



La gerarchia delle memorie

Prof. Alberto Borghese
Dipartimento di Scienze dell'Informazione
borgnese@dsi.unimi.it

Università degli Studi di Milano



Sommario

Caratteristiche di un sistema di memoria

Scrittura di una memoria

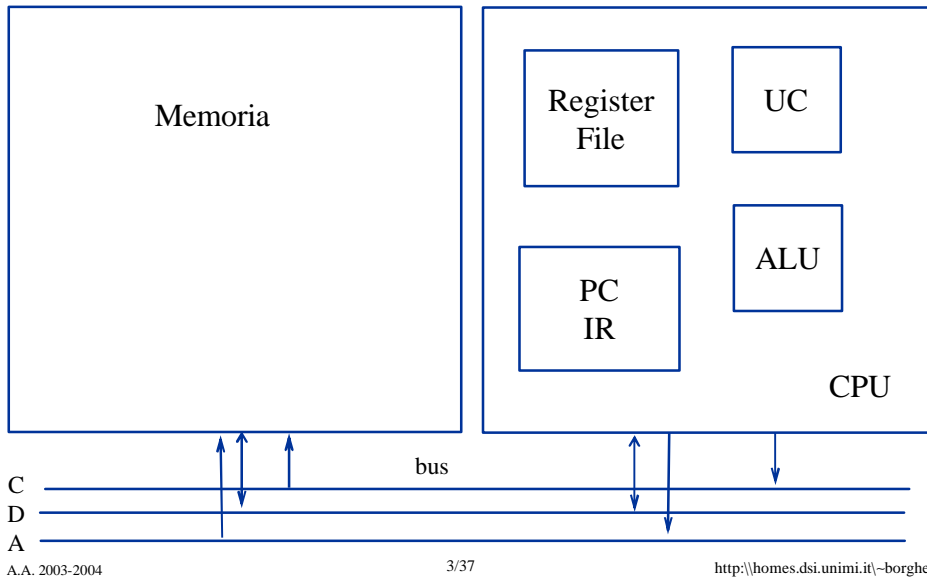
Principio di funzionamento di una memoria cache

Cache a mappatura diretta

Il campo tag di una cache



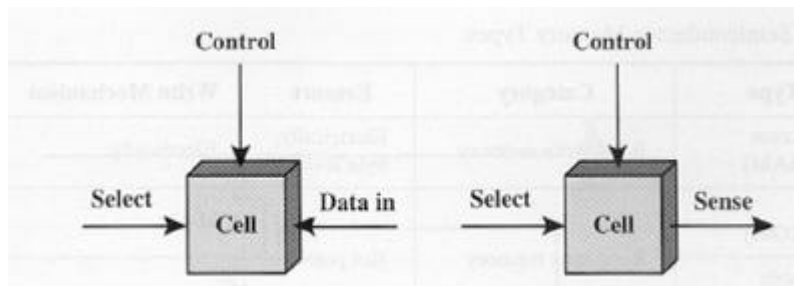
Gli attori principali di un'architettura



Cella di memoria



La memoria è suddivisa in celle, ciascuna delle quali assume un valore binario stabile.
Si può scrivere il valore 0/1 in una cella.
Si può leggere il valore di ciascuna cella.



Quale struttura di memoria abbiamo già incontrato?



Caratteristiche della memoria



Posizione della memoria:

- Processore
- Interna (cache)
- Esterna (cache + Memoria Principale)
- Disco

Metodo di accesso:

