

好的，我们来一起提炼这份PPT "CS253_L1112.pdf" 中的知识点，并进行分类和图片讲解，尽量引用PPT内容。

模块一：内存概述与历史

1.1 内存架构 (Memory Architecture)

- **定义与目标：** "Memory architecture describes the methods used to implement electronic computer data storage in a manner that is a combination of the fastest, most reliable, most durable, and least expensive way to store and retrieve information." (PPT L11, S2) 架构描述了以速度、可靠性、持久性和成本的最佳组合来实现电子计算机数据存储的方法。
- **本质：** "Memory a means of 'remembering '" (PPT L11, S2) 内存是一种“记忆”的手段。

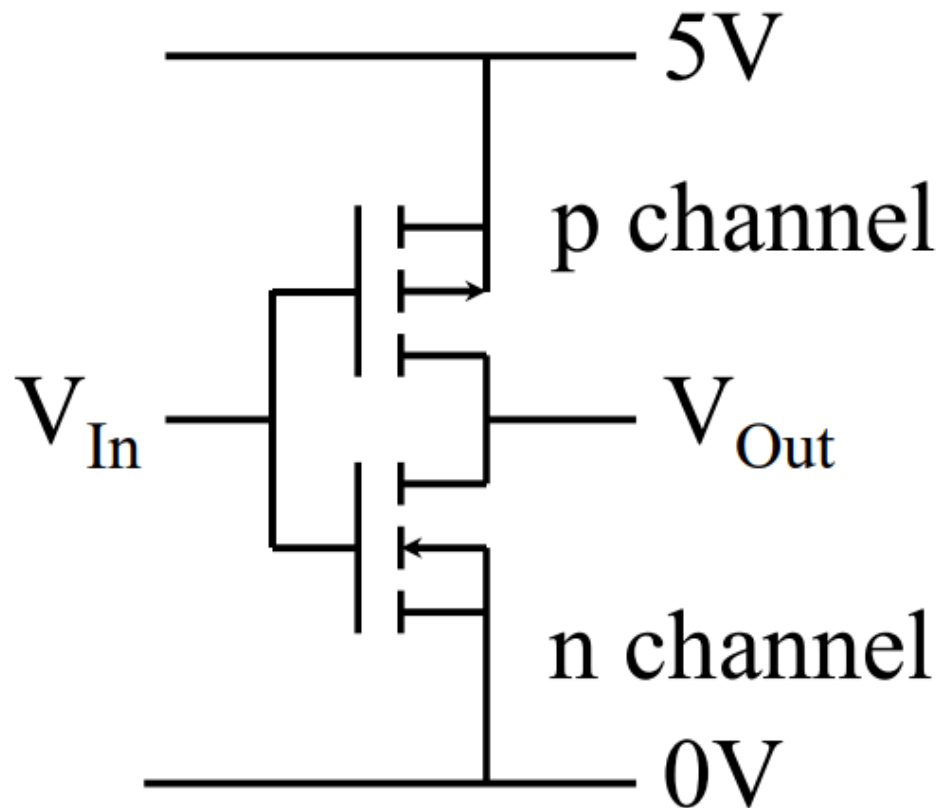
1.2 存储介质的演变 (图片与表格内容解读)

- **Cuneiform script (楔形文字泥板 - 苏美尔, 土耳其) 3100BC:** 存储 "100's bytes", 历史 "5000years"。 (PPT L11, S3)
- **Stone Tablet (石碑) 5000BC:** 存储 "100's bytes"。 (PPT L11, S3)
- **Papyrus (莎草纸) 3000BC:** 存储 "1K bytes"。 (PPT L11, S3)
- **Vellum (牛皮纸/羊皮纸) 800AD:** 存储 "10'sKbytes"。 (PPT L11, S3)
- **Paper book (纸质书) 1455AD:** 存储 "10's bytes" (PPT此处数据可能偏低)。 (PPT L11, S3)
- **Tape (磁带 - C90):** 存储 "660 kBytes"。 (PPT L11, S3)
- **LP (黑胶唱片):** "40mins Audio"。 (PPT L11, S3)
- **VHS (家用录像系统磁带) 1976:**
存储 "6 GB 200min"。 (PPT L11, S3)
 - **图片解读 (PPT L11, S3):** 该幻灯片通过图片和文字展示了从古老的泥板、石碑到近代的纸张、磁带等存储介质的演进，直观地体现了存储容量和技术在数千年间的巨大发展。

1.3 现代内存基础 - CMOS 反相器 (CMOS Inverters)

- **CMOS 定义：** "CMOS-Complementary Metal Oxide Semiconductor" (PPT L11, S4) 互补金属氧化物半导体。
- **电路结构：**
由一个P沟道 (p channel) MOSFET和一个N沟道 (n channel) MOSFET组成。 (PPT L11, S4)
 - 当输入 V_{in} 为高电平 (5V) 时，P沟道截止，N沟道导通，输出 V_{Out} 为低电平 (0V)。

- 当输入 V_{in} 为低电平 (0V) 时, P沟道导通, N沟道截止, 输出 V_{Out} 为高电平 (5V)。



- 功耗特点:

"Note: there is only dissipation during switching" (PPT L11, S4) 只有在开关状态切换的瞬间才有功耗, 静态时功耗很低。

- 图片解读 (CMOS Inverters - PPT L11, S4):

- 左侧图示为一个基本的CMOS反相器电路图, 清晰标注了P沟道和N沟道MOSFET, 输入 V_{in} , 输出 V_{out} , 以及电源(5V)和地(0V)。
- 右侧图示了输入电压(V_{in})和输出电压(V_{out})的理想反相关系曲线, 表明输入高则输出低, 输入低则输出高。

模块二: 内存层级 (Memory Hierarchy)

2.1 内存层级结构

- **目的:** 结合不同存储技术的优点, 平衡速度、容量和成本。
- 典型层级 (从快/小/贵到慢/大/便宜):

(PPT L11, S5)

1. **Registers (寄存器)**
2. **Cache (缓存)**
3. **Main Memory (主内存)**
4. **Disk Cache (磁盘缓存)**
5. **Magnetic Disk (磁盘)**
6. **Magnetic Tape (磁带) / Optical Disk (光盘)**

- 图片解读 (Memory Hierarchy - 三角形图 - PPT L11, S5):

- 此图以金字塔形式展示内存层级。塔尖是寄存器，速度最快、容量最小、成本最高。塔底是磁带/光盘，速度最慢、容量最大、成本最低。CPU直接与靠近塔尖的层级交互，数据根据访问频率和局部性原理在不同层级间移动。

2.2 内存层级特性对比表 (PPT L11, S6)

- **特性维度:** "Location, Capacity, Access Speed, Volatility" (PPT L11, S6)
- 寄存器 (Registers):
 - 位置 (Location): "On CPU chip"
 - 容量 (Capacity): "Bytes"
 - 访问速度 (Access Speed): "Instruction <3 clock cycles"
 - 易失性 (Volatility): "Yes"
- 缓存 (Cache):
 - 位置 (Location): "On or off CPU chip"
 - 容量 (Capacity): "M-bytes"
 - 访问速度 (Access Speed): "10nS" (纳秒)
 - 易失性 (Volatility): "Yes"
- 主内存 (Main Memory):
 - 位置 (Location): "Motherboard"
 - 容量 (Capacity): "G-bytes"
 - 访问速度 (Access Speed): "80nS"
 - 易失性 (Volatility): "Yes"
- 磁盘 (Disk):
 - 位置 (Location): "In PC Case"
 - 容量 (Capacity): "T-bytes"
 - 访问速度 (Access Speed): "10mS" (毫秒) - "the time shown is the latency (time taken to read the first byte)." (PPT L11, S6)
 - 易失性 (Volatility): "No"
- DVD:
 - 位置 (Location): "In/Out PC case"
 - 容量 (Capacity): "**G-bytes**"
 - 访问速度 (Access Speed): "100mS" - "**DVD and Disk** can transfer data fast (Mbits/sec), the time shown is the latency." (PPT L11, S6)
 - 易失性 (Volatility): "No"

模块三：主内存技术 - SRAM 与 DRAM

3.1 静态RAM (Static Ram - SRAM) (PPT L11, S7)

- **构造基础:** "Construction based on Flip Flops" (基于触发器)。
- **单元尺寸:** "Each memory cell uses up a lot of real estate" (每个存储单元占用较多芯片面积)。
- **特性:**

- "Memory is volatile, but does not need refreshing" (内存是**易失的**，但**不需要刷新**)。
- "Much faster than **Dynamic Ram**, used for caches, 10nS" (比动态RAM快得多，用于缓存，速度约10纳秒)。
- SRAM 存储单元 (Static RAM Cell - PPT L11, S8):
 - **功能**: "Stores 1 bit of information" (存储1位信息)。
 - **结构**: "Requires up to 8 FETs per bit" (每个位需要多达8个场效应晶体管)。通常标准6T SRAM单元使用6个晶体管。
 - **特性**: "Volatile" (易失的), "No-refresh" (无需刷新), "Very fast 10nS" (非常快，10纳秒)。
 - **访问**: "RAM=Random Access Memory" (随机存取存储器)。通过X地址和Y地址线选定单元。
 - 图片解读 (Static RAM Cell - PPT L11, S8):
 - 展示了一个典型的SRAM存储单元的电路示意图。由交叉耦合的反相器（形成触发器）和访问控制晶体管组成。图中可见两条数据线 (Data, Data非)，两条地址线 (X Address, Y Address) 用于选择单元，以及电源(+5V)和地(0V)。

3.2 动态RAM (Dynamic Ram - DRAM) (PPT L11, S7)

- **构造基础**: "Construction based on charge stored on the gate of a FET." (基于存储在场效应晶体管栅极上的电荷)。更准确地说是存储在微型电容器上，FET用于访问控制。
- **密度**: "Much higher densities possible (larger memory)" (可以实现更高的存储密度，即更大的内存容量)。
- **成本**: "Lower Cost" (成本较低)。
- **刷新**: "But requires extra refresh control signal." (但需要额外的刷新控制信号)。
- **速度与功耗**: "Typically 60nS, no dissipation when not accessed." (典型速度60纳秒，不访问时无功耗)。
- **DRAM 存储单元** (Simplified Dynamic Memory 1 bit - PPT L11, S9):
 - **原理**: "Information stored as **charge on capacitor**." (信息作为**电荷**存储在电容器上)。
 - **刷新需求**: "Needs refreshing (rereading) to keep charge on capacitor." (需要刷新（重新读取）以保持电容器上的电荷)。
 - **结构**: "Requires apx 1 FET per bit." (每个位大约需要1个场效应晶体管和一个电容)。
 - **特性**: "**Volatile**" (易失的), "Slower 100nS to read or write" (读写较慢，约100纳秒)。
 - **组件**: "Sense Amplifier" (读出放大器) 用于检测微弱的电荷信号，"Buffer" (缓冲器)。
 - 图片解读 (Simplified Dynamic Memory 1 bit - PPT L11, S9):
 - 展示了一个最简化的DRAM存储单元。包含一个晶体管(FET)和一个电容器。通过行选通线(Row)和列选通线(Col)以及读(R)/写(W)控制信号来访问。A0和A1是地址位。信息（1位）存储在电容器C中。读出放大器用于在读取时放大存储的电荷信号。
- 4位DRAM示例 (Simplified Dynamic Memory 4 bit - PPT L11, S10):
 - **图片解读**: 展示了如何将多个1位DRAM单元组织起来形成多位存储。图中显示了4个DRAM单元，通过行选择(Row select R0, R1)和列选择(Column select A1)来访问特定的位。RAS (Row Address Select) 和 CAS (Column Address Select) 是行地址选通和列地址选通信号，通常RAS先行，CAS后行。

3.3 SRAM 与 DRAM 对比 (PPT L11, S11)

- 芯片面积与密度：
 - SRAM (6T): "**200nm x 500nm**". 在1平方毫米硅片上约可集成 "10M" 个晶体管, 对应 "Apx. 2Mbytes" (这里应为bits或更小容量的bytes, 原文Mbytes可能有误, 因为10M晶体管/8晶体管每位 = 1.25M位, 约0.15MB。若按6T算, 约0.2MB)。
 - DRAM: "**3.6um x 0.6um**" (这里原文可能是笔误, 应为更小的尺寸如36nm x 60nm, 或者um的单位使用有误, 因为3.6um x 0.6um的面积远大于SRAM单元)。假设原文数据有误, PPT结论是: "For the same silicon DRAM gives x3 the memory." (在相同的硅片面积上, **DRAM提供的存储容量大约是SRAM的3倍**)。
 - 图片解读 (Compare SRAM and DRAM - PPT L11, S11):
 - 左侧是SRAM单元的显微照片或版图, 显示了其**相对复杂**的结构 (6个晶体管)。
 - 右侧是DRAM单元的示意图或版图, 显示了其**相对简单**的结构 (1个晶体管和1个电容)。
 - 通过对比, 直观展示了**DRAM在集成密度上的优势**。

