



# Valutazione delle prestazioni

Prof. Alberto Borghese  
Dipartimento di Scienze dell'Informazione  
[borgnese@dsi.unimi.it](mailto:borgnese@dsi.unimi.it)

Università degli Studi di Milano



## Sommario

Cosa vuol dire valutare le prestazioni?

Benchmark

Miglioramento delle prestazioni

Valutazione delle prestazioni del sistema di memoria



## Perché valutare le prestazioni?



- Misura/Valutazione quantitativa delle prestazioni (velocità...).
- Fare scelte intelligenti (e.g. installare nuovo hardware o nuovo sw).
- Orientarsi nell'acquisto di nuovo hw.
- Fatturazione delle prestazioni.

### La misura delle prestazioni è il tempo.

Prestazioni<sub>X</sub> > prestazioni<sub>Y</sub> => tempo<sub>X</sub> < tempo<sub>Y</sub> =>  
1/tempo<sub>X</sub> > 1/tempo<sub>Y</sub>

Prestazioni<sub>X</sub> = (1+n/100)xPrestazioni<sub>Y</sub>.  
Tempo<sub>Y</sub>=(1+n/100)xTempo<sub>X</sub> =>  $n=100 \times (\text{Tempo}_Y - \text{Tempo}_X) / \text{Tempo}_X$   
**ΔT**

### *Le prestazioni migliorano perché:*

- Incrementano le prestazioni.
- Diminuisce il tempo di esecuzione.



## Criteri (metrica) di valutazione orientati all'utente



Velocità di esecuzione + quantità di informazione elaborata.

Il criterio di valutazione dipende dall'utilizzo del calcolatore!

- 1) Utilizzo personale -> **tempo di esecuzione.**
- 2) Utilizzo come server -> **throughput.**

### Throughput:

Ammontare di lavori svolti in un dato tempo.  
(accessi a banche dati, programmi, transazioni commerciali...).

### Domande:

Un processore più veloce cosa influenza?  
Più processori dedicati, cosa modificano?



## Esempio



$$p_A = 2 \quad (t_A = 0.5) \quad p_B = 1.5 \quad (t_B = 0.666\dots)$$

### Valutazione in termini di tempo di prestazioni:

$$p_B / p_A = 0.75 \quad \text{B ha prestazioni pari al 75\% di A.}$$

$$\text{Variazione delle prestazioni percentuale: } (p_B - p_A) / p_A = -0.25\%$$

### Utilizzando il tempo di esecuzione:

$$(1/t_B) / (1/t_A) = t_A / t_B = 75\%$$

$$\text{Variazione delle prestazioni percentuale: } (1/t_B - 1/t_A) / (1/t_A) = t_A * (1/t_B - 1/t_A) = 1/2 * (-1/2) = -0.25\%$$

### Valutazione in termini di tempo di esecuzione:

$$t_B/t_A = (2/3) / (1/2) = 4/3 = 1.3333\dots \text{ B richiede il 133\% del tempo di A per eseguire il programma.}$$

$$\text{Variazione delle prestazioni percentuale: } (t_B - t_A) / t_A = (2/3 - 1/2) / (1/2) = 1/3 \Rightarrow 33.3\dots\%$$

B richiede 33.3...% in più per l'esecuzione del programma.



## Criteri (metrica) di valutazione orientati alla macchina



**Tempo di risposta** rappresenta la latenza per il completamento di un lavoro includendo accessi a disco, accessi a memoria, attività di I/O, ...

**Tempo di CPU** rappresenta il tempo speso dalla CPU per eseguire il programma dato: **non** include il tempo di attesa per I/O o per l'esecuzione di altri programmi. Comprende il **tempo utente di CPU** (tempo speso dalla CPU per eseguire le linee di codice che stanno nel nostro programma) + **tempo di CPU di sistema** (speso dal sistema operativo per eseguire i compiti richiesti dal programma)

Comando time Unix: 90.7u 12.9s 2:39 65%

Funzione Profile in Windows.



## Unità di misura delle prestazioni (CPI)



*Tempo di CPU =*

$$\begin{aligned} & \text{Numero\_cicli\_clock} * \text{Durata\_clock} = \\ & \text{Numero\_cicli\_clock} / \text{Frequenza\_clock}. \end{aligned}$$

*Determinazione del numero di cicli di clock:*

$$\begin{aligned} \text{Cicli di clock per istruzione (CPI)} = \\ \text{Cicli\_clock\_CPU\_programma} / \text{Numero\_istruzioni} \end{aligned}$$

*Quindi:*

$$T_{\text{CPU}} = \text{CPI} * \text{Numero\_Istruzioni} * T_{\text{clock}}$$



## Esempio



Tempo di esecuzione del programma: 1.2s

Numero di istruzioni: 400k.

