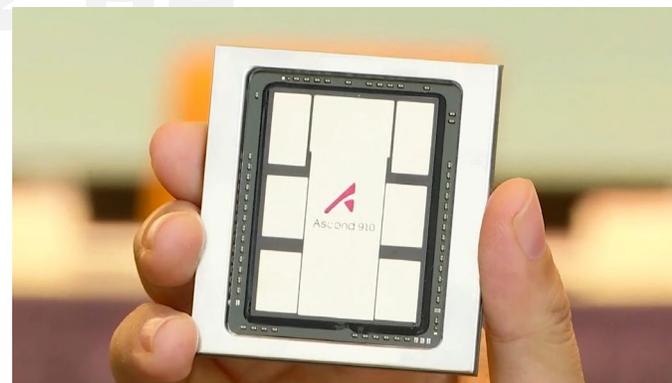
新器件AI芯片



Nvidia  
GPU



Qualcomm Cloud AI



华为昇腾

# 目录

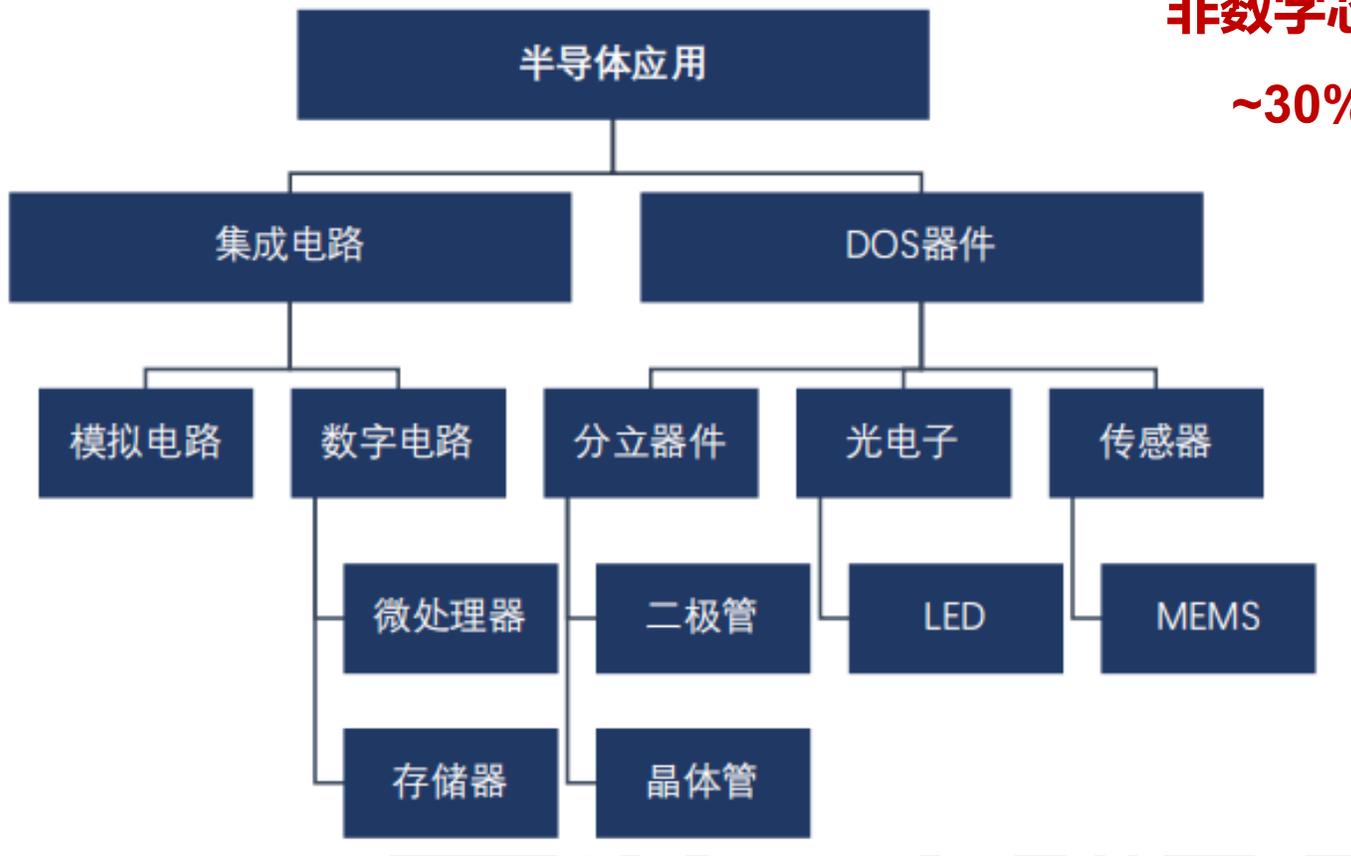
CONTENTS



01. 课程简介与体系结构概念
02. 智能芯片历史与发展趋势
03. 智能芯片产业国内外现状
04. 新兴技术与前沿发展趋势

# 智能芯片产业按半导体应用分类

- 集成电路可分为模拟电路和数字电路，DOS器件分为光电子、传感器等



数字电路的市场份额占70%

非数字芯片

~30%

光电芯片 (9%)

分立器件 (5%)

传感器 (3%)

模拟芯片 (13%)

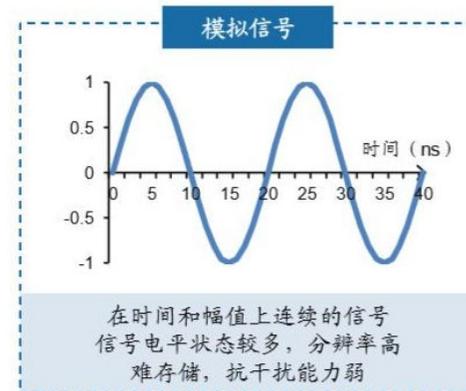
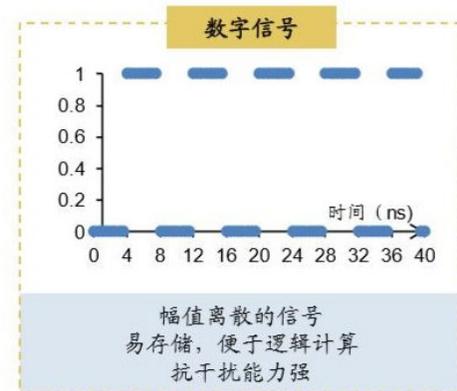
微处理器 (16%)

存储器 (28%)

数字芯片

~70%

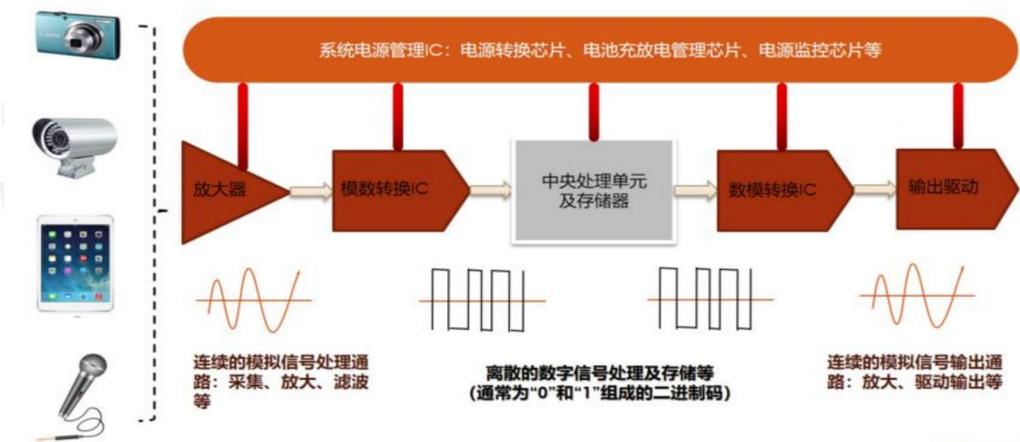
逻辑芯片 (26%)