3.4 DRAM 操作周期 (Dynamic Ram Cycle - PPT L11, S12)

1. "A row is selected" (选择一行)。
 2. "Data is read out of each cell" (从该行的每个单元读出数据)。
 3. "If RAM is in read mode the column line is activated and the data written back to the memory cell." (如果RAM处于读模式, 列线被激活, **数据被写回存储单元** - 这是破坏性读出后的恢复)。
 4. "If RAM is in write mode the column line is activated and new data is written to the memory cell." (如果RAM处于写模式, 列线被激活, **新数据**被写入存储单元)。
 5. **刷新**: "Each memory cell must be re-written every 2mS (refreshed) so that charge is not lost from the capacitor in the memory cell, the RAS and CAS lines are used to achieve this." (每个存储单元必须**每2毫秒**被重写 (刷新) 一次, 以防止电荷从电容器中丢失, RAS和CAS线用于实现此目的)。
- **信号**: "RAS=Row Address Select, CAS=Column Address Select" (行地址选通, 列地址选通)。

3.5 DRAM 读时序图 (Timing Diagram D-Ram (Read) - PPT L11, S13)

- **信号线**: CAS (列地址选通), RAS (行地址选通), Address (地址总线, 分时复用行地址A0-A8和列地址A9-17), R/W (读写控制线), Gate (可能是输出使能或数据有效指示)。
- 过程:
 1. RAS信号先有效 (变低), 此时地址总线上提供行地址。
 2. 稍后CAS信号有效 (变低), 此时地址总线上提供列地址。
 3. R/W线保持**在读状态**。
 4. 经过一段时间的延迟后 (TRD=100nS), 数据总线上出现有效数据 (Valid Data)。
- **图片解读**: 该图清晰地展示了DRAM读操作中各个控制信号和地址/数据信号的时间先后关系。RAS先行, 然后是CAS。地址总线在RAS有效期间传输行地址, 在CAS有效期间传输列地址。数据在CAS有效一段时间后才出现在数据总线上。

3.6 DRAM 观察与特性 (Dynamic Ram, Observations - PPT L11, S14)

- 微处理器访问DRAM的步骤：
 1. "Split the address in two halves." (将地址分成两半，**行地址和列地址**)。
 2. "Select Read/Write." (选择**读或写**操作)。
 3. "Refresh each row every 2mS" (**每2毫秒刷新**每一行)。
- **解决方案**: "Dedicated chips are available to support the uP with this task." (有专门的芯片 (DRAM控制器) 来辅助微处理器完成这些任务)。
- **预充电与速度**: "Dynamic ram precharges a cell before it can be re-read, speed of ram is thus the total cycle time TRD which is about twice the access time (Typically 70nS)." (DRAM在重新读取前需要对单元进行预充电，因此RAM的速度是总周期时间TRD，大约是访问时间的两倍，典型值为70纳秒)。
- **等待状态**: "Due to low speed uP has to introduce wait states when accessing dram." (由于速度较慢，微处理器在访问DRAM时不得不引入等待状态)。
- **高密度**: "Very high density available e.g. 16Mbyte per device." (可实现非常高的密度，例如每个器件16MB)。
- **低成本**: "Low cost."

3.7 静态内存物理形态 (Static Memory (realised with gates) - PPT L11, S15-S17)

- **图片解读 (PPT L11, S15-S17)**: 这三张幻灯片展示了不同封装形式的静态内存芯片的物理外观照片，如DIP (双列直插封装) 芯片。这些图片帮助学生认识内存芯片的实体形态。

3.8 内存模块类型 (DIMS, SIMS, DIPS - PPT L11, S18)

- **SIMMs**: "Single Inline Memory Modules" (单列直插内存模块)。
- **DIMMs**: "Double Inline Memory Modules" (双列直插内存模块)。
- **DIP**:

"Dual in line package" (双列直插封装，通常指单个芯片)。

 - **图片解读**: 展示了SIMM和DIMM内存条以及DIP封装芯片的图片，这些是PC中常见的内存形态。

模块四：内存管理与缓存策略 (Lecture 12)

4.1 改善内存性能的问题 (Improving memory - PPT L12, S2)

- 核心问题:

"Is there a way of combining the Speed of SRAM with the high storage capacity of DRAM?" (没有办法将SRAM的速度与DRAM的高存储容量结合起来?)

 - **图片解读**: 用一个跷跷板的图示，一边是**SRAM (速度快但容量小)**，另一边是**DRAM (容量大但速度慢)**，中间是CPU，形象地表达了在速度和容量之间的权衡。

4.2 局部性原理 (Locality of reference - PPT L12, S3)

- **时间局部性 (Temporal locality):** "If an item is referenced, it will tend to be referenced again soon." (如果一个数据项被访问, 那么它很可能在不久的将来再次被访问)。
- **空间局部性 (Spatial locality):** "If an item is referenced, nearby items will tend to be referenced soon." (如果一个数据项被访问, 那么它附近的 (地址相邻的) 数据项也很可能在不久的将来被访问)。
- "Hello world" 示例分析 (PPT L12, S4):
 - "Bytes stored in data segment, spatial locality." (存储在数据段的字节 (如字符串 "Hello, world.\$") **具有空间局部性, 因为它们是连续存储的**。
 - "Code in loops temporal (and spatial) locality." (循环中的代码具有**时间局部性 (因为会重复执行)** 和**空间局部性** (因为循环体内的指令通常是连续的))。
 - **图片解读:** 展示了一段简单的汇编代码 (打印 "Hello, world."), 并指出了其中数据存储和循环代码体现的局部性原理。

4.3 缓存 (Cache) 的概念与工作原理 (PPT L12, S5 & S7)

- 基本思想:

"The code and data currently being used by the uP are copied from D-RAM (...) to S-RAM (...)" (微处理器当前正在使用的代码和数据从DRAM复制到SRAM (缓存) 中)。

 - 示例 (P100): DRAM (128Mbyte, 180nS) -> SRAM (256Kbyte, 45nS)。
 - 示例 (Intel i7): DRAM (8Gbyte, 180nS) -> SRAM (8Mbyte, 45nS)。
- 工作流程:
 1. "When the uP sends out an address to the cache controller it checks to see if the data is in cache or main memory." (当微处理器发出地址给缓存控制器时, 控制器检查数据是在缓存中还是在主内存中)。
 2. **缓存命中 (Hit):** "If it is in cache the data is sent and there are no waits to the uP." (如果数据在缓存中, 则直接发送给微处理器, 没有等待状态)。
 3. **缓存未命中 (Miss):** "If it is not in the cache the data is read from main memory and sent to the uP (slowly) a copy is kept on file in the cache for future use." (如果数据不在缓存中, 则从主内存读取 (速度较慢) 并发送给微处理器, 同时在缓存中保留一份副本以备将来使用)。
- **命中率 (Hit rate):** "The percentage of S-RAM to (DRAM+SRAM) reads is known as the hit rate. The hit rate is a figure of merit for the effectiveness of the cache. Hit rates on a PC are typically 90%." (SRAM读取次数占总读取次数的百分比称为命中率, 是衡量缓存有效性的指标, PC上的命中率通常为90%)。

4.4 典型缓存参数 (Typical Cache Parameters - PPT L12, S6 & S8)

- **块 (行) 大小 (Block (line) size):** "4-128 bytes"
- **命中时间 (Hit Time):** "1-4 Clock cycles (1 is best)"
- **未命中惩罚 (Miss Penalty):** "8-32 clock cycles" (从未命中到数据从主存取回所需额外时间)
- **未命中率 (Miss rate):** "1%-20%"
- **缓存大小 (Cache size):** "32KB-64MB" (PPT L12, S6) 或 "1KB-1MB" (PPT L12, S8, 可能是早期数据或L1缓存数据)。

4.5 典型PC RAM系统 (Typical PC RAM system, x86 - PPT L12, S9)

- **组件:** 80386 uP (微处理器), Cache Control (缓存控制器), Cache (32K SRAM), DRAM Control (DRAM控制器), Main memory (DRAM 1Mx32)。
- **总线:** 32 Data Bus (32位数据总线), 32 Address Bus (32位地址总线)。
- **信号:**
RAS, CAS (用于DRAM控制)。
 - **图片解读:** 该图展示了一个基于80386处理器的典型PC内存系统框图。CPU通过地址总线 and 数据总线与缓存控制器和DRAM控制器交互。缓存控制器管理SRAM缓存, DRAM控制器管理主内存DRAM。

4.6 缓存组织方式 (Organisation of the cache - PPT L12, S10)

- **缓存目录 (Cache directory):** "The cache controller contains a cache directory." (缓存控制器包含一个缓存目录, 用于存储缓存中数据块与主存地址的映射关系)。
- **直接映射缓存 (Direct Mapped Cache):**
 - "The cache is organised to contain a copy of a single page of memory." (缓存被组织成包含内存单个页面的副本)。
 - "A particular line from main memory will always map to the same line in cache." (主内存中的特定行总是映射到缓存中的同一行)。
 - **抖动 (Thrashing):** "If processor switches rapidly between pages the effect is known as thrashing this will slow the system down." (如果处理器在不同页面 (但映射到缓存同一位置) 之间快速切换, 会导致缓存内容频繁替换, 命中率下降, 称为抖动)。
- **二路组相联缓存 (Two-Way Set Associative):**
 - "Uses two separate caches to allow fast switching between pages and reduce chance of thrashing." (使用两个独立的缓存 (或者说一个缓存被分成两路), 允许在页面间快速切换, 减少抖动几率)。
- **全相联缓存 (Fully Associative):**
 - "Each d-ram address and associated data is stored in cache memory as it is used." (每个DRAM地址及其关联数据在使用时都可以存储在缓存中的任何位置)。
 - "Could be slow since you need to look through the directory." (查找可能较慢, 因为需要搜索整个目录)。
 - "How do you discard old data in cache?" (当缓存满时如何替换旧数据是个问题)。

4.7 直接映射缓存详解 (Direct Mapped Cache - PPT L12, S11-S13)

- **示例系统 (80386):** 32位地址总线 (4GB可寻址空间), 32位数据总线 (每4字节为一行 line)。
- **缓存参数示例:** 8192行, 总共32KB。行被组织成块 (blocks), 每块8行, 即32字节/块。 (PPT L12, S11中 "256 bytes" 应为笔误, 32bytes/block * 1024 blocks = 32KB)。
- **页面划分:** "The cache controller treats the 4GByte memory as $2^{32}/32768=131072$ pages of 32Kbytes each." (缓存控制器将4GB内存视为131072个32KB的页面)。
- **映射关系:** "A particular line from main memory will always map to the same line in cache."
- **缓存目录 (Tag):** "The cache directory will store the page location (main memory) for each line in cache. This information is known as the tag." (缓存目录存储缓存中每一行对应的主存页面位置, 即标签)。

- 地址划分 (Accessing memory - PPT L12, S12):
 - 微处理器输出32位地址。
 - "A15-A32: Set the page in main memory" (高17位用于与缓存目录中的标签比较, 确定页号)。
 - "A5-A14: Set the block number in cache" (中间10位用于选择缓存中的块号/组号)。
 - "A2-A4: Set the line number in the block" (接着3位用于选择块内的行号)。
 - "A0-A1: Specific bytes in the line" (最低2位用于选择行内的字节)。
- 图片解读 (Direct Mapped Cache - PPT L12, S13):
 - 左侧是主内存, 被划分为多个页面(Page 0, Page 1, ...)。每个页面又包含多个块(Blk0, ...)。
 - 中间是缓存, 其大小等于主内存的一个页面大小 (32KB)。缓存也被划分为与页面内块对应的块。
 - 右侧是缓存目录, 存储每个缓存块对应的源自主内存的页号 (Tags) 和有效位 (validity bit)。
 - 当访问内存时, 主存地址的一部分用于索引缓存行, 另一部分(高位)与该缓存行目录中的Tag比较。如果匹配且有效位为1, 则命中。
- 抖动示例 (Thrashing - PPT L12, S14):
 - "If a program accesses page 0 block X and then page 1 block X as part of a repeating loop. Each memory access will result in the DRAM access and an update of the cache." (如果程序在循环中交替访问映射到缓存同一位置的不同主存页面的块X, 例如地址0和地址32768, 就会导致缓存不断被替换, 发生抖动)。
 - **图片解读:** 图示了主存中的Page 0的Blk0和Page 1的Blk0都映射到Cache中的同一个Blk0位置。如果程序交替访问这两个主存块, Cache中的内容会不断被覆盖, 导致命中率降低。

4.8 组相联缓存 (Set Associative - PPT L12, S15-S16)