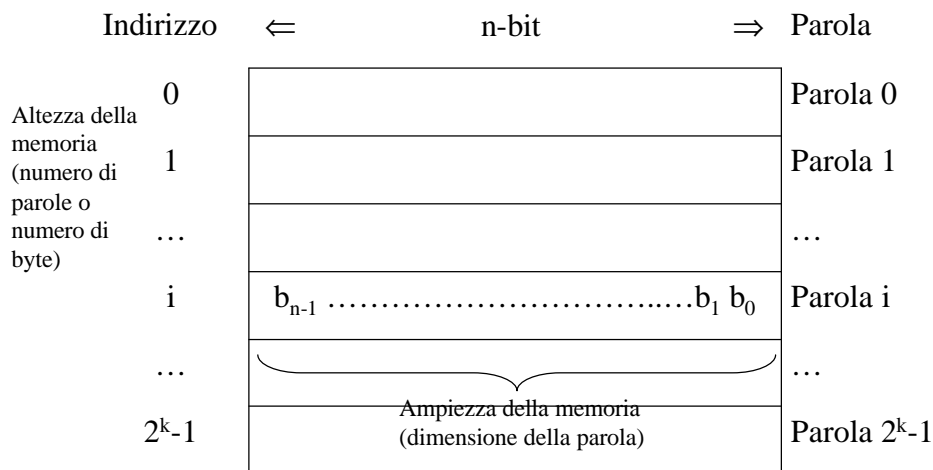
- Sequenziale (e.g. Nastri).
- Diretto (posizionamento + attesa, e.g. Dischi).
- Random Access (circuiti di lettura / scrittura HW, tempo indipendente dalla posizione e dalla storia, e.g. Cache e Memoria principale).
- Associativa (Random Access, il contenuto viene recuperato a partire da un sottoinsieme incompleto dello stesso).

Caratteristiche fisiche:

- Nelle memorie volatili (E.g. Cache), l'informazione sparisce quando si toglie l'alimentatore (memorie a semiconduttore).
- Nelle memorie non-volatili, l'informazione è duratura (un esempio di memoria volatile è la memoria magnetica dei dischi e dei nastri). Esistono memorie a semiconduttore non-volatili (ROM).



Capacità della memoria



Capacità della memoria $C = \#_parole \times dim_parola$



Misura della capacità di una memoria



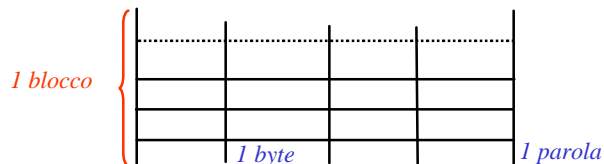
Parola di memoria. E' l'unità naturale in cui la memoria viene organizzata (e.g. 32 bit in MIPS).

Unità indirizzabile. E' il minimo numero di *unità contigue* indirizzabili. In quasi tutti i sistemi si tratta del byte.

Spazio dell'indirizzamento della memoria è: $N = \log_2 \text{Capacità}$

Unità di trasferimento:

- Blocco
- Parola



Prestazione di una memoria



Tempo di accesso.

Per una memoria a random-access, è il tempo richiesto per eseguire un'operazione di lettura / scrittura. Più precisamente è il tempo che intercorre dall'istante in cui l'indirizzo si presenta alla porta di lettura della memoria all'istante in cui il dato diventa disponibile.

Per una memoria non random-access, è il tempo richiesto per posizionare il dispositivo di lettura / scrittura in corrispondenza dell'informazione da leggere / scrivere.

Memory cycle time: si applica ad una memoria random-access. E' il tempo di accesso più il tempo necessario perchè possa avvenire un secondo accesso a memoria.

Transfer rate: quantità di informazione trasferita nell'unità di tempo si misura in:

Random-access memory = $1 / \text{Memory_cycle_time}$.

Not random-access memory = $1 / [T_A + N / R]$ con R trasfer rate.

Numero_parole / s.



Sommario



Caratteristiche di un sistema di memoria

Struttura di una memoria

Principio di funzionamento di una memoria cache

Cache a mappatura diretta

Il campo tag di una cache



Principio di progettazione di una memoria



Quanta memoria?

Quanto deve essere veloce?

Quanto deve costare?

Maggiore è la velocità di accesso, maggiore il costo per bit.

Maggiore è la capacità, minore il costo per bit.

Maggiore è la capacità, maggiore è il tempo di accesso.



Memorie piccole e veloci.

Memorie grandi e lente.



