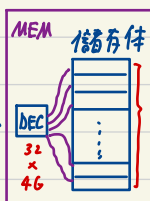
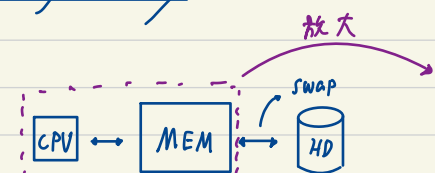


Memory hierarchy



$2^{22} = 46$ 個位置

Note:

Decoder 有几个输出, 就有几个 AND GATE
故以左圖為例, 有 46 個 AND GATE
MEM size 越大, decoder 越大
解碼時間越長
CPU 存取速度越慢

根據馬·鈕曼原則

當要執行某個 executable file 時

需要先將其 swap in 到 MEM 中, 再由 CPU 去抓指令來執行

但 swap 速度很慢 (是 disk I/O operation)

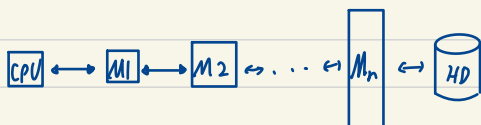
而當 MEM 越小, swap 次數會增加

① 站在減少 swap 次數上, MEM 應越大越好.

② 站在減少 CPU 存取時間上, MEM 應越小越好.

∴ 使用 MEM Hierarchy 方法:

靠近 CPU 的越小, 靠近 HD 的越大



∴ 設 M_1, M_2, \dots, M_n 為資料構成之集合

則: $M_1 \subseteq M_2 \subseteq M_3 \dots \subseteq M_{n-1} \subseteq M_n$

CPU 只能存取 M_1 的資料, 若欲存取資料不在 M_1 的話

需向 M_2 - 一直向後找, 找到再向前傳至 CPU

設程式執行時, 每一指令, 資料執行機率均等

∴ M_1 很小 ∴ miss rate 高 資料傳輸慢

但事實上, 有 locality 的現象 每一指令, 資料執行機率不均等

∴ 可把執行機率高指令, 資料放至 M_1

降低 miss rate 至很低, 則存取速度約等同 M_1 速度

① Principle of Locality

Def: 程式執行時, 任一時間只能存取一小部份的地址空間

1. Temporal Locality:

程式中某些部份被存取後很快又會被存取到

```
for i=1 to n do {  
    * loop body *  
}
```

2. Spatial Locality:

程式中某些部份被存取後, 它位址附近的部份很快又會被存取到

Sequential Execution

Array

Ex. - 一段 C code 同時具 spatial locality & temporal locality

```
void clear ( int A[], int n ) {
```

```
    int i;
```

```
    for ( i=0; i < n; i++ ) {
```

```
        A[i] = 0;
```

```
    }
```

```
}
```

$A[i] = 0$ 指令的擷取具 temporal locality

$A[i]$ 資料的擷取具 spatial locality

② MEM 建構技術

① SRAM: 資料存在 Latch, 存 1 bit 需要 6 顆電晶体



速度較快, 成本高

static: 資料不隨時間而變

② DRAM: 資料存在電容器中, 存 1 bit 只需 1 顆電晶体

會有寄生電容效應 (充放電太久, \therefore 有資料速度慢)



dynamic: \therefore 電容器會漏電, 需要在漏電量沒很多時補回, 稱作 refresh

如果不 refresh, 存的資料值會改變

③ Magnetic Disk: 磁碟上有磁性物質

而讀寫頭會有線圈, 可通電流產生感應磁場

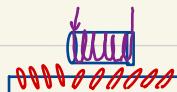
使磁碟上磁性物質沿一方向排列

\therefore 改變電流方向即可改變 0, 1

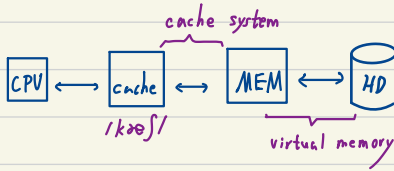
磁碟



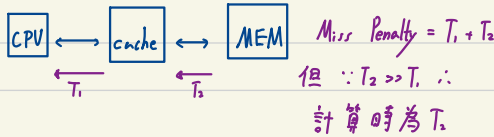
\Rightarrow



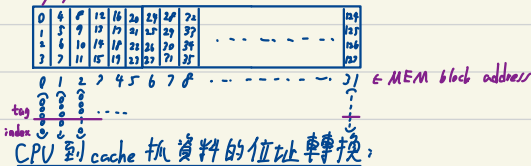
但磁化過程很慢，故速度慢



5. miss penalty, 設 CPU 存取上層資料時 miss 後, 需往更上層找資料
找到了傳回一個 block 和傳回 CPU 的時間



MEM block address % 6



CPU $\rightarrow 38 \rightarrow \lfloor \frac{38}{4} \rfloor = 9$ $\rightarrow 9 \% 6 = 3 \rightarrow$ 到 cache index 找资料