- **解决抖动:** "On a PC the thrashing problem can be solved by splitting the cache into two." (通过将缓存分成两路或多路来解决抖动问题)。
- 二路组相联 (Two-Way Set Associative):
 - "Each cache has its own directory." (或者说, 每个“组”对应多个可以存放数据块的“路”, 每路有自己的Tag)。
 - **LRU (Least Recently Used) 位:** "An additional bit in the tag directory identifies which cache contains the most up to date version in the case both caches contain the same page. (LRU least recently used bit)." (PPT这里的描述可能不完全准确, LRU用于替换策略, 当一个组满了需要替换时, 选择最久未使用的那一路。如果两路都缓存了来自不同页但映射到同一组的数据块, 它们可以共存)。
 - 图片解读 (Two-Way Set Associative - PPT L12, S16):
 - 缓存被逻辑上分成了Cache A和Cache B两路 (或者一个缓存有多个“槽位”属于一个“组”)。主存的一个块可以映射到这个“组”中的任何一路。
 - 缓存目录需要存储每路的Tag、有效位, 以及用于替换策略的LRU位。
 - "Sometimes Tags are stored in special TAG ram" (标签有时存储在专门的TAG RAM中)。

4.9 全相联缓存 (Fully Associative - PPT L12, S17)

- **映射方式:** "A line (4bytes) can be written to any location in cache." (主存中的一个行可以被加载到缓存中的任何位置)。
- **Tag大小:** "Each tag in the cache directory must be 30 bits." (因为地址中用于索引的部分减少了, 所以Tag需要包含更多的地址信息)。
- **优点:** "This cache can hold the same line from many different pages of memory." (可以同时缓存来自不同主存页面的、地址不冲突的行, 抖动最小)。
- **缺点 (Problems):**
 1. "you have to search through the cache directory to find your data." (查找目录的开销大, 因为需要并行比较所有Tag)。
 2. "How do you discard old stuff when the cache is full?" (替换策略复杂)。
 - "Random discard : Uniform" (随机替换)
 - "Least Recently Used: Tricky to keep count of number of reads" (LRU实现复杂)
- **应用:** "Not used on PC's (1995)." (在1995年的PC上不常用, 因为硬件成本高)。

4.10 缓存未命中的原因 (What causes a cache to miss? - PPT L12, S18)

- **强制性未命中 (Compulsory / Cold start miss):** "First read of each block must be from D-RAM, cold start miss." (对每个数据块的第一次读取必然是未命中, 因为此时缓存是空的或不包含该数据)。
- **容量性未命中 (Capacity miss):** "If the number of blocks in regular use exceeds the cache size, then thrashing can occur." (如果程序活跃使用的数据块总大小超过了缓存的容量, 即使是全相联缓存也会因为空间不足而发生未命中和抖动)。
- **(补充)冲突性未命中 (Conflict miss):** 在直接映射或组相联缓存中, 即使缓存总容量足够, 如果多个活跃的数据块映射到同一个缓存位置 (直接映射) 或同一个组 (组相联, 且组已满), 也会导致未命中。

4.11 缓存的效果 - 阿姆达尔定律 (What effect does a cache have? - PPT L12, S19)

- **阿姆达尔定律 (Amdahl's Law):** "The improvement to be gained from using some faster mode of execution is limited by the fraction of the time the faster mode can be used." (通过使用某种更快的执行方式所能获得的改进, 受限于该更快方式能被使用的时间比例)。
- **公式:**
$$\text{Speedup} = 1 / ((1 - p) + p/s)$$
 - **s** is the speedup due to the improved resources (改进资源的加速比, 例如缓存比主存快10倍, $s=10$)。
 - **p** is the proportion of execution time that the part benefiting from improved resources originally occupied (受益于改进资源的部分最初占用的执行时间比例, 例如命中率为95%, $p=0.95$)。

- **示例计算：** "The addition of a cache improves the speed of memory by 10 (from 100ns to 10n) so $s=10$, a typical program has a hit rate of 95% so $p=0.95$. The addition of the cache makes memory x6.89 faster than DRAM" ($1 / ((1-0.95) + 0.95/10) = 1 / (0.05 + 0.095) = 1 / 0.145 \approx 6.89$)。

4.12 估算有效访问时间 (Estimating Effective Access Time - PPT L12, S20)

- 参数：
 - 缓存访问时间 (cache access time) = 10 ns
 - 未命中率 (miss rate) = 5% (即 $p_{\text{miss}} = 0.05$, 所以命中率 $p_{\text{hit}} = 1 - p_{\text{miss}} = 0.95$)
 - DRAM访问时间 (dynamic ram speed) = 100 ns
- 公式：

EAT (Effective Access Time) = (Hit Rate * Cache Access Time) + (Miss Rate * Main Memory Access Time)

 - PPT中的公式: "EAT = $(1 - p) \times 10 + p(100)$ " (这里的p是未命中率)
 - 按PPT公式计算: $EAT = (1 - 0.05) \times 10 + 0.05 \times 100 = 0.95 \times 10 + 0.05 \times 100 = 9.5 + 5 = 14.5$ ns
- **加速比：** " $100\text{ns}/14.5\text{ns}=6.89$ Speedup (Same as Amdahl's Law slide earlier)" (与阿姆达尔定律计算结果一致)。

4.13 写策略 (Write Policies - PPT L12, S21)

- **一致性问题：** "D-RAM and Cache must always be the same." (DRAM和缓存中的数据必须保持一致)。
- **DMA问题：** "The DMA could cause problems with cache." (DMA直接内存访问可能导致缓存与主存数据不一致，因为DMA可能绕过CPU直接修改主存)。
- 写通 (Write through):
 - "Every write to cache is also made to main memory." (每次对缓存的写操作也同时写入主内存)。
 - "Write is slowed down to D-RAM speed." (写操作速度受限于DRAM速度)。
 - "Reading is at cache speed." (读操作速度是缓存速度)。
- 写回 (Write back):
 - "Data is written to the cache only." (数据仅写入缓存)。
 - "When the block is replaced the old block is transferred to main memory." (当缓存中的某个块被替换时，如果该块被修改过 (通常用一个“脏位”标记)，则旧块的内容会被写回主内存)。
 - "Uses less memory [bandwidth]." (使用较少的内存带宽，因为不是每次写都访问主存)。

4.14 多级缓存 (Multilevel caches - PPT L12, S22)

- **概念：** "On modern PC's there are two or more levels of cache." (现代PC通常有两级或更多级的缓存)。
- L1 Cache (Primary / Level 1 cache):

- "is built onto the uP chip, typically 1K to 32K." (集成在微处理器芯片内部, 典型大小1KB到32KB)。
- "These caches can have zero wait states." (可以实现零等待状态访问)。
- "e.g. mov ax,[bx] is done in the minimum number of clock cycles."
- **L2 Cache (Secondary):** "2Mbyte, 500MHz" (示例大小和速度)。
- **主内存:** "Main memory up to 64Gbytes" (示例大小)。

模块五：直接内存访问 (DMA)

5.1 DMA 概念 (DMA on 8086 - PPT L12, S23)

- **定义:** "Many devices are able to write directly to memory. If this can be done it frees up valuable CPU time. The technique is known as DMA or Direct Memory Access." (许多设备能够直接写入内存, 这可以释放宝贵的CPU时间, 该技术称为DMA)。
- **应用:** "DMA is used on devices such as frame grabbers, and sound cards." (DMA用于帧捕获卡、声卡等设备)。

5.2 DMA 控制器与总线 (DMA Controllers - PPT L12, S24)

- **组件:** uP (微处理器), Memory (内存), Latch (锁存器, 用于地址), Disk Drive (磁盘驱动器, 作为DMA设备示例), DMA (DMA控制器)。
- **信号线:**
 - AD0-AD15 (复用的地址/数据总线)
 - ALE (地址锁存使能)
 - R/W (读写控制)
 - DMA RW (DMA读写控制)
 - DACK (DMA应答)
 - DREQ (DMA请求)
 - HLDA (总线保持应答)
 - HOLD (总线保持请求)
 - **图片解读:** 该图展示了DMA控制器如何协调外部设备 (如磁盘驱动器) 和内存之间的直接数据传输。DMA控制器在获得总线控制权后, 可以直接控制地址和数据总线, 使数据在设备和内存间传输, 而无需CPU介入每个字节的传输。

5.3 DMA 事件序列 (DMA Sequence of events - PPT L12, S25)

1. "DMA controller connects uP to memory via a set of switches." (DMA控制器通过一组“开关”连接微处理器和内存, 这里的开关指总线控制权的切换)。
2. "To read the disk a series of commands are sent to the smart disk controller to obtain the data to be placed in main memory." (为读取磁盘, 向智能磁盘控制器发送一系列命令以获取要放入主内存的数据)。
3. "When the device is ready DREQ is set high." (当设备准备好数据时, DREQ (DMA请求)信号置高)。
4. "The HOLD (hold request) is then sent to the uP." (HOLD (总线保持请求)信号发送给微处理器, 请求总线控制权)。

5. "When the uP is ready HLDA (hold acknowledge) is brought high." (当微处理器准备好释放总线时, HLDA (总线保持应答)信号置高)。
6. "The DMA controller switches the data and address line over to the smart drive." (DMA控制器将数据和地址总线的控制权切换给智能驱动器 (或DMA控制器自身控制总线进行传输))。
7. "DMA sends a DACK0 (DMA acknowledge) signal." (DMA控制器发送DACK0 (DMA应答)信号给设备, 通知其可以开始传输)。
8. "The device writes to the memory." (设备将数据写入内存)。
9. "HOLD is then released and the uP regains access to the memory." (HOLD信号被释放, 微处理器重新获得对内存的访问权)。

模块六：只读存储器 (ROM)

6.1 ROM 概述 (ROMS - PPT L12, S26)

- **特性:** "ROMS are non-volatile memories, normally they contain BIOS and Boot Sequence Code." (ROM是非易失性存储器, 通常包含BIOS和引导序列代码)。
- **类型:**
 - **Mask Programmed ROM (掩模型ROM):** "Programmed during manufacture, cannot be altered." (制造过程中编程, 不可更改)。
 - **PROMS (可编程ROM):** "One time programmable ROMS, programmed by blowing fusible links in the device." (一次性可编程ROM, 通过烧断设备中的熔丝进行编程)。
 - **EPROMS (可擦除可编程ROM):** "Programmed electrically by leaving charge on the gate of a FET. Erased using UV light that discharges the gate." (通过在FET栅极上留下电荷进行电编程, 使用紫外线擦除以释放栅极电荷)。
 - **EEPROMS (电可擦除可编程ROM):** "Programmed and re-programmable electrically." (可电编程和电重编程)。
- **速度与应用:** "ROMS are quite slow, so sometimes the contents are copied into a SRAM after the machine has been powered up." (ROM速度较慢, 有时其内容在机器加电后会被复制到SRAM中执行, 即所谓的Shadow RAM)。

6.2 EPROM 原理 (EPROM - PPT L12, S27)

- **编程:** "A high voltage placed on the insulated gate causes electrons to tunnel leaving the plate biased positive." (施加在高压到绝缘栅极上, 导致电子隧穿效应, 使浮动栅极带正电荷, 从而改变晶体管的阈值电压, 实现存储'0'或'1')。
- **擦除:** "Shining light (UV) on the gate discharges the capacitor." (用紫外线照射栅极会使电容器 (浮动栅) 放电, 恢复到初始状态)。
- **EEPROMS:**

"Electrically erasable programmable read only memory." (电可擦除可编程只读存储器, 与EPROM不同, 它可以用电信号擦除, 无需紫外线)。

 - **图片解读 (EPROM - PPT L12, S27):**
 - 展示了一个EPROM存储单元的简化截面图。包含源极(Source)、漏极(Drain)、控制栅(Set/Reset)、浮动栅(Insulated Gate/Plate)和衬底(p)。编程时通过高压使电子隧穿到浮动栅。擦除时用紫外光照射, 使浮动栅上的电荷泄放。

6.3 2716 EPROM 管脚图 (Pin-outs of a 2716 EPROM - PPT L12, S28)

- **型号含义:** " $2^{12}=2048\text{Bytes}=16384\text{Bits}$, Hence 27-16" (2716表示容量为16K bits, 即2K Bytes。需要11根地址线 A0-A10 来寻址 $2^{11}=2048$ 个字节, PPT中标注到A11, 可能是指2732或更大容量的兼容管脚, 或者2716本身有A11但内部未使用)。
- **主要管脚:**
 - A0-A10/A11: 地址线 (Address lines)
 - D0-D7: 数据线 (Data lines)
 - OE#: 输出使能 (Output enable, 低有效)
 - CE# (或 CS#): 片选 (Chip select, 低有效)
 - Vprg: 编程电压 (Programming voltage, 编程时使用)
 - +5V: 电源
 - 0V: 地
 - **图片解读:** 展示了2716 EPROM芯片的典型DIP封装管脚排列图, 标明了各个管脚的功能。