# 智能芯片产业按半导体应用分类

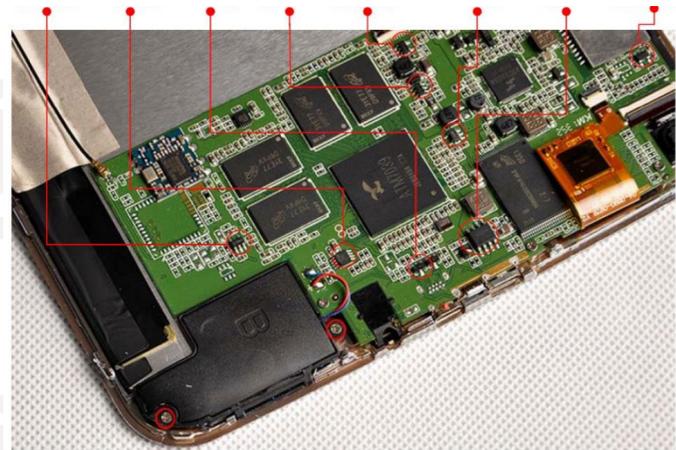
- 集成电路可分为模拟电路和数字电路，DOS器件分为光电子、传感器等

Real World Entry



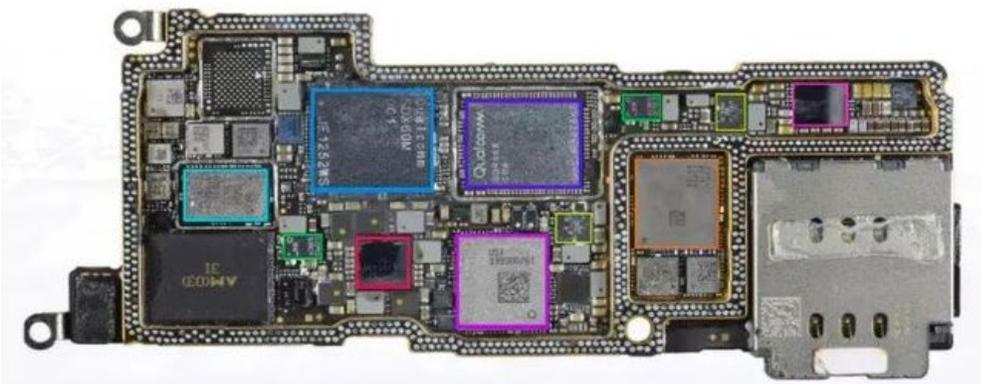
Real World Action

功率放大、电源管理、时钟生成、比较器、射频滤波、接口、数模转换、线性稳压等



模拟芯片  
应用实例

CPU、GPU、存储器芯片、可编程逻辑芯片、MCU、DSP、NPU等



数字芯片应用实例

传感芯片

声、光、电、热、磁、压力、气体、震动、速度、湿度、惯性、流量、电磁波等

光电芯片

激光器芯片、半导体发光芯片等

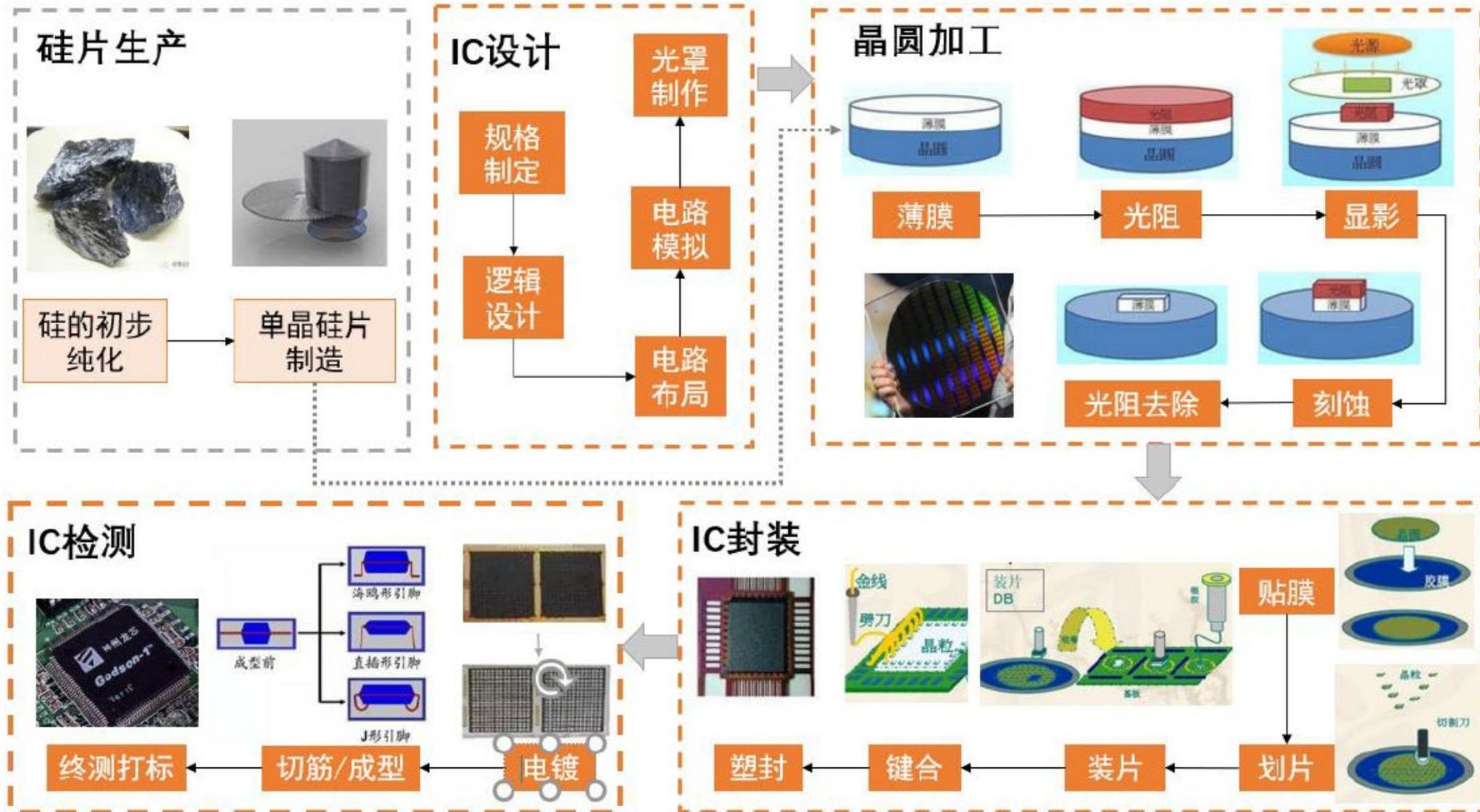
分立器件

电阻、电容、电感、振荡器、晶体管、功率器件等

# 智能芯片产业现状 – 产业链极长、关联几乎所有工业门类

## • 国际分工合作的庞大产业链生态

中国与世界先进水平差距较大



### 硅片生产企业

- 信越化学 (日本)
- 三菱住友 (日本)
- 环球晶圆 (台湾)

### 晶圆加工企业

- 台积电 (台湾)
- 三星 (韩国)
- 格芯 (美国)

### 芯片设计企业

- Intel (美国)
- Qualcomm (美国)
- 海思半导体 (中国)

### 芯片封测企业

- 日月光 (台湾)
- 安靠 (美国)
- 长电 (中国)

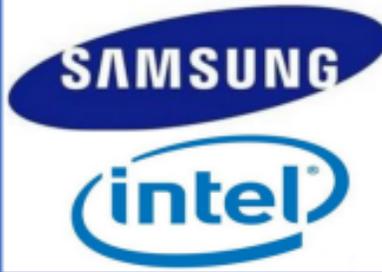
# 智能芯片产业现状 – 产业链极长、关联几乎所有工业门类

## • 国际分工合作的庞大产业链生态



# 智能芯片产业的三种运作模式

- IDM (垂直整合)、Fabless (纯设计) 和 Foundry (晶圆加工)

IDM模式

<p>集芯片设计、制造、封测于一身。早期多数集成电路企业采用的模式，目前仅有极少数企业能够维持</p>
<p>设计、制造等环节协同优化，有助于充分发掘技术潜力；能有条件率先实验并推行新的半导体技术</p>
<p>公司规模庞大，管理成本较高；运营费用较高，资本回报率偏低</p>

Fabless模式

<p>只负责芯片的电路设计与销售；将生产、测试、封装等环节外包</p>
<p>资产较轻，初始投资规模小，创业难度相对较小；企业运行费用较低，转型相对灵活</p>
<p>与IDM相比无法与工艺协同优化，因此难以完成指标严苛的设计；与Foundry相比需要承担各种市场风险</p>

Foundry模式

<p>不负责芯片设计，只负责制造或封测；可以同时为多家设计公司提供服务，但受制于公司间的竞争关系</p>
<p>不承担由于市场调研不准、产品设计缺陷等决策风险。</p>
<p>投资规模较大，维持生产线正常运作费用较高；需要持续投入维持工艺水平，一旦落后追赶难度较大。</p>

典型厂商

基本特点

主要优势

主要劣势

早期企业都是IDM运营模式（垂直整合），这种模式涵盖设计、制造、封测等整个芯片生产流程，这类企业一般具有规模庞大、技术全面、积累深厚的特点，如Intel、三星等

随着专注于晶圆加工的台积电的出现，演化出Fabless和Foundry模式，专攻设计或者制造，各司其职

# 目录

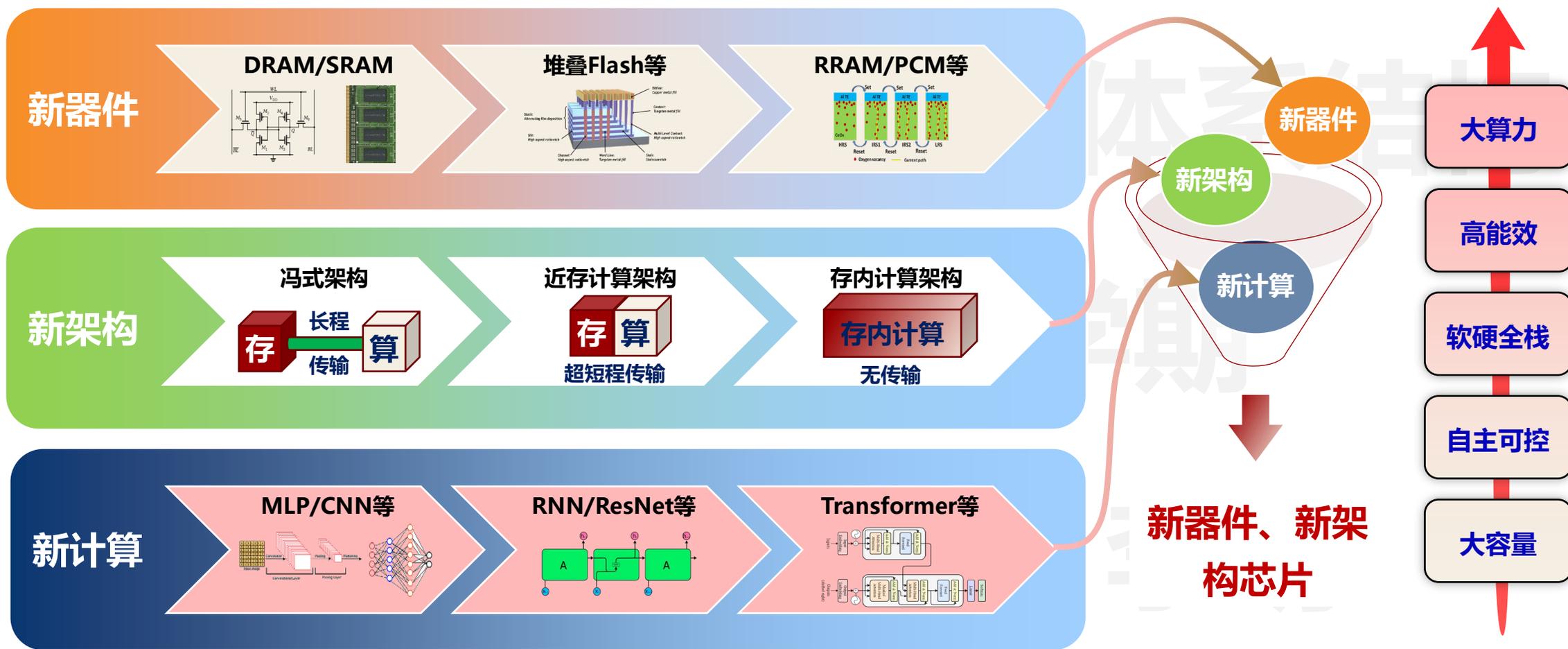
CONTENTS



01. 课程简介与体系结构概念
02. 智能芯片历史与发展趋势
03. 智能芯片产业国内外现状
04. 新兴技术与前沿发展趋势

# 融合新器件、新架构、新计算是后摩尔时代体系结构的发展趋势

- 融合新器件、新架构、新计算是突破后摩尔时代大算力、高效能瓶颈的重大关键技术领域



# AI大模型网络结构、参数规模与算力需求快速演进

- 以AI大模型为代表的新一代人工智能系统对高性能AI芯片提出了新的要求



历史时期	算力需求	翻倍间隔
前深度学习时代 1952 – 2010	30 KOPS – 200 TOPS	21.3月
深度学习时代 2010 – 2022	700 TOPS – 2 EOPS	5.7月
大模型时代 2016 – 2022	1 ZOPS – 1 YOPS	9.9月