Principi di località



I programmi riutilizzano dati e istruzioni che hanno usato di recente.

Regola pratica: un programma spende circa il **90%** del suo tempo di esecuzione per solo il **10%** del suo codice.

Basandosi sul passato recente del programma, è possibile predire con ragionevole accuratezza quali dati e istruzioni userà nel prossimo futuro.

Località temporale: elementi ai quali si è fatto riferimento di recente saranno utilizzati ancora nel prossimo futuro.

Località spaziale: elementi i cui indirizzi sono vicini, tendono ad essere referenziati in tempi molto ravvicinati.

Si possono organizzare programmi e dati in modo da sfruttare al massimo il principio di località (e.g. scrittura di blocchi di dati nei dischi, salti locali...).

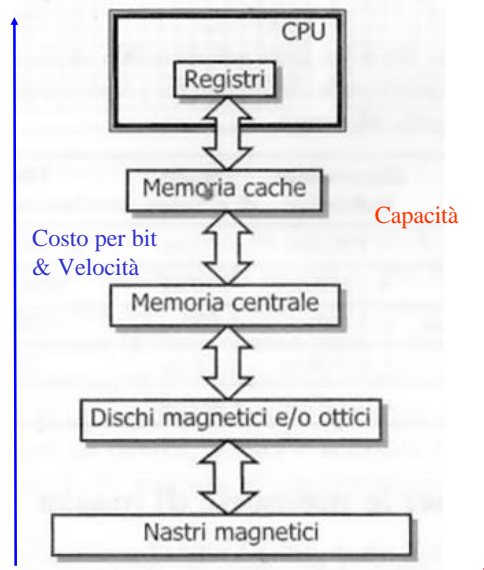


Gerarchia di memorie



Nel livello superiore troviamo un sottoinsieme dei dati del livello inferiore.

Ciascun livello vede il livello inferiore.





Gerarchia di memorie - caratteristiche



Livello	Dimensioni indicative	Tempo di Accesso	Velocità di Trasferimento (Mbyte/s)
Registri	< 1 Kbyte	< 0,01 ns	400,000 (4 byte)
Cache Primaria (Pentium 4, 3Ghz) ExecTrace cache	8kbyte 12kbyte	0.16ns	192,000 (32/64 byte)
Cache Secondaria (Pentium 4, 3Ghz)	256-512 kbyte	0.3ns	96,000 (32 byte in parallelo)
Memoria centrale (RAM, DRAM)	< 4 Gbyte	< 3-5 ns (233Mhz / 320Mhz)	1,600 – 3,000 di picco (DDSRAM – doppia lettura)
Bus PCI Bus PCI 64	133 Mhz 100 Mhz	8 byte 64 byte	133 (8 byte) 1064 (64 byte)
Dischi	> 50 Gbyte	< 10ms	< 200 (Seagate Cheetah SCSI)
Nastri	> 100 Gbyte	> 100ms	1



Tassonomia del funzionamento



HIT Successo nel tentativo di accesso ad un dato: è presente al livello superiore della gerarchia.

MISS Fallimento del tentativo di accesso al livello superiore della gerarchia => il dato o l'indirizzo devono essere cercati al livello inferiore.

HIT_RATE Percentuale dei tentativi di accesso ai livelli superiori della gerarchia che hanno avuto successo.

$$\text{HIT_RATE} = \text{Numero_successi} / \text{Numero_accessi_memoria}$$

MISS_RATE Percentuale dei tentativi di accesso ai livelli superiori della gerarchia che sono falliti