6.4 ROM 扩展 (Combining ROMs (or other memory) - PPT L12, S29)

- **方法:**

"Using the chip select CS line and a demultiplexer it is possible to to combine many memory devices in memory space." (使用片选(CS)线和地址译码器 (如解复用器) 可以将多个存储设备组合在内存空间中, 以扩展总存储容量或位宽)。

 - **图片解读:** 图中展示了使用三片2716 EPROM芯片和一个地址译码器 (由A12, A13驱动产生Y0, Y1, Y2, Y3等片选信号) 来扩展存储空间例子。A0-A11连接到所有芯片的地址线, 数据线也并联。通过A12和A13的不同组合, 译码器会选中其中一片EPROM, 使其CS有效, 从而CPU可以访问该芯片的数据。

模块七: 虚拟内存 (Virtual Memory)

7.1 虚拟内存概念 (Virtual Memory - PPT L12, S30)

- **定义:**
 - "The address space mapped by the CPU is called the virtual memory, on a typical PC this is 64Gbytes." (CPU映射的地址空间称为虚拟内存, 例如在典型PC上可能是64GB, 这里指逻辑地址空间大小)。
 - "The main memory is the physical memory available, e.g. 4Gbytes." (主内存是实际可用的物理内存, 例如4GB)。
- **机制:** "A virtual memory system uses an external memory to store data. Pages of the external memory are switched into main memory." (虚拟内存系统使用外部存储 (通常是硬盘) 来存储数据, 外部存储的页面被换入主内存中)。
- **与缓存的相似性:** "Same methods as caching." (使用了与缓存类似的方法, 如页面替换算法)。
- **安全问题:** "Note: on older systems this could produce security problems. Virtual memory should be cleared before use." (在旧系统中可能产生安全问题, 虚拟内存使用前应被清除)。

7.2 虚拟内存与性能 (Virtual Memory - PPT L12, S31)

- **硬盘作扩展内存:** "If there is insufficient Main memory then the Hard Disk can be used as a form of extended memory. This memory space is known a virtual memory." (如果主内存不足, 硬盘可以用作一种扩展内存, 这个内存空间称为虚拟内存)。
- **未命中惩罚:** "The miss penalty for using virtual memory can be as high as 100,000 cycles." (使用虚拟内存 (即发生页错误, 需要从硬盘调页) 的未命中惩罚可能高达100,000个时钟周期, 非常高昂)。
- **性能影响:** "If you open a number of applications in Windows you can fill memory. The result is that the operating system starts to use virtual memory. Things start to really slow down!" (如果在Windows中打开多个应用程序填满了物理内存, 操作系统开始使用虚拟内存, 系统性能会显著下降)。

希望这份提炼对您有帮助!

好的, 这是您要求的选择题、填空题及其对应的答案和解析, 以Markdown格式呈现。

选择题

【初级 - 基础回忆】

1. 根据PPT L11 S2, 内存架构的目标是实现电子计算机数据存储方式的组合, 使其达到最佳的:
 - A. 仅速度和成本
 - B. 仅可靠性和持久性
 - C. 速度、可靠性、持久性和成本
 - D. 仅容量和速度
2. CMOS反相器在什么情况下会产生主要的功耗? (PPT L11, S4)
 - A. 输入为高电平时
 - B. 输入为低电平时
 - C. 仅在开关状态切换的瞬间
 - D. 始终有较大功耗
3. 在内存层级结构中, 访问速度最快、容量最小的通常是: (PPT L11, S5)
 - A. 主内存 (Main Memory)
 - B. 缓存 (Cache)
 - C. 寄存器 (Registers)
 - D. 磁盘 (Disk)
4. 以下哪种存储器是易失性的, 但不需要刷新? (PPT L11, S7)
 - A. DRAM
 - B. SRAM
 - C. ROM
 - D. EEPROM
5. DRAM存储单元将信息存储为: (PPT L11, S9)
 - A. 触发器的状态
 - B. 磁畴的方向
 - C. 电容器上的电荷
 - D. 熔丝的通断

6. 在DRAM的操作中，RAS信号通常代表：(PPT L11, S12)
- A. 列地址选通 (Column Address Select)
 - B. 行地址选通 (Row Address Select)
 - C. 读访问信号 (Read Access Signal)
 - D. 刷新激活信号 (Refresh Active Signal)
7. “时间局部性”原理指的是：(PPT L12, S3)
- A. 如果一个数据项被访问，它附近的数据项也很可能被访问
 - B. 如果一个数据项被访问，它很可能在不久的将来再次被访问
 - C. 数据在内存中存储的时间越长，访问越快
 - D. CPU访问内存特定区域的时间是固定的
8. 缓存的“命中率”是指：(PPT L12, S5)
- A. CPU在缓存中找到所需数据的次数
 - B. CPU在主存中找到所需数据的次数
 - C. SRAM读取次数占总 (SRAM+DRAM) 读取次数的百分比
 - D. 缓存中数据被更新的频率
9. 在直接映射缓存中，如果处理器在映射到缓存同一位置的不同主存页面之间快速切换，可能会发生什么现象？(PPT L12, S10)
- A. 缓存命中率显著提高
 - B. 抖动 (Thrashing)
 - C. 缓存一致性问题
 - D. 写冲突
10. 阿姆达尔定律主要用于评估什么？(PPT L12, S19)
- A. 存储器的物理尺寸
 - B. 使用某种更快执行方式所能获得的系统性能改进上限
 - C. CPU的时钟频率
 - D. 数据总线的宽度
11. DMA的主要目的是：(PPT L12, S23)
- A. 提高CPU的运算速度
 - B. 允许外设直接访问内存，释放CPU资源
 - C. 增加主内存的容量
 - D. 管理多级缓存
12. 以下哪种ROM类型可以通过紫外线擦除数据？(PPT L12, S26)
- A. Mask ROM
 - B. PROM
 - C. EPROM
 - D. EEPROM

【中级 - 理解与比较】

1. 对比SRAM和DRAM，以下说法正确的是：(PPT L11, S7, S11)
- A. SRAM单元结构比DRAM简单，集成度更高
 - B. DRAM速度比SRAM快，常用于缓存
 - C. SRAM不需要刷新，DRAM需要定期刷新
 - D. SRAM成本比DRAM低，容量更大
2. DRAM的访问时间 (Access Time) 和周期时间 (Cycle Time / TRD) 的关系通常是：(PPT L11, S14)
- A. 访问时间远大于周期时间
 - B. 访问时间约等于周期时间
 - C. 周期时间大约是访问时间的两倍，因为需要预充电
 - D. 两者没有直接关系

3. 缓存的“写回 (Write back)”策略与“写通 (Write through)”策略相比, 其主要优点是: (PPT L12, S21)
- A. 实现更简单
 - B. 读取速度更快
 - C. 能更好地保证缓存与主存的一致性
 - D. 减少了对主内存总线的访问次数, 可能提高写性能
4. 在二路组相联缓存 (Two-Way Set Associative Cache) 中, 当一个“组”已满且需要替换数据块时, LRU算法会选择替换哪个数据块? (PPT L12, S15)
- A. 随机选择一个
 - B. 最近刚被访问过的数据块
 - C. 最久未被访问过的数据块
 - D. 地址最小的数据块
5. 虚拟内存系统中, 当发生“页错误 (Page Fault)”时, 其“未命中惩罚”与缓存未命中相比通常是: (PPT L12, S31)
- A. 更小
 - B. 相当
 - C. 大得多
 - D. 不确定, 取决于硬盘速度
6. 在DMA传输过程中, 当外部设备准备好数据传输时, 它会向DMA控制器发送什么信号? (PPT L12, S25)
- A. HOLD
 - B. HLDA
 - C. DREQ
 - D. DACK

【高级 - 分析与应用】

1. 假设一个直接映射缓存大小为32KB, 块大小为32字节。主存地址为32位。那么用于在缓存目录中进行Tag比较的地址位数是多少? (基于PPT L12, S11-S12的逻辑推断)
- A. 10位
 - B. 12位
 - C. 17位
 - D. 15位
2. 如果一个系统使用了多级缓存 (L1, L2), 并且CPU发出的访存请求在L1缓存中未命中, 但在L2缓存中命中, 这个过程通常被称为:
- A. L1命中, L2命中
 - B. L1未命中, L2未命中
 - C. L1未命中, L2命中
 - D. 强制性未命中
3. 考虑一个具有缓存的系统, 缓存访问时间为10ns, 主存访问时间为100ns。如果缓存命中率为90%, 那么有效平均访问时间 (EAT) 是多少? (PPT L12, S20)
- A. 10ns
 - B. 19ns
 - C. 90ns
 - D. 100ns
4. 在使用地址译码器扩展ROM存储空间时, 如果使用2根地址线 (如A12, A13) 作为译码器的输入, 最多可以产生多少个独立的片选信号来选择不同的ROM芯片? (基于PPT L12, S29的逻辑)
- A. 2个
 - B. 3个

- C. 4个
- D. 8个

填空题

【初级 - 基础回忆】

1. 根据PPT L11 S3, 楔形文字泥板大约能存储 __ 字节的数据。
2. CMOS电路的主要特点是只有在 __ 时才有功耗。
3. 内存层级结构中, 位于CPU芯片内部, 按字节存取, 访问速度最快的是 __。
4. SRAM的构造基础是 __, 而DRAM的构造基础是存储在 __ 上的电荷。
5. DRAM单元由于电荷会泄漏, 所以需要周期性的 __ 操作。
6. DRAM读操作时, 地址总线会分时复用传输 __ 地址和 __ 地址。
7. 程序访问内存时, 如果一个数据项被访问后, 其附近的数据项也很可能被访问, 这体现了 __ 局部性原理。
8. 缓存控制器中用于存储缓存块与主存地址映射关系的部分称为 __。
9. 当CPU访问的数据不在缓存中, 需要从主存读取, 这种情况称为缓存 __。
10. DMA是 __ 的缩写。
11. ROM是非易失性存储器, 通常包含计算机的 __ 和引导序列代码。
12. CPU映射的地址空间称为 __ 内存, 而实际可用的物理内存称为物理内存。

【中级 - 理解与比较】

1. 内存层级中, 从顶层到底层, 存储器的访问速度逐渐 __, 容量逐渐 __, 每字节成本逐渐 __。
2. 一个典型的SRAM存储单元比一个典型的DRAM存储单元需要 __ (更多/更少) 的晶体管。
3. DRAM的刷新操作是通过 __ 和 __ 信号线来协调完成的。
4. 缓存的“块 (line)”是缓存与主存之间数据交换的 __ 单位。
5. 直接映射缓存的主要缺点是容易发生 __ 现象, 可以通过 __ 缓存或全相联缓存来缓解。
6. 在DMA传输过程中, CPU通过向DMA控制器发出 __ 信号来放弃总线控制权, 并通过接收 __ 信号得知总线已被释放。
7. EPROM可以通过 __ 进行编程, 通过 __ 进行擦除。
8. 使用虚拟内存技术时, 如果频繁地在主存和硬盘之间交换页面, 会导致系统性能 __。

【高级 - 分析与应用】

1. 一个32位地址总线的系统, 其可寻址的虚拟内存空间最大为 __ GB。
 2. 在直接映射缓存中, 主存地址通常被划分为三部分: __、__ 和块内偏移地址。
 3. 如果缓存的命中时间为 H , 未命中惩罚为 M (指从主存读取的额外时间), 未命中率为 m , 则有效平均访问时间 (EAT) 的计算公式可以表示为 $EAT =$ __。
 4. 在多级缓存系统中, L1缓存通常集成在 __ 内部, 其访问延迟通常可以做到 __ 等待状态。
 5. 扩展ROM存储空间时, 除了使用地址线 and 数据线, 还需要利用每个ROM芯片的 __ 信号和地址 __ 器来实现片选。
-
-

选择题 - 答案与解析

【初级 - 基础回忆】

1. C. 速度、可靠性、持久性和成本
 - **解析：** PPT L11, S2明确指出内存架构的目标是这四个因素的最佳组合。
2. C. 仅在开关状态切换的瞬间
 - **解析：** PPT L11, S4指出 "Note: there is only dissipation during switching".
3. C. 寄存器 (Registers)
 - **解析：** PPT L11, S5和S6的表格显示寄存器位于层级顶部，速度最快，容量最小。
4. B. SRAM
 - **解析：** PPT L11, S7描述SRAM "Memory is volatile, but does not need refreshing". DRAM易失且需刷新，ROM和EEPROM非易失。
5. C. 电容器上的电荷
 - **解析：** PPT L11, S9明确指出DRAM "Information stored as charge on capacitor".
6. B. 行地址选通 (Row Address Select)
 - **解析：** PPT L11, S12定义 "RAS=Row Address Select".
7. B. 如果一个数据项被访问，它很可能在不久的将来再次被访问
 - **解析：** PPT L12, S3对时间局部性的定义。
8. C. SRAM读取次数占总 (SRAM+DRAM) 读取次数的百分比
 - **解析：** PPT L12, S5对命中率的定义 (虽然原文是 "The percentage of S-RAM to (DRAM+SRAM) reads", 更准确的理解是SRAM命中次数占总访问次数的百分比)。
9. B. 抖动 (Thrashing)
 - **解析：** PPT L12, S10和S14描述了这种现象。
10. B. 使用某种更快执行方式所能获得的系统性能改进上限
 - **解析：** PPT L12, S19对阿姆达尔定律的描述。
11. B. 允许外设直接访问内存，释放CPU资源
 - **解析：** PPT L12, S23的核心思想 "it frees up valuable CPU time".
12. C. EPROM
 - **解析：** PPT L12, S26描述EPROM "Erased using UV light".

