

# Architetture a Confronto

# Prestazioni dei calcolatori (1)

Gli elementi che consentono di calcolare

$$T_{esecuzione} = T_{clock} \cdot \left( \sum_{i=0}^n N_i \cdot CPI_i \right)$$

**$T_{clock}$  periodo di clock della macchina:** dipende dalla tecnologia e dalla organizzazione del sistema.

**$CPI_i$  numero di clock per istruzione di tipo  $i$ .** Cioè il numero di clock occorrenti affinché venga eseguita l'istruzione di tipo  $i$ . Se il set di istruzioni contiene istruzioni di tipo semplice, CPI assumerà valori ridotti; se invece le istruzioni eseguono operazioni più complesse si avranno CPI elevati.

**$N_i$  numero di istruzioni di tipo  $i$**  (somme, sottrazioni, salti, ecc.). Dipende dal set di istruzioni e dalla qualità dei compilatori (per qualità dei compilatori si intende la capacità dei compilatori di ottimizzare il programma, ad es. eliminare le istruzioni non necessarie). Se il set di istruzioni contiene istruzioni di tipo semplice avremo un  $N_i$  più elevato: caso opposto per istruzioni di tipo più complesso che permetteranno di ridurre  $N_i$ . Infatti con un set di istruzioni di tipo semplice, occorreranno più istruzioni per svolgere le stesse funzioni.

# Prestazioni dei calcolatori (2)

**IPS** (Instructions **P**er **S**econd) è una misura che serve per valutare le prestazioni di un sistema; indica quante istruzioni vengono eseguite al secondo. Viene comunemente riportata come **MIPS** (**M**illion **I**nstructions **P**er **S**econd):

$$MIPS = \frac{\text{Frequenza del clock}}{10^6 \cdot CPI}$$

Questa misura è affetta da una forte approssimazione:

- Non tiene conto delle possibilità offerte dal set di istruzioni (quanto più il set di istruzioni di una certa macchina è potente tanto più è possibile ridurre la lunghezza dei programmi che vengono eseguiti su di essa),
- Non tiene conto delle percentuali delle diverse istruzioni all'interno di programmi reali (una macchina più veloce in un programma può essere più lenta in un altro).
- Non tiene conto dell'ampiezza dei bus, della presenza di cache o altre tecniche (che si vedranno nel seguito) che ottimizzano i tempi di esecuzione.

Una stima di IPS viene solitamente ottenuta attraverso il **benchmark Dhrystone**, che contiene solo operazioni in aritmetica intera. Un benchmark è un insieme di programmi di prova rappresentativi di una particolare classe di applicazioni.

**FLOPS** (**F**loating point **P**er **S**econd) È una misura che è significativa se si intende utilizzare un programma dove la maggior parte delle operazioni sono di tipo floating point. Tipicamente riportata come MFLOPS, GLOPS, o TFLOPS.

**LINPACK benchmarks:** Misurano il numero di FLOPS a 64 bit (principalmente moltiplicazioni e somme) per risolvere sistemi di equazioni lineari densi  $n \times n$ . Un'implementazione parallela chiamata HPL, scritta in C, viene utilizzata per tenere aggiornata la lista dei TOP500 calcolatori (<https://www.top500.org>).

# MIPS a confronto

Processore	Anno	(Dhrystone) MIPS	Frequenza (MHz)
Pencil and Paper	1892	$1.19 \times 10^{-8}$	-
Zuse 1	1938	$4.24 \times 10^{-8}$	-
ENIAC	1946	0.00289	-
IBM System/370 model 158	1972	0.64	9
Intel 8080	1974	0.29	2
Motorola 68000	1979	1	8
Intel 386DX	1988	9	25
Intel 486DX	1992	54	66
PowerPC 600s (G2)	1994	35	33
ARM 7500FE	1996	36	40
Intel Pentium Pro	1996	541	200
Intel Pentium III	1999	1354	500
AMD Athlon	2000	3561	1200
Pentium 4 Extreme Edition	2003	9726	3200
ARM Cortex A8	2005	2000	1000
Xbox360 IBM "Xenon" Triple Core	2005	6400	3200
IBM Cell All SPEs	2006	12096	3200
AMD Athlon FX-57	2005	12000	2800
AMD Athlon FX-60 (Dual Core)	2006	18938	2600
Intel Core 2 Extreme QX6700	2006	57063	3330
Intel Core i7 Extreme 965EE	2008	76383	3200
AMD Phenom II X6 1090T	2010	68200	3200
Intel Core i7 Extreme Edition 980X	2010	147600	3300
Intel Core i7 5960X	2014	238310	3000
Intel Core i7 6950X	2016	317900	3000
AMD Ryzen 7 1800X	2017	304510	3600

