



Esercizi sulle prestazioni delle memorie cache

Prof. Alberto Borghese
Dott. Massimo Marchi
Dipartimento di Scienze dell'Informazione
alberto.borghese@unimi.it
Università degli Studi di Milano

A.A. 2012-2013 1/14 http:\\borghese.di.unimi.it\



Sommario

Impatto delle memorie cache
Progettazione delle memorie cache

A.A. 2011-2012 2/34 http:\\homes.dsi.unimi.it/~borghese



Tempo medio di accesso alla memoria



TEMPO DI ACCESSO. E' legato al tempo di accesso del livello inferiore di memoria.

TEMPO DI TRASFERIMENTO. E' legato alla larghezza di banda del canale di comunicazione tra i due livelli di memoria (e.g. bus).

Il tempo medio di accesso alla memoria sarà:

$$T_{\text{medio}} = \text{HIT_RATE} * \text{HIT_TIME} + \text{MISS_RATE} * \text{MISS_TIME} =$$

$$\text{HIT_TIME} * \text{HIT_RATE} + \text{MISS_RATE} * (\text{HIT_TIME} + \text{MISS_PENALTY}) =$$

$$\text{HIT_TIME} * (\text{HIT_RATE} + \text{MISS_RATE}) + \text{MISS_RATE} * \text{MISS_PENALTY} =$$

$$\text{HIT_TIME} + \text{MISS_RATE} * \text{MISS_PENALTY}$$

Hit – massimizzo hit_rate

Miss – Minimizzo miss_penalty (costo delle miss -> stalli)

A.A. 2011-2012

3/34

<http://homes.dsi.unimi.it/~borgnese>



Impatto di una memoria cache



Il tempo di CPU è composto dal tempo richiesto dalla CPU per eseguire il programma e dal tempo che la CPU trascorre in attesa di risposta dal sottosistema di memoria.

$$T_{\text{CPU}} = (\#\text{Cicli della CPU in esecuzione} + \#\text{Cicli di stallo}) * T_{\text{Clock}}$$

Ipotesi:

- Tutti gli stalli di memoria sono dovuti al fallimento di accesso alla cache.
- I cicli di clock utilizzati per un accesso alla cache riuscito (HIT) sono inclusi nei cicli di clock della CPU in esecuzione.

$$\text{CPI} = \frac{\text{Numero_Cicli_Clock_per_Istruzione}}{\text{Numero di Cicli di Clock / Numero di istruzioni.}}$$

A.A. 2011-2012

4/34

<http://homes.dsi.unimi.it/~borgnese>



Impatto di una memoria cache



$$\# \text{ Cicli_clock_stallo} = \# \text{ Accessi_Memoria} * \text{MISS_RATE} * \text{MISS_PENALTY}$$

$$\text{Tempo}_{\text{CPU Programma}} = (\# \text{ Cicli_clock} + \# \text{ Cicli_clock_stallo}) * T_{\text{clock}} = \\ \# \text{ Istruzioni} * \text{CPI}_{\text{exec}} * T_{\text{clock}} + \# \text{ Cicli_clock_stallo} * T_{\text{clock}}$$

$$\text{CPI}_{\text{con_cache}} = \text{CPI}_{\text{exec}} + \# \text{ Cicli_clock_stallo} / \# \text{ Istruzioni} = \\ \text{CPI}_{\text{exec}} + (\# \text{ Accessi_memoria} / \# \text{ Istruzioni}) * \text{MISS_RATE} * \text{MISS_PENALTY}$$

Caso ideale: (100% HIT, 0% MISS): $\text{CPI}_{\text{con_cache}} = \text{CPI}_{\text{exec}}$

Caso senza cache: (100% MISS): $\text{CPI}_{\text{senza_cache}} = \\ \text{CPI}_{\text{exec}} + (\# \text{ Accessi_memoria} / \# \text{ Istruzioni}) * \text{MISS_PENALTY}$



Esercizio su cache



Si consideri il VAX-11/780. La MISS_PENALTY è di 6 cicli di clock, mentre tutte le istruzioni impiegano 8.5 cicli di clock se si ignorano i MISS (stalli della memoria). Ipotizzando un MISS_RATE dell'11% e che vi siano in media 2 riferimenti alla memoria per ogni istruzione,

⇒ Qual è l'impatto sulle prestazioni quando viene inserita la cache reale rispetto ad una cache ideale?