【中级 - 理解与比较】

1. C. SRAM不需要刷新，DRAM需要定期刷新
 - **解析：** 这是两者关键区别之一 (PPT L11, S7)。SRAM单元复杂，集成度低；SRAM速度快，常用于缓存；DRAM成本低，容量大。
2. C. 周期时间大约是访问时间的两倍，因为需要预充电
 - **解析：** PPT L11, S14提到 "speed of ram is thus the total cycle time TRD which is about twice the access time".
3. D. 减少了对主内存总线的访问次数，可能提高写性能
 - **解析：** 写回策略只在块被替换且被修改过时才写回主存，减少了总线流量 (PPT L12, S21)。写通每次都写主存，一致性更容易保证但可能慢。

4. C. 最久未被访问过的数据块

- **解析：** LRU (Least Recently Used) 算法的含义 (PPT L12, S15)。

5. C. 大得多

- **解析：** PPT L12, S31提到虚拟内存的未命中惩罚 "can be as high as 100,000 cycles", 远大于缓存未命中的几十个周期。

6. C. DREQ

- **解析：** PPT L12, S25描述 "When the device is ready DREQ is set high."。

【高级 - 分析与应用】

1. C. 17位

- **解析：**
 - 缓存大小 = 32KB = 2^{15} 字节。
 - 块大小 = 32 字节 = 2^5 字节。
 - 缓存行数 = 缓存大小 / 块大小 = $2^{15} / 2^5 = 2^{10}$ 行。
 - 因此, 需要10位作为索引 (Index) 来选择缓存中的行。
 - 块内偏移需要5位 (因为块大小是 2^5 字节)。
 - Tag位数 = 主存地址总位数 - 索引位数 - 块内偏移位数 = $32 - 10 - 5 = 17$ 位。
 - PPT L12 S12的例子是32KB缓存, 32位地址, A15-A32 (高18位) 用于Tag, 但它的块大小和行组织与此题假设不同。此题按标准计算方法。如果严格按PPT L12, S12的地址划分 (A5-A14为块号, A2-A4为行号, A0-A1为字节号), 那么块号10位, 行号3位, 字节号2位, 总共15位用于索引和偏移, 则Tag = $32 - 15 = 17$ 位。

2. C. L1未命中, L2命中

- **解析：** 这是多级缓存工作的典型场景。

3. B. 19ns

- **解析：** $EAT = (\text{命中率} \times \text{缓存访问时间}) + (\text{未命中率} \times \text{主存访问时间})$
 $EAT = (0.90 \times 10\text{ns}) + (0.10 \times 100\text{ns}) = 9\text{ns} + 10\text{ns} = 19\text{ns}$ 。(PPT L12, S20的公式是 $(1 - p_{\text{miss}})T_{\text{cache}} + p_{\text{miss}}T_{\text{main}}$, 这里的 p_{miss} 是0.10)

4. C. 4个

- **解析：** 2根地址线可以产生 $2^2 = 4$ 种不同的组合, 每种组合可以对应一个片选信号。

填空题 - 答案与解析

【初级 - 基础回忆】

1. 根据PPT L11 S3, 楔形文字泥板大约能存储 **100's (数百)** 字节的数据。
2. CMOS电路的主要特点是只有在 **开关状态切换的瞬间 (switching)** 时才有功耗。
3. 内存层级结构中, 位于CPU芯片内部, 按字节存取, 访问速度最快的是 **寄存器 (Registers)**。
4. SRAM的构造基础是 **触发器 (Flip Flops)**, 而DRAM的构造基础是存储在 **电容器 (capacitor)** 上的电荷。
5. DRAM单元由于电荷会泄漏, 所以需要周期性的 **刷新 (refreshing)** 操作。
6. DRAM读操作时, 地址总线会分时复用传输 **行 (Row)** 地址和 **列 (Column)** 地址。

7. 程序访问内存时，如果一个数据项被访问后，其附近的数据项也很可能被访问，这体现了 **空间 (Spatial)** 局部性原理。
8. 缓存控制器中用于存储缓存块与主存地址映射关系的部分称为 **缓存目录 (Cache directory / Tag store)**。
9. 当CPU访问的数据不在缓存中，需要从主存读取，这种情况称为缓存 **未命中 (Miss)**。
10. DMA是 **直接内存访问 (Direct Memory Access)** 的缩写。
11. ROM是非易失性存储器，通常包含计算机的 **BIOS (基本输入输出系统)** 和引导序列代码。
12. CPU映射的地址空间称为 **虚拟 (virtual)** 内存，而实际可用的物理内存称为物理内存。

【中级 - 理解与比较】

1. 内存层级中，从顶层到底层，存储器的访问速度逐渐 **减慢 (slower)**，容量逐渐 **增大 (larger)**，每字节成本逐渐 **降低 (lower)**。
2. 一个典型的SRAM存储单元比一个典型的DRAM存储单元需要 **更多 (more)** 的晶体管。(SRAM通常6T，DRAM通常1T1C)
3. DRAM的刷新操作是通过 **RAS (行地址选通)** 和 **CAS (列地址选通)** 信号线来协调完成的。(更准确地说是通过RAS-only refresh或CAS-before-RAS refresh等周期)
4. 缓存的“块 (line)”是缓存与主存之间数据交换的 **基本 (fundamental / smallest)** 单位。
5. 直接映射缓存的主要缺点是容易发生 **抖动 (Thrashing)** 现象，可以通过 **组相联 (Set Associative)** 缓存或全相联缓存来缓解。
6. 在DMA传输过程中，CPU通过向DMA控制器发出 **HOLD (总线保持请求)** 信号来放弃总线控制权，并通过接收 **HLDA (总线保持应答)** 信号得知总线已被释放。
7. EPROM可以通过 **电信号 (electrically)** 进行编程，通过 **紫外线 (UV light)** 进行擦除。
8. 使用虚拟内存技术时，如果频繁地在主存和硬盘之间交换页面，会导致系统性能 **显著下降 (really slow down / degrade significantly)**。

【高级 - 分析与应用】

1. 一个32位地址总线的系统，其可寻址的虚拟内存空间最大为 **4 GB**。 $(2^{32} \text{ bytes} = 4 \text{ GB})$
2. 在直接映射缓存中，主存地址通常被划分为三部分：**标签 (Tag)**、**索引 (Index / Line number / Set number)** 和块内偏移地址。
3. 如果缓存的命中时间为 H ，未命中惩罚为 M (指从主存读取的额外时间，即 $T_{\text{main}} - T_{\text{cache}}$)，未命中率为 m ，则有效平均访问时间 (EAT) 的计算公式可以表示为 $EAT = H + m \times M$ 或者 $H \times (1-m) + (H+M) \times m$ 或者 $(1-m) \times H + m \times T_{\text{main}}$ (其中 $T_{\text{main}} = H + M$)。PPT L12 S20的公式是 $(1 - p_{\text{miss}}) \times T_{\text{cache}} + p_{\text{miss}} \times T_{\text{main}}$ ，所以是 $(1-m) \times H + m \times (\text{主存访问时间})$ 。
4. 在多级缓存系统中，L1缓存通常集成在 **CPU芯片 (uP chip)** 内部，其访问延迟通常可以做到 **零 (zero)** 等待状态。
5. 扩展ROM存储空间时，除了使用地址线 and 数据线，还需要利用每个ROM芯片的 **片选 (Chip Select / CS)** 信号和地址 **译码 (decoder / demultiplexer)** 器来实现片选。

Thinking...

Deconstructing User Needs

I've started by analyzing the user's request for a detailed explanation of the IVT, acknowledging the 'abstract' nature of the provided information. My focus is on going beyond a simple recall and building a step-by-step, analogy-driven breakdown to foster deep understanding, rather than superficial knowledge. This will inform the response's structure.

Refining Structure & Analogy

I've reviewed the provided text and am now focusing on structuring the explanation. I'll start with the 'why' and build up using clear definitions and analogies. The phone book analogy seems promising. I'll incorporate textual visual aids, followed by detailed examples.

Perfecting Example Breakdown

I've now detailed the core concepts, constructing the analogies and outlining the structure: Why, What, Where, How. I'm generating concrete examples, including their calculation, with visual aids. This emphasizes the non-abstract nature by using simple examples to drive home user understanding.

好的，这是您要求的解答题题目、凝练的标准答案以及解析，以Markdown格式呈现，并通过二级标题进行分离。

解答题题目

【初级 - 概念与原理阐述】

1. 简述内存层级结构 (Memory Hierarchy) 的设计目标以及至少三个典型的层级名称，并说明它们在速度和容量上的相对关系。(参考模块二：内存层级)
2. 什么是SRAM和DRAM？请分别说明它们存储一位信息的基本原理、是否需要刷新以及一个主要应用场景。(参考模块三：主内存技术 - SRAM 与 DRAM)
3. 解释什么是“时间局部性 (Temporal locality)”和“空间局部性 (Spatial locality)”原理。这两个原理对于缓存设计有什么重要意义？(参考模块四：内存管理与缓存策略)
4. 简述DMA (直接内存访问) 的工作原理及其主要优点。(参考模块五：直接内存访问)
5. 什么是ROM？列举至少三种不同类型的ROM及其主要特点。(参考模块六：只读存储器)

【中级 - 理解与分析】

1. 请解释DRAM的读操作时序，特别是RAS和CAS信号的作用以及地址是如何分两次提供的。并说明为什么DRAM的周期时间通常大于其访问时间。(参考模块三：主内存技术 - SRAM 与 DRAM)
2. 什么是缓存命中 (Cache Hit) 和缓存未命中 (Cache Miss)？当发生缓存未命中时，系统会执行哪些主要操作？命中率是如何定义的？(参考模块四：内存管理与缓存策略)
3. 比较直接映射缓存 (Direct Mapped Cache)、组相联缓存 (Set-Associative Cache) 和全相联缓存 (Fully Associative Cache) 的主要特点、优缺点以及它们在解决“抖动 (Thrashing)”问题上的不同表现。(参考模块四：内存管理与缓存策略)
4. 解释缓存的“写通 (Write-Through)”和“写回 (Write-Back)”两种写策略的工作机制，并分析它们各自的优缺点，特别是在总线流量和数据一致性方面。(参考模块四：内存管理与缓存策略)
5. 什么是虚拟内存？它试图解决什么问题？简述其基本工作机制，并解释为什么当系统频繁使用虚拟内存（即硬盘作为扩展内存）时，性能会显著下降。(参考模块七：虚拟内存)

【高级 - 设计与计算/综合应用】

1. 假设一个系统的主存地址为32位，缓存大小为64KB，缓存块大小为16字节。

- a) 如果采用直接映射方式, 请计算主存地址中用于Tag、索引 (Index) 和块内偏移 (Offset) 的位数分别是多少?
 - b) 如果一个程序顺序访问从地址 `0x00001000` 开始的连续256字节数据, 对于上述直接映射缓存, 理想情况下会发生多少次强制性未命中 (Compulsory Miss)? 假设缓存初始为空。
2. 一个计算机系统, 其CPU访问主内存 (DRAM) 的时间是80ns。为了提高性能, 引入了一个缓存 (SRAM), 其访问时间是10ns。
- a) 如果通过测试得知, 该缓存在运行某个典型程序时的命中率为95%, 请计算该系统的有效平均访问时间 (EAT)。
 - b) 根据阿姆达尔定律, 如果将缓存视为“改进部分”, 其相对于主存的加速比是多少? 在95%的命中率下, 整个内存访问性能相对于没有缓存时提升了多少倍? (计算Speedup)
3. 请设计一个简化的DMA传输过程的事件序列 (至少包括CPU、DMA控制器、外设和内存之间的关键交互步骤), 描述从外设请求DMA传输数据到内存, 直到传输完成, CPU收回总线控制权的整个流程。
4. 如果要使用多片2KB大小的EPROM芯片 (例如2716, 其地址线为A0-A10) 来构建一个8KB的ROM存储区域, 请简述如何通过地址译码和片选信号来实现。你需要多少片EPROM? 需要使用CPU的哪些高位地址线进行译码?
-
-