Clock: 1Mhz.

Per l'esecuzione del programma, occorrono: #Cicli\_clock =  $10^6 * 1.2$

CPI = #Cicli\_clock / #Istruzioni = 3.

**NB Sulle macchine di oggi il CPI è inferiore ad 1.**

$$T_{\text{medio}} = \frac{\sum t_i}{\text{\#Istruzioni}} = \frac{T_{\text{tot}}}{\text{\#Istruzioni}}$$

$$T_{\text{tot}} = \text{CPI} * T_{\text{clock}} * \text{\#Istruzioni}$$



## Misura delle prestazioni



Tempo esecuzione singola istruzione, ma:

In genere, istruzioni di tipo diverso richiedono quantità diverse di tempo. Esempi:

- la moltiplicazione richiede più tempo dell'addizione
- l'accesso alla memoria richiede più tempo dell'accesso ai registri.

Tempo esecuzione medio (pesato) di un mix di istruzioni:

$$t_{medio} = \frac{\sum_{i=0}^S l_i t_i}{\sum_{i=0}^S l_i}$$



## Misura delle prestazioni mediante CPI



$$T_{CPU} = CPI * \text{Numero\_Istruzioni} * T_{clock}$$

$$t_{medio} = \frac{\sum_{i=0}^S l_i t_i}{\sum_{i=0}^S l_i} \quad CPI_{medio} = \frac{\sum_{i=1}^n (CPI_i * l_i)}{\sum_{i=1}^n l_i} = \sum_{i=1}^n (CPI_i * f_i)$$

- $CPI_i$  numero di cicli di clock per istruzioni di tipi  $i$ .
- $l_i$  Numero di volte che l'istruzione  $i$  viene eseguita nel programma.
- $f_i$  Frequenza con cui l'istruzione  $i$  viene eseguita nel programma.

(  $\sum_{i=1}^n l_i$  rappresenta il numero di istruzioni)

$$T_{CPU} = \sum_{i=1}^n (CPI_i * l_i) * T_{clock}$$



## Esempio



Si consideri un calcolatore in grado di eseguire le istruzioni riportate in tabella:

	Frequenza	cicli di clock
ALU	43%	1
Load	21%	4
Store	12%	4
Branch	12%	2
Jump	12%	2

Calcolare CPI e il tempo di CPU per eseguire un programma composto da 200 istruzioni supponendo di usare una frequenza di clock pari a 500 MHz.

$$CPI = 0,43 * 1 + 0,21 * 4 + 0,12 * 4 + 0,12 * 2 + 0,12 * 2 = 2,23$$

$$T_{CPU} = 200 * 2,23 * 2_{ns} = 892_{ns}$$



## MIPS = milioni di istruzioni per secondo



$$MIPS = (\text{numero\_istruzioni} / 10^6) / \text{tempo\_esecuzione}$$

$$MIPS = \text{frequenza\_clock} / (CPI * 10^6) = 1/t_{clock} * 1 / (CPI * 10^6) = 1 / (t_{medio} * 10^6)$$

$$t_{clock} * CPI = t_{medio}$$

### Problemi:

- dipende dall'insieme di istruzioni, quindi è difficile confrontare computer con diversi insiemi di istruzioni;
- varia a seconda del programma considerato;
- può variare in modo inversamente proporzionale alle prestazioni!
- valore di picco, scelgo il mix di istruzioni per massimizzare il MIPS misurato (fuorviante).

**Esempio:** macchina con hardware opzionale per virgola mobile. Le istruzioni in virgola mobile richiedono più cicli di clock rispetto a quelle che lavorano con interi, quindi i programmi che usano l'hardware opzionale per la virgola mobile in luogo delle routine software per tali operazioni impiegano meno tempo ma hanno un MIPS più basso. L'implementazione software delle istruzioni in virgola mobile esegue semplici istruzioni, con il risultato di avere un elevato MIPS, ma ne esegue talmente tante da avere un più elevato tempo di esecuzione!!



## Sommario



Cosa vuol dire valutare le prestazioni?