⇒ Qual è l'impatto sulle prestazioni tra il caso di cache reale e senza inserimento della cache?




Soluzione esercizio su cache

Dati di ingresso: MISS_PENALTY=6 $CPI_{exec}=8.5$ MISS_RATE=0,11
#Accessi_memoria/#Istruzioni = 2

$$CPI_{con_cache} = 8,5 + 0,11 * (2 * 6) = 9,82$$

$$CPI_{con_cache_ideale} = 8,5 + 0 * (2 * 6)$$

$$CPI_{senza_cache} = 8,5 + 1 * (2 * 6) = 20,5$$

Perdita in prestazioni (speed-up): $CPI_{con_cache_ideale} / CPI_{con_cache} \Rightarrow 8,5 / 9,82 = 0,865$

Guadagno in prestazioni (speed-up): $CPI_{senza_cache} / CPI_{con_cache} \Rightarrow 20,5 / 9,82 = 2,087$

A.A. 2011-2012 7/34 http://homes.dsi.unimi.it/~borgnese




Sommarrio

Impatto delle memorie cache
Progettazione delle memorie cache

A.A. 2011-2012 8/34 http://homes.dsi.unimi.it/~borgnese



Esercizio 1



Si progetti una cache di 16KByte a 8 vie per un sistema con indirizzamento al byte di 32bit, bus dati a 32 bit, bus indirizzi a 30bit, e word di 4 byte. L'ampiezza della cache è di 1 word. In quale cella della cache e con quale tag viene memorizzato il dato letto dall'indirizzo 0x40404040 della memoria principale (lw \$t0, 0x40404040)? Quanta memoria è necessaria per implementare la cache?

Una cache a 8 vie è una cache in grado di memorizzare 8 blocchi diversi per ogni indice. 16Kbyte corrispondono a $16 / 8 = 2\text{KByte}$ per via. La cache è organizzata a matrice. L'ampiezza è di 4 Byte e l'altezza sarà quindi $2\text{KByte} / 4 \text{ Byte} = 512 = \text{Numero di linee}$. Servono quindi 11 bit ($\log_2(2\text{KByte})$) per indirizzare i Byte all'interno della cache. Il tag risulta quindi pari a $30-11 = 19$ bit. Gli 11 bit di indirizzamento all'interno della cache vengono suddivisi come segue. L'ampiezza della linea di ciascuna via della cache è pari a 4 Byte, servono quindi 2 bit. Gli altri $11-2 = 9$ bit servono per indirizzare la linea ($2^9 = 512$ linee)

L'indirizzo a 30 bit viene quindi suddiviso come segue:

	tag	indice	unused
bit	29-11	10-2	1-0

L'indirizzo **0x40404040 = 01.00 0000 0100 0000 0100 0.000 0100 00.00₂**, corrisponde all'indice **000 0100 00 = 16** ed al tag = **000 0000 1000 0000 1000 = 2056**.

Lo spazio fisicamente necessario per implementare la cache è pari a:
19 bit di tag + 32 bit di dato + 1 bit di validate + 1 dirty bit per banco = 53 bit per linea per banco.

Il numero totale di bit sarà quindi:
 $53\text{bit} * 512 \text{ linee} * 8 \text{ vie} = 217,088 \text{ bit} (27,136 \text{ Byte})$.

La memoria quindi a fronte di una capacità di memoria di 16KByte, richiede un totale di 26,5 Kbyte.

