



La struttura delle memorie cache

Prof. Alberto Borghese
Dipartimento di Informatica
alberto.borghese@unimi.it

Università degli Studi di Milano

Riferimento Patterson: 5.2, 5.4, 5.9, B.9



Sommario

Memory miss

SRAM

DRAM

Trasferimento dati



Gerarchia di memorie

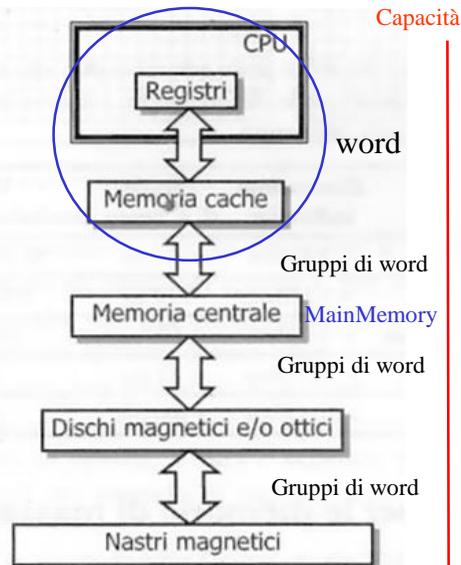


Livelli multipli di memorie con diverse dimensioni e velocità.

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

Ciascun livello vede il livello inferiore e viceversa.

Cache (memoria nascosta)



Gestione dei fallimenti di una cache



La gestione avviene tra CPU e MMU.

Hit – è quello che vorremmo ottenere, il funzionamento della CPU non viene alterato.

Miss – **in lettura** devo aspettare che il dato sia stato caricato nella linea di cache e sia pronto -> eccezione particolare della CPU che crea uno **stallo** della CPU e non un flush.

Nelle CPU super-scalari si sfrutta l'esecuzione fuori ordine per nascondere questa latenza.

Passi da eseguire in caso di Miss (miss penalty):

- 1) Bloccare tutte le istruzioni nella pipeline (blocco dei registri di pipeline per uno o più cicli di clock)
- 2) Scaricare in MM la linea di cache interessata (micro-blocco).
- 3) Richiedere che la MM legga il micro-blocco contenente il dato da leggere e lo porti fuori nel MDR.
- 4) Trasferire il micro-blocco in cache, aggiornare i campi validità e tag.
- 5) Riavviare l'esecuzione della pipeline.

NB Il programma non può continuare!!



Controllo mediante FSM



STATI del controllore (S):

- “Idle”: non ci sono richieste alla cache
- “Compare”: identificazione di Hit / Miss
- Scaricamento della linea da cache
- Trasferimento della linea nella MM

$\langle S, I, Y, f(S,D), g(S), S_0 \rangle$

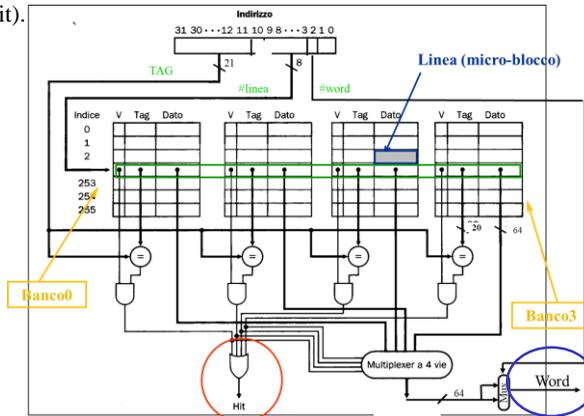
Meccanismo per decidere se una linea va scaricata: **dirty bit** = 1 se la linea è stata modificata dalla CPU (bit di validità + dirty bit).

INPUT al controllore (I):

- Read/Write
- DirtyBit
- MM Ready

OUTPUT del controllore (Y):

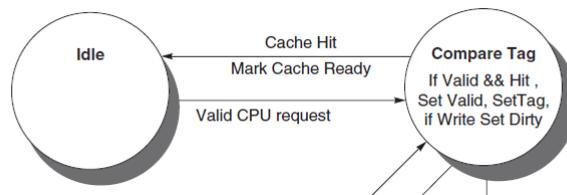
- Hit/miss (il dato presente in uscita è quello richiesto dalla CPU).



A.A. 2020-2021



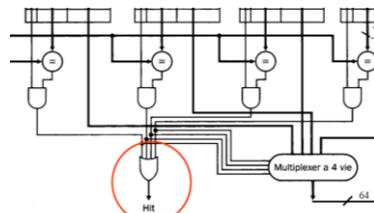
Gestione di una Hit



A seguito di una richiesta di lettura/scrittura, il controllore della cache si sposta nello stato in cui viene controllato se il dato è presente in cache ed è valido.

Se la condizione è verificata si ha una hit -> il dato in uscita dalla cache è quello richiesto dalla CPU.

Il controllore della cache ritorna nello stato idle.



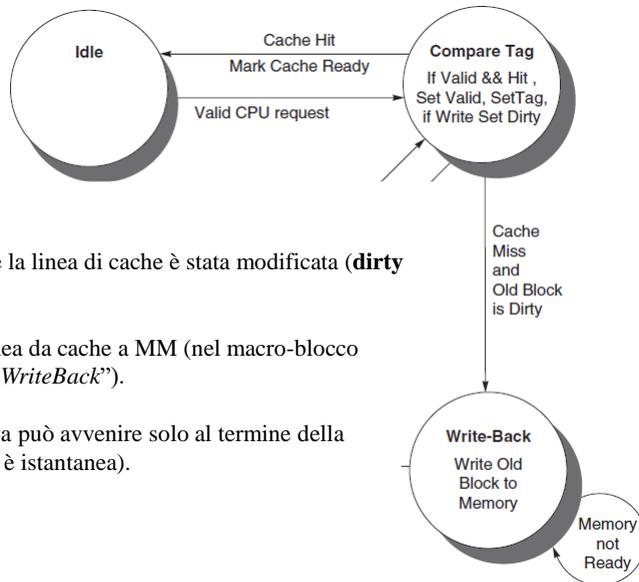
A.A. 2020-2021

6/52

<http://horghese.di.unimi.it/>



Scaricamento di una linea di cache in MM



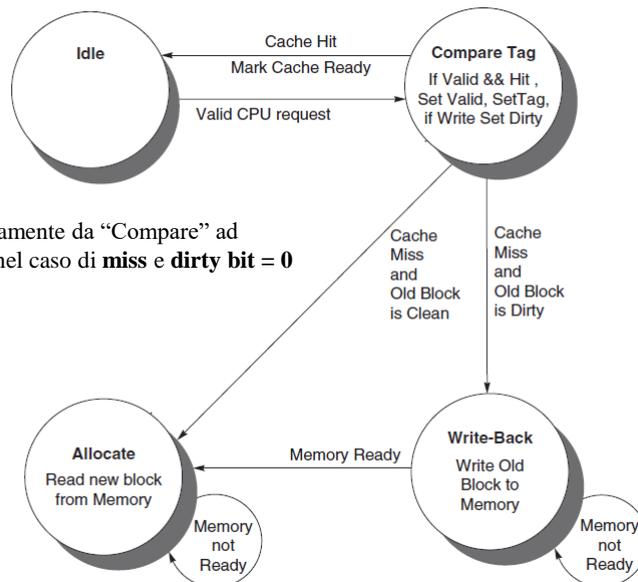
Se si verifica una **miss** e la linea di cache è stata modificata (**dirty bit = 1**),

Va prima trasferita la linea da cache a MM (nel macro-blocco identificato dal TAG – “*WriteBack*”).

La transizione successiva può avvenire solo al termine della scrittura della MM (non è istantanea).



