



Le memorie Cache

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

Università degli Studi di Milano



Sommario

Circuito di lettura / scrittura di una cache a mappatura diretta

Memorie associative

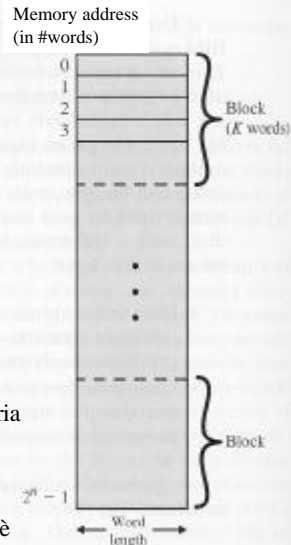
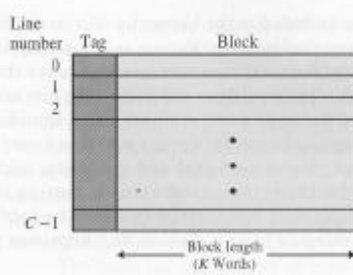
Memorie n-associative



Mappatura diretta di una 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 Memoria Principale è scritto nella corrispondente linea di cache.



Come leggere / scrivere su cache



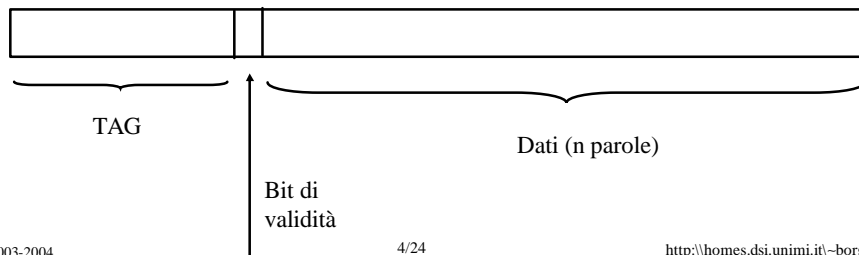
Individuare la linea della cache dalla quale leggere / scrivere (operazione analoga all'indirizzamento del register file).

Confrontare il campo tag con il blocco di Memoria Principale in cui risiede il dato.

Controllare il bit di validità.

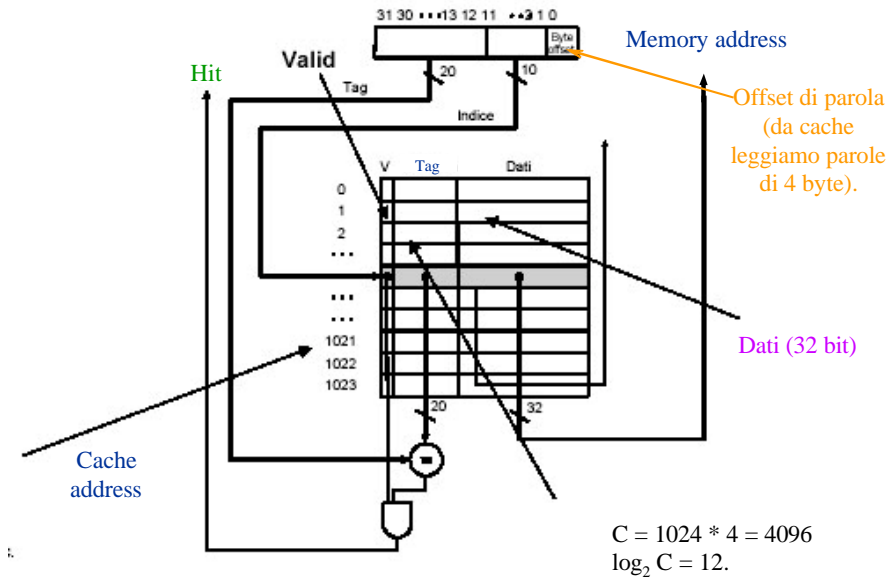
Leggere (scrivere) il dato.