# Classificazione dei calcolatori

(parametri dimensionali al 2018)

## Computer usa e getta

Cartoline d'auguri musicali, RFID

*centesimi €, alcuni MIPS*



## Sistemi embedded

in SmartCard, Orologi, Automobili, Elettrodomestici, Hi-Fi, Lettori audio/video, giocattoli, apparati medicali

*1..50€, 1..1000 MIPS*



## SmartPhone e Tablet

Applicazioni mobili di ogni genere

*100..1000€, 1..500 GFLOPS (CPU)*



# Classificazione dei calcolatori (2)

## Console da gioco

Videogiochi *100..500€, 250..6000 GFLOPS (GPU)*



## Personal Computer (PC)

Desktop o Portatili

*500..2000€, 50..1000 GFLOPS (CPU)*



## Server (e Workstation)

Gestione archivi, centralizzazione servizi, multiutenza, HPC (High Performance Computing), Deep Learning.

*5..100K€, 1..500 TFLOPS (CPU + GPU)*

# Classificazione dei calcolatori (3)

## Cluster (raggruppamento) di più Server



Per scalare verso l'alto prestazioni (in genere il livello di concorrenza) di singoli Server.

*10K€..1M€, 1..500 TFLOPS*

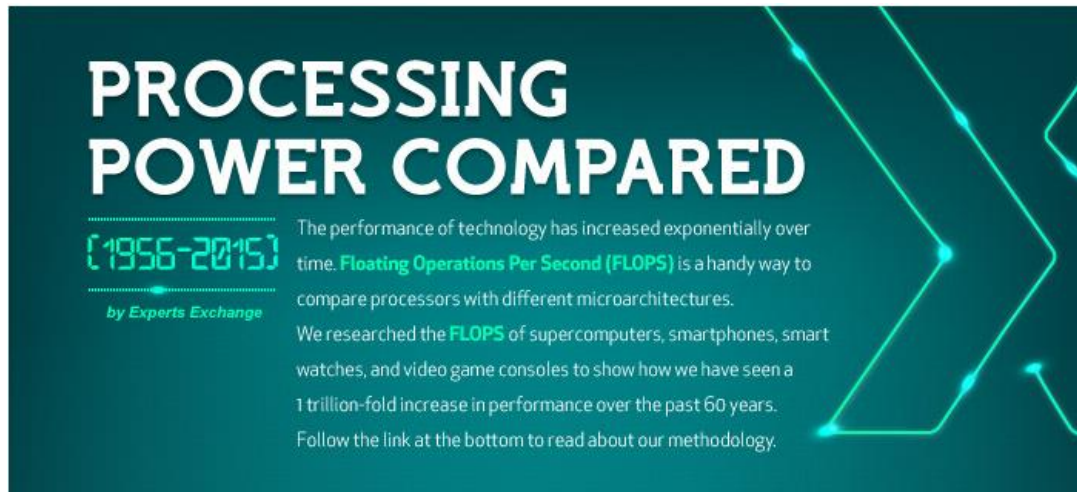
## Supercalcolatore

Elevata potenza di calcolo per applicazioni scientifiche

*50M€..250M€, 500 TFLOPS..93 PFLOPS*



# Confronto Potenze di Calcolo



<http://pages.experts-exchange.com/processing-power-compared>  
(valido al 20/02/2018)



# Come misurare il miglioramento delle prestazioni

Poter valutare le prestazioni di un calcolatore permette di calcolare:

- La bontà di un investimento
- L'adeguatezza del calcolatore al tipo di applicazione

$$\text{Prestazione } P = \frac{1}{T_{\text{esecuzione}}}$$

$$\text{Speed up} = \frac{P_A - P_B}{P_B} = \frac{T_B - T_A}{T_A} \quad (\text{del sistema A rispetto a B})$$

La **legge di Amdhal** afferma che il miglioramento delle prestazioni che si ottiene in un sistema di elaborazione accelerando un qualsiasi sottoinsieme del calcolatore è proporzionale alla percentuale di tempo per cui quel sottoinsieme è utilizzato.

$$T_{\text{es. finale}} = \frac{p \cdot T_{\text{es. iniziale}}}{a} + (1 - p) \cdot T_{\text{es. iniziale}}$$

Dove  $p$  è la percentuale di tempo utilizzata dal sottoinsieme migliorato, e  $a$  è il fattore di accelerazione.

**Esempio:** Modificando l'unità che esegue le operazioni in virgola mobile (unità floating point) si riduce a 1/10 il tempo necessario per eseguire le stesse. Si supponga che del tempo di esecuzione solo il 40% venga utilizzato per eseguire operazioni in virgola mobile.

$$T_{\text{es. finale}} = \frac{0.4 \cdot T_{\text{es. iniziale}}}{10} + (1 - 0.4) \cdot T_{\text{es. iniziale}} = 0.64 \cdot T_{\text{es. iniziale}}$$

$$\text{Speed up} = \frac{T_{\text{es. iniziale}} - 0.64 \cdot T_{\text{es. iniziale}}}{0.64 \cdot T_{\text{es. iniziale}}} = 0.56 = 56\%$$

# RISC VS CISC

La velocità di esecuzione di un'istruzione all'interno della CPU determina in larga misura la velocità della CPU ed è da sempre oggetto di discussione tra due correnti di pensiero.

**CISC** (**Complex Instruction Set**): i sostenitori di questa idea ritengono che l'Instruction Set di un calcolatore debba contenere quante più istruzioni possibili, anche se ognuna di queste richiede più cicli di data path, poiché ciò permette di creare macchine più potenti (es. 8086..Pentium).

**RISC** (**Reduced Instruction Set**): i sostenitori di questa idea ritengono che ogni istruzione dell'Instruction Set debba essere eseguita in un solo ciclo (o comunque pochi cicli) di data path: saranno necessarie più istruzioni RISC per ottenere lo stesso risultato di una istruzione CISC, ma il sistema risulterà comunque più veloce poiché non sarà più necessario interpretare le istruzioni (es. PowerPC).

Le macchine CISC hanno dominato il mercato negli anni '70 e '80 mentre attualmente la tecnologia è fortemente orientata verso soluzioni RISC.

La superiorità di una soluzione rispetto all'altra è comunque un fatto relativo che dipende da fattori di mercato e fattori tecnologici: un radicale innovazione tecnologico potrebbe muovere nuovamente l'ago della bilancia a favore delle macchine CISC.

# Principi di progettazione RISC

Tutte le istruzioni del livello ISA vengono eseguite direttamente dall'hardware: in particolare l'unità di controllo all'interno della CPU è **cablata** (realizzata con un circuito digitale sequenziale ovvero una macchina a stati) e non **microprogrammata** (istruzioni associate a microprogrammi composti da micro-istruzioni). Eliminando il livello di interpretazione aumenta la velocità della maggior parte delle istruzioni.

Ottimizzare la velocità con la quale le istruzioni vengono mandate al primo stadio di esecuzione: anche in presenza di operazioni complesse la velocità del processore è determinata dal numero di istruzioni “iniziate” per secondo (**MIPS: Millions of Instructions Per Second** ). Nella maggior parte dei processori moderni si fa ampio uso di pipelining e parallelismo.

Le istruzioni dovrebbero essere facilmente decodificabili: la velocità di esecuzione delle istruzioni dipende anche dal tempo necessario a identificare le risorse necessarie a eseguirle. Questo processo può essere velocizzato utilizzando istruzioni con **struttura regolare e lunghezza fissa**.