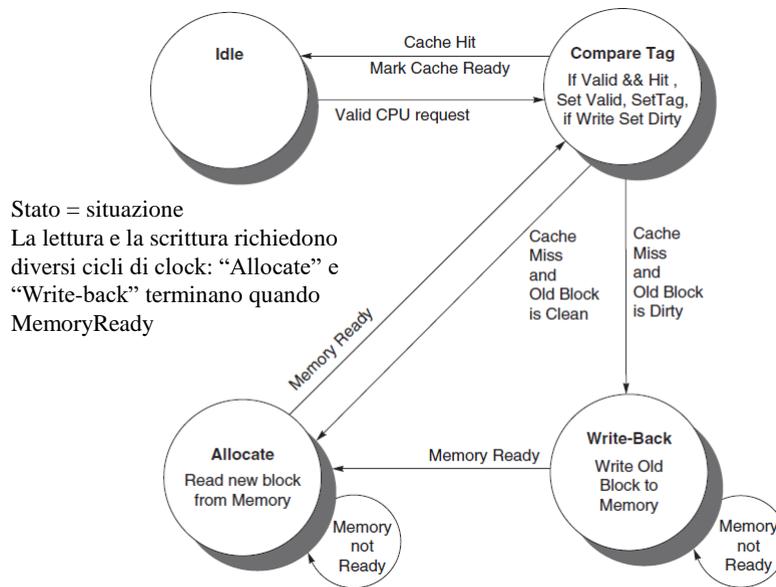
Letture di una linea di cache dalla MM



Passo direttamente da “Compare” ad “Allocate” nel caso di **miss** e **dirty bit = 0**



Controllore della cache



Sommario

Memory miss

SRAM

DRAM

Trasferimento dati



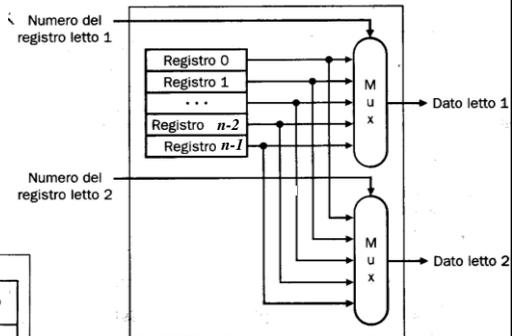
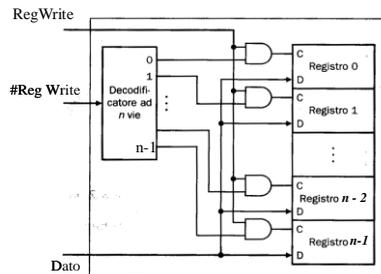
Register file



Il tempo di lettura dipende dal cammino critico dei Mux.

Il tempo di scrittura dipende dal cammino critico del Decoder.

Numero_registro = selettore.



Selezione - #registro

Letture - sempre disponibile in uscita (dopo tempo di commutazione del MUX)

Scrittura - segnale esplicito (in AND con il clock in caso di memoria sincrona).



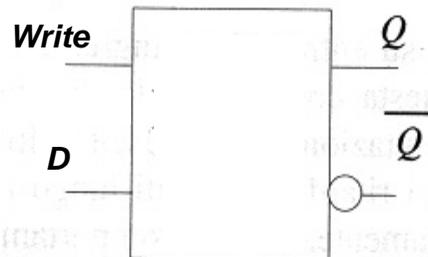
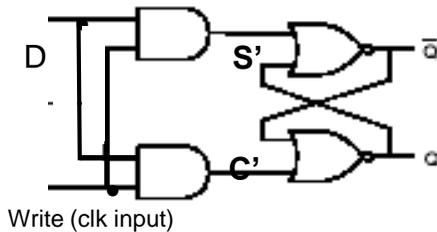
Cella SRAM



E' trasparente quando Write = 1

Se Write = 1 $Q_{t+1} = D$

Se Write = 0 $Q_{t+1} = Q_t$



Letture - sempre disponibile in uscita

Scrittura - segnale esplicito («apertura porta di accesso al latch»)



SRAM oggi



Static RAM. Bistabili (4-6 **transistor** per cella, ottimizzazione della cella rispetto al latch, cf. Register File).

Tempo di accesso uguale per ogni dato.

Informazione stabile (non ci sono disturbi da parte di quello che succede intorno)

Non c'è bisogno di rinfrescare il contenuto della memoria (refresh)

Sono memorie volatili (dipendono dall'alimentazione)

Con la tecnologia CMOS, consumano energia solo quando commutano..
Poca energia viene consumata per mantenere il dato (**standby**).

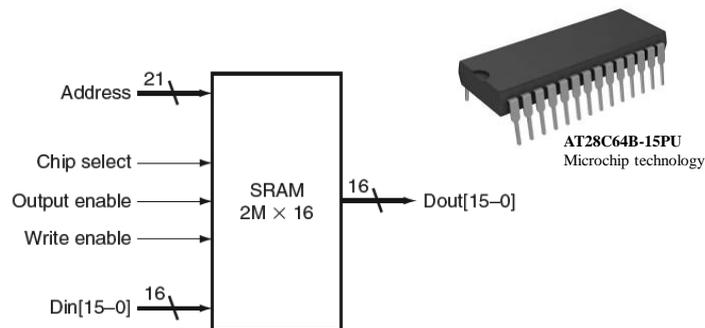
1 porta di lettura / scrittura (o leggo o scrivo o non faccio nulla: 2 segnali)

Cache -> inserita nella CPU -> Produzione delle SRAM è principalmente da parte dei produttori di CPU.

Memorie veloci per dispositivi embedded.



Un chip di SRAM



Altezza (2 Mbit) x Ampiezza (16 bit) – Oggi 1-4 bit

Altezza (2 Mbit) → 21 bit di indirizzamento

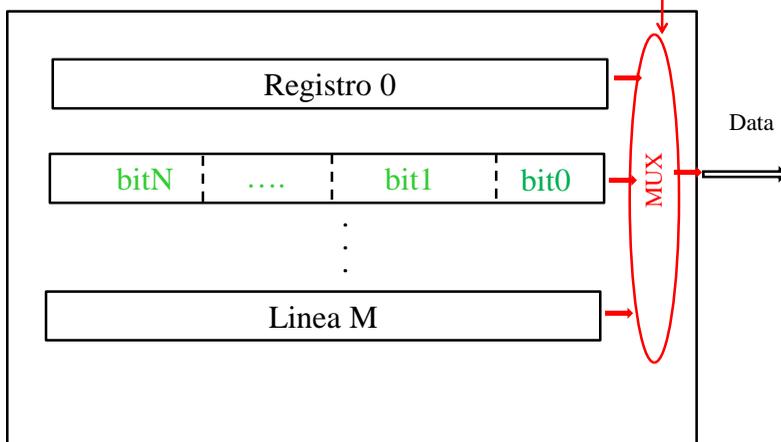
Ampiezza (16 bit) → Bus dati su 16 bit.

Chip select -> Abilitazione del chip (per risparmiare energia).

Write enable (abilitazione scrittura), Output enable (abilitazione lettura!).