解答题答案与解析

【初级 - 概念与原理阐述】

1. 简述内存层级结构 (Memory Hierarchy) 的设计目标以及至少三个典型的层级名称, 并说明它们在速度和容量上的相对关系。
- 标准答案:
 - **设计目标:** 内存层级结构的设计目标是在**成本可接受**的前提下, 提供一个尽可能**接近最快存储器速度**的、具有**大容量存储空间**的存储系统。它通过组合不同速度、容量和成本的存储技术来实现这一目标, 利用程序访问的局部性原理来提高平均访问速度。
 - **典型层级名称及关系 (至少三个, 从快/小到慢/大):**
 1. **寄存器 (Registers):** 位于CPU内部, 速度最快, 容量最小 (字节级)。
 2. **缓存 (Cache):** 位于CPU与主存之间 (或CPU内部), 速度较快, 容量较小 (KB到MB级)。
 3. **主内存 (Main Memory / DRAM):** 计算机的主要工作内存, 速度中等, 容量较大 (GB级)。
(可选更多层级: 磁盘缓存、磁盘、磁带/光盘)
 - **相对关系:** 越靠近CPU的层级 (如寄存器、缓存), 其访问速度越快, 但单位存储成本越高, 因此容量也越小。越远离CPU的层级 (如主内存、磁盘), 其访问速度越慢, 单位存储成本越低, 因此容量可以做得很大。
 - **解析:**

内存层级结构是计算机体系结构中的核心概念。没有任何一种单一的存储技术能够同时满足高速、大容量和低成本的要求。因此, 设计者采用分层的方式, 将少量昂贵的高速存储器 (如SRAM用于缓存) 放置在靠近CPU的地方, 用于存放当前最频繁访问的数据和指令, 而将大量廉价的、速度相对较慢但容量大的存储器 (如DRAM作为主存, 磁盘作为辅存) 放置在较远的层级。数据根据访问的局部性原理在这些层级之间动态调度, 使得CPU在大多数情况下能够从高速层级获取数据, 从而提高整体性能。

2. 什么是SRAM和DRAM? 请分别说明它们存储一位信息的基本原理、是否需要刷新以及一个主要应用场景。

○ 标准答案:

- SRAM (Static Random Access Memory - 静态随机存取存储器):
 - **基本原理:** 基于**触发器 (Flip-Flops)** 存储一位信息。一个触发器由多个晶体管 (通常是6个, 形成交叉耦合的反相器) 构成, 只要供电, 触发器就能保持其状态 (0或1)。
 - **是否需要刷新:** **不需要刷新。**
 - **主要应用场景:** CPU的**高速缓存 (Cache)**, 寄存器文件等对速度要求极高的场合。
- DRAM (Dynamic Random Access Memory - 动态随机存取存储器):
 - **基本原理:** 基于存储在微型**电容器 (Capacitor)** 上的**电荷**来存储一位信息。电容器充电表示1, 放电表示0。一个晶体管用作访问开关。
 - **是否需要刷新:** **需要周期性刷新。**因为电容器上的电荷会随时间逐渐泄漏, 如果不刷新, 存储的信息就会丢失。
 - **主要应用场景:** 计算机的**主内存 (Main Memory)**。

○ 解析:

SRAM和DRAM是构成现代计算机内存系统的两种主要半导体存储技术。

SRAM的“静态”指的是只要电源接通, 它就能保持存储的数据, 无需刷新操作。这是因为它使用触发器锁存数据, 结构相对复杂, 每个存储单元占用更多晶体管和芯片面积, 因此成本较高, 集成度较低, 但优点是速度非常快。

DRAM的“动态”指的是它存储的数据需要不断地被“刷新” (重新读取并写回)。这是因为它使用电容器存储电荷, 而电荷会自然泄漏。DRAM单元结构简单 (通常一个晶体管一个电容), 集成度高, 成本较低, 容量可以做得很大, 但速度相对SRAM慢, 并且需要额外的刷新电路和刷新周期, 这会占用一部分内存带宽。

3. 解释什么是“时间局部性 (Temporal locality)”和“空间局部性 (Spatial locality)”原理。这两个原理对于缓存设计有什么重要意义?

○ 标准答案:

- **时间局部性 (Temporal locality):** 指的是如果一个内存位置 (数据或指令) 当前被访问, 那么它在不久的将来很可能被再次访问。例如, 循环中的指令和变量。
- **空间局部性 (Spatial locality):** 指的是如果一个内存位置被访问, 那么与它地址相邻的内存位置 (它附近的数据或指令) 在不久的将来也很可能被访问。例如, 顺序执行的指令、数组元素、数据结构中的相邻字段。
- 对缓存设计的重要意义:
 1. **缓存存在的理论基础:** 局部性原理是缓存之所以有效的基础。因为程序访问具有局部性, 所以可以将当前和近期访问过的数据及其附近的数据从慢速主存调入快速缓存中。
 2. 提高命中率:
 - 利用时间局部性: 当一个数据被调入缓存后, 由于它很可能被再次访问, 后续访问就可以直接从缓存中快速获取, 从而提高命中率。
 - 利用空间局部性: 当发生缓存未命中时, 不仅仅将被访问的那个数据字调入缓存, 而是将包含该字的整个数据块 (Cache Line/Block) 一起调入。因为程序很可能接下来会访问这个块中的其他数据, 这样可以预取数据, 提高后续访问的命中率。

3. **指导替换算法：** 缓存替换算法（如LRU - 最近最少使用）的设计也基于局部性原理，倾向于替换掉那些近期最不可能被再次访问的数据块。

○ **解析：**

程序在执行时，对内存的访问往往不是完全随机的，而是呈现出一定的规律性，这就是局部性原理。

时间局部性意味着“用过的东西可能马上还要用”。缓存通过保留最近访问过的数据来利用这一点。

空间局部性意味着“用过的东西旁边的好东西可能也要用”。缓存通过一次加载一个数据块（包含多个相邻字）来利用这一点。

如果程序没有局部性（例如，完全随机地访问内存），那么缓存的效率会大大降低，因为刚调入缓存的数据可能马上就不会再用了，或者旁边的数据也不会被用到。幸运的是，大多数实际程序都表现出良好的局部性。

4. **简述DMA (直接内存访问) 的工作原理及其主要优点。**

○ **标准答案：**

■ **工作原理：**

DMA是一种允许外部设备（如硬盘控制器、网卡、声卡）与系统主内存之间直接传输数据，而无需CPU介入每个字节或字传输的技术。其过程大致如下：

1. CPU初始化DMA控制器，设置传输的源地址、目标地址、数据长度以及传输方向（读/写）。
2. 当外设准备好数据传输时，向DMA控制器发出DMA请求 (DREQ)。
3. DMA控制器向CPU发出总线请求 (HOLD)。
4. CPU在完成当前总线周期后，释放总线控制权，并向DMA控制器发送总线应答 (HLDA)。
5. DMA控制器获得总线控制权后，向外设发出DMA应答 (DACK)，并直接控制地址总线、数据总线和控制总线，在内存和外设之间传输数据。
6. 数据传输完成后，DMA控制器释放总线请求 (撤销HOLD)，CPU重新获得总线控制权。

■ **主要优点：**

1. **提高数据传输效率：** 对于大量数据的传输，DMA比CPU控制的程序I/O或中断I/O方式快得多，因为它消除了CPU逐个字节搬运数据的开销。
2. **释放CPU资源：** 在DMA传输期间，CPU可以并行执行其他任务，无需等待数据传输完成，从而提高了CPU的利用率和系统整体性能。

○ **解析：**

想象一下CPU是大老板，内存是仓库，外设是供应商。

没有DMA时，供应商（外设）要把货（数据）送到仓库（内存），需要大老板（CPU）亲自一件一件地搬。大老板很忙，这种搬运工作效率低，而且占用了大老板做更重要决策的时间。

有了DMA控制器（一个能干的仓库管理员），大老板只需要告诉管理员：“把供应商的这批货（指定源、目标、数量）搬到仓库的某个位置”。然后大老板就可以去忙别的了。仓库管理员（DMA控制器）会协调供应商，直接把货搬到仓库，搬完了再通知大老板。

这样，数据传输速度快了，大老板也解放出来了。

5. **什么是ROM? 列举至少三种不同类型的ROM及其主要特点。**

○ **标准答案：**

- **ROM (Read-Only Memory - 只读存储器)：** 是一种非易失性存储器，即在断电后其存储的内容不会丢失。它的内容通常在制造时或通过特定编程过程写入，在正常工作时主要用于读取，不能轻易修改（或完全不能修改）。

- 至少三种不同类型的ROM及其主要特点:

1. Mask Programmed ROM (掩模型ROM):

- **特点:** 内容在芯片制造过程中的掩模阶段就被固定写入, 用户无法更改。成本低 (大批量生产时), 可靠性高。

2. PROM (Programmable ROM - 可编程ROM):

- **特点:** 用户可以使用专门的PROM编程器进行一次性编程 (通常通过烧断内部熔丝实现)。编程后内容不可更改。适用于小批量生产或原型开发。

3. EPROM (Erasable Programmable ROM - 可擦除可编程ROM):

- **特点:** 可以通过电信号进行编程。存储的内容可以通过强紫外线照射芯片特定窗口进行擦除, 擦除后可以重新编程。允许多次编程, 但擦除过程不方便。
(可选: EEPROM (Electrically Erasable Programmable ROM - 电可擦除可编程ROM): 可以通过电信号进行编程和擦除, 无需紫外线, 擦除更方便, 可以按字节或块进行擦除和重写。Flash Memory也是一种EEPROM的变体。)

- **解析:**

ROM的主要特性是“只读”(在正常工作状态下)和“非易失”。它非常适合存储那些不应被用户轻易修改且在系统启动时就需要存在的关键程序和数据, 如计算机的BIOS (基本输入输出系统)、固件程序、引导加载程序等。

不同类型的ROM提供了不同程度的灵活性和成本。掩模型ROM适用于最终产品的大规模生产。PROM允许用户自行编程一次。EPROM提供了可重复编程的能力, 但擦除不便。

EEPROM和Flash Memory则提供了更方便的电擦除和重写能力, 使得固件更新成为可能。

【中级 - 理解与分析】

1. 请解释DRAM的读操作时序, 特别是RAS和CAS信号的作用以及地址是如何分两次提供的。并说明为什么DRAM的周期时间通常大于其访问时间。

- **标准答案:**

- DRAM读操作时序与信号作用:

1. **地址分两次提供:** DRAM为了减少芯片引脚数量 (从而降低成本和封装尺寸), 采用了地址多路复用技术。完整的内存地址被分为**行地址**和**列地址**两部分。
2. **RAS (Row Address Strobe/Select - 行地址选通):** 当RAS信号变为有效 (通常是下降沿触发) 时, DRAM芯片锁存地址总线上的信号作为**行地址**。这个行地址用于选择DRAM内部存储阵列中的特定一行。
3. **CAS (Column Address Strobe/Select - 列地址选通):** 在RAS有效之后一段时间, CAS信号变为有效 (通常也是下降沿触发)。此时, DRAM芯片锁存地址总线上的信号作为**列地址**。这个列地址用于从已被选中的行中选择特定的存储单元 (或一组单元)。
4. **数据输出:** 在CAS信号有效, 并且读写控制信号 (R/W) 指示为读操作后, 经过一定的延迟 (访问时间), 所选单元的数据会出现在数据输出引脚上。

- 为什么DRAM周期时间大于访问时间:

- **访问时间 (Access Time, t_{AC}):** 通常指从CAS有效开始, 到数据稳定出现在数据输出引脚上所需的时间。
- **周期时间 (Cycle Time, t_{RC} 或 TRD 在PPT中指总周期):** 指的是完成一次完整的读 (或写) 操作并使DRAM准备好进行下一次操作所需的最短时间。

- DRAM的读操作是**破坏性的**，即读取单元内容时，存储在该单元电容器上的电荷会被消耗掉。因此，在数据被读出后，必须立即将数据**写回 (Restore/Precharge)** 到原来的存储单元中，以恢复其电荷状态。这个写回和预充电（为下一次访问做准备）的过程需要额外的时间。
- 因此，周期时间 = 访问时间 + 预充电/恢复时间。所以周期时间通常大于（甚至是访问时间两倍左右，如PPT L11 S14所述）访问时间。CPU不能在刚获取到数据后立即开始下一次对同一bank的访问，必须等待整个周期完成。

○ **解析：**

DRAM的地址复用是为了节省引脚。想象一个大表格（存储阵列），你需要先告诉它你要第几行（RAS + 行地址），然后再告诉它你要这一行里的第几列（CAS + 列地址）。

“破坏性读出”是DRAM的一个固有特性。读取就像是把一个小水桶里的水倒出来看看有多少，倒出来后水桶就空了，所以必须马上把同样多的水再倒回去，不然下次来看就没水了。这个“倒回去再准备好接下一次水”的过程使得完成一次操作的总时间（周期时间）比仅仅“看到水有多少”（访问时间）要长。