MEM
byte
address

MEM
block
address

Cache
Index

$q \div b = 1 \Rightarrow$ 比 2^k 大 (找不到就 cache miss)

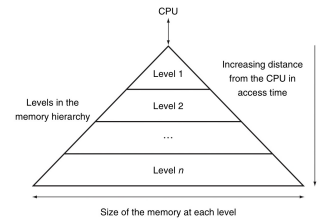


FIGURE 5.3 This diagram shows the structure of a memory hierarchy: as the distance from the processor increases, so does the size. This structure, with the appropriate operating mechanisms, allows the processor to have an access time that is determined primarily by level 1 of the hierarchy and yet have a memory as large as level n . Maintaining this illusion is the subject of this chapter. Although the local disk is normally the bottom of the hierarchy, some systems use tape or a file server over a local area network as the next levels of the hierarchy.

共需做2次除法(除法速度慢)

0. $\boxed{CPU} = 0100110$
 MEM Block Address
 ↓
 01001
 Tag Index

⇒ 不需做除法, 只需做 bit 切割即可

cache miss 處理

e.g.
 CPU → 0010100
 Tag Index offset

設 cache miss, CPU 會 stall

Cache Controller 會要求 MEM 把5号内容讀到 cache

讀取完成後, Controller 才會通知 CPU

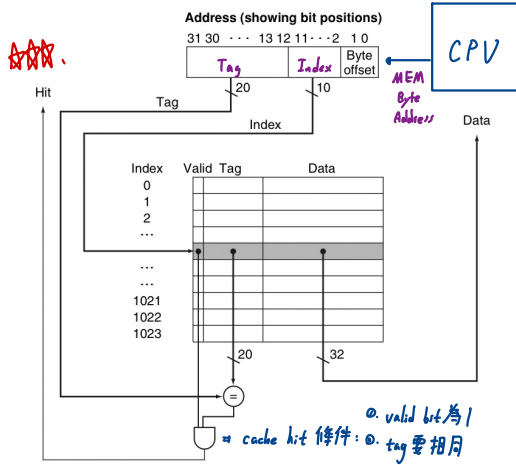
之後, CPU 照 cache index, 到對應编号 offset 抓資料

Memory Address 切割方式:

0. 看 MEM 共有几个 blocks, 設為 n 个 blocks $\Rightarrow \lg n$ 个 bit

0. 看 cache 共有几个 blocks, 設為 m 个 blocks $\Rightarrow \lg m$ 个 bit

Memory Address: $\boxed{\lg n - \lg m} \boxed{\lg m} \boxed{\quad}$
 $\lg n$



其中, valid bit 用作區別 cache block 內資料是否有效

若為0, 表示不是有效資料

(除非 CPU 用到前面 program 資料)

FIGURE 5.10 For this cache, the lower portion of the address is used to select a cache entry consisting of a data word and a tag. This cache holds 1024 words or 4 KiB. We assume 32-bit addresses in this chapter. The tag from the cache is compared against the upper portion of the address to determine whether the entry in the cache corresponds to the requested address. Because the cache has 2^{10} (or 1024) words and a block size of one word, 10 bits are used to index the cache, leaving $32 - 10 = 20$ bits to be compared against the tag. If the tag and upper 20 bits of the address are equal and the valid bit is on, then the request hits in the cache, and the word is supplied to the processor. Otherwise, a miss occurs.

● 計算題

给定: 0. cache size: cache 可存資料量

0. block size

0. address length

求: 0. block 總數: $\frac{\text{cache size}}{\text{block size}}$

Tag | Index | offset

0. 切位址求 tag bit 數: $\text{Address length} - \lg(\text{block size}) - \lg(\text{block 總數})$

$$\textcircled{3}. (1 + \text{tag bit} + \lg(\text{block size})) \times \text{block 个数}$$