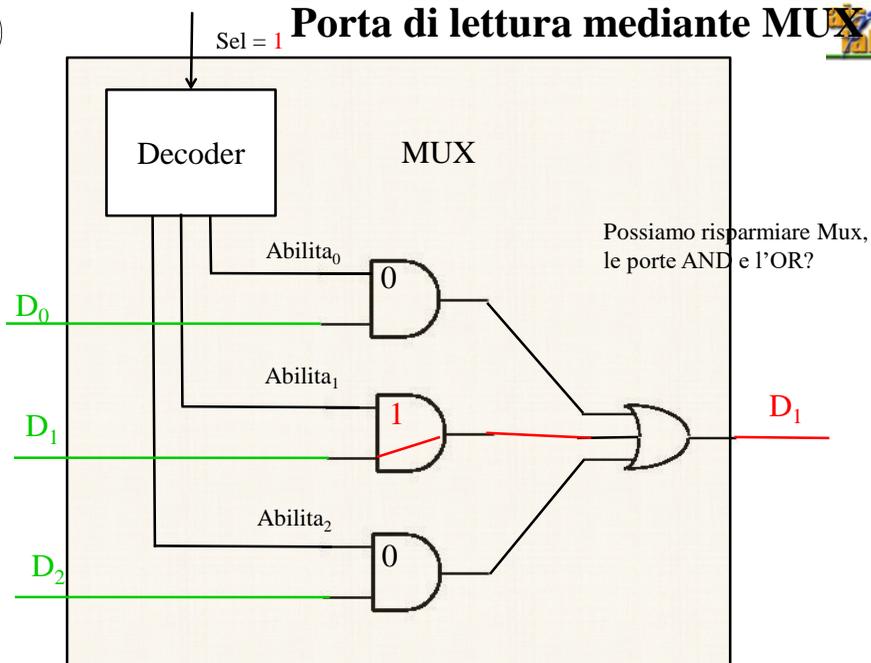
Lettura da una linea di memoria



Selezione uno dei registri = porto in uscita l'uscita di tutti i bit del registro selezionato.
Il mux è pratico per 32 registri, ma non per migliaia di linee come nelle cache L1-L3.



Porta di lettura mediante MUX



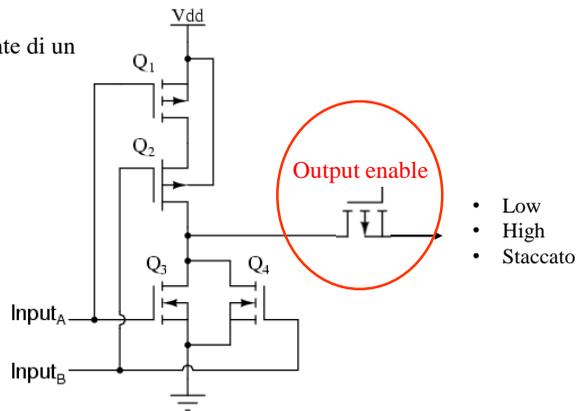


Porte three-state

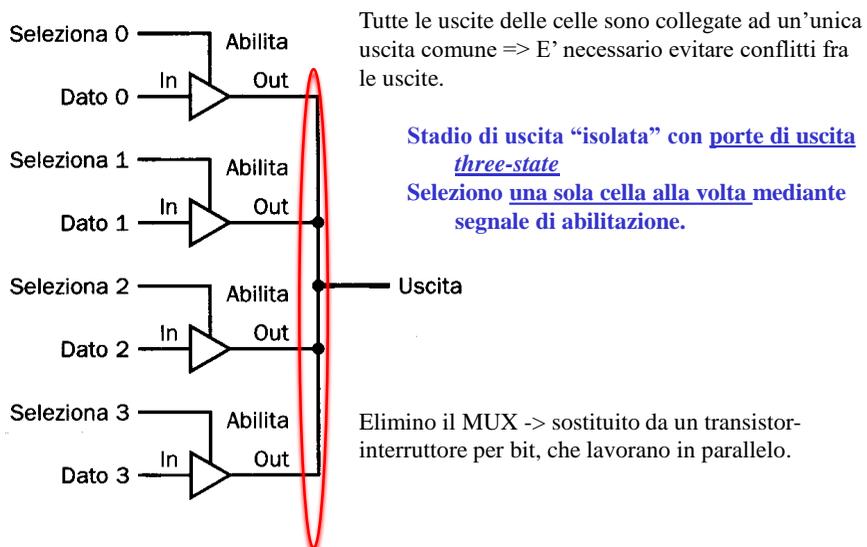


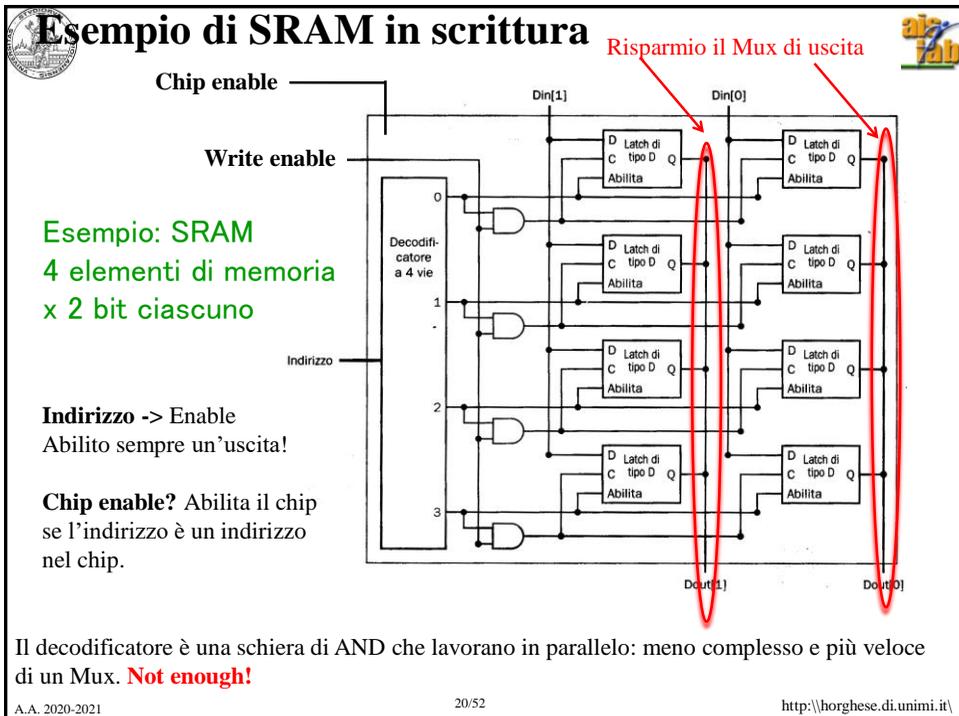
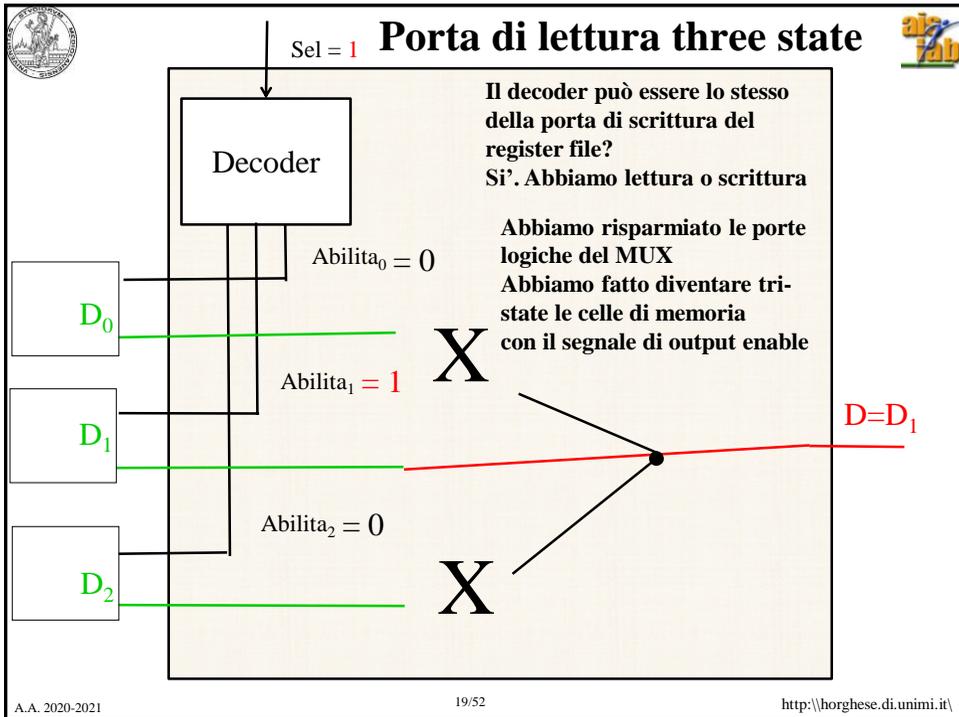
- **Stato 0** – spento – connesso a massa (0V)
- **Stato 1** – acceso – connesso all'alimentazione (+Vcc)
- **Stato disabilitato** – disconnesso – non connesso a nulla, tensione floating (alta impedenza). Interruttore di accensione del circuito.

