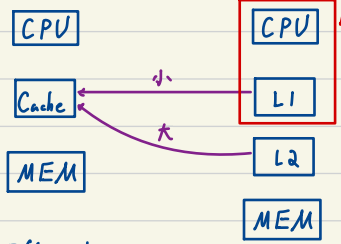


Multi-Level Cache: 减少 miss penalty

1-level cache

2-level cache



同一个 chip 上

- L1 采用 split cache \Rightarrow 增加频宽, 减少 hit time
- L2 采用 combined cache \Rightarrow 减少 miss rate

- L1 采用 write through / write back 皆可
- L2 采用 write back

Cache 只能 medium size

\therefore 太大 \Rightarrow hit time 长

L1 越小越好 \Rightarrow 减少 hit time, clock cycle time

太小 \Rightarrow miss rate 高

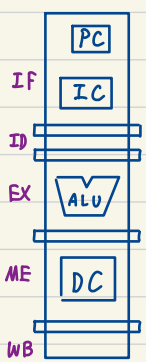
L2 越大越好 \Rightarrow L1 miss 可到 L2 找

减少 miss penalty

Multi-Level Cache 效能分析

Split Cache:

Combined Cache:



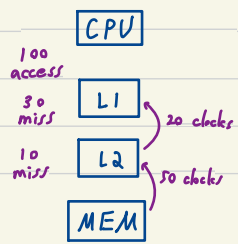
IC Extra CPI = $\frac{\# \text{ of I-C access per instruction}}{\text{CPI}} \times \sum_{i=1}^k \text{Miss Rate}_i \times \text{Miss Penalty}_i$

DC Extra CPI = $\frac{\# \text{ of D-C access per instruction}}{\text{CPI}} \times \sum_{i=1}^n \text{Miss Rate}_i \times \text{Miss Penalty}_i$

	MR	MP
$1 + 1w/sw\%$	MR_1	MP_1
	MR_2	MP_2
	\vdots	\vdots
	\vdots	\vdots
	MR_n	MP_n

of access per instruction

Example: $1w/sw = 2\%$



Total stall cycles: $30 \times 20 + 10 \times 50$

Stall cycles per access: $\frac{30}{100} \times 20 + \frac{10}{100} \times 50 = \sum_{i=1}^n MR_i \times MP_i$

Stall cycles per instruction: $1.2 \cdot \left[\frac{30}{100} \times 20 + \frac{10}{100} \times 50 \right] = (1 + 1w/sw\%) \left(\sum_{i=1}^n MR_i \times MP_i \right)$
 $= (1.2 \times \frac{30}{100}) \times 20 + (1.2 \times \frac{10}{100}) \times 50$
 $= \sum_{i=1}^n \text{miss rate per instruction}_i \times \text{miss penalty}_i$

Original:

CPU	L1	MEM	1x memory access time 1x penalty time
MRP _i	0.02	X	
MP _i	400	X	

增加 L2:

CPU	L1	L2	MEM
MRP _i	0.02	0.005	X
MP _i	20	400	X

<Ex P.48> base CPI: 1.0

clock rate: 4GHz

memory access time: 100ns = 400 clock cycles

L1 miss rate per instr: 2%

if L2 exists, L1 miss penalty: 5ns = 20 clock cycles

Original effective CPI: $1.0 + 0.02 \times 400 = 9$

增加 L2 : $1.0 + 0.02 \times 20 + 0.005 \times 400 = 3.4$

Note: 1. 没给完美 CPI 设为 1

2. 给 miss rate, 没给 lw/sw ratio, 设为 0

3. penalty 时间给定绝对时间, 没给 clock rate & clock cycle time
1x L1 cache hit time 1x clock cycle time

Global & Local Miss Rate

1000 access	CPU	GMR	LMR
50 miss	L1	$\frac{50}{1000}$	$\frac{50}{1000}$
20 miss	L2	$\frac{20}{1000}$	$\frac{20}{50}$
5 miss	L3	$\frac{5}{1000}$	$\frac{5}{20}$
	MEM		

Global Miss Rate: 1x CPU 为观察点

Local Miss Rate: 以上一层缓存存储为观察点

$$MRP_i = (1 + lw/sw\%) \times GMR_i$$

$$= (1 + lw/sw\%) \times \prod_{k=1}^i LMR_k$$

Average memory access time (AMAT)

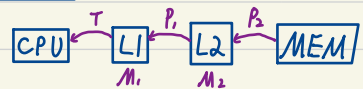
CPU	cache	MEM
hit time	T	
miss rate	M	
miss penalty	P	

$$AMAT = T \times (1 - M) + (T + P) \times M$$

$$= T - T \times M + T \times M + M \times P$$

$$= T + M \times P$$

2-level



$$AMAT = T \times (1 - M_1) + (T + P_1) \times (M_1 - M_2)$$

$$+ (T + P_1 + P_2) \times M_2$$

$$= T + M_1 \times P_1 + M_2 \times P_2$$

n-level cache:

General:

$$T_{L1} + \sum_{i=1}^n MR_i \times MP_i$$

Note: 增加 associativity 不一定會降低 AMAT

$$\therefore AMAT = \text{hit time} + MR \times MP$$

降低 MR, 但 hit time 上升

<Ex p.54>:

		1000 access	60 miss	30 miss	5 miss	
CPU	L1	L2	L3	MEM		
GMR	6%	3%	0.5%			
LMR	6%	50%	$\frac{5}{30}$			
hit time	1	5	10			
MP	5	10	100			

$$AMAT = 1 + 0.06 \times 5 + 0.03 \times 10 + 0.005 \times 100$$

$$= 2.1 \text{ clock cycles}$$

$$\text{Extra CPI} = (1 + 1w/sw \text{ ratio}) \times \sum MR_i \times MP_i$$

$$= 1.5 \times 1.1 = 1.65$$

Note: AMAT per instruction

$$= (1 + 1w/sw\%) \times AMAT$$

$$= (1 + 1w/sw\%) \times (\text{hit time} + MR \times MP)$$

$$= (1 + 1w/sw\%) \times \text{hit time} + MRP \times MP$$