



Le memorie Cache associative

Prof. Alberto Borghese
Dipartimento di Scienze dell'Informazione
alberto.borghese@unimi.it

Università degli Studi di Milano

Riferimento Patterson: 5.2, 5.3



Sommario

Memorie associative

Memorie n-associative

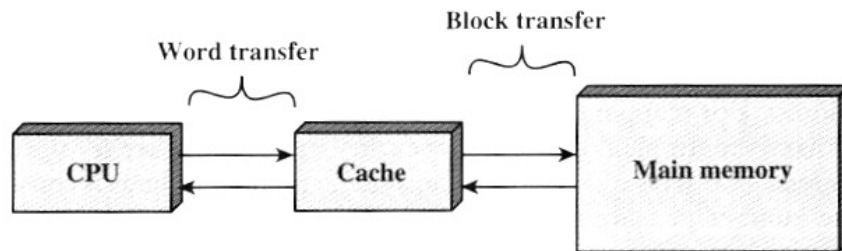


Principio di funzionamento di una cache



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.

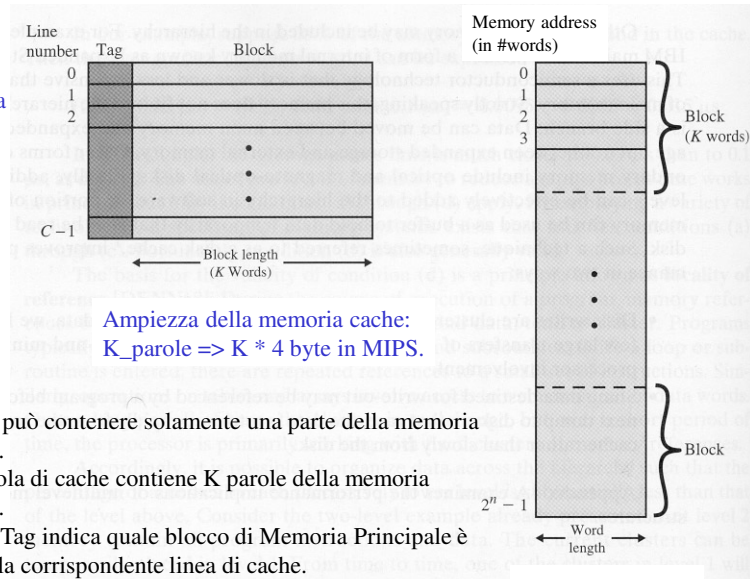
La cache contiene una copia di parte del contenuto della memoria principale. Di che cosa?



Mappatura diretta di una cache



Altezza della memoria cache: # di linee



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

- 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.



Problemi con le cache a mappatura diretta

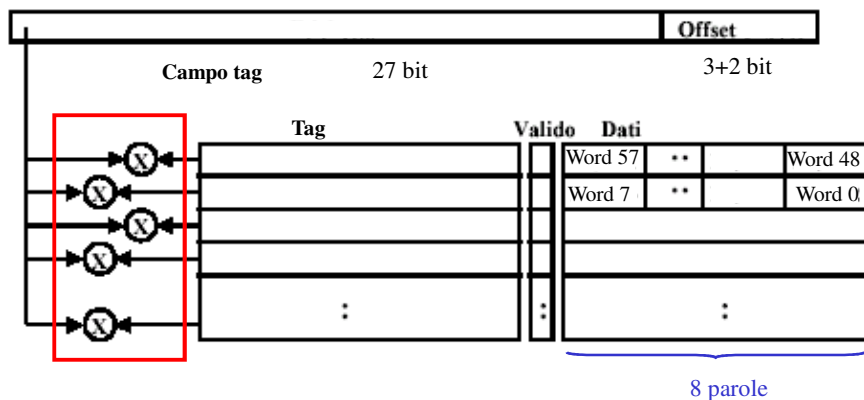


- Riempimento non ottimale (a macchia di leopardo).
- MISS per accesso alla stessa linea di cache con dati appartenenti a blocchi diversi di RAM
- Memoria associativa: il contenuto viene recuperato fornendo degli elementi associati al contenuto (e.g. ricerca di testo, ricerca attraverso ontologie WEB).
- Nelle memorie associative si utilizza una parte dell'indirizzo per recuperare il dato.