Solo le istruzioni Load e Store dovrebbero contenere indirizzi di memoria: in questo modo le rimanenti istruzioni utilizzeranno solo operandi contenuti nei registri eliminando i tempi morti dovuti ai ritardi nella lettura dei dati dalla memoria. Inoltre le operazioni di Load e Store potranno essere parallelizzate in quanto rappresenteranno operazioni a se stanti.

Disporre di molti registri: i registri sono una risorsa fondamentale per ridurre i tempi di accesso alla memoria.

# Come aumentare le prestazioni

I **progressi tecnologici** (aumento del **numero di porte** per unità di superficie, maggiori **velocità di commutazione** delle porte, minor consumo di corrente e quindi **dissipazione di calore**, ...) costituiscono una fonte primaria per il miglioramento delle prestazioni delle CPU.

In ogni caso, **indipendentemente** dal cambiamento di tecnologia, esistono diversi altri modi per **migliorare le prestazioni** di una microarchitettura, e ognuno di questi deve essere valutato considerando il **trade-off prestazioni/costi**:

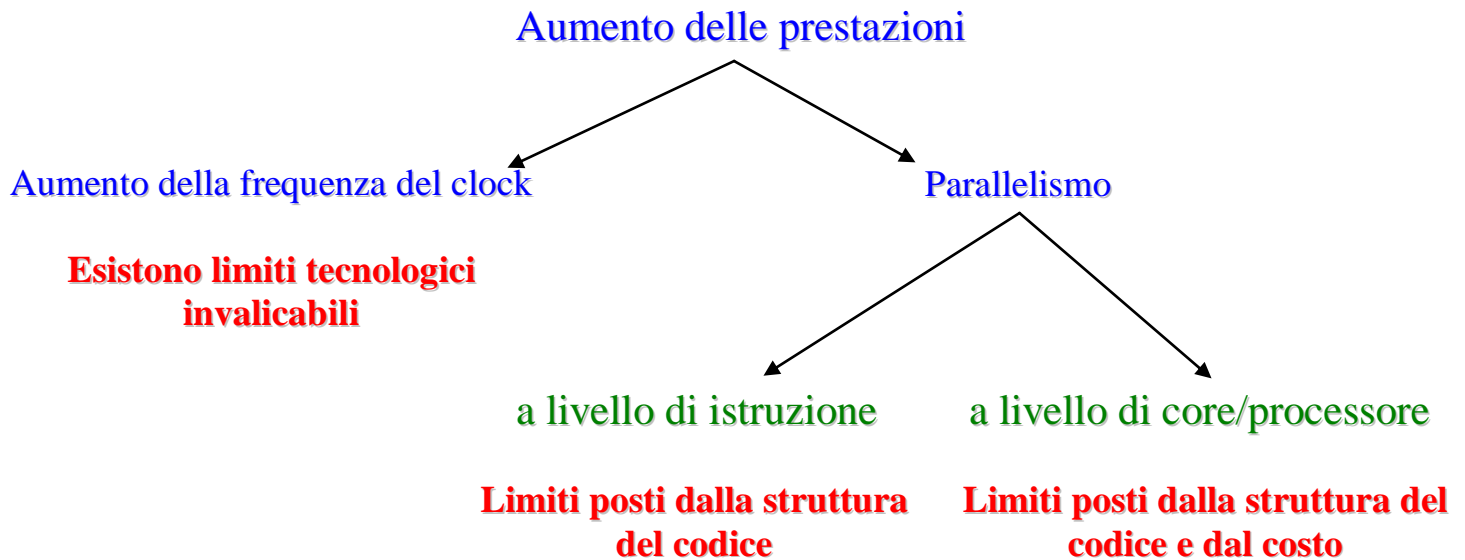
- **Ridurre il numero di cicli** (**path di esecuzione**) necessari per implementare le istruzioni ISA (Livello della microarchitettura).
- Semplificare l'organizzazione per consentire di **aumentare la frequenza di clock**. Limiti fisici.
- **Parallelismo**: sovrapporre l'esecuzione delle istruzioni (architetture superscalari, pipelining e CPU in parallelo)

Analizzeremo inoltre alcune tecniche utilizzate dalle CPU di recente fabbricazione per **migliorare ulteriormente** le prestazioni:

- **Predizione di salto**
- **Esecuzione fuori ordine**
- **Esecuzione speculativa**

# Parallelismo

Per migliorare le prestazioni di un calcolatore si possono seguire più strade:



Nel **parallelismo a livello di istruzione** più istruzioni vengono eseguite contemporaneamente all'interno della stessa CPU tramite tecniche di pipelining e processori superscalari

Nel **parallelismo a livello di core o di processore** più core/CPU cooperano per la soluzione dello stesso problema.