代表性AI大模型	参数量	算力需求
GPT-4	~1.5万亿个	~2.7 YOPS
GPT-3	~1746亿个	~314 ZOPS
GPT-3 Small	~1.25亿个	~224 EOPS

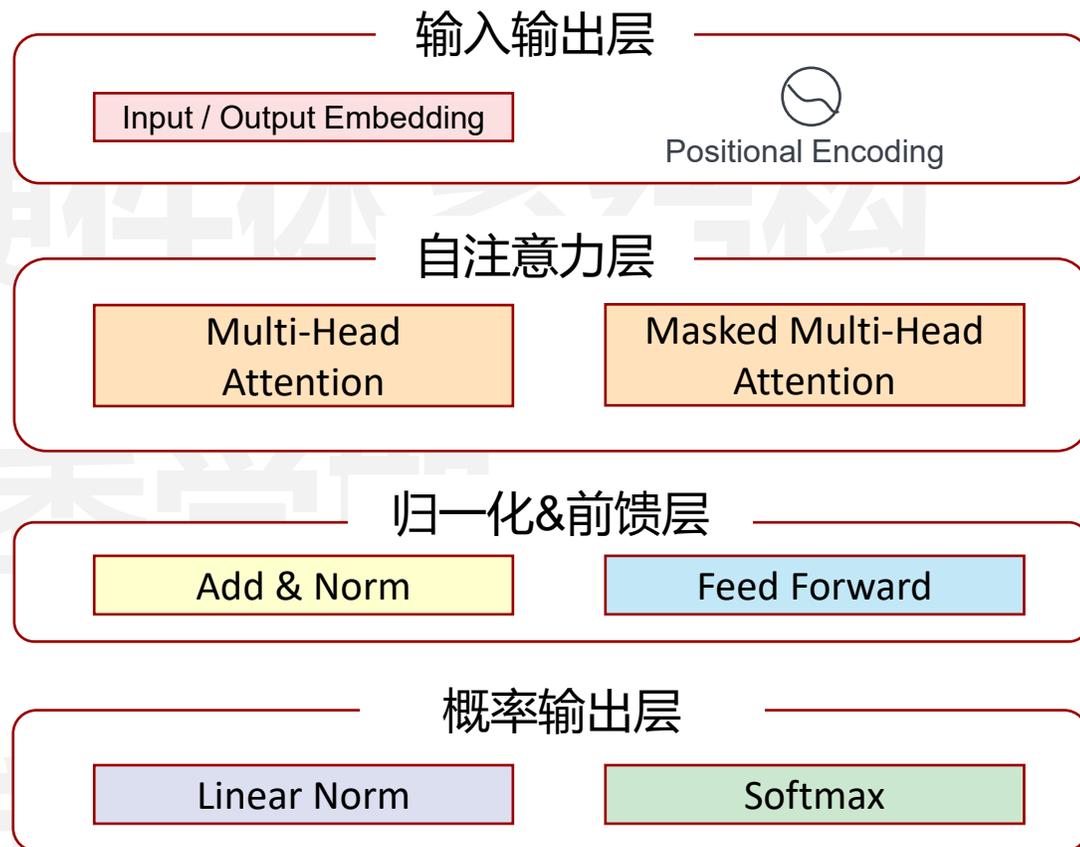
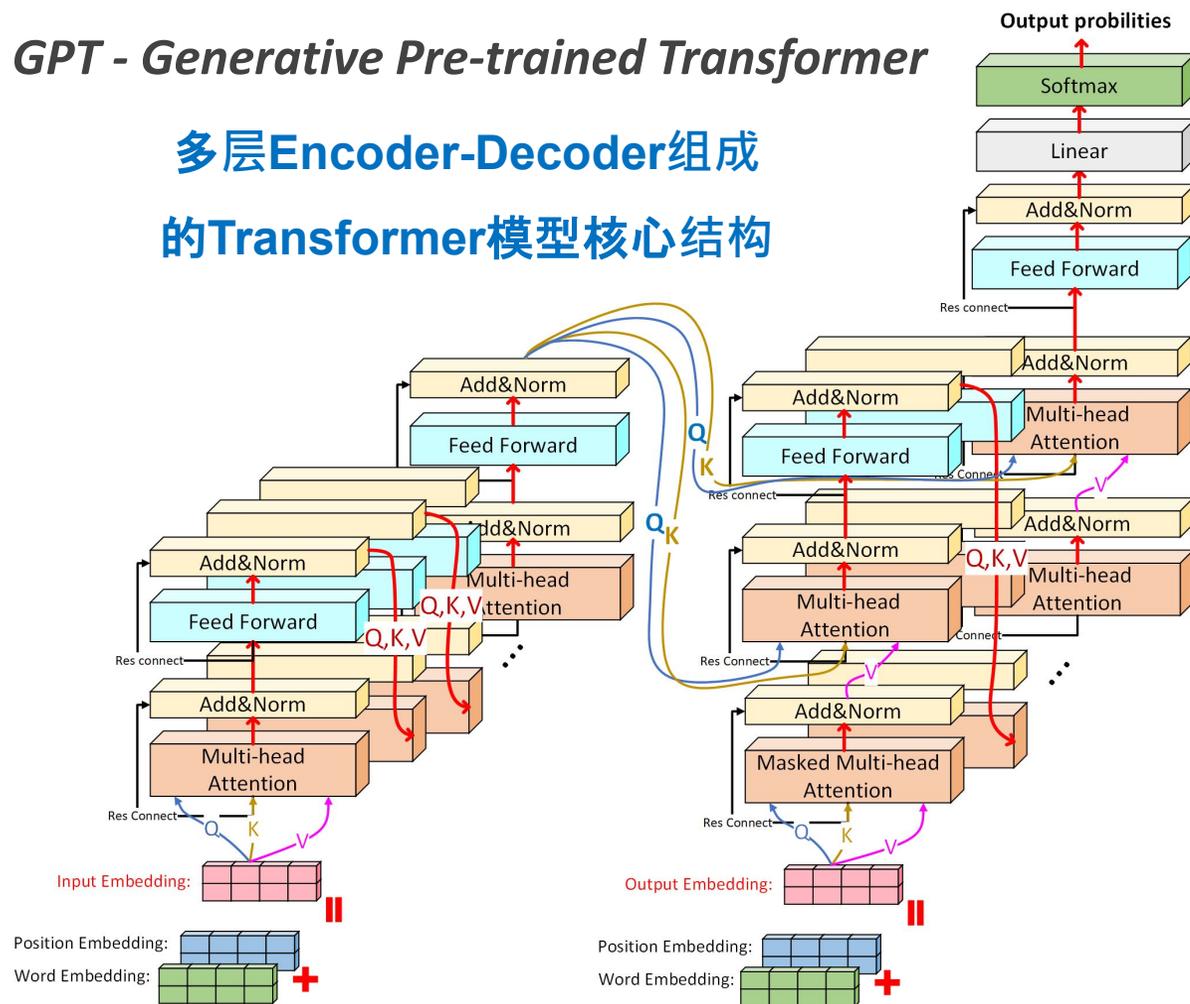
芯片性能成为支撑智能系统从量变产生质变的基石

# 当前AI大模型以Transformer为骨干网络 (以GPT为例)

- Decoder-Encoder层数、Token数量、掩码Mask尺寸、特征矩阵尺寸急剧增大

## GPT - Generative Pre-trained Transformer

多层Encoder-Decoder组成的Transformer模型核心结构



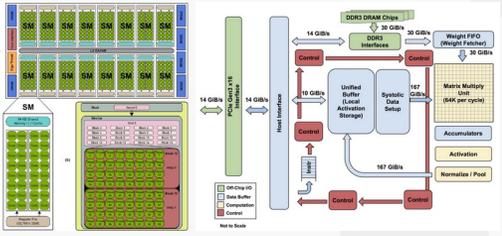
微软/OpenAI于2023年6月提出了LongNet, 将Transformer的Token数提高到了10亿级别

# 当前冯氏芯片对AI大模型的支持

- 按技术路径可分为通用AI芯片、定制AI芯片、可重构AI芯片、神经形态AI芯片等

## 冯氏AI芯片技术路线图及其发展现状

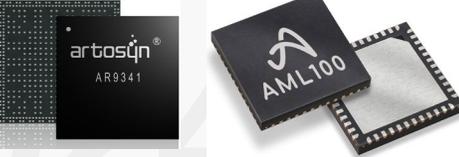
### 通用AI芯片 (云、数据中心、边缘)