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

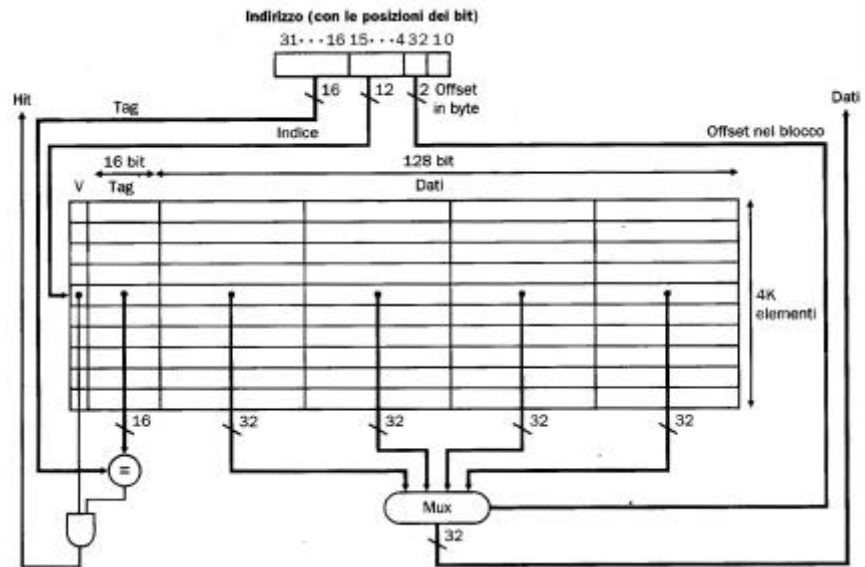




Circuito di lettura/scrittura cache



Cache a blocchi, funzionamento





Esercizi



Sia data una cache a corrispondenza diretta contenente 64Kbyte di dati e avente blocchi di 1 parola. Assumendo che gli indirizzi siano di 32 bit quale è il numero totale di bit richiesto per l'implementazione della cache?

Supponendo che il MIPS abbia una cache di 512byte, indicare cosa succede nei campi della cache quando vengono eseguite le seguenti istruzioni:

lw \$t1, 0x0000(\$t0) \$t0 = 1kbyte = 1,024 byte

lw \$t1, 0x0000(\$t0) \$t0 = 0

lw \$t1, 0x0202(\$t0) \$t0 = 1kbyte = 1,024 byte

lw \$t1, 0x0001(\$t0) \$t0 = 0

lw \$t1, 0x0201(\$t0) \$t0 = 1kbyte = 1,024 byte



Sommario



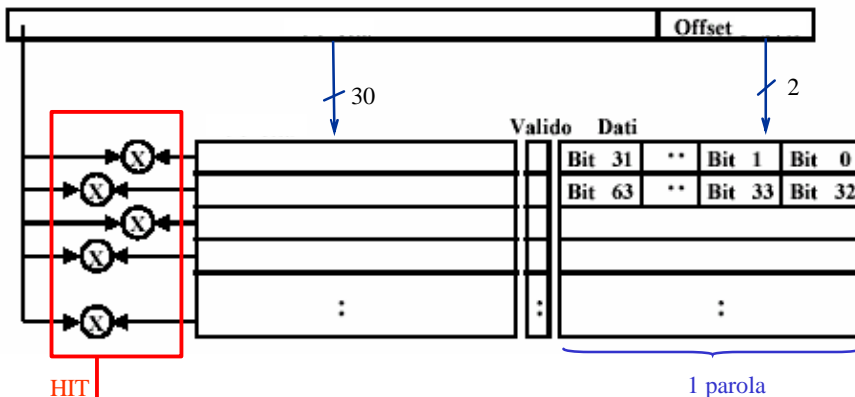
Circuito di lettura / scrittura di una cache a mappatura diretta

Memorie associative

Memorie n-associative



Memorie associative



Consentono di caricare un blocco di Memoria Principale in una qualsiasi linea di cache.

E' una memoria completamente associativa.

Tramite comparatori individuo in quale blocco si trova il mio dato.