# Aumento frequenza di clock

Nel ventennio 1980-2000 le frequenze di clock sono aumentate di **oltre 1000 volte**: da circa 1-10 MHz (IBM XT) a oltre 3 GHz (Pentium 4). La potenza di calcolo è aumentata **ben più di 1000 volte**, e ciò dimostra che aumentare la frequenza di clock **è uno solo** dei modi per rendere un calcolatore più potente.

Nel 2004 con il Pentium 4 sono state raggiunte frequenze di circa **3.8 GHz**, che negli anni successivi sono state superate solo di poco.

Nel futuro la **frequenza di clock non potrà aumentare significativamente** (a meno di scoperte scientifiche straordinarie); infatti siamo oramai giunti nei pressi dei **limiti fisici**:

- le alte frequenze **creano disturbi** e aumentano il **calore** da dissipare
- ci sono **ritardi nella propagazione del segnale**
- **bus skew** (i segnali su linee diverse viaggiano a velocità diverse): problemi di sincronizzazioni.
- nelle architetture superscalari, ci sono anche **limitazioni dovute alla suddivisione delle operazioni** (meglio chiarito nel seguito)

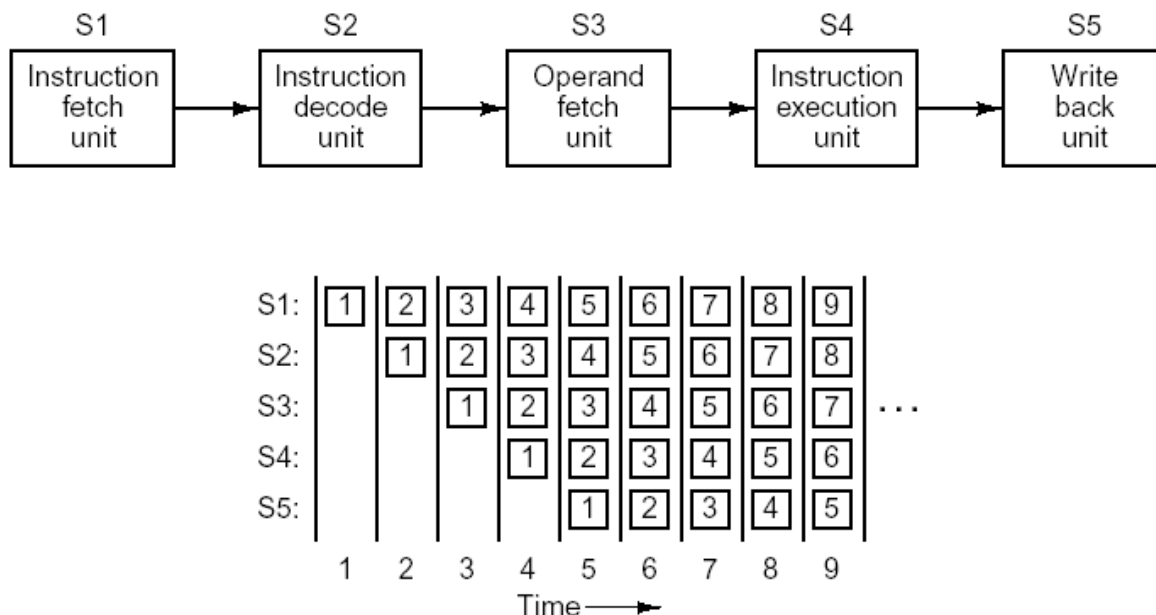
## Esempio:

*in un nanosecondo (frequenza 1 GHz) che distanza può percorrere un impulso digitale ?*

Anche se questo viaggiasse alla **velocità della luce** (**300.000** Km/sec; in realtà la propagazione di elettroni nei conduttori non è elevata come quella nel vuoto) il segnale potrebbe propagarsi in un **ns** di soli  $3 \cdot 10^8 \cdot 10^{-9}$  metri ovvero di circa **33 centimetri** (ovvero circa della lunghezza di una scheda madre !)