**Benchmark**

Miglioramento delle prestazioni

Valutazione delle prestazioni del sistema di memoria



## Misure & Problemi



$MIPS_{relativi} = \frac{tempo_{CPU}}{tempo_{CPU_{ref}}} * MIPS_{CPU_{ref}}$  La  $CPU_{ref}$  è VAX-11/780. Problema: evoluzione dei sistemi.

MFLOPS per i super computer. Problema: misure di picco.

MIPS di picco e sostenuti. Problema: poco significative.

**Benchmarks = Programmi per valutare le prestazioni.**

Benchmarks: Whetstone, 1976; Drystone, 1984.

Kernel benchmark. Loop Livermore, Linpack, 1980. Problema: polarizzazione del risultato.

Benchmark con programmi piccoli (10-100 linee, 1980). Problema: mal si adattano alle strutture gerarchiche di memoria.



## Indici SPEC ('89, '92, '95)



<http://www.spec.org/>. The Standard Performance Evaluation Corporation (SPEC) is a non-profit corporation formed to establish, maintain and endorse a standardized set of relevant benchmarks that can be applied to the newest generation of high-performance computers. SPEC develops benchmark suites and also reviews and publishes submitted results from our [member organizations](#) and other benchmark licensees.

Insieme di programmi test.  
Condizioni diverse: singolo / multiplo processore / time sharing.  
Benchmark specifici per valutare S.O. e I/O.

SPEC'95 -> SPECint, SPECfp, base Sun SPARCstation 10/40.

### ***Benchmark particolari:***

SDM (Systems Development Multitasking).  
SFS (System-level File Server).  
SPECchpc96. Elaborazioni scientifiche ad alto livello.

***Orientamento:*** Benchmark specifici.



## Esempio benchmark SPEC95



### *Elaborazione intera:*

- |            |  |
|------------|--|
| 1) Go      | Intelligenza artificiale   |
| 2) m88ksim | Simulatore chip Motorola 88K; esecuzione di un programma.                                |
| 3) gcc     | Compilatore Gnu C che genera codice SPARC.   |
| 4) compres | Compressione e decompressione di un file in memoria.                                     |
| 5) li      | Interprete lisp  |
| 6) ijpeg   | Compressione e decompressione di immagini grafiche.                                      |
| 7) perl    | Manipolazione di stringhe e numeri primi nel linguaggio di programmazione dedicato Perl. |
| 8) vortex  | Programma di gestione di una base di dati.   |



## Esempio becnhmark SPEC95



### Elaborazione virgola mobile:

- |            |   |
|------------|---|
| 1) Tomcatv | Programma per generazione di griglie.   |
| 2) Swim    | Modello per acqua poco profonda con griglia 513 x 513.  |
| 3) Su2cor  | Fisica quantistica: simulazione MonteCarlo.   |
| 4) Hydro2D | Astrofisica: equazione idrodinamiche di Naiver Stokes.  |
| 5) Mgrid   | Risolutore multi-griglia in campo di potenziale 3D.   |
| 6) Applu   | Equazioni alle differenze parziali paraboliche/ellittiche.                                    |
| 7) Turb3D  | Simulazione di turbolenza isotropica ed omogenea in un cubo.                                  |
| 8) Apsi    | Risoluzione di problemi di temperatura, velocità del vento e diffusione di agenti inquinanti. |
| 9) Fpppp   | Chimica quantistica.  |
| 10) Wave5  | Fisica dei plasmi: simulazione di particelle elettromagnetiche.                               |

URL: <http://www.spec.org/>

A.A. 2007-2008

17/41

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



## SPEC CPU200 CINT2000



Benchmark	Language	Category	Full Descriptions
164.gzip	C	Compression	<a href="#">HTML</a> <a href="#">Text</a>
175.vpr	C	FPGA Circuit Placement and Routing	<a href="#">HTML</a> <a href="#">Text</a>
176.gcc	C	C Programming Language Compiler	<a href="#">HTML</a> <a href="#">Text</a>
181.mcf	C	Combinatorial Optimization	<a href="#">HTML</a> <a href="#">Text</a>
186.crafty	C	Game Playing: Chess	<a href="#">HTML</a> <a href="#">Text</a>
197.parser	C	Word Processing	<a href="#">HTML</a> <a href="#">Text</a>
252.eon	C++	Computer Visualization	<a href="#">HTML</a> <a href="#">Text</a>
253.perlbnk	C	PERL Programming Language	<a href="#">HTML</a> <a href="#">Text</a>
254.gap	C	Group Theory, Interpreter	<a href="#">HTML</a> <a href="#">Text</a>
255.vortex	C	Object-oriented Database	<a href="#">HTML</a> <a href="#">Text</a>
256.bzip2	C	Compression	<a href="#">HTML</a> <a href="#">Text</a>
300.twolf	C	Place and Route Simulator	<a href="#">HTML</a> <a href="#">Text</a>