I circuiti a valle vedono l'equivalente di un "filo tagliato".



Memoria three-state (soluzione HW)







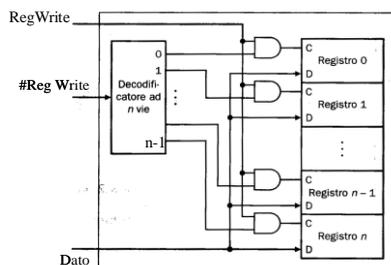
SRAM a matrice



Una **SRAM 4M x 8** avrebbe bisogno di un decoder a 22 bit -> 4M linee, con 4M porte AND ciascuna con 22 bit in ingresso.

C'è un limite elettrico al numero di linee che si possono collegare assieme.

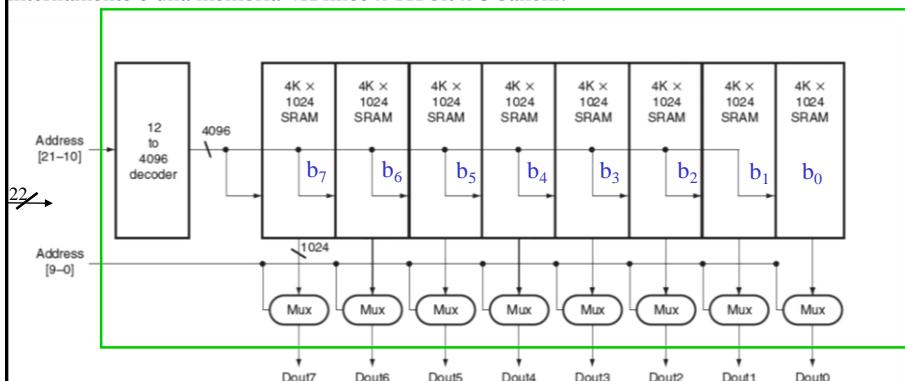
È più conveniente costruire una matrice e separare la lettura delle linee dalla lettura delle colonne (estrazione di una linea «lunga» dalla memoria – cf. Memoria cache).



SRAM a matrice



SRAM 4M x 8. Ad ogni bit assegno un chip (banco) di 4K linee x 1024 bit per linea (4M bit => internamente è una memoria 4K linee x 1K bit x 8 banche).



Il decodificatore sarà a 12 bit ($\log_2 4K$) per selezionare una delle 4096 linee (ciascuna di 1024 bit). Ciascuna linea fornisce 1024 bit.

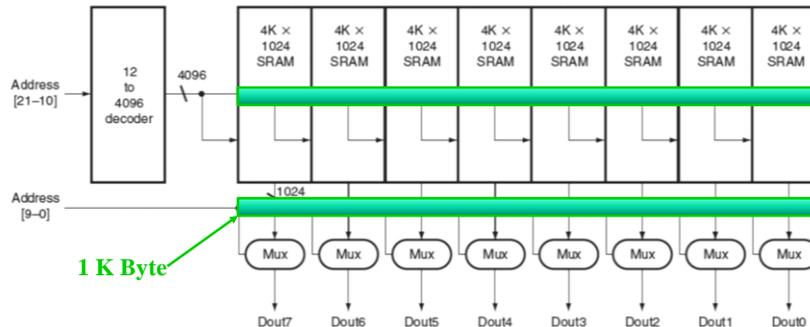
La stessa linea sarà selezionata in ogni banco di memoria → porto fuori 1K x 8 = 8 Kbyte.

Seleziono 1 bit da ogni banco (1 Byte) con il Mux (controllato dai 10 bit meno significativi, $\log_2 1K$).

Nell'approccio non a matrice avrei avuto bisogno di un decodificatore a 22 bit ($\log_2 4M$).



Modalità burst



Leggo prima una intera linea e poi attraverso i mux di uscita leggo la parola (1 Byte)
Lettura gerarchica in 2 passi.

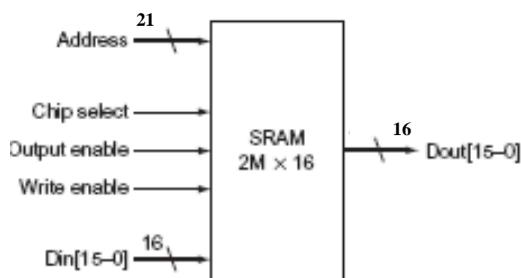
Synchronous Static RAM (SSRAM) -> trasferisco **burst**.

Burst: indirizzo iniziale (e.g. indirizzo del primo elemento di un vettore) + lunghezza del burst (numero byte consecutivi. In questo caso al massimo saranno 1K elementi pari alla lunghezza della linea di uscita).

- Non viene richiesto di specificare un nuovo indirizzo per ogni byte -> sono adiacenti.
- E' una tecnica potente per portare fuori da una SRAM e leggere **blocchi di memoria**.



Chip di SRAM (2M x 16)



Tempo di accesso:
da Address a Dout.

Altezza (2 Mbit) x Ampiezza (16 bit) – Oggi 1-4 bit

Altezza (2 Mbit) → 21 bit di indirizzamento
Ampiezza (16 bit) → Bus dati su 16 bit.

Chip select -> **Abilitazione del chip (per risparmiare energia, indirizzo nel range del chip).**

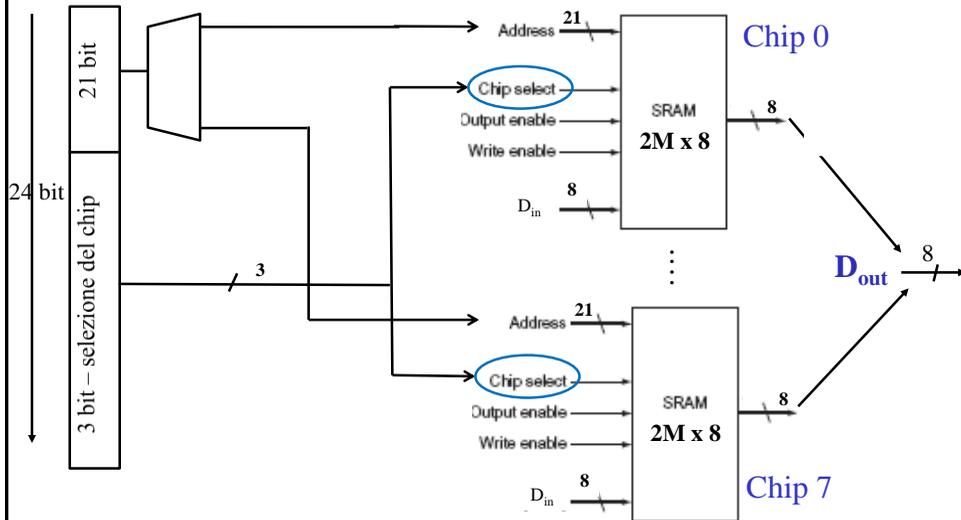
Write enable (abilitazione scrittura), Output enable (abilitazione lettura!).