$$\text{MISS_RATE} = \text{Numero_fall.} / \text{Numero_accessi_memoria}$$

$$\text{HIT_RATE} + \text{MISS_RATE} = 1$$



Sommario



Caratteristiche di un sistema di memoria

Struttura di una memoria

Principio di funzionamento di una memoria cache

Cache a mappatura diretta

Il campo tag di una cache

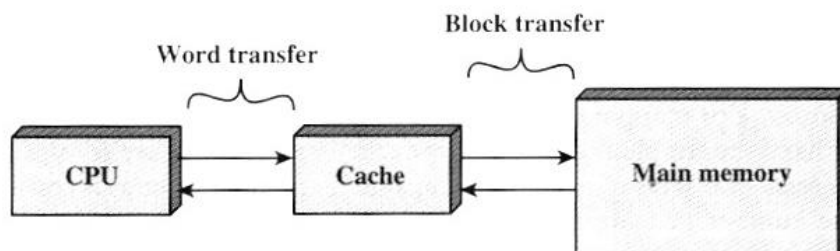


Principio di funzionamento di una cache



Scopo: fornire alla CPU una velocità di trasferimento pari a quella della memoria più veloce con una capacità pari a quella della memoria più grande.

Una cache “disaccoppia” i dati utilizzati dal processore da quelli memorizzati nella Memoria Principale.



Word transfer: Data transfer or Instruction transfer. In MIPS = 1 parola.

La cache contiene una copia di parte del contenuto della memoria principale. Di che cosa?



Sottosistema di memoria

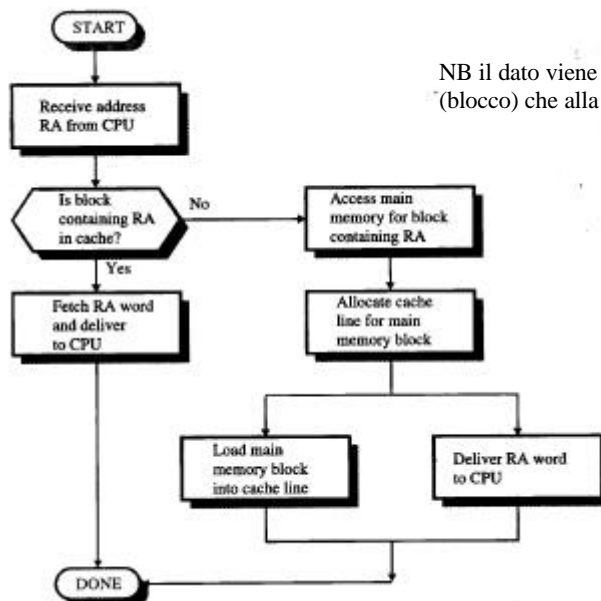


Porta nella cache primaria i dati richiesti mentre il binomio processore-memoria sta lavorando.

- 1) Controlla se una parola è in cache (Hit).
- 2) Porta una parola (e quelle vicine) in cache, prelevandole dal livello inferiore (Miss):



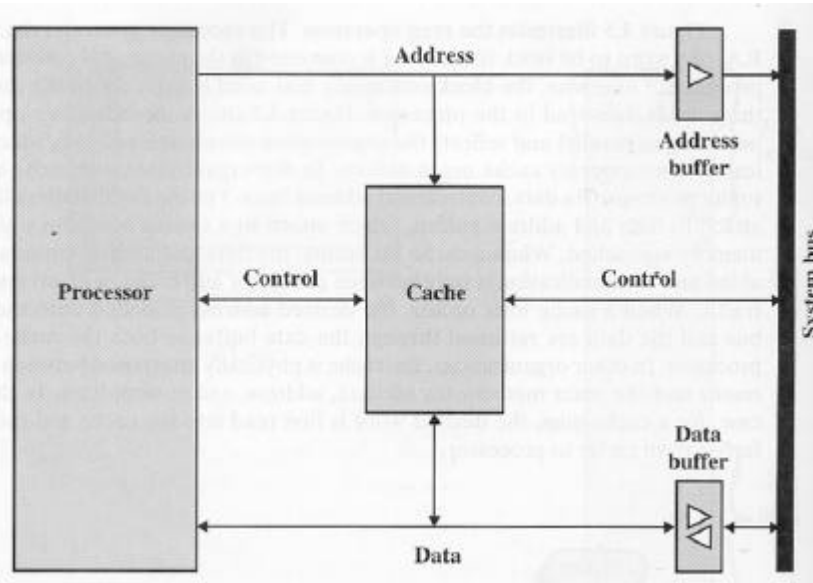
Cache operation



NB il dato viene trasferito sia alla Cache (blocco) che alla CPU (parola)