# Pipelining (1)

Con la tecnica del pipelining l'esecuzione di ogni istruzione viene suddivisa in più fasi (chiamate **stadi**) ognuna delle quali viene gestita da un hardware dedicato.

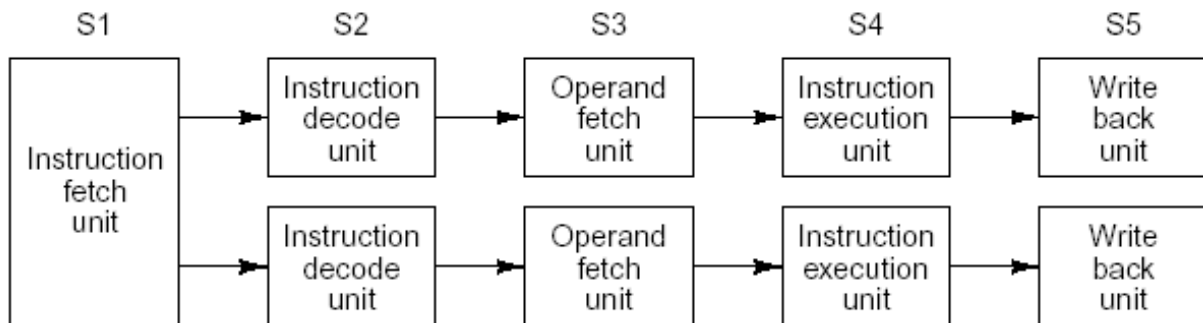


- Se il ciclo di clock della macchina è di 2 nsec sono necessari 10 nsec per completare l'esecuzione della prima istruzione (nessun risparmio rispetto all'assenza di pipelining), **ma una volta riempita la pipeline, si completerà una istruzione ogni 2 nsec.**
- È necessario garantire che le istruzioni **non siano in conflitto** ossia che non dipendano l'una dall'altra. In caso contrario sarà necessario attendere il completamento dell'istruzione precedente prima di avviare la successiva.
  - 1)  $x = a + b;$
  - 2)  $y = x + c;$

L'istruzione 2 ha come operando il risultato dell'istruzione 1. Finché la pipeline per 1 non termina 2 non può essere terminata.
- Le pipeline vengono normalmente utilizzate su macchine RISC tuttavia anche Intel adotta questa tecnica a partire dal processore 486 che è dotato di una pipeline a 5 stadi.

# Pipelining (2)

È anche possibile avere più di una pipeline:



- Vengono lette due istruzioni alla volta che vengono eseguite su pipeline diverse.
- Anche in questo caso è necessario gestire eventuali conflitti tra le istruzioni.
- Il processore Intel Pentium utilizza due pipeline: la prima (**pipeline u**) esegue qualsiasi istruzione, la seconda (**pipeline v**) esegue solo istruzioni semplici su interi.