Memoria SRAM come insieme di chip



Capacità totale = 8 Blocchi * 2 MByte => Capacità di 16 MByte



Principio di località: i dati all'interno di un chip verranno utilizzati negli istanti di tempo successivi. Abilitazione e disabilitazione sono eventi relativamente poco frequenti.

it\



Esempio



Memoria cache ad accesso diretto con linee di 4 word e word di 4 Byte.

Costruzione mediante 32 chip di SRAM da 2 M x 16

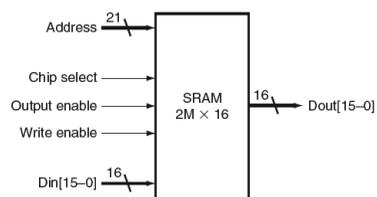
“Tessellazione della cache con i chip”

1w \$t1, 172(\$0)

Line number	Tag	Block
0		
1		
2		
		⋮
N-1		

Ampiezza = 4x4 Byte = 16 Byte

Chip di SRAM



Ampiezza 2 Byte



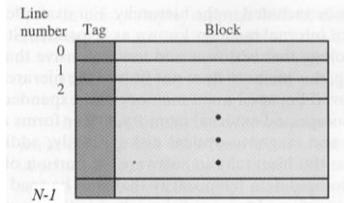
Dimensionamento

Memoria cache ad accesso diretto con linee di 4 word e word di 4 Byte (dim_linea = 16 Byte).

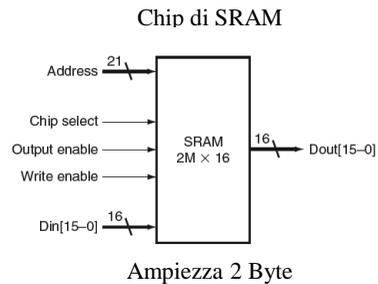
Costruzione mediante 32 chip di SRAM da 2 M x 16 =

32 chip x 2M linee x 2Byte = 128 Mbyte in totale

1w \$t1, 172(\$0)



Ampiezza = 4x4 Byte = 16 Byte



Servono 8 chip di SRAM affiancati per realizzare una linea di cache (8 chip x 2 Byte = 16 Byte).

Ho a disposizione 32 / 8 = 4 schiere di chip → 2M linee x 4 = 8 M linee

Cache sarà di **8M linee x 16 Byte = 128 Mbyte**

I chip saranno disposti in una matrice 4 x 8.



Indirizzamento

Cache sarà di 8M x 16 Byte = 128 Mbyte → 27 bit di indirizzamento

I chip saranno disposti in una matrice di 4 x 8

(4 schiere di 8 chip)





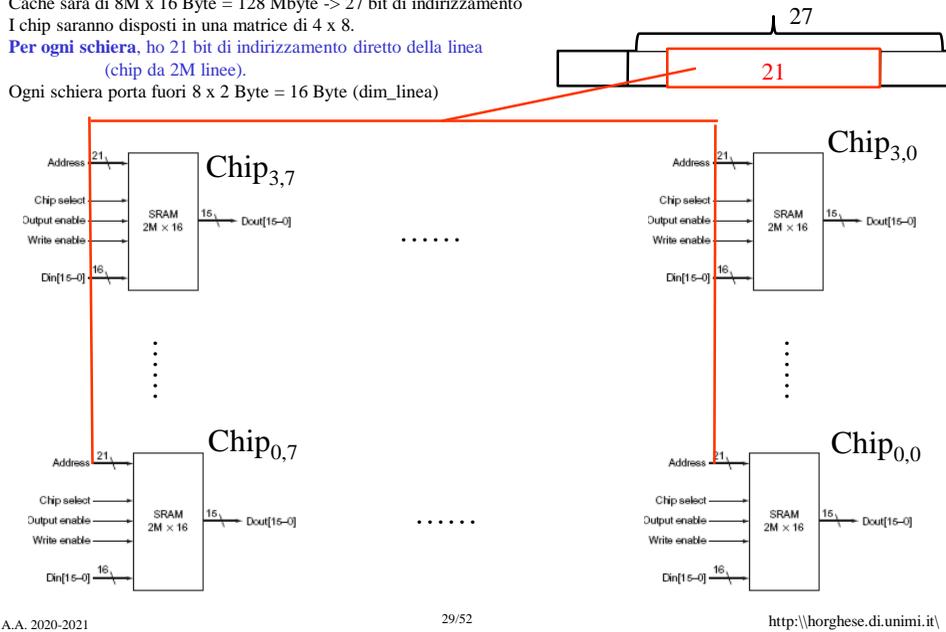
Selezione della linea



Cache sarà di $8M \times 16 \text{ Byte} = 128 \text{ Mbyte} \rightarrow 27$ bit di indirizzamento
I chip saranno disposti in una matrice di 4×8 .

Per ogni schiera, ho 21 bit di indirizzamento diretto della linea
(chip da 2M linee).

Ogni schiera porta fuori $8 \times 2 \text{ Byte} = 16 \text{ Byte}$ (dim_linea)