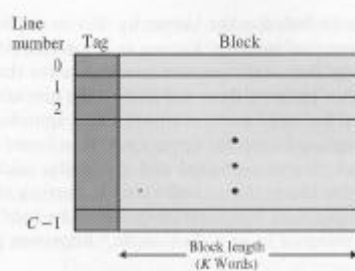
Struttura di indirizzamento della memoria



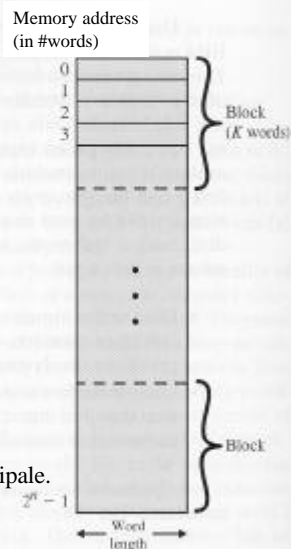
Contenuto della cache



Altezza della memoria cache: # di linee



Ampiezza della memoria cache: $K_{\text{parole}} \Rightarrow K * 4 \text{ byte in MIPS.}$



- La cache può contenere solamente una parte della memoria principale.
- Ogni parola di cache contiene K parole della memoria principale.
- Il campo Tag indica quale blocco di RAM è scritto nella corrispondente linea di cache.



Principio di funzionamento della cache



1w \$t0, 0(\$s4)

$X_n \rightarrow \$t0$

Supponiamo che la linea di cache abbia ampiezza pari ad una parola.

<i>Indirizzo Cache</i>	<i>Prima del caricamento di X_n</i>	<i>Dopo il caricamento X_n</i>
000	X_4	X_4
001	X_1	X_1
010	X_{n-2}	X_{n-2}
011		X_n
100	X_2	X_2
101	X_3	X_3
110		
111	X_{n-1}	X_{n-1}

Problemi:

- Corrispondenza tra 0(\$s4) e l'indirizzo di cache.

Come si fa a sapere se un dato è già in cache?

E se un dato è presente, come si fa a recuperare un dato nella cache?



Sommario



Caratteristiche di un sistema di memoria

Struttura di una memoria

Principio di funzionamento di una memoria cache

Cache a mappatura diretta

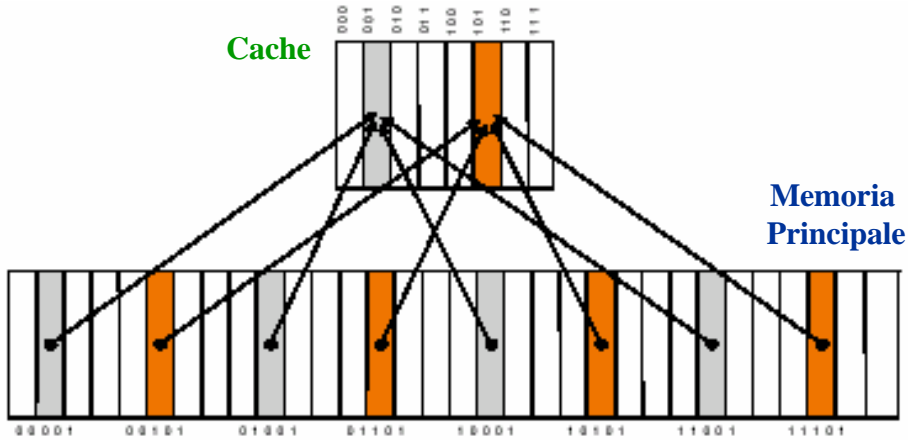
Il campo tag di una cache



Corrispondenza diretta (direct mapped)



Ad ogni indirizzo di Memoria Principale corrisponde un indirizzo di cache.