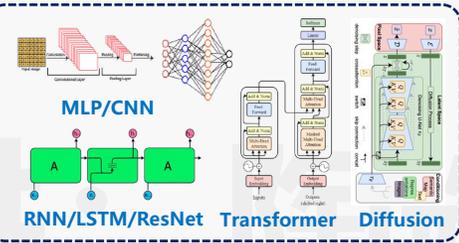
**通用GPU**      **通用TPU/NPU等**

- 支持通用计算的指令集架构
- 支持高计算精度浮点数运算
- 规模大、算力大、编程性佳

### 定制AI芯片 (边缘、终端设备)



**视觉AI芯片**      **语音AI芯片**

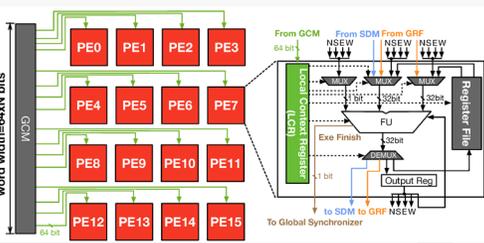


- 支持特定应用的几种AI模型
- 大部分采用定点数运算精度
- 中小规模、有限的可编程性

### 可重构AI加速器 (云、边、端)

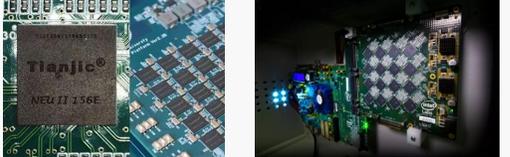


**FPGA加速卡**      **可重构阵列AI芯片**

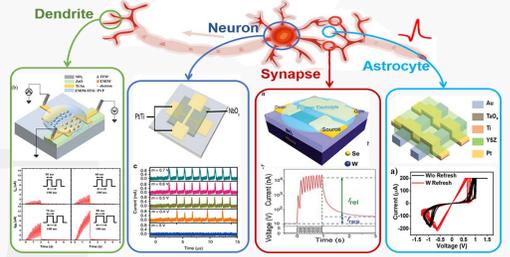


- 快速灵活、硬件逻辑高可编程性
- 不受数据类型限制、吞吐率较高
- 即插即用、云边端应用均可胜任

### 神经形态AI芯片 (终端设备)

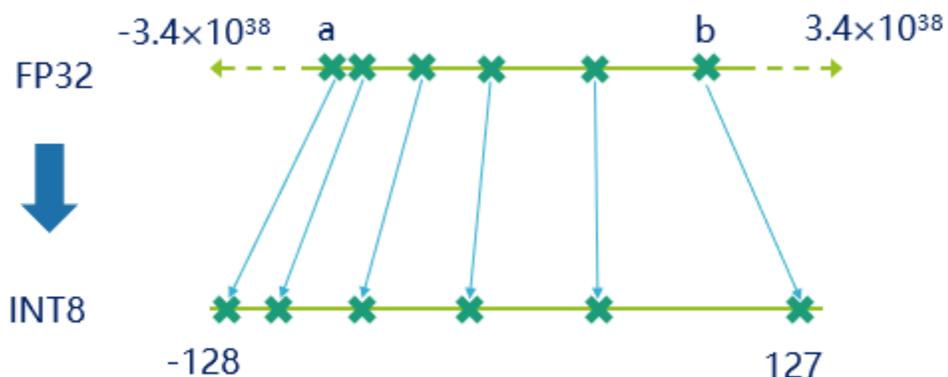


**SNN芯片**      **神经形态芯片**

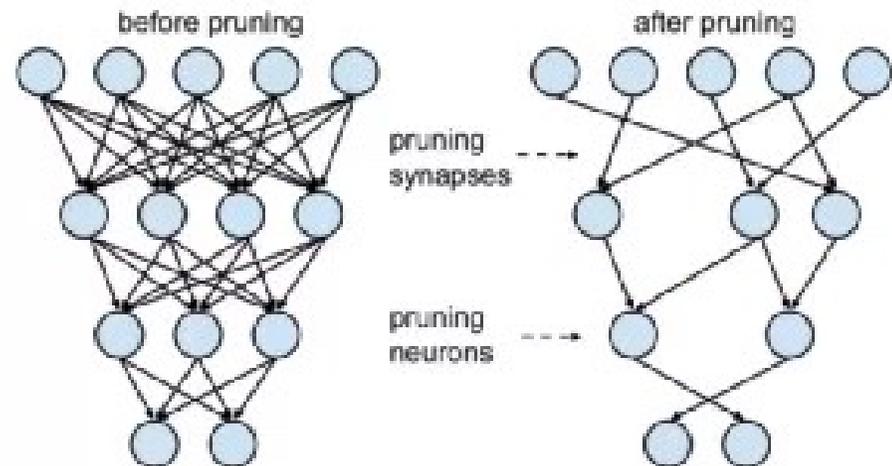


- 利用电路或器件模拟生物神经元
- 大量模拟神经元相连构成接近于人脑神经系统的类脑智能芯片

- 面对复杂应用，单纯的硬件设计已经不足以支持性能需求，需要软硬件协同设计

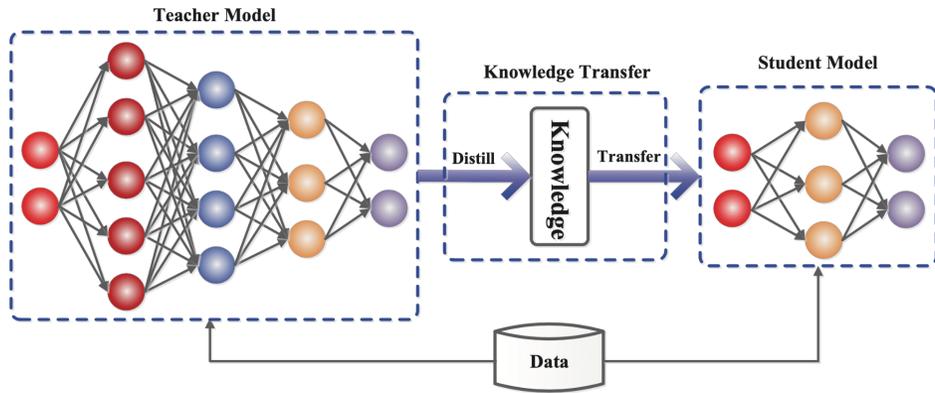