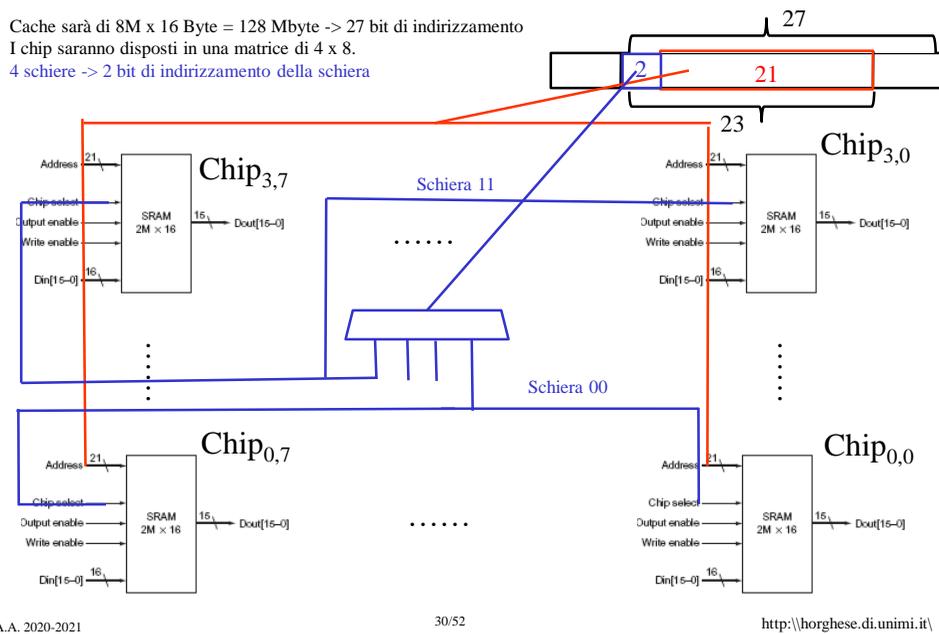
Selezione della schiera di chip

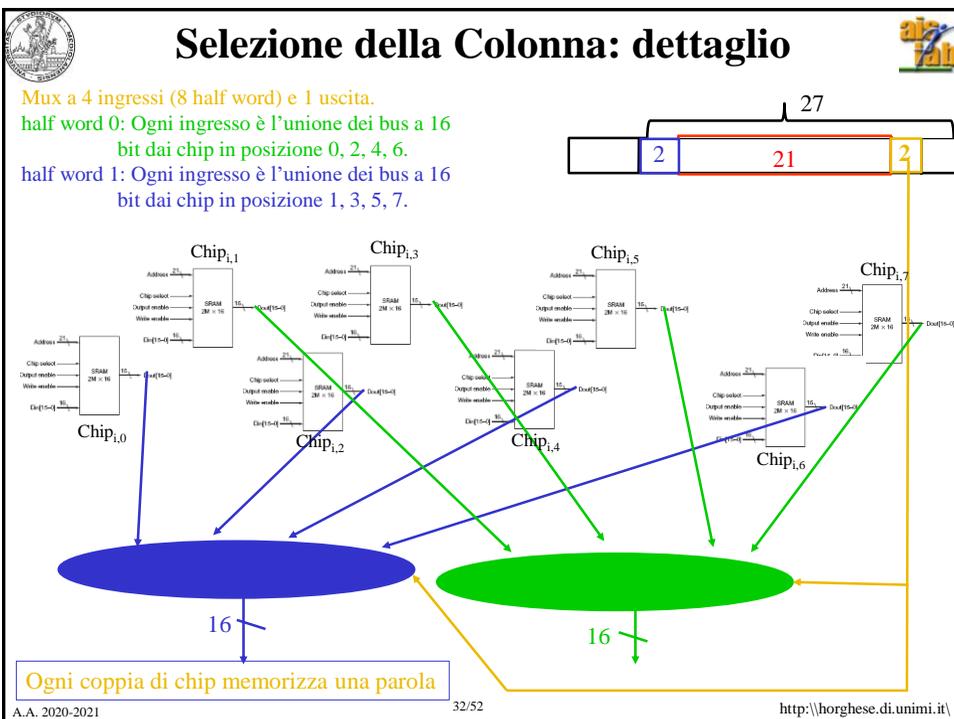
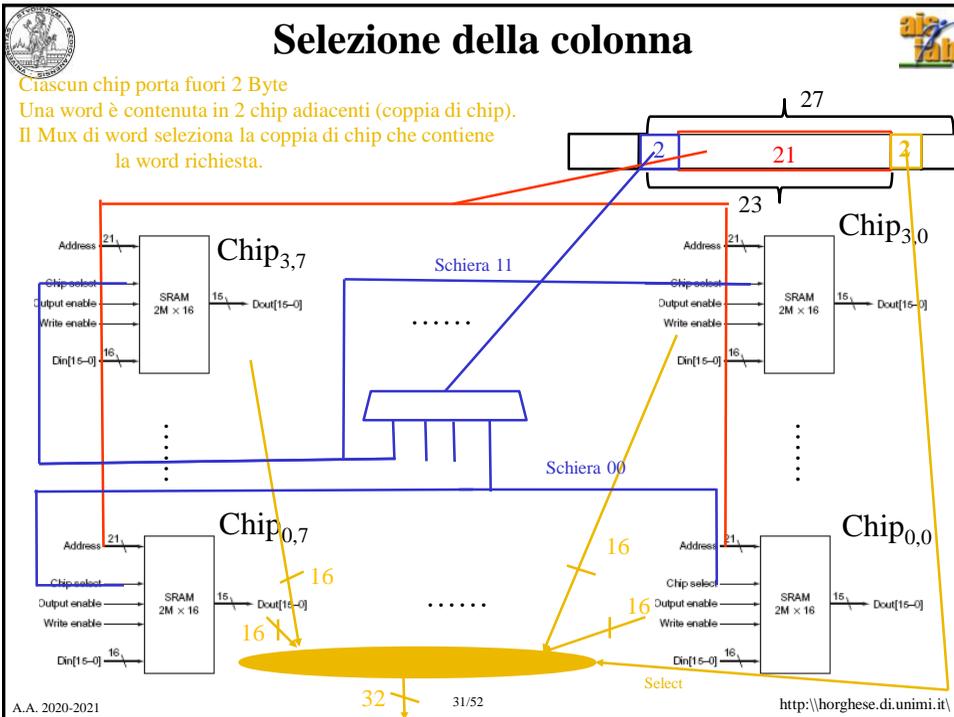


Cache sarà di $8M \times 16 \text{ Byte} = 128 \text{ Mbyte} \rightarrow 27$ bit di indirizzamento

I chip saranno disposti in una matrice di 4×8 .

4 schiere \rightarrow 2 bit di indirizzamento della schiera







Posizione della parola

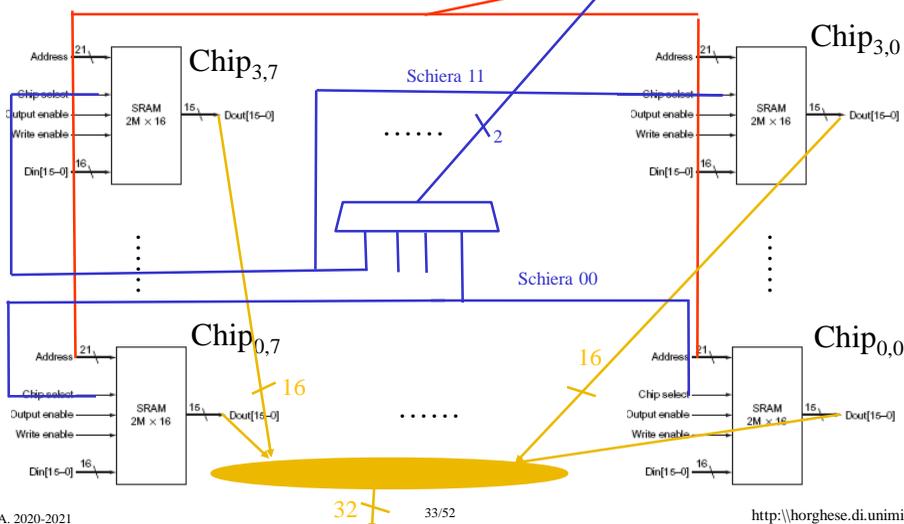


hw St1, 172(S0)

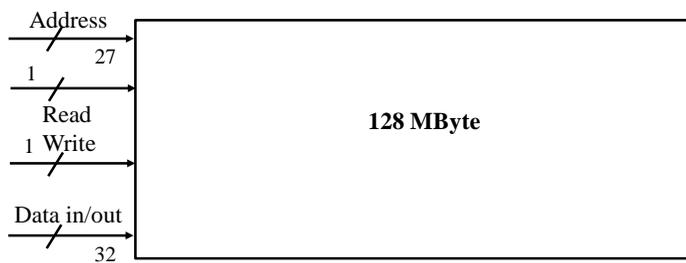
Schiera 0 dei chip (indirizzo < 32 MByte)

Linea 10 all'interno del chip (172 / 16 Byte = 10, R = 12)

Parola 3 => Chip 6 e 7 (12 Byte / 4 Byte / word = 3).



Indirizzamento fisico e logico





Sommario



Memory miss

SRAM

DRAM

Trasferimento dati



Memorie DRAM

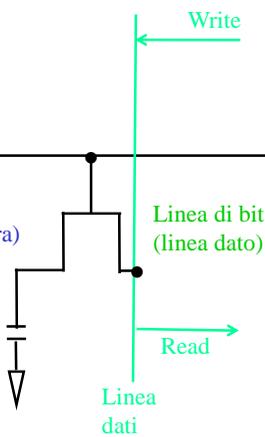


Dynamic RAM. 1 transistor per bit (contro 4-6 transistor della SRAM).

