

Calcolatori Elettronici

La memoria gerarchica

Introduzione

La memoria - generalità

Funzioni:

- **Supporto alla CPU:** deve fornire dati ed istruzioni il più rapidamente possibile
 - **Archiviazione:** deve consentire di archiviare dati e programmi, garantendone la conservazione e la reperibilità anche dopo elevati periodi di tempo
- Esigenze diverse:
- **velocità** per il supporto alla CPU
 - **non volatilità e capacità** per l'archiviazione

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 4

L. Tarantino - a.a. 2004/2005

E' necessario combinare più unità di memoria

- ⌌ *Capacità, velocità e costo NON sono criteri indipendenti:*
 - la velocità si ottiene con capacità ridotta e tecnologie costose
 - la capacità si ottiene con tempi lunghi e tecnologie poco costose
- ⌌ *Le tecnologie elettroniche sono utilizzate per le memorie più veloci*
- ⌌ *Le tecnologie magnetico/optiche sono utilizzate per le memorie più capaci*

Tecnologia ed organizzazione (1)

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 3

L. Tarantino - a.a. 2004/2005

E' necessario combinare più unità di memoria

- ⌌ *Alcuni parametri:*
 - *tempo di accesso* (quanto passa tra la richiesta e la relativa risposta?)
 - *tempo di trasferimento* (quantitativo byte al secondo si possono trasferire?)
- ⌌ *Capacità*
 - quanti byte può contenere?
- ⌌ *Costo*
 - dati capacità e costo totale, qual è il costo per bit?

Caratterizzazione di una memoria

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 6

L. Tarantino - a.a. 2004/2005

Tecnologia ed organizzazione (3)

Osservazioni preliminari:

- la funzione di "supporto alla CPU" è svolta dalla memoria volatile
- i tempi di funzionamento della memoria principale sono comunque sensibilmente maggiori dei tempi di funzionamento della CPU
- se la memoria principale fosse più piccola i tempi di funzionamento si ridurrebbero (più piccolo è più veloce)

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 5

L. Tarantino - a.a. 2004/2005

Tecnologia ed organizzazione (2)

Obiettivo:

Ottenere un sistema di memoria complessivo che

- fornisca all'utente una quantità di memoria pari a quella disponibile nella tecnologia più lenta
- consentendo allo stesso tempo (*quasi sempre*) una velocità di accesso pari a quella garantita dalla tecnologia più veloce

E' necessario combinare più unità di memoria

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 8

L. Tarantino - a.a. 2004/2005

n il comportamento medio dei programmi è prevedibile, in base ad alcuni principi:
 n occorre poter prevedere con una certa previsione il comportamento dei programmi eseguiti

— la regola del 90/10

— il principio di località del riferimento

n località temporale
 n località spaziale

Tecnologia ed organizzazione (5)

Come si individuano i dati da mettere nella memoria veloce?

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 7

L. Tarantino - a.a. 2004/2005

n sistema di memoria (volatile) formato da due unità, una piccola e veloce, e una grande ma più lenta
 n i dati usati più frequentemente si potrebbero tenere nella memoria veloce, in modo che la maggior parte delle volte il tempo di accesso sia quello della memoria veloce
 n solo quando un dato non viene trovato nella memoria veloce lo si va a cercare nella memoria più lenta

Tecnologia ed organizzazione (4)

Combinare più unità di memoria volatile

Esempio (scenario semplificato)

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 10
L. Tarantho - a.a. 2004/2005

In altre parole

se un'istruzione (o un dato) è stato utilizzata in un istante di tempo T, c'è un'elevata probabilità che la stessa istruzione (lo stesso dato) venga utilizzato in istanti di tempo "prossimi" a T

Località temporale

i programmi tendono a riutilizzare dati ed istruzioni che hanno usato recentemente

Esempio

```

while i < 100 {
    somma = somma + A[i]
    ..
}
    
```

ad ogni iterazione viene utilizzata questa istruzione

ad ogni iterazione del ciclo la variabile somma viene utilizzata

Il principio di località del riferimento (1)

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 9
L. Tarantho - a.a. 2004/2005

Principio empirico di progetto:

- investire risorse il dove viene speso il tempo
- Cosa accade durante l'esecuzione di un programma?
- Regola del 90/10

il 90% del tempo viene speso sul 10% del codice (cicli...)

La regola del 90/10

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 12
L. Tarantho - a.a. 2004/2005

Capacità
+ piccolo = + veloce



Tecnologia
+ veloce = + costoso

cerchiamo di mettere il "famoso" 10% in una memoria piccola e costruita con tecnologia "veloce"

Sfruttiamo la regola del 90/10

- » per migliorare le prestazioni posso cercare di abbassare il tempo necessario per gli accessi a memoria relativi al 10% del codice responsabile del 90% del tempo di esecuzione
- » il tempo di accesso in memoria dipende dalla tecnologia usata per la realizzazione della memoria e dalla sua capacità

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 11
L. Tarantho - a.a. 2004/2005

» **In altre parole**

se in un istante di tempo T è stata utilizzata una parola di memoria di indirizzo X, c'è un'elevata probabilità che in istanti di tempo "prossimi" a T vengano usate parole di memoria vicine a X