**模型量化**：将高精度的权重量化为低精度的权重，以一定的精度损失为代价换取更小的存储和计算开销

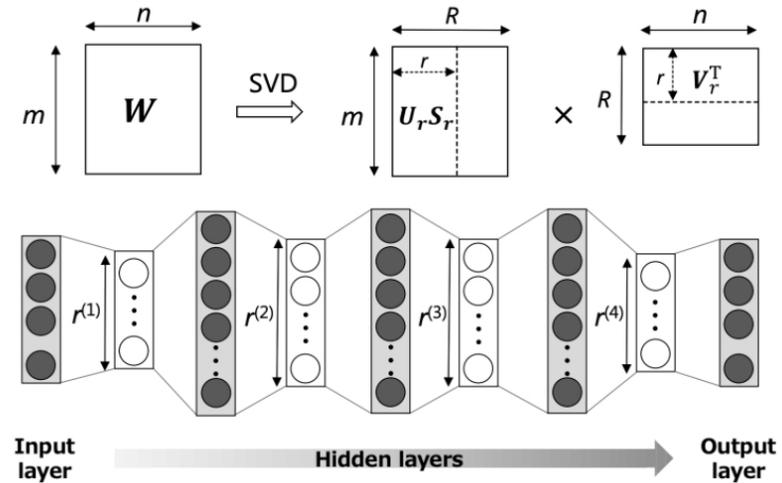


**模型剪枝**：将神经网络中重要性较小的神经元和权重删除，减少计算量，加速神经网络推理

- 面对复杂应用，单纯的硬件设计已经不足以支持性能需求，需要软硬件协同设计

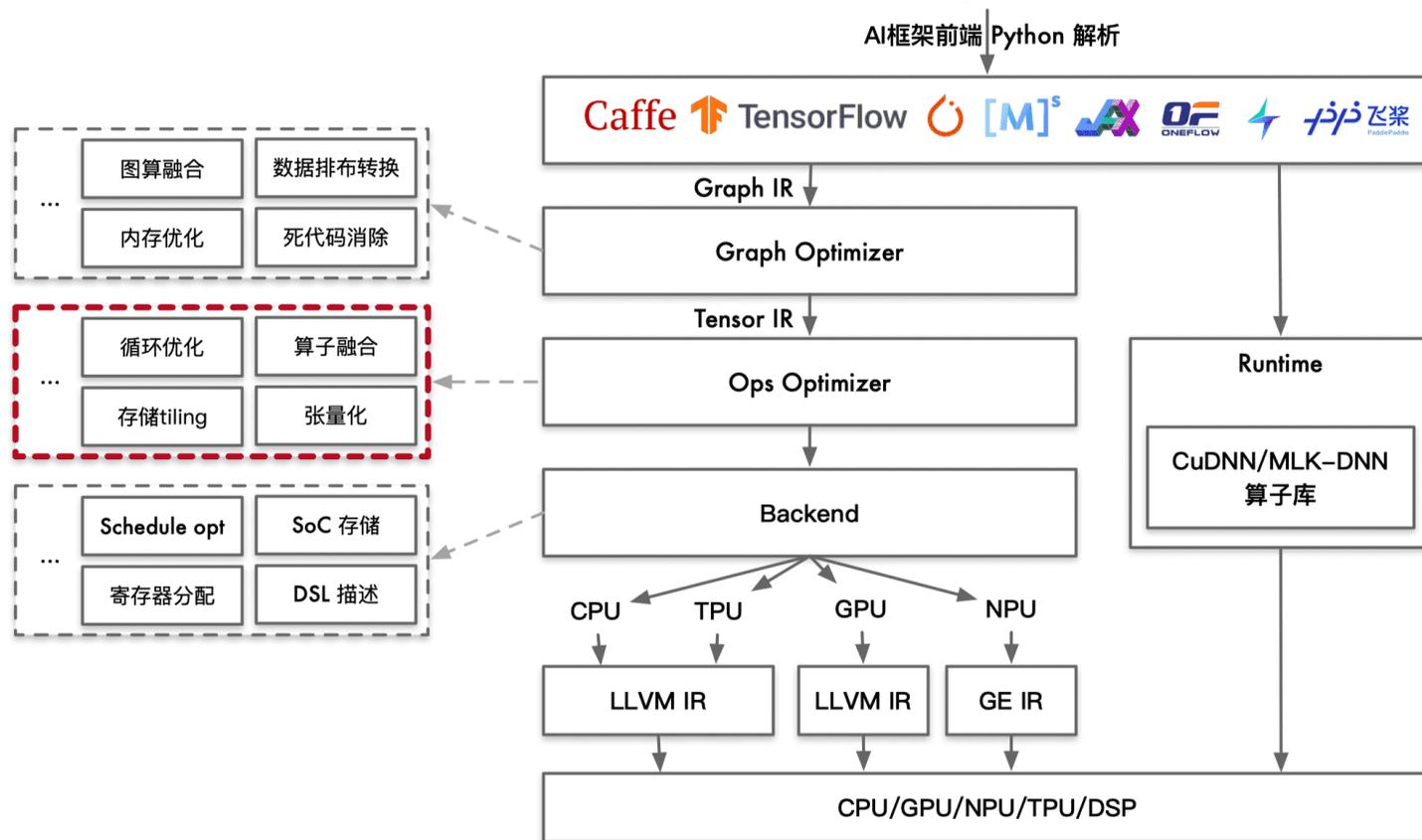


**知识蒸馏：**将规模较大的模型作为 teacher model 训练一个较小的 student model，在尽可能保证性能的情况下减小模型规模



**低秩分解：**将大规模权重分解为两个小规模权重矩阵相乘（SVD），减小矩阵向量乘的计算量

## • 编译层面优化



本系结构

在程序编译过程中

对算子、存储tiling

和寄存器分配等等

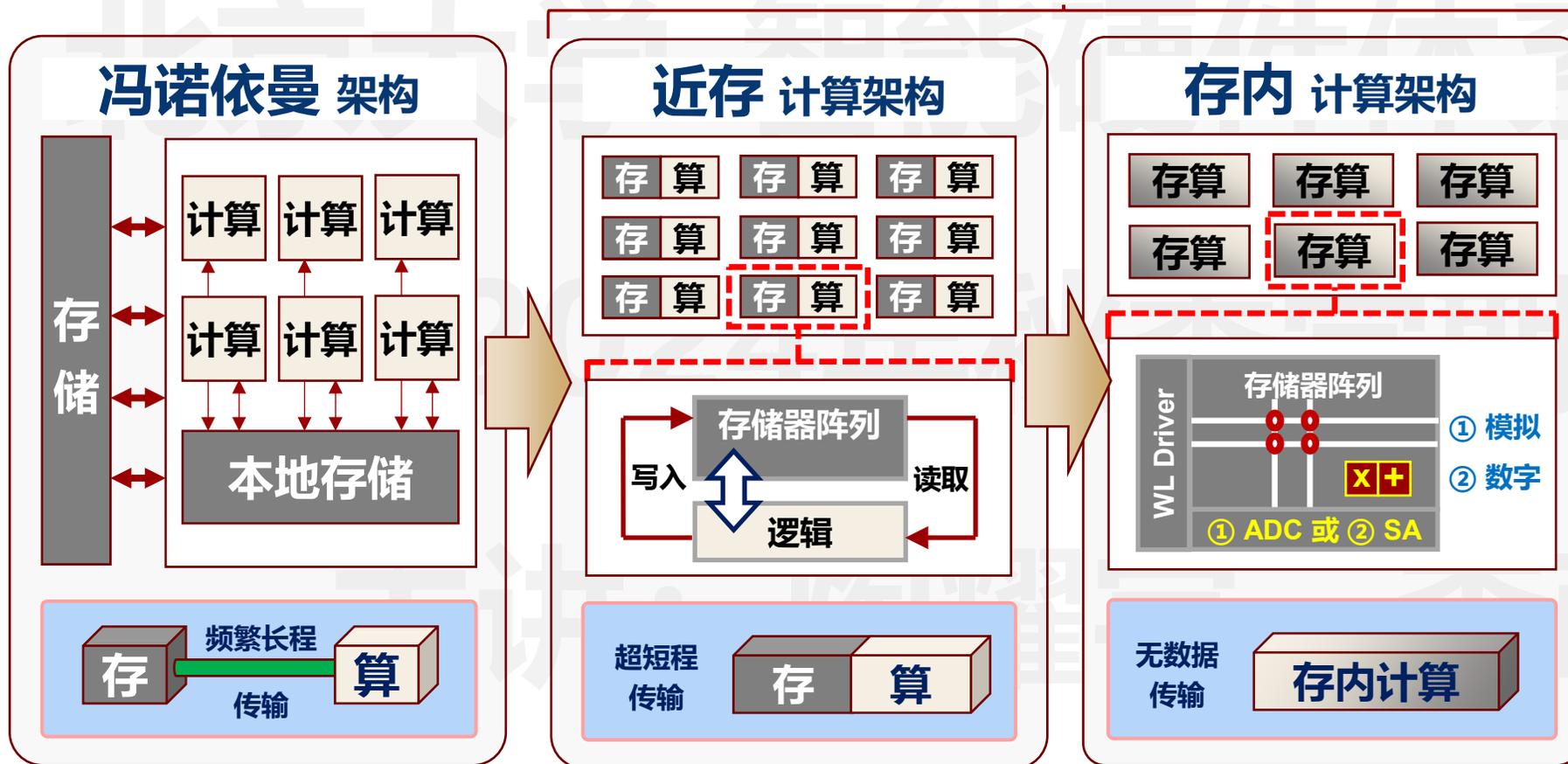
方面进行优化

李萌

# 代表性智能芯片新兴技术之二 – 新架构：存算一体

- 存算一体技术成为后摩尔时代打破算力瓶颈的重要路径

算力提升、能效提升 → 存算一体技术

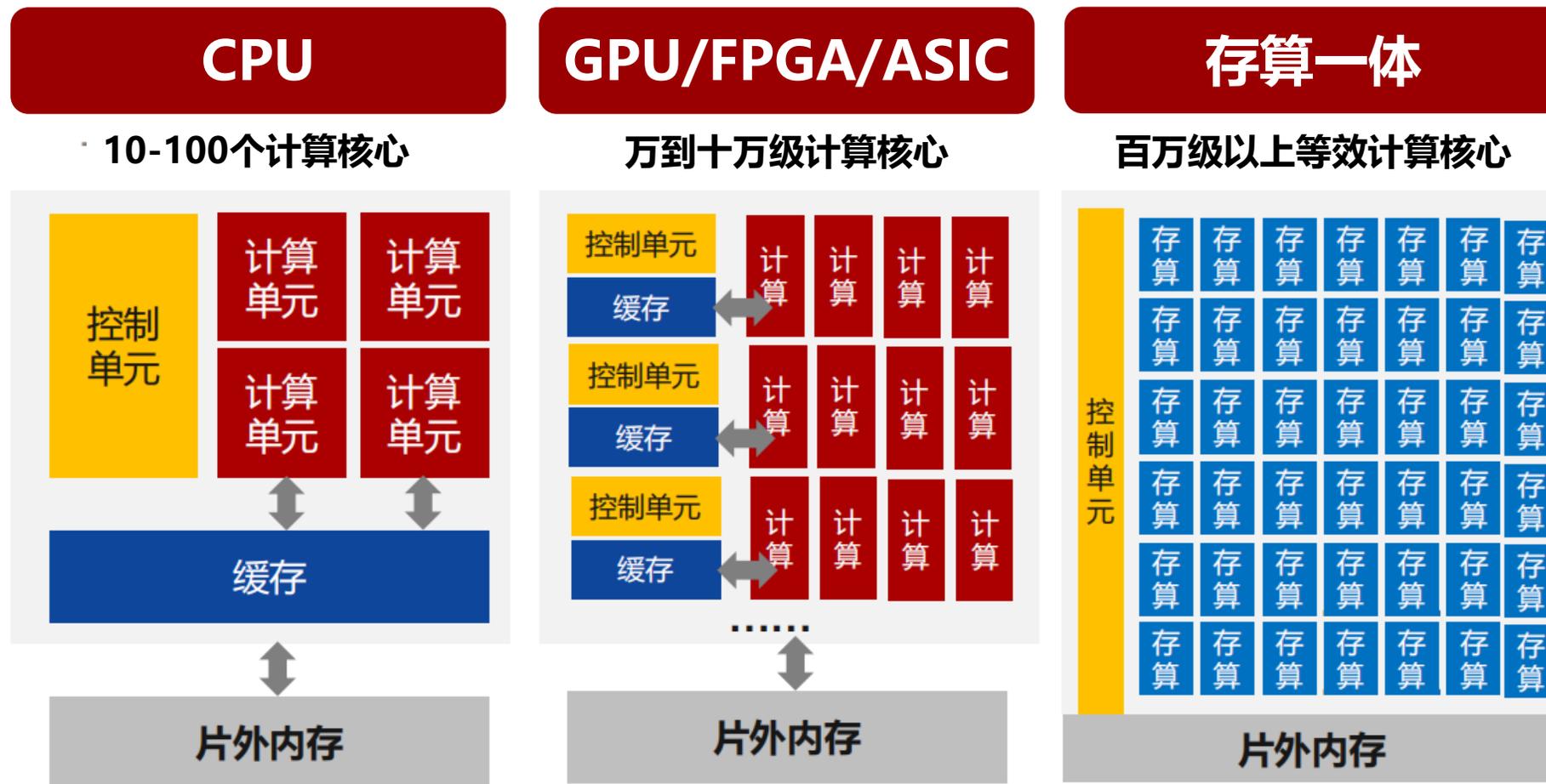


天然优势

-  大算力
-  低功耗
-  低延时

# 存算一体成为打破AI大模型推理算力极具潜力的技术路径

- 存算一体提供比GPU等冯氏芯片高多个数量级的并发度，有效支撑AI大模型推理

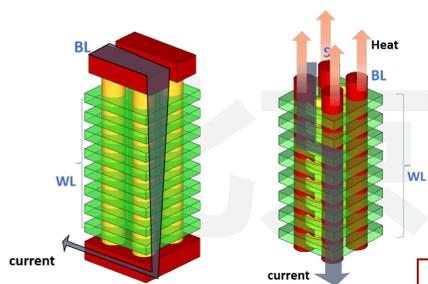


存算一体如何对AI大模型进行有效支持?