Occorre quindi sostituire il meccanismo di accesso diretto tramite numero di linea, mediante un meccanismo associativo che determini il numero della linea.



Memorie associative



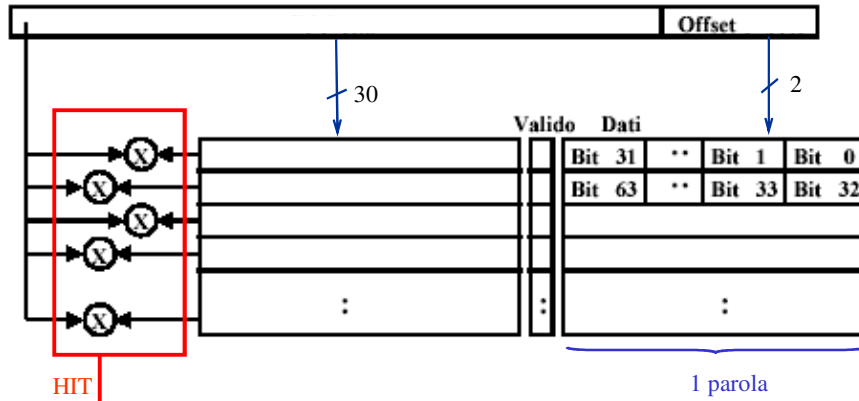
Consentono di caricare un blocco di Memoria Principale in una qualsiasi linea di cache. E' una memoria completamente associativa.

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?



Memorie associative



HIT

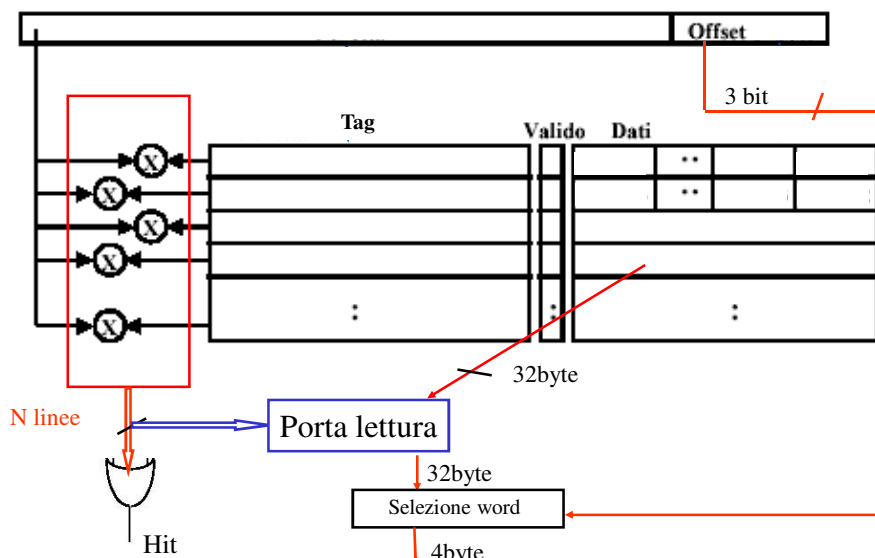
Consentono di caricare un blocco di Memoria Principale in una qualsiasi linea di cache.
 E' una memoria completamente associativa.

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?