A.A. 2011-2012

9/34

http://homes.dsi.unimi.it/~borgnese



Esercizio 2



Si progetti una cache di 128 byte a corrispondenza diretta organizzata in linee di 4 word con word di 2 byte per un sistema con indirizzamento al byte a 16 bit, bus indirizzi a 16 bit. Si consideri il tempo di accesso della cache in lettura/scrittura 2ns ed il tempo di accesso in lettura alla DRAM di 60ns. Si consideri anche che la cache lavora in modalità "write-back" e che il tempo per scaricare una linea di cache sia di 100ns. Data la sequenza di accessi 0,2,8,10,132,2 dire quanto vale il tempo di accesso totale nell'ipotesi che i bit di validate siano tutti a 0 all'inizio.

La cache ha un'ampiezza della linea pari a $4 \text{ word} * 2 \text{ Byte} / \text{Word} = 8 \text{ Byte}$.

Il numero di bit per indirizzare i byte all'interno della cache è quindi pari a $\log_2(128) = 7$ ed il numero di bit di tag sarà quindi: $16 - 7 = 9$.

Il numero di bit necessari ad indirizzare i Byte all'interno di una linea è pari a:
 $\log(\text{ampiezza linea in Byte}) = 3 = (\log_2(8))$.

Il numero di linee sarà quindi pari a: 7 (#bit per indirizzare i Byte dentro la cache) - $3 = 4$. Questo numero si può ottenere anche dividendo la capacità della cache per la capacità della linea, cioè $128 \text{ Byte} / 8 \text{ Byte} = 16 \text{ linee} \Rightarrow 3 \text{ bit}$.

L'indirizzo viene quindi così suddiviso:

	tag (9 bit)	indice (4 bit)	offset (2 bit)	unused
bit	15-7	6-3	2-1	0

Lo spazio occorrente per implementare questa cache è quindi:

$(9 \text{ bit di tag} + 1 \text{ bit di validate} + 1 \text{ dirty bit} * 64 \text{ bit per linea}) * 16 \text{ linee} = 1184 \text{ bit} (148 \text{ Byte})$

A.A. 2011-2012

10/34

http://homes.dsi.unimi.it/~borgnese



Esercizio 2



Si progetti una cache di 128 byte a corrispondenza diretta organizzata in linee di 4 word con word di 2 byte per un sistema con indirizzamento al byte a 16 bit, bus indirizzi a 16 bit. Si consideri il tempo di accesso della cache in lettura/scrittura 2ns ed il tempo di accesso in lettura alla DRAM di 60ns. Si consideri anche che la cache lavora in modalità "write-back" e che il tempo per scaricare una linea di cache sia di 100ns. Data la sequenza di accessi 0 (load), 2 (store), 8 (load), 10 (load), 132 (load), 4 (store) dire quanto vale il tempo di accesso totale nell'ipotesi che i bit di validate siano tutti a 0 all'inizio.

La cache ha un'ampiezza della linea pari a $4 \text{ word} * 2 \text{ Byte / Word} = 8 \text{ Byte}$.

Il numero di bit per indirizzare i byte all'interno della cache è quindi pari a $\log_2(128) = 7$ ed il numero di bit di tag sarà quindi: $16 - 7 = 9$.

Il numero di bit necessari ad indirizzare i Byte all'interno di una linea è pari a:
 $\log(\text{ampiezza linea in Byte}) = 3 = (\log_2(8))$.

Il numero di linee sarà quindi pari a: 7 (#bit per indirizzare i Byte dentro la cache) - $3 = 4$.

Questo numero si può ottenere anche dividendo la capacità della cache per la capacità della linea, cioè $128 \text{ Byte} / 8 \text{ Byte} = 16$ linee $\Rightarrow 3$ bit.