» **Località spaziale**

dati vicini in memoria tendono ad essere utilizzati in istanti di tempo ravvicinati

» **Esempi**

- istruzioni vicine
- componenti di array durante scansioni sequenziali

Il principio di località del riferimento (2)

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 14
L. Tarantho - a.a. 2004/2005

quantità di informazione
 trasferita nell'unità di tempo

<i>Dimensioni</i>	< 1 MB	< 100	< 4 GB
<i>Tempo accesso</i> (nanosec)	5.000 +	10.000 +	1.000 +
<i>Banda passante</i> (MB/sec)	↓ quantità di informazione trasferita nell'unità di tempo		

```

        graph LR
            CPU[CPU] --- Cache[Cache]
            Cache --- Memoria[Memoria principale]
            
```

Un sistema con due livelli di memoria

Calcolatori Elettronici - Introduzione alla memoria gerarchica - Slide 13
L. Tarantho - a.a. 2004/2005

cerchiamo di mettere il "famoso" 10% in una memoria
 piccola e costruita con tecnologia "veloce"

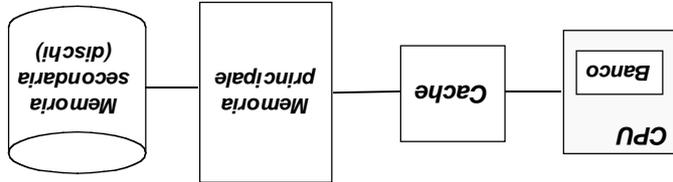
- E' difficile, richiederebbe un'analisi complessa del codice
- **Strutiamo il principio di località**
 - mettiamo i dati usati recentemente nella memoria piccola e veloce (*loc. temporale*)
 - quando spostiamo un dato nella memoria piccola e veloce portiamoci dietro anche le parole vicine (*loc. spaziale*)

Combiniamo 90/10 e località

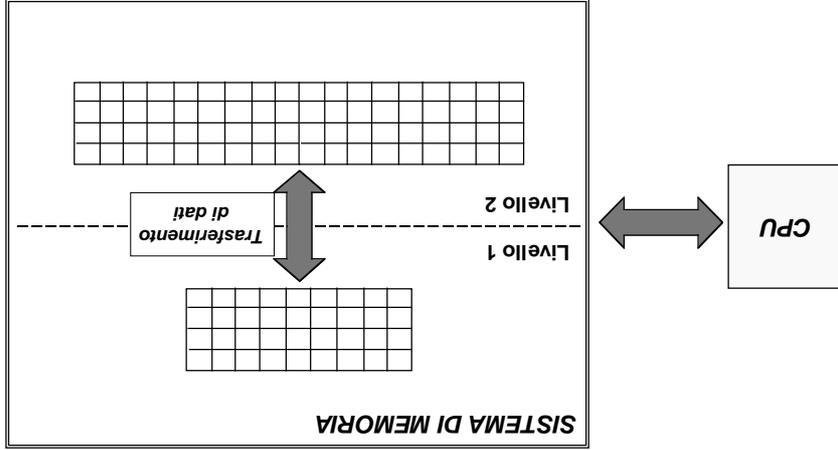
- *svantaggio*: la carica tende a dissiparsi
 - la memoria deve essere “*rinfrascata*” ogni tanto (riletta e riscritta tutta prima che i valori dei bit vengano persi) *REFRESH*
 - le letture tendono a dissipare la carica, quindi ogni lettura è seguita da una scrittura
 - *vantaggio*: il bit è rappresentato dalla capacità tra gate e canale di un mostet (il singolo bit occupa meno spazio, inoltre è meno costoso)
 - *svantaggio*: la carica tende a dissiparsi
 - *vantaggio*: il bit è rappresentato dalla capacità tra gate e canale di un mostet (il singolo bit occupa meno spazio, inoltre è meno costoso)
- La *cache* è realizzata in tecnologia SRAM (RAM statica)
- La *memoria principale* è realizzata in tecnologia DRAM (RAM Dinamica)
- La *memoria secondaria* è realizzata mediante dischi magnetici

Aspetti tecnologici

Dimensioni	Tempo accesso (nanosec)	Banda passante (MB/sec)
< 1KB	< 0,5	20.000 ÷ 100.000
< 1 MB	< 10	5.000 ÷ 10.000
< 4 GB	< 100	1.000 ÷ 5.000
> 20 GB	10.000.000	20 ÷ 40



Il sistema di memoria gerarchica



Un sistema con due livelli di memoria (2)

1) non sempre la parola cercata sarà in cache
 — se c'è è "accesso di tipo hit" (colpito)
 — se non c'è è "accesso di tipo miss" (mancato)
 2) se si verifica un miss devo cercare la parola nel livello di memoria seguente

Il blocco è la più piccola quantità di informazione spostata tra due livelli di memoria contigui

1) per la località
 — mantenere in cache i dati usati più recentemente
 — spostare "blocchi" di parole e non parole singole
 2) facciamo un modo che la maggior parte delle volte i dati e le istruzioni cercate si trovano nella cache
 3) un programma non accede a tutto il suo codice con uguale probabilità

Dove sono dati e istruzioni