Lettura di una memoria associativa



N linee

Porta lettura

Selezione word

Hit



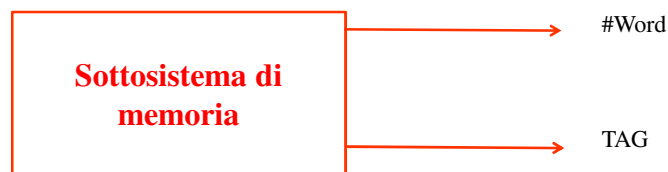
Alcuni dettagli



- L'uscita del comparatore di una linea va in AND con il bit di validità di quella linea. Il segnale diventa quindi = 1 quando il dato è presente (stesso TAG) ed è valido.
- Le uscite dagli N comparatori, ciascuno associato ad una linea diversa, possono avere al massimo un "1".
- Se colleghiamo le uscite degli N comparatori ad un encoder, otteniamo il #Linea contenente il dato che è il segnale di selezione che cercavamo.



Letture del dato



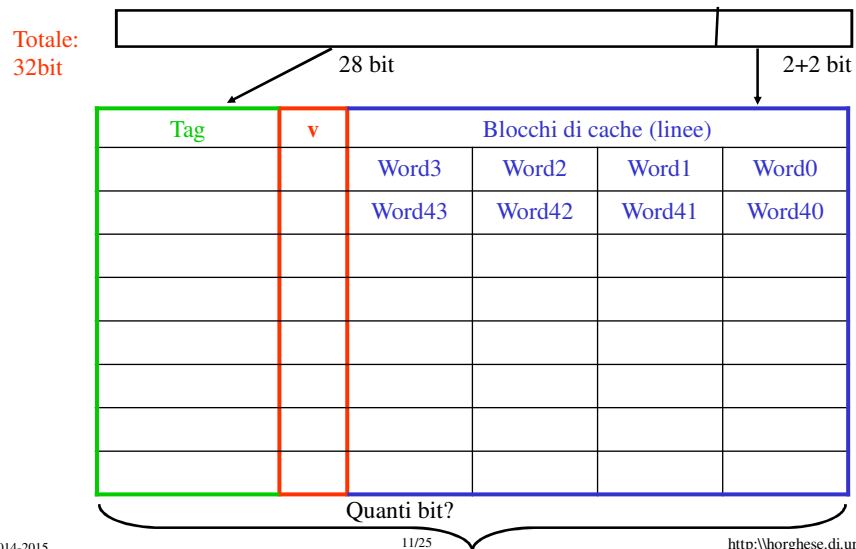
Parte di selezione del dato: #Word – La parola è cercata in tutte le linee

Parte di controllo: TAG



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 indirzzamento: $(s + w)$ 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 indirzzabili: $2^{(s+w)}$ unità ($2^{(s+w)}$ byte in MIPS).

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

Numero totale di macro-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

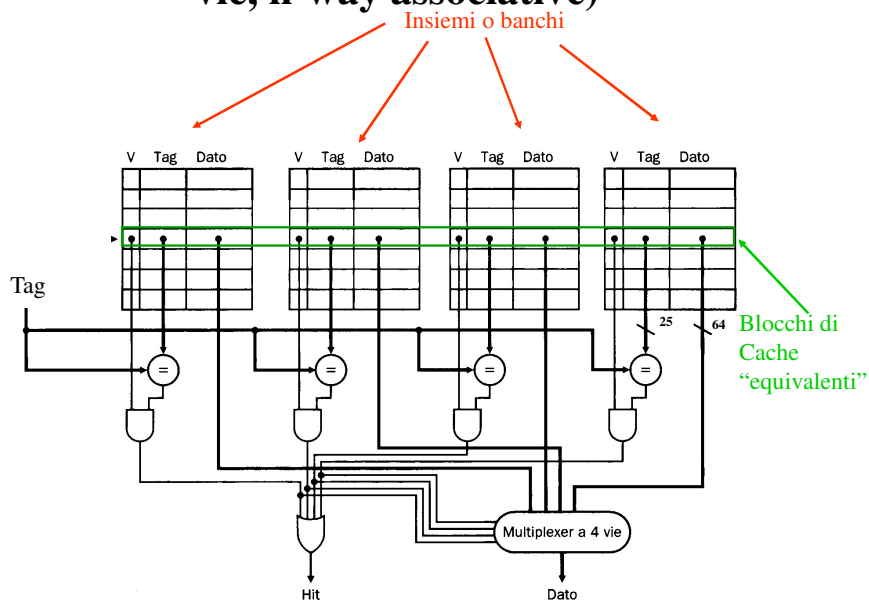


Memorie associative

Memorie n-associative



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





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.