A.A. 2007-2008

18/41

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



## Sommario



Cosa vuol dire valutare le prestazioni

Benchmark

Miglioramento delle prestazioni

Valutazione delle prestazioni del sistema di memoria



## Miglioramento delle prestazioni



- Riduzione del numero di cicli di clock.
- Diminuzione del periodo di clock (aumentare la frequenza).

$\text{Tempo esecuzione} = \text{Numero\_Cicli\_clock} * \text{Durata del clock}$

CPI rappresenta il tempo di esecuzione medio delle istruzioni.

Miglioramenti dell'architettura per ridurre il CPI.

Miglioramento del compilatore per ridurre il CPI.

Ridurre la durata del clock (aumentarne la frequenza).

### Espressione dei risultati

Il tempo totale di esecuzione dipende da diverse caratteristiche: dischi, sottosistema di I/O, sottosistema grafico ....

Per questo motivo occorre menzionare la **configurazione** del sistema.



## Valutazione delle prestazioni, coerenza



	Calcolatore I	Calcolatore II
Istruzione A (s)	1	10
Istruzione B (s)	1000	100
Istruzione C (s)	10	100

Qual è più veloce? Dipende dal peso dei programmi.

Media pesata, tempo medio di esecuzione di un'istruzione:

$$T = 1/n \sum_{i=0}^n n_i t_i$$

Programma 1: 1000 istruzioni A, 1 istruzione B, 10 istruzioni C.

$$T_I = 1/1011 * (1000*1 + 1*1000 + 10*10) = 2100/1011 \approx 2$$

$$T_{II} = 1/1011 * (1000*10 + 1*100 + 10*100) = 11100/1011 \approx 10$$

Programma 2: 100 istruzione A, 10 istruzioni B, 10 istruzioni C.

$$T_I = 1/120 * (100*1 + 10*1000 + 10*10) = 10200/120 \approx 85$$

$$T_{II} = 1/120 * (100*10 + 10*100 + 10*100) = 3000/120 \approx 25$$



## Come rendere più veloci i calcolatori



Rendere veloce il caso più comune.

Si deve favorire il caso più frequente a discapito del più raro.

Il caso più frequente è spesso il più semplice e può essere quindi reso più veloce del caso infrequente.

### *Legge di Amdahl*

Il miglioramento delle prestazioni globali ottenuto con un miglioramento particolare (e.g. un'istruzione), dipende dalla frazione di tempo in cui il miglioramento era eseguito.

Esempio: Pentium e PentiumPro: a fronte di un raddoppio della frequenza di clock che è passata da 100 a 200 Mhz, si è registrato un aumento delle prestazioni misurate tramite SpecInt di 1,7 volte e di 1,4 volte misurate in SpecFloat.



## Speed-up



Il miglioramento globale proporzionale al miglioramento di una parte del sistema?

### Speed up (accelerazione):

$\text{prestazioni\_intero\_lavoro\_con\_miglioramento} / \text{prestazioni\_senza}$

Oppure

$\text{tempo\_intero\_lavoro\_senza\_miglioramento} / \text{tempo\_con\_miglioramento}$ .



## Speed-up - esempio



Consideriamo un calcolatore (CALC1) con ALU ed una FP\_ALU.  
Consideriamo un secondo calcolatore (CALC2) in cui la ALU è stata velocizzata (2x).

Consideriamo un'applicazione che prevede un 90% di istruzioni in aritmetica intera. Di quanto è lo speed-up?