## Pipelining e Frequenza della CPU

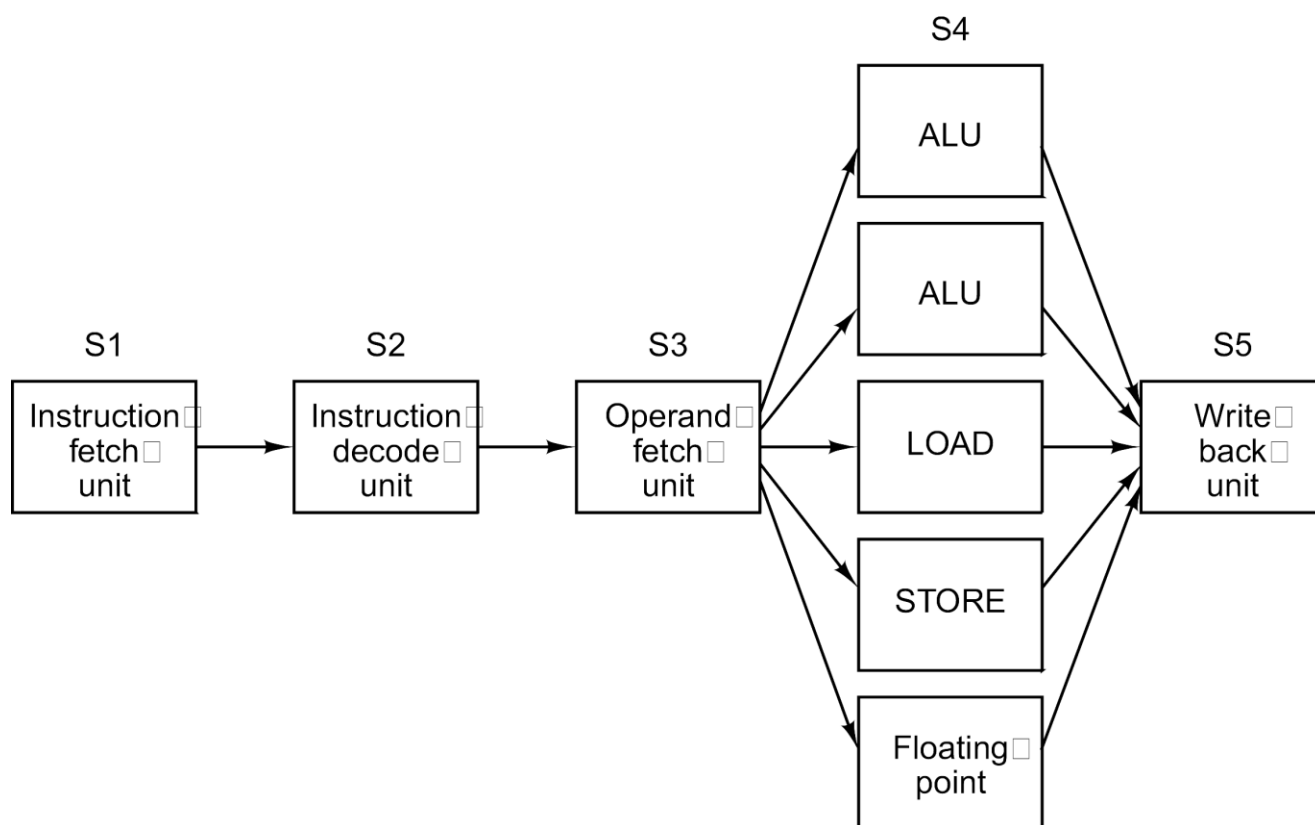
La frequenza di lavoro massima del sistema dipende dal tempo impiegato dalla fase più onerosa. Vedi:

<https://software.intel.com/en-us/blogs/2014/02/19/why-has-cpu-frequency-ceased-to-grow>



# Architetture superscalari

Una soluzione alternativa è quella di avere una sola pipeline dotata di unità funzionali multiple (**architettura superscalare**)



Questa soluzione porta a un aumento delle prestazioni se il tempo di esecuzione delle unità funzionali dello stadio S4 è superiore al tempo di esecuzione degli altri stadi. In questo caso la maggior velocità dei rimanenti stadi sarebbe compensata dal parallelismo.

Una soluzione simile a quella rappresentata in figura è implementata nel processore Intel Pentium II.

# Predizione di salto

Le architetture con **pipeline** (praticamente **tutte** le architetture **moderne**) funzionano molto bene se il codice viene eseguito **sequenzialmente senza salti**. Infatti, a seguito di un salto che determina un **cambiamento** nella sequenza di **istruzioni** da eseguire, una parte del **lavoro già eseguito** (stadi S1, S2, ed S3 nell'esempio precedente) viene **buttato** !