Il segnale di Hit si genera come AND (comparatore_output, Valido)

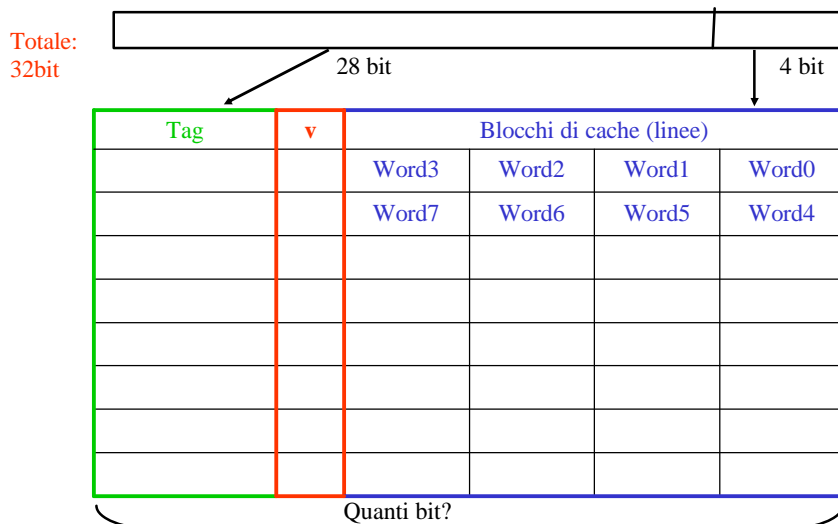
Dove scrivo il blocco?



Accesso alle memorie associative



Posso accedere alla memoria attraverso l'indirizzo completo modulo la dimensione del blocco di cache (lunghezza della linea di cache).





Tassonomia



Spazio di indirizzamento: $(s + w)$ bit: somma della dimensione del campo tag + somma della dimensione dell'offset all'interno della parola. Spazio misurato in word o byte (come nel caso del MIPS).

Numero di unità indirizzabili: $2^{(s+w)}$ unità ($2^{(s+w)}$ byte in MIPS).

Dimensione del blocco = dimensione della linea di cache = 2^w parole o byte.

Numero totale di blocchi della memoria principale: 2^s .

Dimensioni del campo tag: s bit.

Viene aumentato il numero di Hit ma con un appesantimento notevole della circuiteria.



Sommario

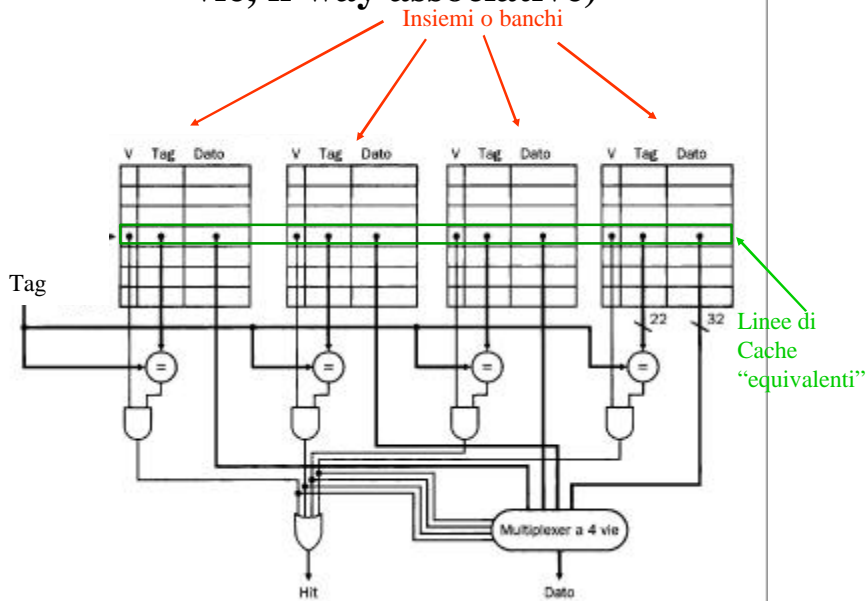


Circuito di lettura / scrittura di una cache a mappatura diretta

Memorie associative

Memorie n-associative

Memorie n-associative (o associative a n-vie, n-way associative)



A.A. 2003-2004

13/24

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

Memorie n-associative

n-associative o set associative o a n vie.

La memoria è suddivisa in n insiemi, o banchi, ciascuno di k linee, posti in parallelo.

Blocco (linea di cache): #parole (byte) lette/scritte contemporaneamente in cache, "parola" della cache.

Insieme (banco): cache elementare.

Cache: è l'insieme dei banchi più i circuiti che li gestiscono.