Linea di indirizzamento

Pass transistor
(transistor di lettura / scrittura)

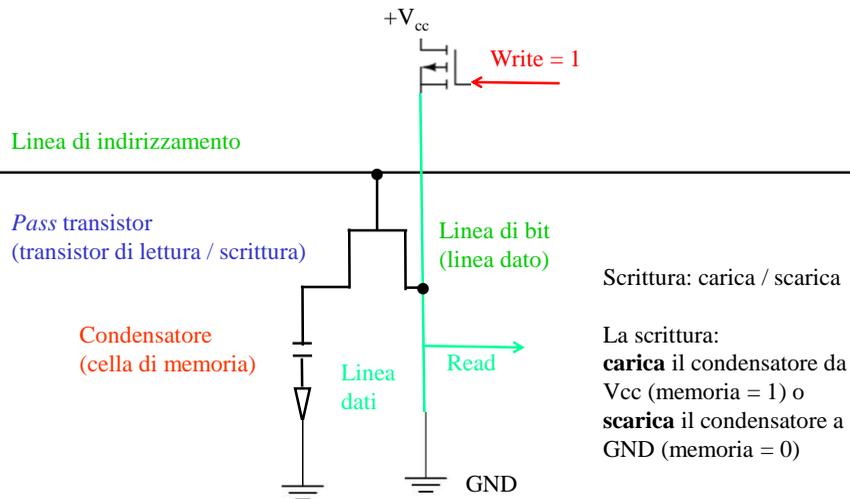
Condensatore
(cella di memoria)



1 Pass transistor + 1 condensatore.



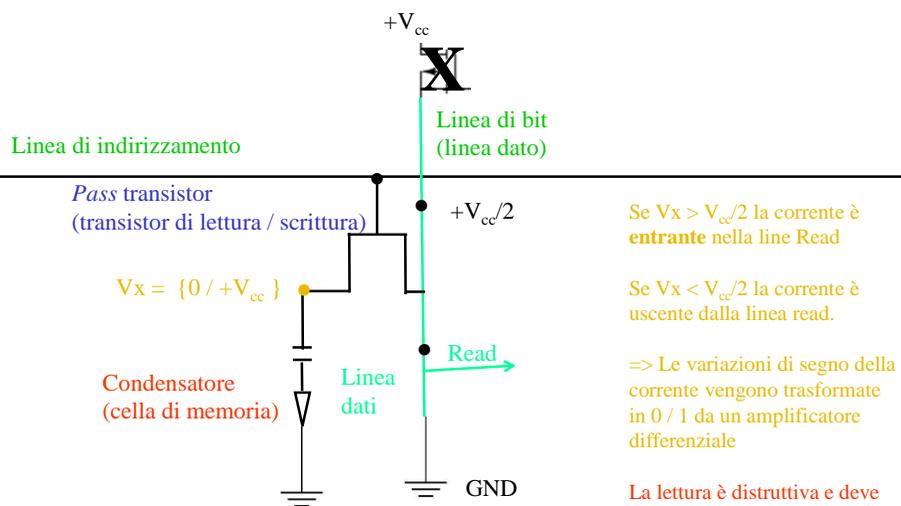
Memorie DRAM - scrittura



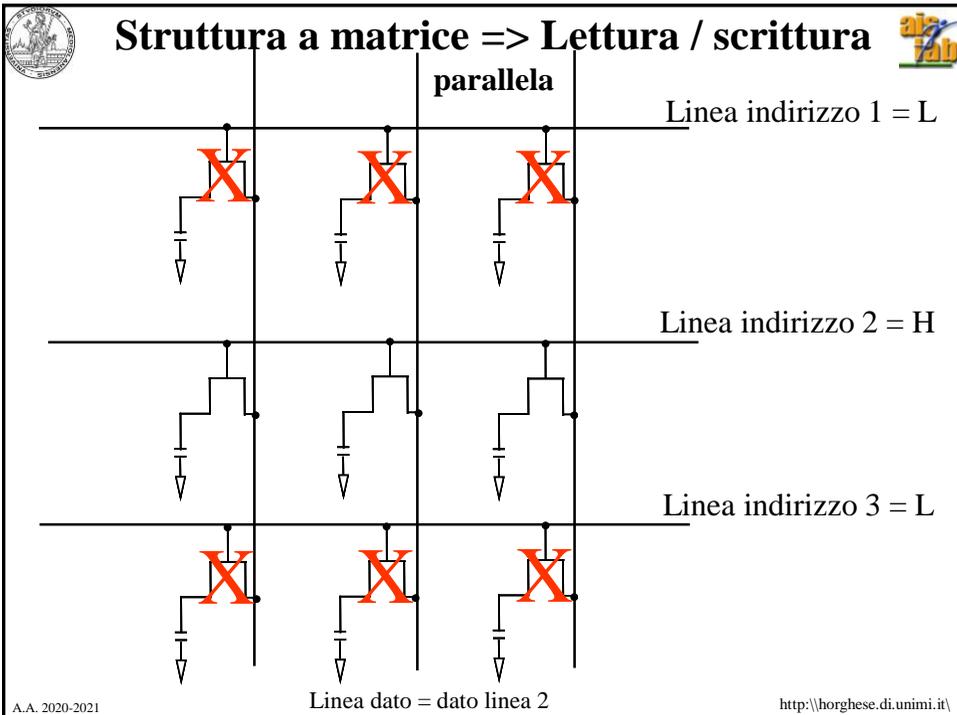
1 Pass transistor + 1 condensatore.



Memorie DRAM - lettura



1 Pass transistor + 1 condensatore.



I problemi delle DRAM

I condensatori sono componenti passivi -> richiedono tempo per caricarsi e scaricarsi.

La lettura è distruttiva.

I condensatori si scaricano (qualche millisecondo)

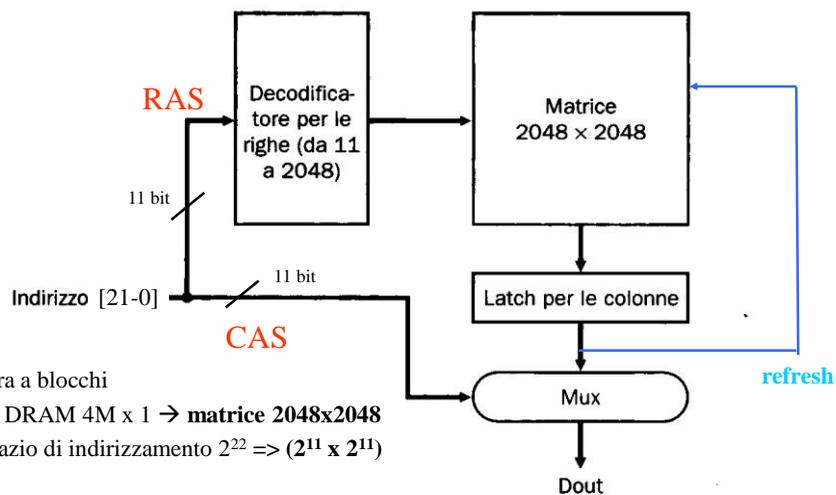
Refresh gestito autonomamente dal controllore della memoria mediante ciclo lettura/scrittura

Cosa leggo/scrivo? Quale/i bit?

A.A. 2020-2021 40/52 <http://horghese.di.unimi.it/>



Struttura a 2 livelli (matrice) di una DRAM



- Struttura a blocchi
 - es. DRAM 4M x 1 → matrice 2048x2048
 - Spazio di indirizzamento $2^{22} \Rightarrow (2^{11} \times 2^{11})$

•Accesso:

selezione riga (RAS) + selezione colonna (CAS)

Efficiente per il refresh (refresh di riga) – (35-70ms, tempo di carica dei condensatori).

Utilizzo per la Memoria Principale.

ese.di.unimi.it/

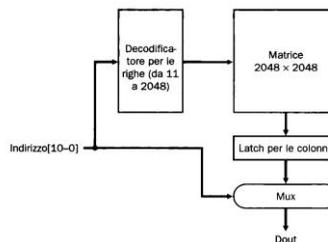


Osservazioni



Latch di Colonna (MDR)

- Refresh
- Trasferimento in modalità burst (solo MUX di uscita, viene nascosta la latenza di lettura)
- Modalità burst: indirizzo iniziale + lunghezza del burst.
- Aiuta il blocking



SDRAM (Synchronous Dynamic RAM) sono memorie sincronizzate. I latch sono sincronizzati.