现有AI大模型推理基本上基于GPU/FPGA/ASIC等冯氏芯片

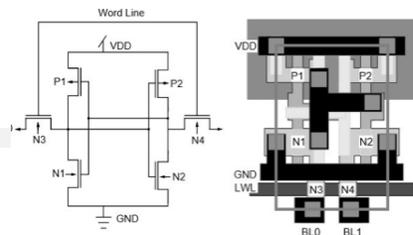
# 代表性智能芯片新兴技术之三 – 新器件：存储-计算融合器件

## 未来存储器介质材料的创新



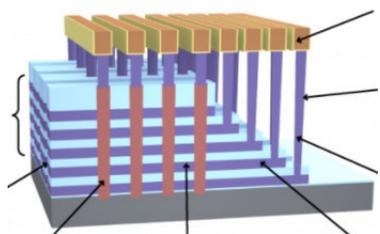
DRAM

**优点：工艺成熟、密度高**  
**缺点：速度低、刷新、只近存**  
**非易失性：否**  
**适合场景：冯氏架构过渡**



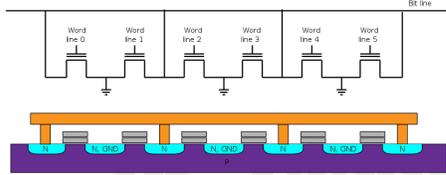
SRAM

**优点：工艺成熟、IP化应用**  
**缺点：能效低、密度低**  
**非易失性：否**  
**适合场景：端侧、边缘中小算力**



SSD/Nand Flash

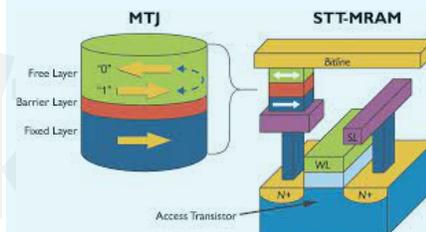
**优点：工艺成熟、容量大、成本低**  
**缺点：速度低、只能近存**  
**非易失性：是**  
**适合场景：云端大容量**



Nor Flash

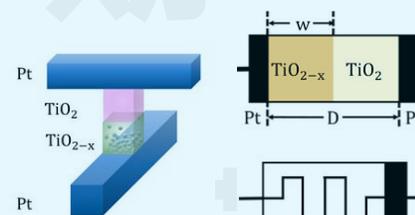
**优点：工艺成熟、密度高、成本低**  
**缺点：对PVT变化敏感、能效低**  
**非易失性：是**  
**适合场景：端侧、边缘低成本**

## 新兴器件



磁器件  
(MRAM)

**优点：能效、速度、密度高**  
**缺点：与CMOS大规模集成难**  
**非易失性：是**  
**适合场景：端侧、边缘中小算力**

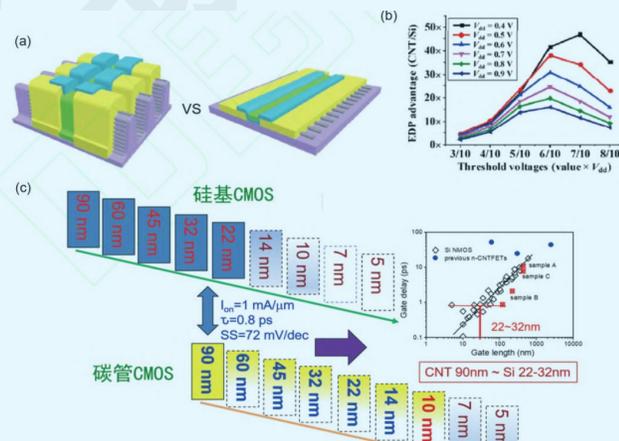
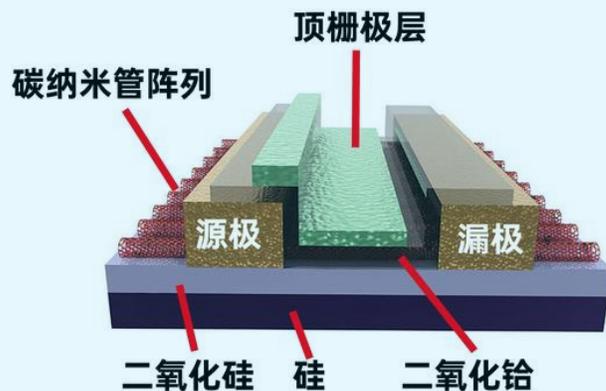
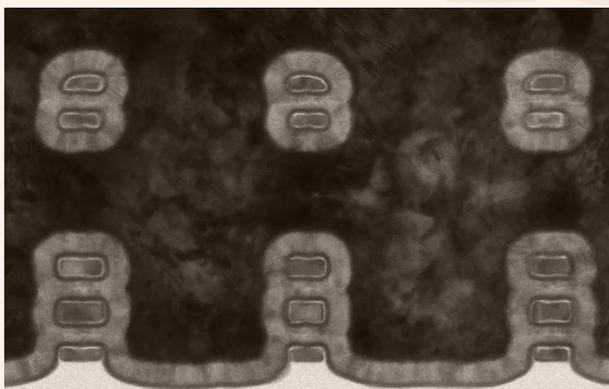
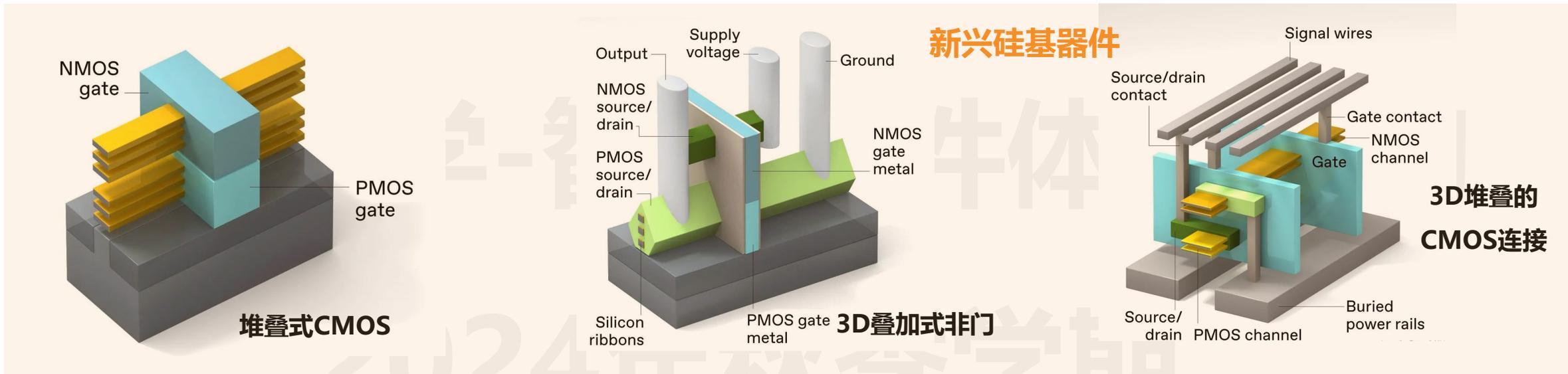


忆阻器  
(RRAM/PCM)