Capacità della cache: #parole = #Insiemi * (#blocchi / insieme) * (#parole / blocco).

La corrispondenza tra Memoria Principale e linea di un banco è a mappatura diretta.
La corrispondenza tra Memoria Principale e banco è associativa.

Per cercare un dato non devo più analizzare tutte le linee di una cache, ma un'unica linea per ogni banco.

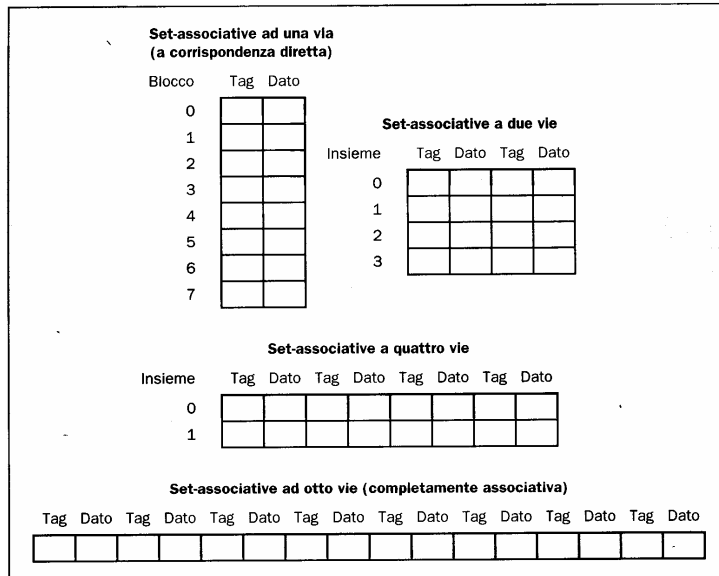
A.A. 2003-2004

14/24

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



Dalle cache a mappatura diretta alle cache associative



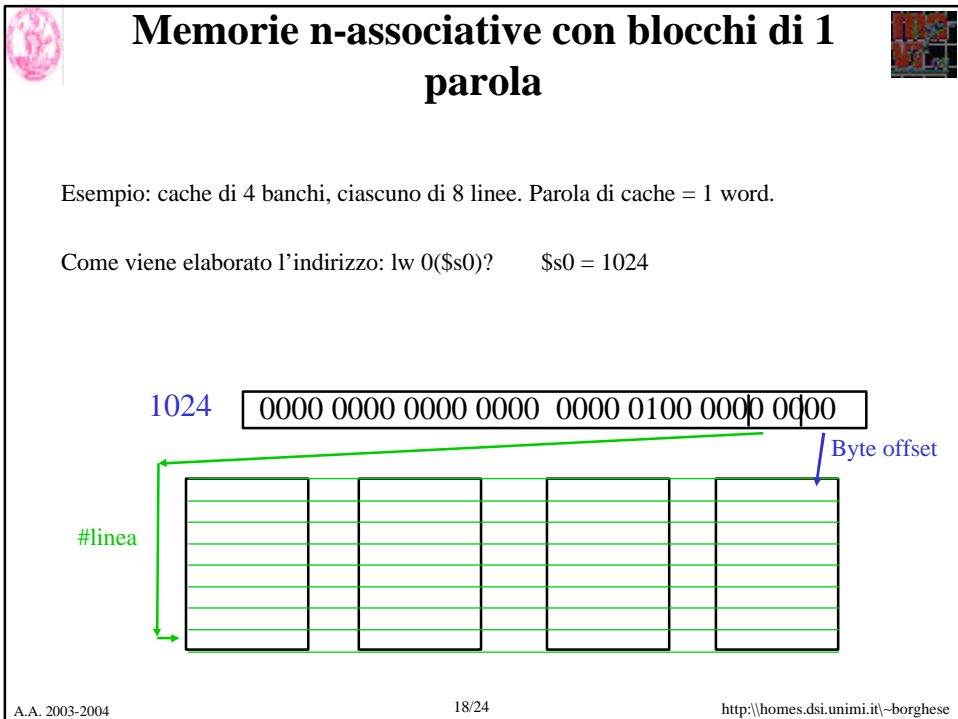
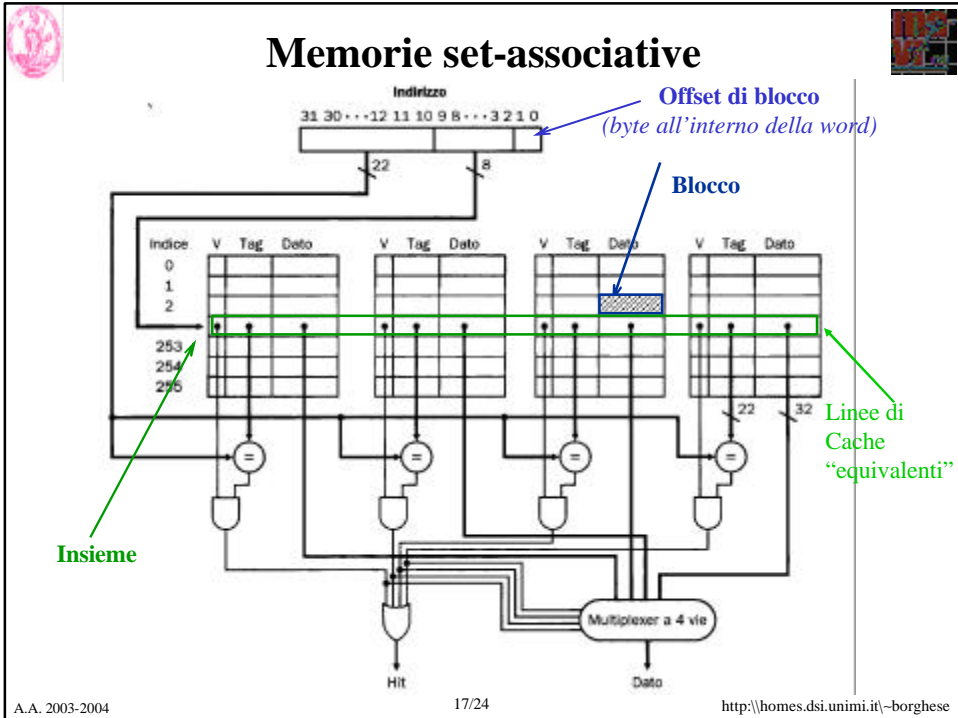
Accesso a cache ad n-way