ISTRUZIONI INTERE			ISTRUZIONI TOTALI	
Calcolatore	T_EXEC	Speedup_m	T_EXEC	Speedup
CALC1	90	1.0	100	1.0
CALC2	45	2.0	55	1.82

Speed-up =  $100/55 = 1.818\dots$



## Corollario della legge di Amdhal



Se un miglioramento è utilizzabile solo per una frazione del tempo di esecuzione complessivo ( $F_m$ ), allora non è possibile accelerare l'esecuzione più del reciproco di uno meno tale frazione:

$$\text{Speedup}_{\text{globale}} < 1/(1-F_m).$$

Definizioni:

1. **Frazione migliorato** ( $F_m \leq 1$ ), ovvero la frazione del tempo di calcolo della macchina originale che può essere modificato per avvantaggiarsi dei miglioramenti. Nell'esempio precedente la frazione è 0.90.

$$\begin{aligned} T_m &= F_m * T_{\text{old}} \\ T_{\text{nm}} &= (1 - F_m) * T_{\text{old}} \end{aligned}$$

2. **Speedup migliorato** ( $S_m \geq 1$ ), ovvero il miglioramento ottenuto dal modo di esecuzione più veloce. Nel precedente esempio questo valore viene fornito nella colonna chiamata Speedup\_migliorato (pari a 2).



## Dimostrazione



$$\begin{aligned} T_{\text{old}} &= T_{\text{old}} * (1 - F_m) + T_{\text{old}} * F_m \\ T_{\text{new}} &= T_{\text{nm}} + T_m = T_{\text{old}} * (1 - F_m) + T_{\text{old}} * F_m / S_m \end{aligned}$$

$$0.1 \quad \quad \quad 0.9 / 2$$

$$T_{\text{new}} = T_{\text{old}} * (1 - F_m + F_m / S_m) = T_{\text{old}} * [1 - F_m * (1 - 1 / S_m)]$$

Istruzioni non accelerate

Istruzioni accelerate

$$\text{Speedup}_{\text{globale}} = T_{\text{old}} / T_{\text{new}} = T_{\text{old}} / (T_{\text{old}} * [1 - F_m * (1 - 1 / S_m)]) =$$

$$1 / [1 - F_m + F_m / S_m] < 1 / [1 - F_m] \text{ c.v.d. } (S_m \rightarrow \infty)$$

Istruzioni non accelerate

Se il tempo di esecuzione delle istruzioni accelerate va all' $\infty$  il tempo di esecuzione diventa il tempo di esecuzione delle istruzioni non accelerate soltanto.

$$\text{Esempio precedente: } T_{\text{new}} = 100 * (1 - 0.9 + 0.9/2) = 55$$



## Esempio 2



### Esempio:

Si consideri un miglioramento che consente un funzionamento **10** volte più veloce rispetto alla macchina originaria, ma che sia utilizzabile solo per il **40%** del tempo. Qual è il guadagno complessivo che si ottiene incorporando detto miglioramento?

$$\text{Speedup}_{\text{globale}} = 1 / [1 - F_m + F_m / S_m]$$

$$\text{Frazione}_{\text{migliorato}} = 0.4$$

$$\text{Speedup}_{\text{migliorato}} = 10$$

$$\text{Speedup}_{\text{globale}} = 1.56$$



## Esempio - 3



Supponiamo di potere aumentare la velocità della CPU della nostra macchina di un fattore 5 (senza influenzare le prestazioni di I/O) con un costo 5 volte superiore.

Assumiamo inoltre che la CPU sia utilizzata per il 50% del tempo ed il rimanente sia destinato ad attesa per operazioni di I/O. Se la CPU è un terzo del costo totale del computer è un buon investimento da un punto di vista costo/prestazioni, aumentare di un fattore cinque la velocità della CPU?

$$\text{Speedup}_{\text{globale}} = 1.67 \quad \text{Incremento di costo} = 2.33$$

L'incremento di costo è quindi più grande del miglioramento di prestazioni: la modifica *non* migliora il rapporto costo/prestazioni.



## Esempio – speedup dovuto a vettorializzazione



Si deve valutare un miglioramento di una macchina per l'aggiunta di una modalità vettoriale. La computazione vettoriale è 20 volte più veloce di quella normale. La *percentuale di vettorializzazione* è la porzione del tempo che può essere spesa usando la modalità vettoriale.

- Disegnare un grafico che riporti lo speedup come percentuale della computazione effettuata in modo vettoriale.
- Quale percentuale di vettorializzazione è necessaria per uno speedup di 2?
- Quale per raggiungere la metà dello speedup massimo?

La percentuale di vettorializzazione misurata è del 70%. I progettisti hardware affermano di potere raddoppiare la velocità della parte vettoriale se vengono effettuati significativi investimenti. Il gruppo che si occupa dei compilatori può incrementare la percentuale d'uso della modalità vettoriale.