**优点：算力、能效、密度高**  
**缺点：工艺爬坡成熟中**  
**非易失性：是**  
**适合场景：云边端大算力**

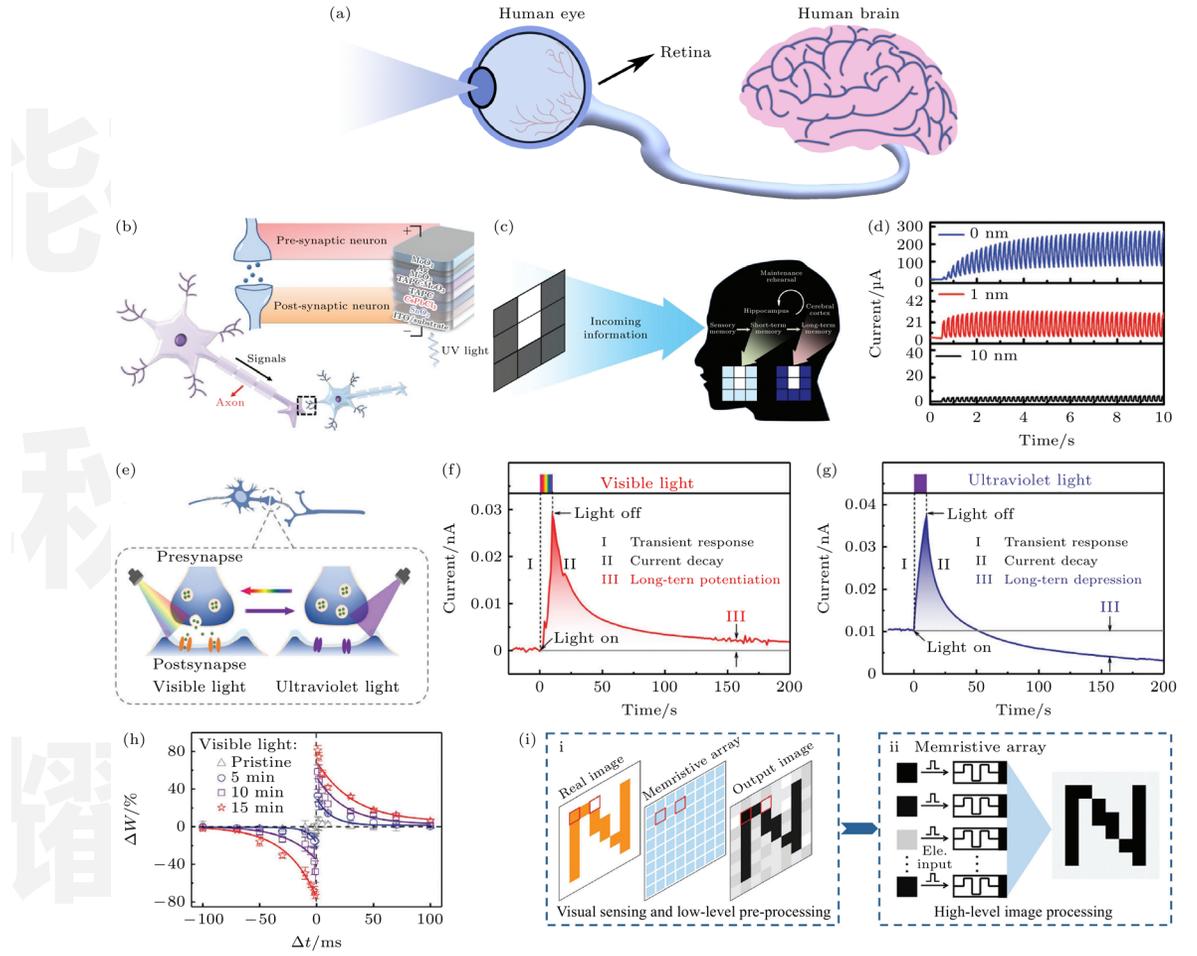
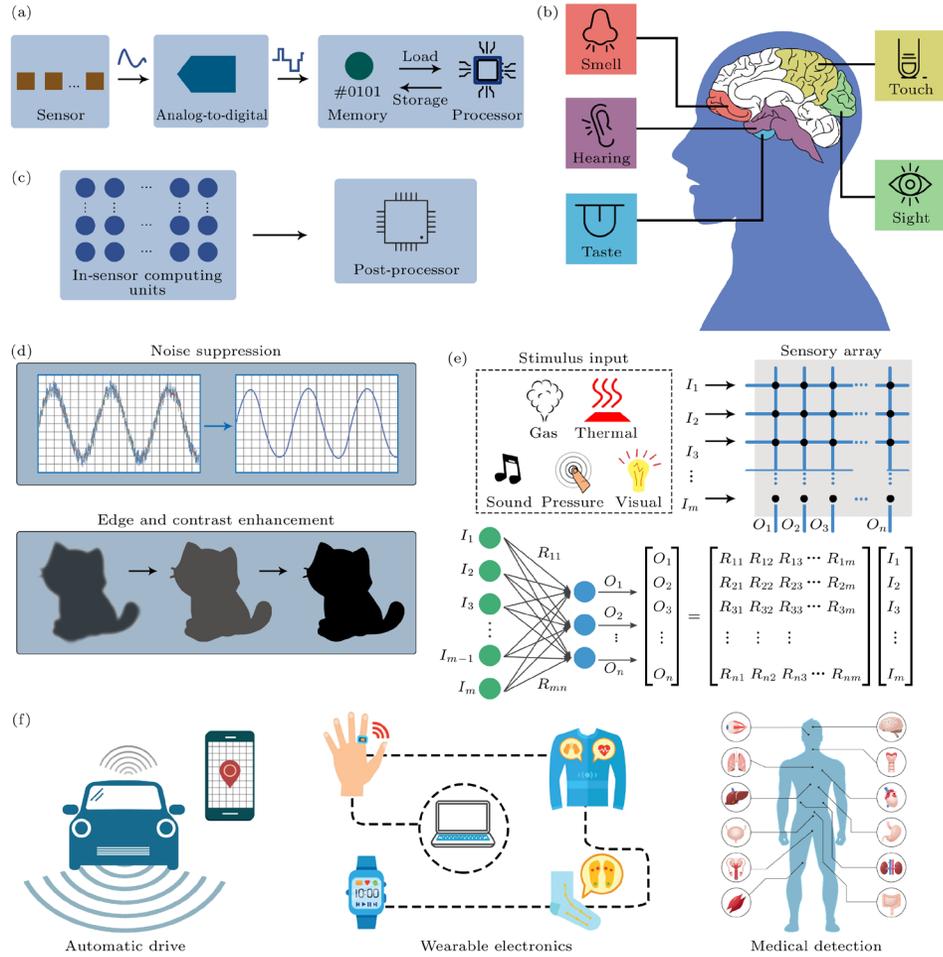
# 代表性智能芯片新兴技术之四 – 新器件：高密度的逻辑器件

## • 未来三维堆叠式晶体管与碳管器件



# 代表性智能芯片新兴技术之五 – 新架构：感存算一体

- 将传感、计算、存储融为一体，大幅降低系统功耗和计算延时，应用前景广阔

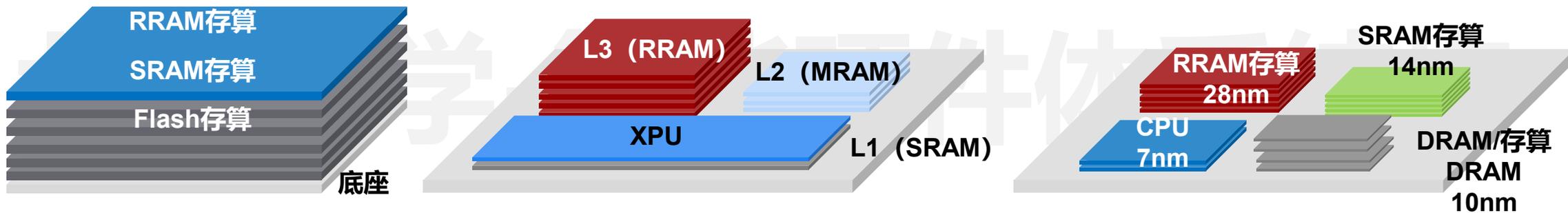


## 视觉感存算一体芯片与硬件系统

# 代表性新兴技术之六 – 新架构：三维异质集成

- 协同先进封装技术，实现多种芯片方案相结合

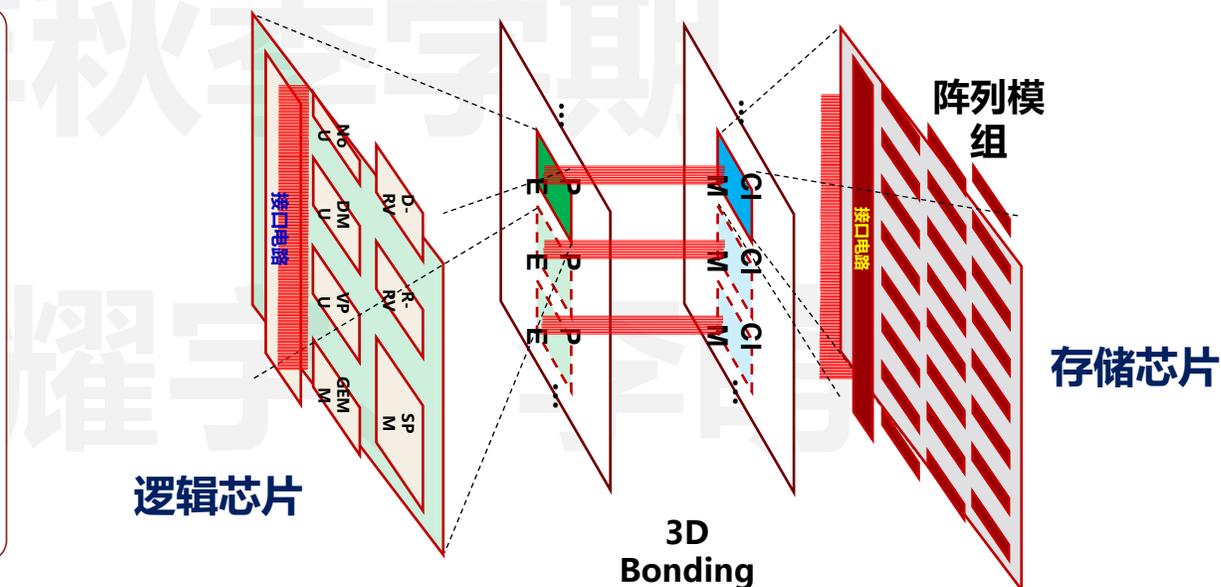
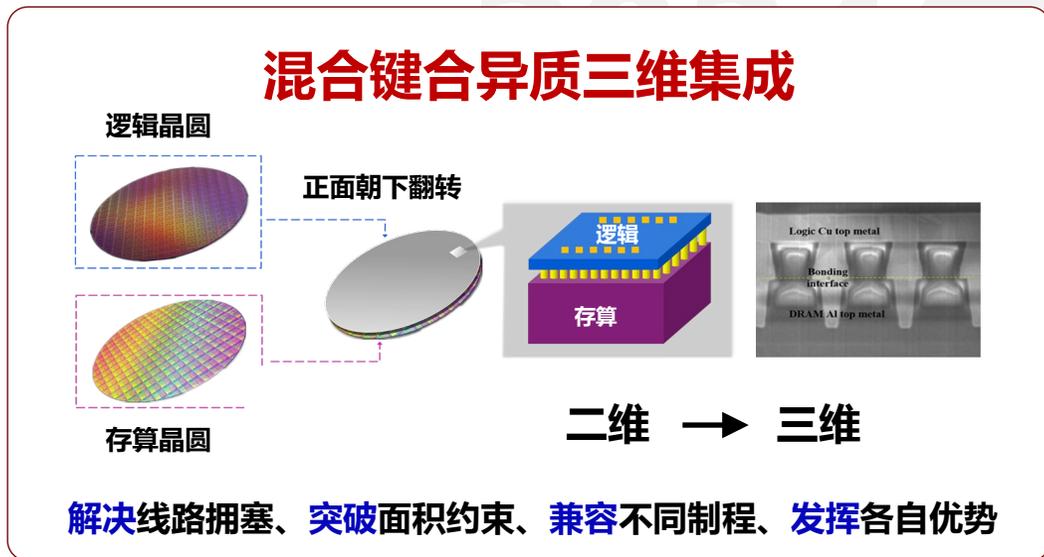
## 先进三维集成芯片示例图