- Esiste un tempo certo per la scrittura e la lettura globale.
- Non c'è bisogno di procedure di conferma (protocolli hand-shaking).
- Il processore può fare altro.

DDR SDRAM (Double Data Rate SDRAM). La sincronizzazione del funzionamento avviene su entrambi i fronti del clock. Velocità di trasferimento di picco di 3,2 Gbyte/s.



Specifiche delle DRAM

Anno introduzione	Dimensione chip	\$ per GiB	Tempo accesso totale a nuova riga / colonna	Tempo medio di accesso alla riga esistente
1980	64 Kibibit	\$ 1 500 000	250 ns	150 ns
1983	256 Kibibit	\$ 500 000	185 ns	100 ns
1985	1 Mebibit	\$ 200 000	135 ns	40 ns
1989	4 Mebibit	\$ 50 000	110 ns	40 ns
1992	16 Mebibit	\$ 15 000	90 ns	30 ns
1996	64 Mebibit	\$ 10 000	60 ns	12 ns
1998	128 Mebibit	\$ 4000	60 ns	10 ns
2000	256 Mebibit	\$ 1000	55 ns	7 ns
2004	512 Mebibit	\$ 250	50 ns	5 ns
2007	1 Gibibit	\$ 50	45 ns	1,25 ns
2010	2 Gibibit	\$ 30	40 ns	1 ns
2012	4 Gibibit	\$ 1	35 ns	0,8 ns

Tempo totale di accesso (dall'arrivo dell'indirizzo al dato in uscita): 35 ns

Tempo di accesso in modalità burst (accesso al buffer di linea di uscita): 0.8 ns

40 volte più veloce!



Esercizio

Memoria cache ad accesso diretto con linee di 8 word e word di 8 Byte.

Costruzione mediante 64 chip di SRAM da 1 M x 16

1w \$1, 1452(\$0)



Organizzazione della memoria



Organizzazione a matrice (cf. cache)
Organizzazione gerarchica.



Sommario



Memory miss
SRAM
DRAM
Trasferimento dati



Gestione dei fallimenti di una cache



La gestione avviene tra CPU e MMU.

Hit – è quello che vorremmo ottenere, il funzionamento della CPU non viene alterato.

Miss – **in lettura** devo aspettare che il dato sia stato caricato nella linea di cache e sia pronto
-> eccezione particolare della CPU che crea uno **stallo** della CPU e non un flush.

Nelle CPU super-scalari si sfrutta l'esecuzione fuori ordine per nascondere questa latenza.

Passi da eseguire in caso di Miss (miss penalty):

- 1) Bloccare tutte le istruzioni nella pipeline (blocco dei registri di pipeline per uno o più cicli di clock)
- 2) Scaricare in MM la linea di cache interessata (micro-blocco).
- 3) Richiedere che la MM legga il micro-blocco contenente il dato da leggere e lo porti fuori nel MDR.
- 4) **Trasferire il micro-blocco in cache, aggiornare i campi validita' e tag.**
- 5) Riavviare l'esecuzione della pipeline.

NB Il programma non può continuare!!



I componenti della Miss penalty



Tempi di accesso:

1 ciclo di clock per inviare l'indirizzo.

15 cicli di clock per ciascuna attivazione della Memoria Principale
(dall'invio dell'indirizzo alla parola in uscita)

1 ciclo di clock per trasferire una parola al livello superiore (cache).

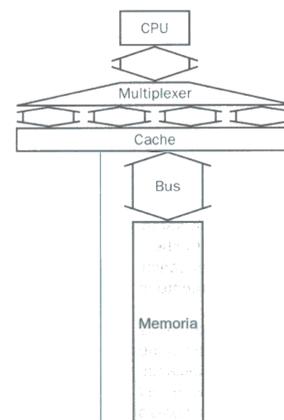
Blocco di cache di 4 parole, blocco di memoria principale di 1 parola:

Miss_Penalty =

$$\begin{aligned}
 &1 \text{ (invio indirizzo)} \\
 &+ \\
 &15 * 4 \text{ (parole) (lettura)} \\
 &+ \\
 &1 * 4 \text{ (parole) (trasferimento a cache)} \\
 &= \\
 &65 \text{ cicli_clock}
 \end{aligned}$$

Obbiettivi:

- Diminuire la penalità di fallimento (miss_penalty).



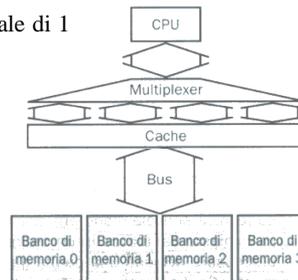


Interleaving



Blocco di cache di 4 parole, blocco di memoria principale di 1 parola:

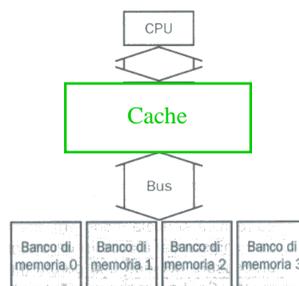
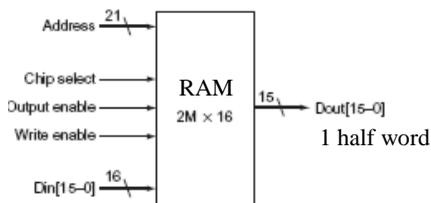
$$\begin{aligned}
 \text{Miss_Penalty} = & \\
 & 1 \text{ (invio indirizzo)} \\
 & + \\
 & 15 * 1 \text{ (insieme di 4 parole) (lettura)} \\
 & + \\
 & 1 * 4 \text{ (parole) (trasferimento a cache)} \\
 = & \\
 & \mathbf{20 \text{ cicli_clock}}
 \end{aligned}$$



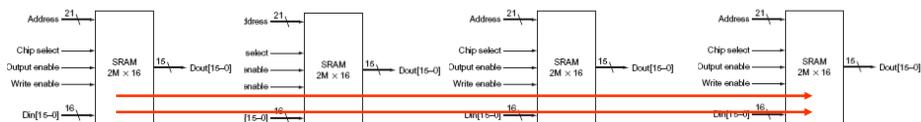
Interleaving (interallacciamento). Banchi che potrebbero essere trasferiti in parallelo alla cache. Si sposa perfettamente con la **struttura gerarchica** della memoria e con la **modalità di trasferimento a burst**.



Osservazioni sull'interleaving



Linea di cache di 4 parole = 16 Byte = 8 half word



Il vettore di MM viene spalmato sulle linee.
Leggo dal MDR in uscita (velocità di 40x)

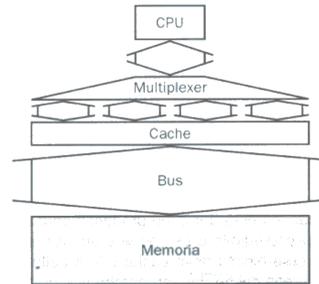


Bus ampio



Blocco di cache di 4 parole, blocco di memoria principale di 1 parola:

$$\begin{aligned} \text{Miss_Penalty} = & \\ & 1 \text{ (invio indirizzo)} \\ & + \\ & 15 * 1 \text{ (insieme di 4 parole) (lettura)} \\ & + \\ & 1 * 1 \text{ (insieme di 4 parole) (trasferimento a cache)} \\ & = \\ & \mathbf{17 \text{ cicli_clock}} \end{aligned}$$



Complessità del bus non giustificata. Si va verso bus sempre piu' stretti (insiemi di bus seriali a 1 bit).



Sommario



Memory miss

SRAM

DRAM

Trasferimento dati