- Quale incremento della percentuale di vettorializzazione sarebbe necessario per ottenere lo stesso guadagno di prestazioni?
- Quale investimento raccomanderebbe?

A.A. 2007-2008

29/41

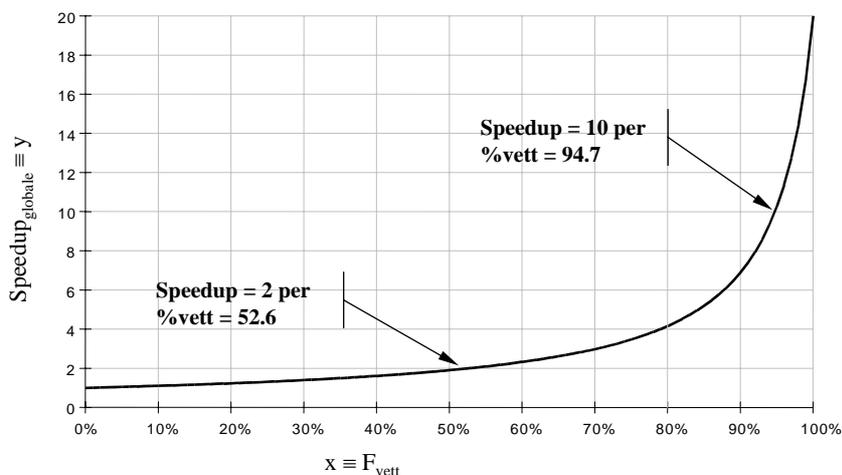
<http://homes.dsi.unimi.it/~borgnese>



## Curva di speed-up



$$\text{Speedup}_{\text{globale}} \equiv y = 1 / [1 - x + x / 20] = 20 / (20 - 19x) \quad x \equiv F_{\text{vett}}$$



A.A. 2007-2008

30/41

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



## Speed-up dovuto a HW



$$\text{Speedup}_{\text{original}} = 1 / [1 - 0.7 + 0.7 / 20] = 1 / (1 - 0.7 * 19 / 20) = 2,9851$$

$$\text{Speedup}_{\text{HW}} = 1 / [1 - 0.7 + 0.7 / 40] = 1 / (1 - 0.7 * 39 / 40) = 3,1496$$

$$\text{Speedup}_{\text{compiler}} = 3,1496 = 1 / [1 - x + x / 20] \rightarrow F_{\text{vettoriale}} = 71,84\%$$



## Sommario



Cosa vuol dire valutare le prestazioni

Benchmark

Miglioramento delle prestazioni

**Valutazione delle prestazioni del sistema di memoria**



## Principio 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.

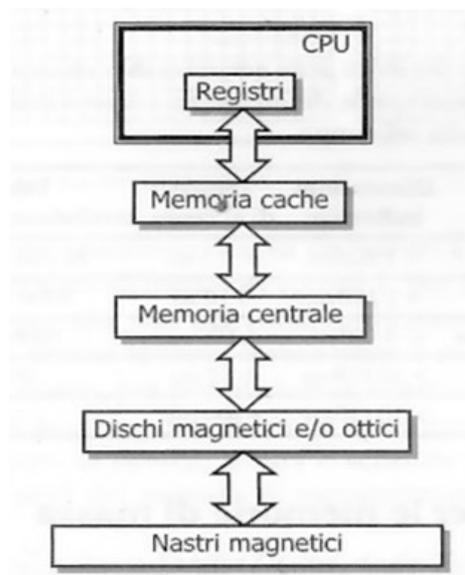
Basandosi sul passato recente del programma, è possibile predire con ragionevole accuratezza quali dati e istruzioni userà nel prossimo futuro.

*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.



## Gerarchia di memorie





## Valutazione prestazioni memoria



Obiettivo principale della gerarchia di memoria è incrementare le prestazioni => diminuire la velocità di accesso sia in caso di HIT che di MISS.

*Cosa succede in caso di MISS?*

**HIT\_TIME** Tempo di accesso al livello superiore (che comprende anche il tempo necessario per determinare se l'accesso ha avuto successo oppure fallisce).

**MISS\_PENALTY** è composto da:

TEMPO DI ACCESSO per accedere alla prima parola del blocco dopo che è stato rilevato il fallimento.

TEMPO DI TRASFERIMENTO per trasferire le altre parole del blocco al livello superiore.