Dalle cache a mappatura diretta alle cache associative



Set-associative ad una via (a corrispondenza diretta)

Blocco	Tag	Dato
0		
1		
2		
3		
4		
5		
6		
7		

Set-associative a due vie

Insieme	Tag	Dato	Tag	Dato
0				
1				
2				
3				

Set-associative a quattro vie

Insieme	Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato
0								
1								

Set-associative ad otto vie (completamente associativa)

Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato



Accesso a cache ad n-vie

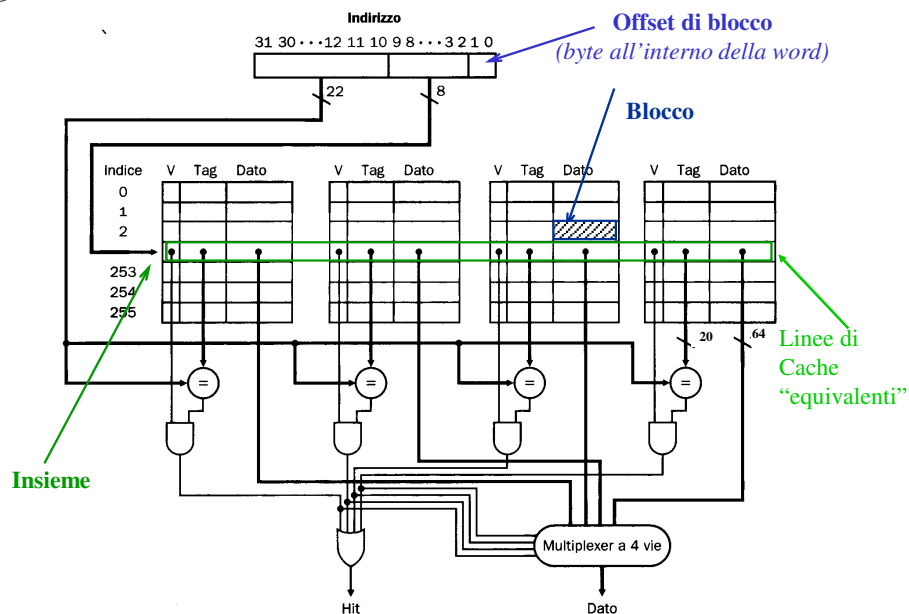
INDICE. Se la parola richiesta è memorizzata in cache, si trova in una particolare linea di uno dei banchi. Questa linea è individuata dall'indice. L'indice è costituito da k bit, dove $k = \log_2(\#linee)$. E' analogo al numero di linea nelle cache a mappatura diretta.

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.



Memorie set-associative



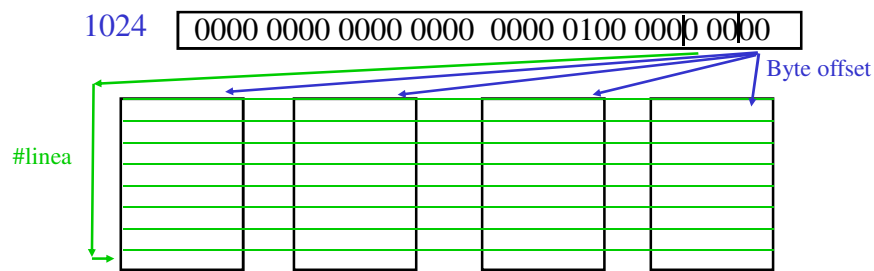


Memorie n-associative con blocchi di 1 parola



Esempio: cache di 4 banche, ciascuno di 8 linee. Parola di cache = 1 word, non c'è offset nel blocco.