2. **什么是缓存命中 (Cache Hit) 和缓存未命中 (Cache Miss)? 当发生缓存未命中时，系统会执行哪些主要操作？命中率是如何定义的？**

○ **标准答案：**

- **缓存命中 (Cache Hit)：** 当CPU需要访问某个内存地址的数据时，如果该数据**已经存在于缓存中**，则称为缓存命中。CPU可以直接从高速的缓存中获取数据。
- **缓存未命中 (Cache Miss)：** 当CPU需要访问某个内存地址的数据时，如果在缓存中**找不到**该数据，则称为缓存未命中。
- **缓存未命中时的主要操作：**
 1. **从主存读取：** CPU（或缓存控制器）必须从较慢的主内存中读取包含所需数据的整个数据块 (Cache Line/Block)。
 2. **数据送往CPU：** 所需的数据（通常是块中的一个字）被发送给CPU以满足其当前请求。
 3. **数据块存入缓存：** 同时，从主存读取的整个数据块会被加载到缓存的某个适当位置，以备将来可能的访问（利用局部性原理）。这可能涉及到缓存替换算法，如果缓存已满，需要选择一个现有的块进行替换。
 4. **更新缓存目录：** 缓存控制器会更新缓存目录 (Tag存储)，记录新调入块的主存地址等信息。
- **命中率 (Hit Rate) 的定义：** 命中率是指CPU访问内存时，在缓存中成功找到所需数据的次数占总访问次数的百分比。

$$\text{命中率} = (\text{缓存命中次数} / (\text{缓存命中次数} + \text{缓存未命中次数})) \times 100\%$$

(PPT L12 S5中提到 "The percentage of S-RAM to (DRAM+SRAM) reads is known as the hit rate", 可以理解为从缓存读取的次数占总的从缓存或主存读取次数的比例。)

○ **解析：**

缓存就像CPU的书桌，主内存是图书馆。

缓存命中： CPU要找一本书（数据），发现书就在书桌上，直接拿来看，很快。

缓存未命中： CPU发现书桌上没有这本书，就得去图书馆（主内存）找。

未命中操作：

1. 图书管理员（缓存控制器）去图书馆把这本书（整个数据块）找出来。
2. 先把CPU急着要看的那一页（一个字）给CPU。
3. 同时，把这本书（整个数据块）放到书桌上（缓存），方便以后再看。如果书桌满了，可能得把一本很久没看的书放回图书馆，给新书腾地方。

4. 在书桌的索引卡（缓存目录）上记下这本书的信息。

命中率： CPU在书桌上找到书的次数占总共找书次数的比例。命中率越高，说明书桌上的书越能满足CPU的需求，CPU去图书馆的次数就越少，工作效率就越高。

3. 比较直接映射缓存 (Direct Mapped Cache)、组相联缓存 (Set-Associative Cache) 和全相联缓存 (Fully Associative Cache) 的主要特点、优缺点以及它们在解决“抖动 (Thrashing)”问题上的不同表现。

○ 标准答案：

■ 直接映射缓存 (Direct Mapped Cache)：

- **特点：** 主内存中的每一个块只能映射到缓存中的一个**固定**的行 (位置)。映射关系由主存块地址的低位 (索引位) 决定。
- **优点：** 实现简单，硬件成本低，查找速度快 (无需比较，直接通过索引定位)。
- **缺点：** 冲突概率高。如果程序频繁访问多个映射到缓存同一行的不同主存块，会导致这些块在缓存中不断被替换，发生严重的抖动。
- **抖动表现：** 对抖动非常敏感，是三种方式中最容易发生抖动的。

■ 全相联缓存 (Fully Associative Cache)：

- **特点：** 主内存中的任何一个块可以映射到缓存中的**任何一个**行 (位置)。
- **优点：** 映射最灵活，冲突概率最低，缓存空间利用率最高，最不容易发生因地址映射冲突导致的抖动。
- **缺点：** 实现非常复杂，硬件成本高。查找时需要将目标块的Tag与缓存中所有行的Tag进行并行比较，比较器电路复杂且功耗大。替换算法也相对复杂。
- **抖动表现：** 只有当缓存容量不足以容纳工作集时才会发生容量性抖动，不会因地址映射冲突导致抖动。

■ 组相联缓存 (Set-Associative Cache)：

- **特点：** 是直接映射和全相联的一种折中。缓存被划分为若干个“组 (Set)”，每个组包含若干个行 (路/Way，例如2路、4路组相联)。主内存中的一个块首先通过直接映射方式映射到缓存中的一个**特定组**，但它可以存储在该组内的**任何一个行 (路)**中。
- **优点：** 相对于直接映射，显著降低了冲突概率，减少了抖动；相对于全相联，硬件实现成本和查找复杂度较低。
- **缺点：** 冲突概率仍高于全相联，但远低于直接映射。硬件复杂度介于两者之间。
- **抖动表现：** 对抖动的抵抗能力介于直接映射和全相联之间。如果一个组内的所有路都被来自不同主存块 (但映射到该组) 的数据占满，并且这些块被频繁交替访问，仍可能发生组内抖动。路数越多，抵抗抖动的能力越强。

○ 解析：

这三种映射方式是缓存组织的核心。

- **直接映射：** 就像给每个主存块分配了一个固定的停车位 (缓存行)。如果两个主存块 (比如你家车和你邻居车) 被分配到同一个停车位，而你们俩又总想同时停车，那就会不停地挪车 (抖动)。优点是找车位快 (直接看号)，管理简单。
- **全相联：** 就像一个大停车场，任何主存块的车都可以停在任何一个空位。冲突最小，但找车时得看遍所有停车位上的车牌号 (Tag比较)，管理复杂。
- **组相联：** 是前两者的平衡。停车场被分成几个区 (组)，你的车 (主存块) 只能停在指定的那个区，但在那个区里你可以随便找个空位 (路) 停。比如2路组相联，就是每个区有两个停车位。这样比直接映射灵活，比全相联管理简单。

4. 解释缓存的“写通 (Write-Through)”和“写回 (Write-Back)”两种写策略的工作机制，并分析它们各自的优缺点，特别是在总线流量和数据一致性方面。

○ 标准答案：

■ 写通 (Write-Through):

- **工作机制：** 当CPU执行写操作时，数据**同时写入缓存和主内存**。
- **优点：**
 - **数据一致性好：** 缓存和主内存中的数据始终保持一致，简化了多处理器系统或有DMA操作时的一致性维护。
 - **实现相对简单：** 控制逻辑较简单。
- **缺点：**
 - **写操作慢：** 每次写操作都必须等待较慢的主内存写操作完成，降低了写性能。
 - **总线流量大：** 频繁的写操作会产生大量的主内存总线流量，可能成为系统瓶颈。

■ 写回 (Write-Back):

- **工作机制：** 当CPU执行写操作时，数据**仅写入缓存**，同时将该缓存块标记为“脏 (dirty)”（表示其内容已与主存不同）。只有当这个“脏”块被从缓存中替换出去时，其内容才会被写回到主内存。
- **优点：**
 - **写操作快：** 写操作通常只需访问高速缓存，速度快。
 - **总线流量小：** 只有在脏块被替换时才需要写回主存，如果一个块在被替换前被多次写入，则只有最后一次的结果需要写回，大大减少了主内存总线流量。
- **缺点：**
 - **数据一致性维护复杂：** 缓存和主内存之间可能存在数据不一致（当缓存中有脏块时）。在多处理器系统或有DMA操作时，需要额外的机制（如缓存一致性协议）来保证数据一致性。
 - **实现相对复杂：** 需要管理“脏位”，并在替换时执行写回操作。
 - **块替换时可能产生延迟：** 如果被替换的块是脏的，写回操作会增加替换的延迟。

○ 解析：

这是处理CPU写操作时如何更新缓存和主存的两种策略。

- **写通：** 就像你改了笔记（缓存），马上去图书馆把原书上的对应内容也改了。好处是你的笔记和书永远一样，别人来看书也没问题。坏处是每次改笔记都要跑一趟图书馆，慢，而且老占着去图书馆的路。
- **写回：** 你改了笔记，先在笔记上做个标记（脏位），表示这页和图书馆的书不一样了。只有当你这本笔记不用了，要换新的笔记时（块替换），你才把所有做过标记的页一次性去图书馆更新到原书上。好处是改笔记快，不用老跑图书馆。坏处是如果别人直接去看图书馆的书，可能看到的是旧内容。

5. 什么是虚拟内存？它试图解决什么问题？简述其基本工作机制，并解释为什么当系统频繁使用虚拟内存（即硬盘作为扩展内存）时，性能会显著下降。

○ 标准答案：

- **虚拟内存 (Virtual Memory):** 是一种内存管理技术, 它为进程提供了一个比物理内存实际容量大得多的逻辑地址空间 (虚拟地址空间)。操作系统负责将程序使用的虚拟地址映射到物理内存地址。
- 试图解决的问题:
 1. **物理内存容量限制:** 允许运行大于物理内存的程序, 或者同时运行多个程序而它们的总需求超过物理内存。
 2. **内存碎片:** 通过分页或分段机制, 可以更有效地管理物理内存, 减少碎片。
 3. **内存保护与共享:** 为每个进程提供独立的地址空间, 实现进程间的内存保护, 并能方便地实现内存共享。
- 基本工作机制:
 1. 程序使用的地址是虚拟地址。
 2. 操作系统将虚拟地址空间和物理内存都划分为固定大小的**页 (Page)** (或可变大小的段)。
 3. 当程序访问一个虚拟地址时, CPU中的**内存管理单元 (MMU)** 会查找**页表 (Page Table)** (或段表), 将虚拟页号转换为物理页号。
 4. 如果所需的虚拟页面当前在物理内存中 (**页命中 Page Hit**), 则MMU计算出物理地址, 访问完成。
 5. 如果所需的虚拟页面当前不在物理内存中, 而是在硬盘的交换空间 (页错误 Page Fault), 则会产生一个缺页中断。操作系统会:
 - 暂停当前进程。
 - 在硬盘上找到该页面。
 - 如果物理内存已满, 选择一个物理页面换出到硬盘 (如果该页被修改过, 则先写回)。
 - 将所需的页面从硬盘调入物理内存。
 - 更新页表。
 - 重新执行导致缺页的指令。

- **性能显著下降的原因:**

硬盘的访问速度 (毫秒级) 远慢于物理内存的访问速度 (纳秒级), 两者之间存在几个数量级的差异。当物理内存不足, 系统频繁地在硬盘和物理内存之间进行页面交换 (称为**颠簸 Thrashing** 或 **页抖动**) 时, CPU大部分时间都在等待慢速的硬盘I/O操作完成, 而不是执行有用的计算任务。这种大量的页面换入换出操作 (未命中惩罚极高) 导致系统整体性能急剧下降, 用户会感觉到系统非常卡顿。

- **解析:**

虚拟内存就像给了你一张无限大的工作台 (虚拟地址空间), 但你实际的物理工作台 (物理内存) 可能很小。

- **解决问题:** 你想同时做好多项目 (运行大程序或多程序), 物理工作台放不下, 虚拟内存让你感觉好像都能放下。
- **工作机制:** 你常用的工具和材料 (活跃页面) 放在物理工作台上。不常用的暂时存到旁边的储藏室 (硬盘交换空间)。你需要某个储藏室的东西 (页错误), 管家 (操作系统) 就去储藏室拿, 如果物理工作台满了, 还得先把一些东西从物理工作台搬回储藏室。

- **性能下降：** 如果你物理工作台太小，而你又老是要用储藏室里的东西，管家就得不间断地在工作台和储藏室之间来回搬运。去储藏室拿东西比在工作台上拿慢得多得多，所以你大部分时间都在等管家，而不是自己干活，效率自然就低了。

【高级 - 设计与计算/综合应用】

1. 假设一个系统的主存地址为32位，缓存大小为64KB，缓存块大小为16字节。

- a) 如果采用直接映射方式，请计算主存地址中用于Tag、索引 (Index) 和块内偏移 (Offset) 的位数分别是多少？
- b) 如果一个程序顺序访问从地址 `0x00001000` 开始的连续256字节数据，对于上述直接映射缓存，理想情况下会发生多少次强制性未命中 (Compulsory Miss)? 假设缓存初始为空。
- 标准答案：

- a) 地址划分：

- 1. 块内偏移 (Offset):

- 块大小 = 16 字节 = 2^4 字节。
 - 所以，块内偏移需要 **4位** (用于寻址块内的16个字节)。

- 2. 索引 (Index):

- 缓存大小 = 64 KB = $2^6 \times 2^{10}$ 字节 = 2^{16} 字节。
 - 缓存行数 (Number of lines/sets for direct mapped) = 缓存大小 / 块大小 = $2^{16} / 2^4 = 2^{12}$ 行。
 - 所以，索引需要 **12位** (用于选择缓存中的 2^{12} 行之一)。