Indirizzi diversi di Memoria Principale corrispondono allo stesso indirizzo di cache.
Quali indirizzi della memoria principale si considerano?



Determinazione della legge di corrispondenza



Le **linee** di una cache indicano blocchi costituiti da **k parole**.
Ciascuna **parola** è costituita da **m byte** (MIPS: 4 byte).

Posso mettere in corrispondenza un blocco di Memoria Principale di $n = k * m$ byte con una linea di cache.

Come ottengo l'indirizzo di cache corrispondente ad un indirizzo di memoria principale?



Come si ottiene l'indirizzo di cache?



Identifico il blocco di Memoria principale a cui appartiene il byte da leggere / scrivere.
Associo il blocco di Memoria principale ad una linea di cache mediante operazione di modulo.

Posizione_byte_cache = indirizzo_Memoria principale (in #byte) *modulo* #byte_cache.

Utilizzo come unità di conteggio il blocco di cache.

Note:

L'operazione di modulo se la dimensione del blocco è multiplo della base di conteggio si ottiene **scartando** alcune delle cifre meno significative.

La divisione per un multiplo di 2 si ottiene come **shift a sx**, le cifre che "escono" rappresentano il **resto**.



Esempio



La cache con linee di ampiezza pari a 4 parole ed altezza di 8 linee:

Il blocco di dati della memoria principale che può essere contenuto in ogni linea di cache, ha dimensioni: $n = 4 \text{ parole} * 4 \text{ byte}$.

La capacità della cache sarà $C = 8 \text{ linee} * n = 8 * 16 = 128 \text{ byte}$.

• $\text{lw } \$t0, 68(\$zero)$ – 68 è il 5° byte della 5ª linea ($68 / 16 = 4$, **13ª linea della cache**).

Indirizzo_cache = Indirizzo_Memoria principale *modulo* Capacità_linea_cache
Indirizzo_cache = $68 / 16 = 4 \rightarrow$ resto 4. {Il resto indica la posizione del primo byte della parola all'interno della linea di cache \rightarrow 5° byte = 1° byte della 2ª parola.}
Il dato viene scritto (assieme al suo blocco) nella linea 5ª della cache.

• $\text{lw } \$t0, 196(\$zero)$ – $196 > 128 \text{ byte} \Rightarrow$ mappiamo il 2° blocco di RAM sulla cache.

196 è il 5° byte della 13ª linea ($196 / 16 = 12$, **13ª linea**).

Indirizzo_blocco_memoria principale \Rightarrow resto ($196 / 128$) = 68.

Indirizzo_cache = Indirizzo_memoria principale (relativo all'inizio del blocco) *modulo* capacità_linea_cache $\rightarrow 68 / 16 = 4 \rightarrow$ resto 4. {Il resto indica la posizione del primo byte della parola all'interno della linea di cache \rightarrow 5° byte.}

Il dato viene scritto (assieme al suo blocco) nella stessa linea 5ª della cache.



Indirizzamento scartando i bit più significativi



Indirizzo cache	Indirizzo decimale RAM	Indirizzo binario RAM
111	112-127, 240-255, 368-383,...	00 0111 0000 - 00 0111 1111
110	96, 224, 352	00 0110 0000 - 00 0110 1111
101	80-95, 208, 336...	00 0101 0000 - 00 0101 1111
100	64-79, 192-207, 320-335, 448, 461,...	00 0100 0000 - 00 0100 1111
011	48-63, 176, 304...	00 0011 0000 - 00 0011 1111
010	32-47, 160, 288...	00 0010 0000 - 00 0010 1111
001	16-31, 44, 272...	00 0001 0000 - 00 0001 1111
000	0-15, 128, 256, 384,...	00 0000 0000 - 00 0000 1111

NB: La capacità della cache è di: $8 * 16 \text{ byte} = 128 \text{ byte} = x000\ 0000 - x111\ 1111$

0000 0000 0000 00(00)

→

0000 0000 0111 11(11)