Come viene elaborato l'indirizzo: lw 0(\$s0)? \$s0 = 1024

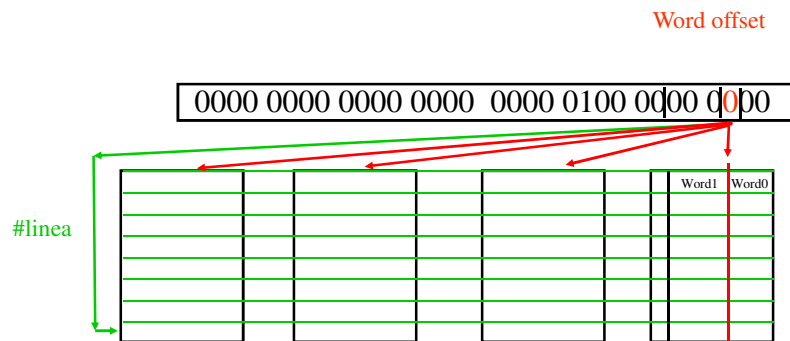


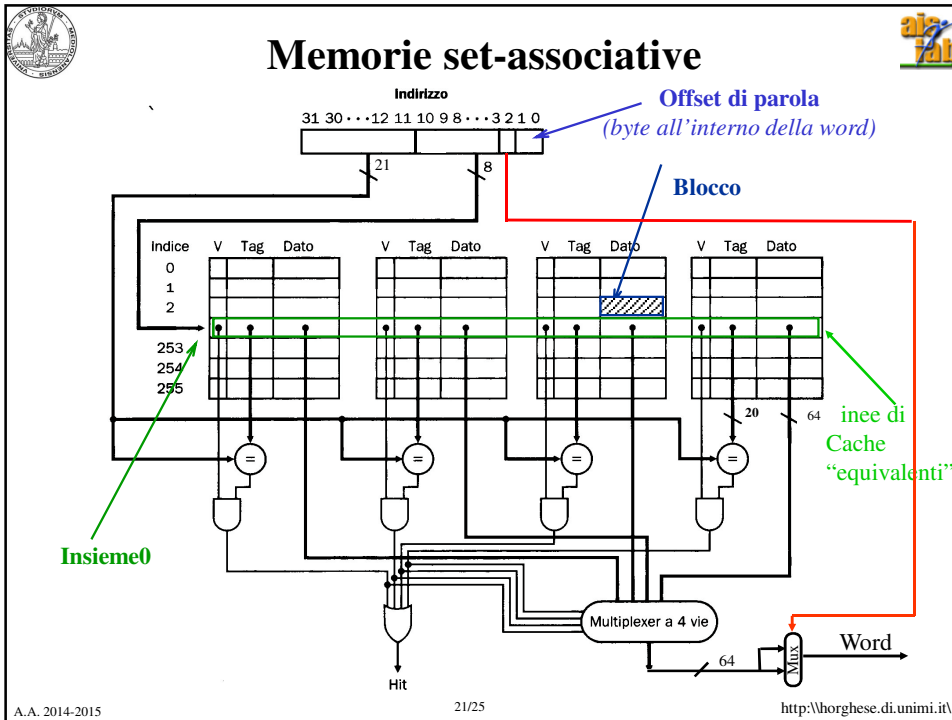
Memorie n-associative con blocchi di 2 parole



Esempio: cache di 4 banche, ciascuno di 8 linee. Parola di cache = 2 word.

Come viene elaborato l'indirizzo: lw 0(\$s0)? \$s0 = 1024





Criteri di sostituzione di un blocco

Dove inserisco il blocco letto dalla RAM?

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!!

22/25

http://whorghese.di.unimi.it/



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



Memorie associative

Memorie n-associative