- 3. 标签 (Tag):

- 主存地址总位数 = 32位。
 - Tag位数 = 主存地址位数 - 索引位数 - 块内偏移位数 = $32 - 12 - 4 = 16$ 位。

- **总结：Tag = 16位, Index = 12位, Offset = 4位。**

- b) 强制性未命中次数：

- 1. 程序访问数据量 = 256 字节。
 - 2. 缓存块大小 = 16 字节。
 - 3. 当发生一次缓存未命中时，会从主存加载一个16字节的数据块到缓存中。
 - 4. 由于是顺序访问，并且缓存初始为空，所以每当访问到一个新的数据块的第一个字节时，就会发生一次强制性未命中。
 - 5. 需要加载的数据块数量 = 总访问字节数 / 每个块的字节数 = $256 / 16 = 16$ 个块。
 - 6. 因此，理想情况下会发生 **16次** 强制性未命中。

- 解析：

- a):

- 块内偏移决定了CPU如何从一个已经加载到缓存的块中找到具体的字节。块大小是16字节，所以需要 $\log_2(16)=4$ 位。
 - 索引决定了主存中的一个块应该映射到缓存中的哪一行。缓存有 $64\text{KB}/16\text{B} = 4096$ (2^{12}) 行，所以需要12位索引。
 - 剩下的高位地址就是Tag，用于区分映射到同一缓存行的不同主存块。 $32 - 12 - 4 = 16$ 位。

■ b):

强制性未命中（也叫冷启动未命中）是指对一个数据块的第一次访问，此时该数据块肯定不在缓存中。由于是顺序访问，每16个字节的数据会填满一个缓存块。256字节的数据会顺序填满 $256/16 = 16$ 个不同的缓存块。因为缓存初始为空，所以对这16个块的第一次访问都会导致强制性未命中。

2. 一个计算机系统，其CPU访问主内存（DRAM）的时间是80ns。为了提高性能，引入了一个缓存（SRAM），其访问时间是10ns。

- a) 如果通过测试得知，该缓存在运行某个典型程序时的命中率为95%，请计算该系统的有效平均访问时间 (EAT)。
- b) 根据阿姆达尔定律，如果将缓存视为“改进部分”，其相对于主存的加速比是多少？在95%的命中率下，整个内存访问性能相对于没有缓存时提升了多少倍？（计算Speedup）
- 标准答案：

■ a) 有效平均访问时间 (EAT):

- 命中率 (H) = 0.95
- 未命中率 (M_rate) = $1 - H = 1 - 0.95 = 0.05$
- 缓存访问时间 (T_cache) = 10 ns
- 主内存访问时间 (T_main) = 80 ns
- $EAT = (H \times T_{\text{cache}}) + (M_{\text{rate}} \times T_{\text{main}})$
- $EAT = (0.95 \times 10 \text{ ns}) + (0.05 \times 80 \text{ ns})$
- $EAT = 9.5 \text{ ns} + 4 \text{ ns}$
- $EAT = 13.5 \text{ ns}$

■ b) 阿姆达尔定律与性能提升:

1. 改进部分的加速比 (s):

- 这里“改进部分”指的是缓存访问相对于主存访问的加速。
- $s = T_{\text{main}} / T_{\text{cache}} = 80 \text{ ns} / 10 \text{ ns} = 8$ 。
- (注意：PPT L12 S19中s的定义是 "speedup due to the improved resources"，即改进后比改进前快多少倍。而PPT L12 S20的计算中，直接用主存时间/缓存时间作为s是不准确的，阿姆达尔定律中的s是针对整个系统中被加速部分而言的。更严谨的说法是，如果访问发生在缓存中，速度是10ns，如果发生在主存是80ns。我们计算的是整体访问时间的加速。)

2. 受益于改进资源的部分最初占用的执行时间比例 (p):

- 这里p是指如果所有访问都通过缓存完成（即100%命中），这部分访问所占的比例。在本题中，更直接的理解是，有95%的访问受益于缓存的快速访问。
- $p = \text{命中率} = 0.95$ 。

3. 整体性能提升 (Speedup):

- 没有缓存时的访问时间 = $T_{\text{main}} = 80 \text{ ns}$ 。
- 有缓存时的有效平均访问时间 = $EAT = 13.5 \text{ ns}$ 。
- $\text{Speedup} = (\text{时间无改进}) / (\text{时间有改进}) = T_{\text{main}} / EAT$
- $\text{Speedup} = 80 \text{ ns} / 13.5 \text{ ns} \approx 5.926 \text{ 倍}$ 。

■ 若严格套用PPT L12 S19的阿姆达尔公式：

$$\text{Speedup} = 1 / ((1 - p) + p/s)$$

这里s是“改进资源的加速比”，即缓存相对于主存快多少倍， $s = 80\text{ns} / 10\text{ns} = 8$ 。

p是“受益于改进资源的部分最初占用的执行时间比例”，这里指命中缓存的比例， $p = 0.95$ 。

$$\text{Speedup} = 1 / ((1 - 0.95) + (0.95 / 8))$$

$$\text{Speedup} = 1 / (0.05 + 0.11875)$$

$$\text{Speedup} = 1 / 0.16875 \approx \mathbf{5.926 \text{ 倍}}。$$

(结果与直接用 $T_{\text{main}} / \text{EAT}$ 计算一致)

○ 解析：

- a)：EAT是考虑了命中和未命中两种情况的加权平均访问时间。95%的情况下访问时间是10ns，5%的情况下访问时间是80ns。

- b)：

- 改进部分的加速比s：当访问命中缓存时，速度是原来的 $80/10 = 8$ 倍。
- p：有95%的内存访问能够享受到这种8倍的加速。
- 整体加速是针对整个内存访问而言的。没有缓存时，所有访问都是80ns。有了缓存后，平均访问时间降到了13.5ns。所以整体性能提升了 $80/13.5$ 倍。阿姆达尔定律提供了一个通用的计算框架，其核心思想是，系统的整体性能提升受限于系统中未能被加速部分所占的比例。

3. 请设计一个简化的DMA传输过程的事件序列（至少包括CPU、DMA控制器、外设和内存之间的关键交互步骤），描述从外设请求DMA传输数据到内存，直到传输完成，CPU收回总线控制权的整个流程。

○ 标准答案：

一个简化的DMA传输数据到内存的事件序列：

1. CPU初始化DMA控制器：

CPU向DMA控制器编程，设置：

- 源地址（外设的某个寄存器或缓冲区地址，虽然DMA控制器通常直接与外设接口交互，但逻辑上指定了数据来源）。
- 目标地址（内存中的起始地址）。
- 传输字节数（或字数）。
- 传输方向（外设到内存）。
- 启动DMA传输模式。

2. 外设准备就绪，请求DMA：当外部设备（如磁盘）有数据准备好传输时，它向DMA控制器发出一个**DMA请求信号 (DREQ)**。

3. **DMA控制器请求总线**：DMA控制器收到DREQ后，向CPU发出一个**总线请求信号 (HOLD)**，请求获得系统总线（地址总线、数据总线、控制总线）的控制权。

4. **CPU响应总线请求**：CPU在完成当前的总线周期（例如，当前指令的内存访问）后，会释放总线，并向DMA控制器发送一个**总线应答信号 (HLDA)**，表示同意DMA控制器接管总线。此时CPU进入暂停状态，不使用总线。

5. **DMA控制器获得总线控制权并通知外设**：DMA控制器检测到HLDA有效后，它就获得了总线的控制权。然后，DMA控制器向请求DMA的外设发送一个**DMA应答信号 (DACK)**，通知该外设可以开始数据传输。

6. 数据传输：

- DMA控制器将要写入内存的地址放到地址总线上。
- 外设将数据放到数据总线上。
- DMA控制器发出内存写控制信号。
- 数据从外设直接写入内存的指定地址。
- DMA控制器内部的地址计数器和字节计数器会相应更新。
- 这个过程会重复进行，直到所有指定字节数的数据传输完毕。

7. **传输完成，DMA控制器释放总线：** 当设定的数据量传输完成后，DMA控制器会撤销对CPU的总线请求信号（释放HOLD）。

8. **CPU重新获得总线控制权：** CPU检测到HOLD信号无效后，会撤销HLDA信号，并重新获得对系统总线的控制权，继续执行之前的任务。

9. **(可选) DMA控制器通知CPU传输完成：** DMA控制器可以通过产生一个中断信号来通知CPU，DMA传输已经完成。

○ **解析：**

这个过程就像CPU（老板）委托DMA控制器（经理）去处理一批货物的入库（外设到内存）。

1. 老板给经理下指令（初始化DMA）。
2. 供应商（外设）告诉经理货备好了（DREQ）。
3. 经理跟老板说：“我要用车道（总线）运货”（HOLD）。
4. 老板说：“好，你用吧，我先歇会儿”（HLDA）。
5. 经理拿到车道使用权，告诉供应商：“可以开始卸货了”（DACK）。
6. 经理指挥把货直接从供应商送到仓库指定位置（数据传输）。
7. 货运完了，经理不再占用车道（释放HOLD）。
8. 老板重新接管车道（CPU恢复总线控制）。
9. 经理可能还会跟老板汇报一下：“货都入库了”（DMA完成中断）。

4. **如果要使用多片2KB大小的EPROM芯片（例如2716，其地址线为A0-A10）来构建一个8KB的ROM存储区域，请简述如何通过地址译码和片选信号来实现。你需要多少片EPROM？需要使用CPU的哪些高位地址线进行译码？**

○ **标准答案：**

■ **所需EPROM芯片数量：**

- 目标ROM存储区域大小 = 8KB。
- 单片EPROM大小 = 2KB。
- 所需芯片数量 = $8\text{KB} / 2\text{KB} = 4$ 片。

■ **地址译码与片选实现：**

1. **芯片内部寻址：** 每片2KB的EPROM（如2716）需要11根地址线（A0-A10）来寻址其内部的 $2^{11} = 2048$ 个字节。因此，CPU地址总线的低11位（A0-A10）应并联到所有4片EPROM的对应地址输入引脚上。

2. **片选信号生成（地址译码）：**

为了能够独立选择这4片EPROM中的某一片，我们需要使用CPU地址总线中高于A10的地址线进行译码，以产生4个独立的片选信号。

- 8KB的存储区域需要寻址 2^{13} 个字节 ($8\text{KB} = 8 \times 1024\text{B} = 2^3 \times 2^{10}\text{B} = 2^{13}\text{B}$)。所以总共需要13根地址线（A0-A12）来覆盖整个8KB区域。
- 其中A0-A10用于片内寻址。

- 剩下的地址线 **A11和A12** 可以用来选择4片芯片中的一片。

3. 译码逻辑：

使用一个2-4译码器（例如74LS139，它包含两个2-4译码器，或者用逻辑门搭建）。

- 将CPU的地址线A11和A12作为译码器的输入。
- 译码器的4个输出 (Y0, Y1, Y2, Y3) 分别连接到4片EPROM的片选输入端 (CS# 或 CE#，通常是低电平有效)。

4. 工作方式：

- 当A12 A11 = 00 时，译码器输出Y0有效，选中第一片EPROM (地址范围 0000-07FFh 相对基址)。
- 当A12 A11 = 01 时，译码器输出Y1有效，选中第二片EPROM (地址范围 0800-0FFFh 相对基址)。
- 当A12 A11 = 10 时，译码器输出Y2有效，选中第三片EPROM (地址范围 1000-17FFh 相对基址)。
- 当A12 A11 = 11 时，译码器输出Y3有效，选中第四片EPROM (地址范围 1800-1FFFh 相对基址)。

5. **数据线：**所有4片EPROM的数据线 (D0-D7) 应并联到CPU的数据总线上。当某一片芯片被选中且处于读状态时，其数据会输出到数据总线。

- **需要使用的CPU高位地址线进行译码： A11 和 A12。**

○ 解析：

想象你有4本书（4片2KB的EPROM），每本书都有2048页（字节），页码从0到2047（对应A0-A10）。你想把这4本书合成一本大书（8KB的ROM区域）。

- 你需要一个方法来区分这4本书。这就是地址译码的作用。
- CPU的低11位地址线 (A0-A10) 告诉你翻到每本书的第几页。
- CPU的更高位地址线 (A11, A12) 用来告诉你要看的是这4本书中的哪一本。
 - A12 A11 = 00 -> 看第1本书
 - A12 A11 = 01 -> 看第2本书
 - A12 A11 = 10 -> 看第3本书
 - A12 A11 = 11 -> 看第4本书
- 译码器就像一个图书管理员，根据A11和A12的指示，只把对应那本书的“允许阅读”灯（片选信号）打开。
- 所有书的相同页码的内容（数据线）都接到同一个地方，但只有被选中的那本书的内容才能被读出来。