三维集成

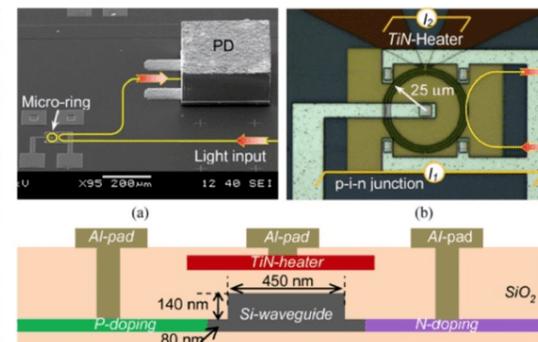
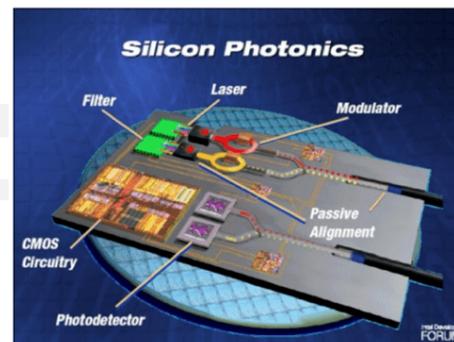
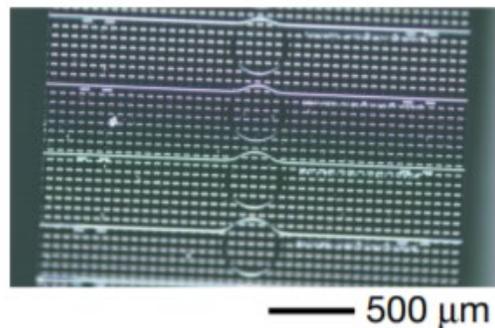
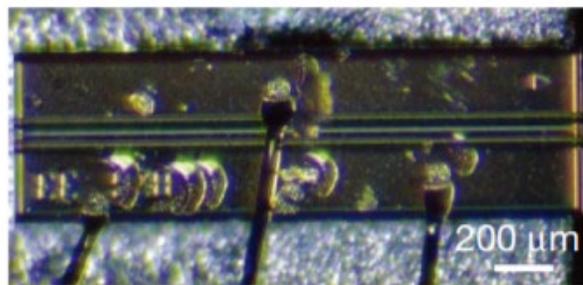
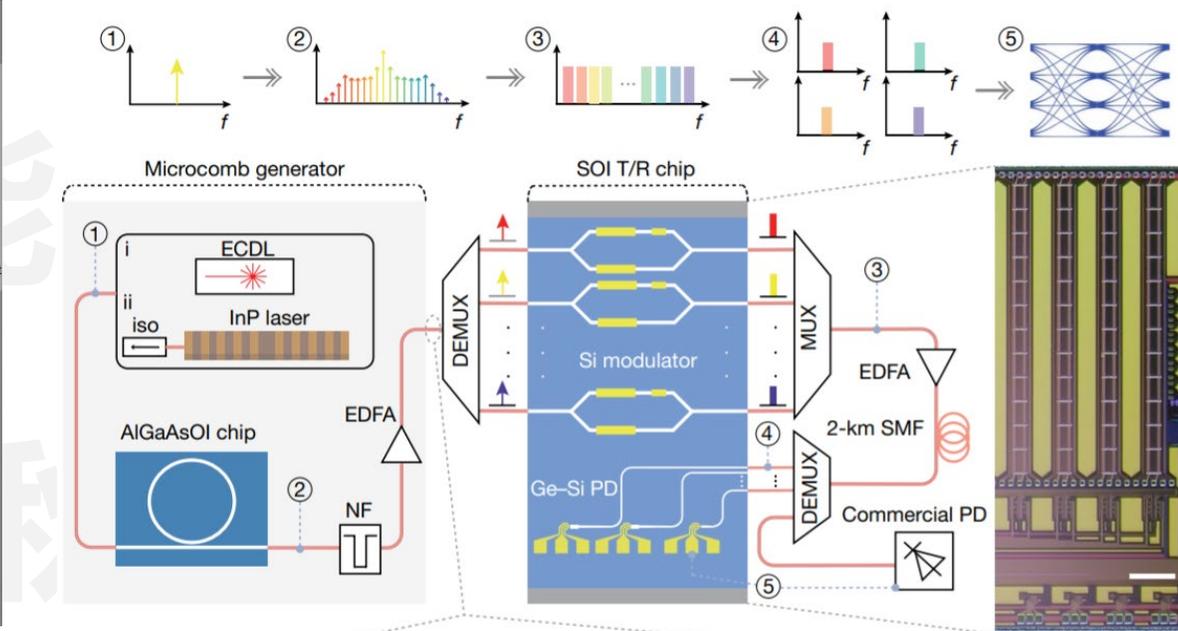
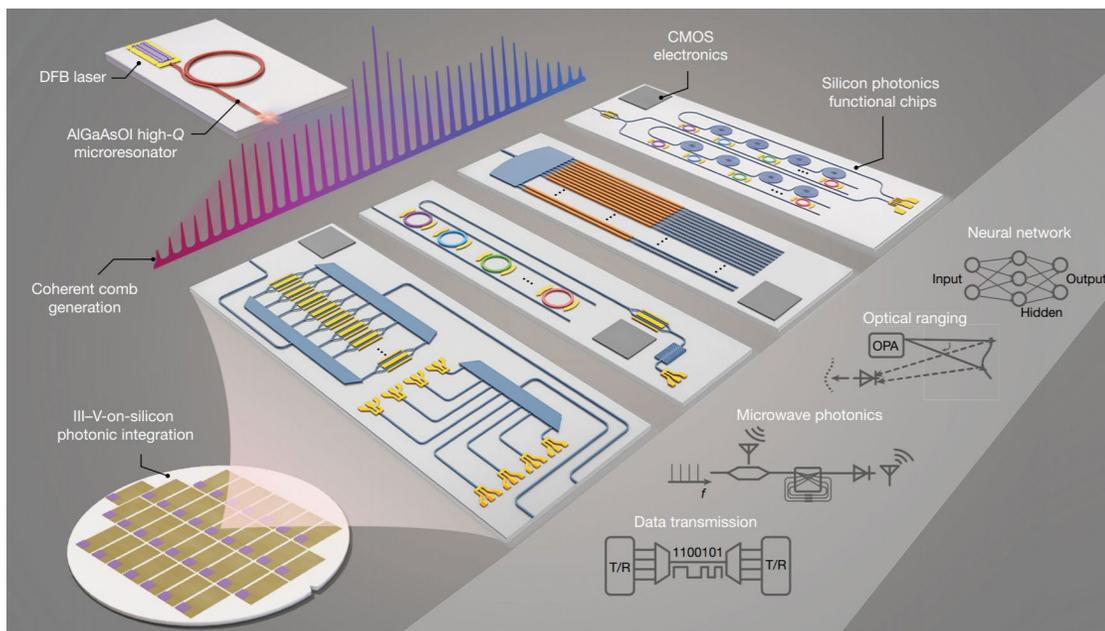
多级存储器堆叠SoC

异构小芯粒封装



# 代表性新兴技术之七 – 新架构：片上光互连技术

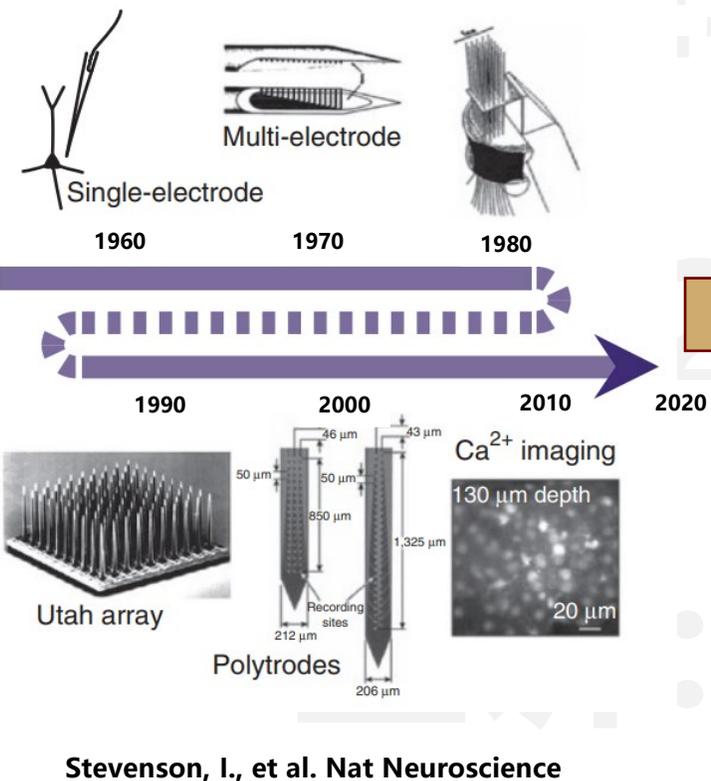
- 片上集成光电子通信系统有望突破信号传递延时的瓶颈，打破金属互连的物理上限



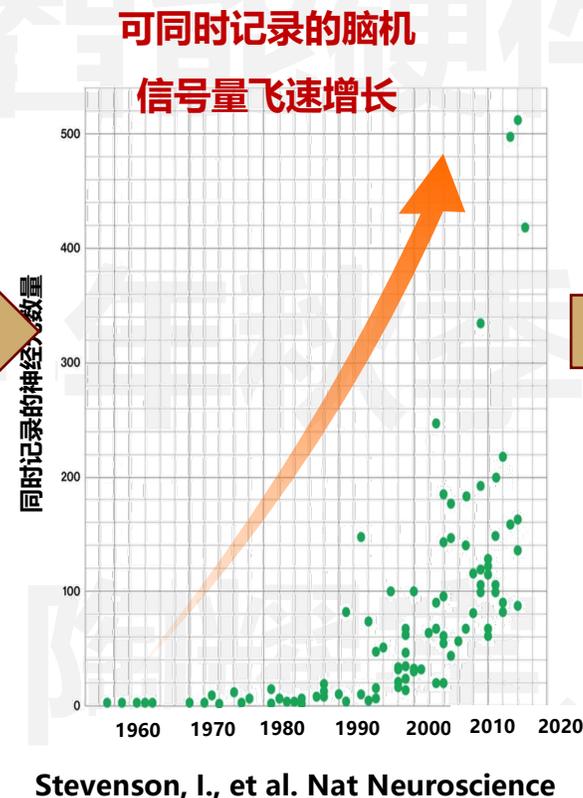
# 代表性新兴技术之八 – 新计算：例如脑机接口芯片与系统

- 为脑机接口服务的芯片与系统将在未来数十年成为人类发展的方向之一

## 脑机数据采集工具的发展



## 脑机数据量的“摩尔定律”



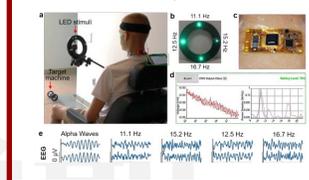
### 脑机打字



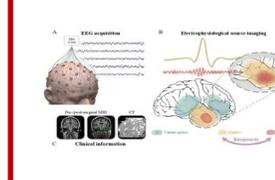
### 脑控无人机



### 脑控座椅



### 脑机癫痫监测



### 脑控机械臂