128 indirizzi diversi (32 parole di 4 byte)

0000 0000 0111 1100 -31				
0000 0000 0111 1000 -30				
0000 0000 0111 0100 -29				
0000 0000 0111 0000 -28				
0000 0000 0110 1100 -27				
0000 0000 0110 1000 -26				
0000 0000 0110 0100 -25				
0000 0000 0110 0000 -24				
0000 0000 0101 1100 -23				
0000 0000 0101 1000 -22				
0000 0000 0101 0100 -21				
0000 0000 0101 0000 -20				
0000 0000 0100 1100 -19				
0000 0000 0100 1000 -18				
0000 0000 0100 0100 -17				
0000 0000 0100 0000 -16				
0000 0000 0011 1100 -15				
0000 0000 0011 1000 -14				
0000 0000 0011 0100 -13				
0000 0000 0011 0000 -12				
0000 0000 0010 1100 -11				
0000 0000 0010 1000 -10				
0000 0000 0010 0100 -9				
0000 0000 0010 0000 -8				
0000 0000 0001 1100 -7				
0000 0000 0001 1000 -6				
0000 0000 0001 0100 -5				
0000 0000 0001 0000 -4				
0000 0000 0000 1100 -3				
0000 0000 0000 1000 -2				
0000 0000 0000 0100 -1				
0000 0000 0000 0000 -0	Byte 0	Byte 1	Byte 2	Byte 3

Indirizzamento di una MP



Parsing dell'indirizzo



0000 0000 0000 00(00)

0000 0000 0111 11(11)

128 indirizzi diversi (32 parole di 4 byte)

- Questi 128 byte vengono messi in corrispondenza con i 128 byte di una cache che è organizzata come 8 x 16: 4 parole per linea.
- Gli ultimi 2 bit a destra indicano i byte interni alla parola. Non vengono considerati nei trasferimenti da / per la memoria. Viene considerata la *Word*.
- I 2 bit seguenti indicano la posizione della parola ricercata all'interno della linea della cache.
- I 3 bit seguenti indicano la linea della cache nella / dalla quale si scrive / legge.
- I bit seguenti indicano il numero del blocco della memoria principale messo in corrispondenza con la cache.



Metodo di parsing dell'indirizzo



0000 0000 0000 00(00)

0000 0000 0111 11(11)

128 indirizzi diversi (32 parole di 4 byte)

Questo vale per:

- Cache di 8 linee.
- Linee di 4 parole.
- Parole di 4 byte.

Prendiamo un indirizzo ed operiamo per divisioni successive nel caso più generale:

Capacità della cache

$$M / [(\#byte / parola) * (\#parole / linea) * \#linee] = \text{Numero_blocco di RAM.}$$

Il resto, R_1 , rappresenta l'offset in byte all'interno della cache a partire dal byte 0.

$$R_1 / [(\#byte / parola) * (\#parole / linea)] = \text{Numero di linee di cache.}$$

Il resto, R_2 , rappresenta l'offset in byte all'interno della linea di cache.

$$R_2 / [(\#byte / parola)] = \text{Numero di parola all'interno della linea di cache.}$$

Il resto, R_3 , rappresenta l'offset in byte all'interno della parola.



Esempio di parsing dell'indirizzo



0000 0000 0000 00(00) → 0000 0000 0111 11(11) 128 indirizzi diversi (32 parole di 4 byte)

La cache con linee di 4 parole (ampiezza) ed altezza di 8 linee:

Il blocco di dati contenuto in ogni linea di cache è di dimensioni: $n = 4 * 4 \text{ byte} = 16 \text{ byte}$.

La capacità della cache è di $8 * 16 \text{ byte} = 128 \text{ byte}$.

`lw $t0, 196($zero)`

$196 / [4 * 4 * 8] = 1$ (2° blocco di RAM) con resto $R_1 = 196 - 1 * 128 = 68$.

Il resto, R_1 , rappresenta l'offset in byte all'interno della cache.