L'indirizzo viene quindi così suddiviso:

	tag (9 bit)	indice (4 bit)	offset (2 bit)	unused
bit	15-7	6-3	2-1	0

Lo spazio occorrente per implementare questa cache è quindi:

$(9 \text{ bit di tag} + 1 \text{ bit di validate} + 1 \text{ dirty bit} * 64 \text{ bit per linea}) * 16 \text{ linee} = 1184 \text{ bit} (148 \text{ Byte})$

A.A. 2011-2012

11/34

http://homes.dsi.unimi.it/~borgnese



Esercizio 2 - accessi



L'accesso a **0 (00000000-0000-00-0)** genera una **miss** perchè la cache è totalmente vuota. il blocco corrispondente con **tag=0** e **indice=0** viene caricato dalla memoria in un solo accesso (64 bit di bus dati). TEMPO: 60ns

L'accesso a **2 (00000000-0000-01-0)** corrisponde allo stesso blocco, **indice=0** quindi abbiamo una **hit** (con dirty bit impostato a 1). TEMPO: 2ns.

L'accesso a **8 (00000000-0001-00-0)** corrisponde al **tag=0**, **indice=1** che non è ancora stato letto quindi è una **miss**. TEMPO: 60ns.

L'accesso a **10 (00000000-0001-01-0)** corrisponde a **tag=0**, **indice=1**; in questo caso è una **hit**. TEMPO: 2ns.

L'accesso a **132 (00000001-0000-10-1)** corrisponde a **tag=1**, **indice=0**; in questo caso è una **miss** (da notare che la linea 0 contiene anche il dirty bit = 1 perchè scritta dalla seconda istruzione. Occorre quindi procedere anche al write-back della linea 0 con tag 0). TEMPO: $60\text{ns} + 100\text{ns}$.

L'accesso a **4 (00000000-0000-10-0)** corrisponde a **tag=0**, **indice=0**, non è ancora presente in cache quindi **miss**. Occorre prima caricare la linea aggiornata dalla memoria e poi operare la scrittura. TEMPO: 60ns.

TEMPO TOTALE = $60 + 2 + 60 + 2 + 160 + 60 = 344 \text{ ns}$.

A.A. 2011-2012

12/34

http://homes.dsi.unimi.it/~borgnese



Esercizio 3



Si progetti una cache a corrispondenza diretta di 64Kbytes, a 4 vie, che utilizzi word di 8 Byte e memorizzi 4 word per ogni linea, bus indirizzi su 16 bit, indirizzamento al byte. Dimensionare tutti i parametri della cache.

La cache ha una capacità per banco di 64KByte / 4 = 16 Kbyte.
 Il numero di bit per indirizzare i byte all'interno della cache è quindi pari a $\log_2(16K) = 14$ bit ed il numero di bit di tag sarà quindi: $16 - 14 = 2$.
 L'ampiezza della linea per banco sarà di: $4 \text{ Word} * 8 \text{ Byte} / \text{Word} = 32 \text{ Byte}$.
 L'altezza della cache sarà perciò di $16 \text{ Kbyte} / 32 \text{ Byte} = 512$ Linee.
 Il numero di bit necessari ad indirizzare i Byte all'interno di una linea è pari a: $\log_2(32) = 5$.
 Il numero di linee sarà quindi pari a: $14 - 5 = 9$.
 La parola all'interno della linea viene selezionata dai bit 3 e 4 dell'indirizzo, mentre i bit 0, 1, 2 non vengono utilizzati ed indirizzano uno degli 8 Byte che costituiscono la parola.

La dimensione totale della cache è pari a:
 $(2 \text{ (bit di TAG)} + 1 \text{ (bit validità)} + 1 \text{ (dirty bit)} + 32 * 8 \text{ (dati)}) * 512 * 4 = 532,480 \text{ bit (66,560 Byte)} = 65\text{KByte}$

Con un tag piccolo, c'è un overhead sull'allocazione di memoria piccolo (1KByte).

La scomposizione dell'indirizzo sarà:

	tag (2 bit)	indice (13 bit)	offset (2 bit)	unused (1 bit)
bit	14-15	13-5	5-4	3

Si noti che in questo caso la RAM può essere copiata tutta in cache, ricopiando i diversi macro blocchi nei diversi banchi della memoria cache.



Sommario



Impatto delle memorie cache
 Progettazione delle memorie cache