**MISS\_TIME = > HIT\_TIME + MISS\_PENALTY**



## Tempo medio di accesso alla memoria



TEMPO DI ACCESSO. E' legato al tempo di accesso del livello inferiore di memoria.

TEMPO DI TRASFERIMENTO. E' legato alla larghezza di banda del canale di comunicazione tra i due livelli di memoria (e.g. bus).

Il tempo medio di accesso alla memoria sarà:

**$T_{\text{medio}} = \text{HIT\_RATE} * \text{HIT\_TIME} + \text{MISS\_RATE} * \text{MISS\_TIME} =$**

**$\text{HIT\_TIME} * \text{HIT\_RATE} + \text{MISS\_RATE} * (\text{HIT\_TIME} + \text{MISS\_PENALTY}) =$**

**$\text{HIT\_TIME} * (\text{HIT\_RATE} + \text{MISS\_RATE}) + \text{MISS\_RATE} * \text{MISS\_PENALTY} =$**

**$\text{HIT\_TIME} + \text{MISS\_RATE} * \text{MISS\_PENALTY}$**



## Impatto di una memoria cache



Il tempo di CPU è composto dal tempo richiesto dalla CPU per eseguire il programma e dal tempo che la CPU trascorre in attesa di risposta dal sottosistema di memoria.

$$T_{CPU} = (\#Cicli \text{ della CPU in esecuzione} + \#Cicli \text{ di stallo}) * T_{Clock}$$

### *Ipotesi:*

- Tutti gli stalli di memoria sono dovuti al fallimento di accesso alla cache.
- I cicli di clock utilizzati per un accesso alla cache riuscito (HIT) sono inclusi nei cicli di clock della CPU in esecuzione.



## Impatto di una memoria cache



$$\#Cicli\_clock\_stallo = \#Accessi\_Memoria * MISS\_RATE * MISS\_PENALTY$$

$$\begin{aligned} \text{Tempo}_{CPU\_Programma} &= (\#Cicli\_clock + \#Cicli\_clock\_stallo) * T_{clock} = \\ &\#Istruzioni * CPI_{exec} * T_{clock} + \#Cicli\_clock\_stallo * T_{clock} \end{aligned}$$

$$\begin{aligned} CPI_{con\_cache} &= CPI_{exec} + \#Cicli\_clock\_stallo / \#Istruzioni = \\ &CPI_{exec} + (\#Accessi\_memoria / \#Istruzioni) * MISS\_RATE * MISS\_PENALTY \end{aligned}$$

**Caso ideale:** (100% HIT, 0% MISS):  $CPI_{con\_cache} = CPI_{exec}$

**Caso senza cache:** (100% MISS):  $CPI_{senza\_cache} =$   
 $CPI_{exec} + (\#Accessi\_memoria / \#Istruzioni) * MISS\_PENALTY$



## Esercizio su cache



Si consideri il VAX-11/780. La MISS\_PENALTY è di 6 cicli di clock, mentre tutte le istruzioni impiegano 8.5 cicli di clock se si ignorano i MISS (stalli della memoria). Ipotizzando un MISS\_RATE dell'11% e che vi siano in media 2 riferimenti alla memoria per ogni istruzione,

⇒ Qual è l'impatto sulle prestazioni quando viene inserita la cache reale rispetto ad una cache ideale?

⇒ Qual è l'impatto sulle prestazioni tra il caso di cache reale e senza inserimento della cache?



## Soluzione esercizio su cache



Dati di ingresso: MISS\_PENALTY=6 CPI<sub>exec</sub>=8.5 MISS\_RATE=0,11  
#Accessi\_memoria/#Istruzioni = 2

$$CPI_{con\_cache} = 8,5 + 0,11 * (2 * 6) = 9,82$$

$$CPI_{con\_cache\_ideale} = 8,5 + 0 * (2 * 6)$$

$$CPI_{senza\_cache} = 8,5 + 1 * (2 * 6) = 20,5$$

$$\text{Perdita in prestazioni (speed-up): } CPI_{con\_cache\_ideale} / CPI_{con\_cache} \Rightarrow 8,5 / 9,82 = 0,865$$

$$\text{Guadagno in prestazioni (speed-up): } CPI_{senza\_cache} / CPI_{con\_cache} \Rightarrow 20,5 / 9,82 = 2,087$$



## Sommario



Cosa vuol dire valutare le prestazioni

Benchmark

Miglioramento delle prestazioni

Valutazione delle prestazioni del sistema di memoria