$68 / [4 * 4] = 4$ (5ª linea della cache) con resto $R_2 = 68 - 4 * 16 = 4$.

Il resto, R_2 , rappresenta l'offset in byte all'interno della linea di cache.

$4 / 4 = 1$ (2ª parola della cache) con resto $R_3 = 4 - 1 * 4 = 0$.

Il resto, R_3 , rappresenta l'offset in byte all'interno della parola.



Esempio di parsing dell'indirizzo



0000 0000 0000 00(00) → 0000 0001 1111 11(11) 512 indirizzi diversi (128 parole, 4 byte)

La cache con linee di 8 parole (ampiezza) ed altezza di 16 linee:

Il blocco di dati contenuto in ogni linea di cache è di dimensioni: $n = 8 * 4 \text{ byte} = 32 \text{ byte}$.

La capacità della cache è di $16 * 32 \text{ byte} = 512 \text{ byte}$.

`lw $t0, 600($zero)`

$600 / [16 * 8 * 4] = 1$ (2° blocco di RAM) con resto $R_1 = 600 - 1 * 512 = 88$.

Il resto, R_1 , rappresenta l'offset in byte all'interno della cache.

$88 / [8 * 4] = 2$ (3ª linea della cache) con resto $R_2 = 88 - 2 * 32 = 24$.

Il resto, R_2 , rappresenta l'offset in byte all'interno della linea di cache.

$24 / 4 = 6$ (7ª parola della linea di cache) con resto $R_3 = 24 - 6 * 4 = 0$.

Il resto, R_3 , rappresenta l'offset in byte all'interno della parola.



Sommario



Caratteristiche di un sistema di memoria

Struttura di una memoria

Principio di funzionamento di una memoria cache

Cache a mappatura diretta

Il campo tag di una cache



Come si può sapere se un dato è presente in cache?



Aggiungiamo a ciascuna delle linee della cache un campo **tag**.

Il tag contiene i bit che costituiscono la parte più significativa dell'indirizzo. E' costituito da K bit.

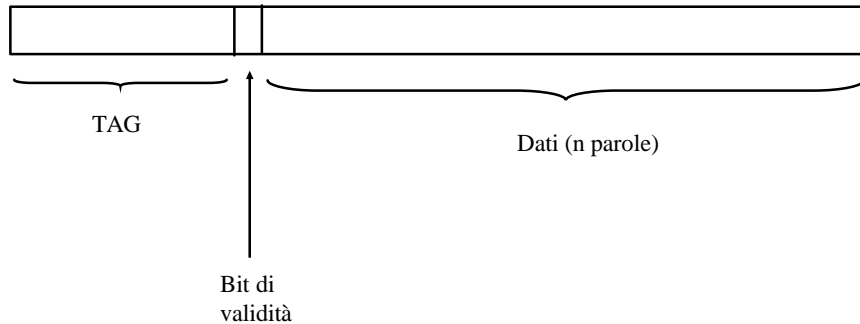
$$K = M - \text{sup}(\log_2 \text{Capacità})$$

Nell'esempio precedente: $K = 32 - \text{sup}(\log_2 128) = 25$ bit.

Occorre inoltre l'informazione data_valid / data_not_valid: **bit di validità**.



Struttura di una linea di cache



Quanti blocchi di RAM avremo delle dimensioni della cache?

Nel caso precedente, avremo linee di cache di lunghezza: $25 + 1 + 4 * 4 * 8 = 134\text{bit}$.



Come leggere / scrivere su cache



Individuare la linea della cache dalla quale leggere / scrivere (operazione analoga all'indirizzamento del register file).

Confrontare il campo tag con il blocco RAM in cui risiede il dato.

Controllare il bit di validità.

Leggere (scrivere) il dato.

Per blocchi più ampi di una parola, occorre individuare una parola tra le k presenti nella linea di cache.



Sommario



Caratteristiche di un sistema di memoria

Struttura di una memoria

Principio di funzionamento di una memoria cache

Cache a mappatura diretta

Il campo tag di una cache