INDICE. Se la parola richiesta è memorizzata in cache, si trova in una particolare linea di uno dei banchi. Questa linea è individuata dall'indice. L'indice è costituito da k bit, dove $k = \log_2(\#linee)$. E' analogo al numero di linea nelle cache a mappatura diretta.

TAG – contiene il blocco della RAM a cui appartiene il dato. Cerca il tag di Memoria Principale all'interno dei TAG associati alla linea individuata in ciascun banco.

L'insieme dei segnali di HIT pilotano anche il MUX che trasferiscono in uscita il contenuto del banco opportuno della cache.





Criteri di sostituzione di un blocco



Dove inserisco il blocco letto dalla memoria?

Soluzione hardware, algoritmo semplice.

LRU – Least recently Used. Viene associato ad ogni blocco un bit di USE.
Efficiente per memorie a 2 vie.

FIFO – Implementazione tramite buffer circolare.

LFU – Least frequently Used. Associa un contatore ad ogni blocco di cache.

RANDOM – Non funziona molto peggio!!



Dove si può posizionare un blocco di RAM in cache?



Corrispondenza diretta: in un'unica posizione.

Memoria ad 1 via.

#posizioni = #linee.

Completamente associative: in n posizioni (n banchi).

Ciascun banco è costituito da 1 linea.

n insiemi o banchi.

N-associative: in m posizioni (m grado di associatività).

Ho m insiemi (banchi)

Ciascun insieme è costituito da n linee.



Come si trova un blocco di RAM in cache?



Corrispondenza diretta: indicizzazione.

Controllo del tag del blocco (1 comparazione).

Associativa: ricerca in tutti gli elementi della cache.

n comparazioni: controllo di tutti i tag.

La memoria virtuale è di questo tipo (tramite la *Page Table*).

N-associativa: ricerca negli m insiemi,

m comparazioni.



Sommario



Circuito di lettura / scrittura di una cache a mappatura diretta

Memorie associative

Memorie n-associative