La figura mostra un esempio di codice (C) e il corrispondente assembler: **sapere se  $i$  sarà uguale a 0** sarebbe di grande aiuto per eseguire il pre-fetching corretto !

if (i == 0)	CMP i,0; compare i to 0
k = 1;	BNE Else; branch to Else if not equal
else	Then: MOV k,1; move 1 to k
k = 2;	BR Next; unconditional branch to Next
	Else: MOV k,2; move 2 to k
	Next:

Per **ottimizzare** le operazioni di **pre-fetching**, le CPU moderne possono utilizzare tecniche di **predizione di salto**, che a fronte di **istruzioni di salto condizionale** cercano di prevedere (*indovinare?*) se il programma salterà oppure no:

- **predizione statica** vengono utilizzati **criteri di buon senso** derivati dallo studio delle abitudini dei programmatori e del comportamento dei compilatori; ad esempio: **assumiamo che tutti i salti condizionali all'indietro vengano eseguiti**.
- **predizione dinamica** vengono mantenute **statistiche** (tabelle interne) circa la **frequenza** con cui i recenti salti condizionali sono stati **eseguiti**. Sulla base di queste statistiche la CPU esegue la predizione. *Esempio:* se nelle ultime  $n$  volte il salto corrispondente a una data istruzione è stato eseguito almeno  $n/2$  volte, allora probabilmente sarà ancora eseguito.

# Esecuzione fuori ordine ed esecuzione speculativa

Gran parte delle CPU moderne sono sia **pipelined** sia **superscalari**. La progettazione di una CPU è più semplice se tutte le **istruzioni vengono eseguite nell'ordine** in cui vengono lette; ciò però non sempre porta a prestazioni ottimali per via delle **dipendenze** esistenti tra le varie operazioni:

*Infatti se un'istruzione A richiede un valore calcolato da un'istruzione precedente B, A non può essere eseguita fino a che l'esecuzione di B non è terminata.*

Nel tentativo di ovviare a questi problemi e di **massimizzare le prestazioni**, alcune CPU consentono di **saltare** temporaneamente **istruzioni** che hanno **dipendenze** (lasciandole in attesa) e di passare ad eseguire istruzioni successive non dipendenti. Questo tipo di tecnica prende il nome di **esecuzione fuori ordine** e deve comunque garantire di ottenere esattamente gli **stessi risultati** dell'esecuzione ordinata.

Un'altra tecnica correlata all'esecuzione fuori ordine, prende il nome di **esecuzione speculativa**: essa consiste nell'**anticipare** il più possibile l'esecuzione di alcune **parti del codice** (**gravose**) prima ancora di sapere se queste serviranno.

L'anticipazione (**hoisting** = atto di alzare / portare in alto il codice) può essere ad esempio relativa a un'operazione floating-point o a una lettura da memoria, ... Ovviamente **nessuna** delle **istruzioni anticipate** deve produrre effetti **irrevocabili**.

In realtà, molto spesso l'**hardware** della CPU **non è in grado** da solo di anticipare istruzioni, se non nei casi banali o se non con il supporto dei compilatori. In alcuni processori di **recente progettazione**, vengono previste **particolari istruzioni** o **direttive** che il compilatore può utilizzare per ottimizzare il comportamento della CPU.

# Bug “architetturali”

A inizio 2018 sono stati resi pubblici da ricercatori dell'Università di Graz e da Google importanti bug di sicurezza (noti come **Meltdown** e **Spectre**) che coinvolgono diverse CPU moderne (Intel, AMD, ARM) che fanno uso di Cache ed Esecuzione Fuori Ordine.



- [https://it.wikipedia.org/wiki/Meltdown\\_\(vulnerabilit%C3%A0\\_di\\_sicurezza\)](https://it.wikipedia.org/wiki/Meltdown_(vulnerabilit%C3%A0_di_sicurezza))
- [https://it.wikipedia.org/wiki/Spectre\\_\(vulnerabilit%C3%A0\\_di\\_sicurezza\)](https://it.wikipedia.org/wiki/Spectre_(vulnerabilit%C3%A0_di_sicurezza))

L'attacco **Meltdown** permette di leggere il contenuto di tutta la memoria fisica **eludendo** i vincoli che non consentono a un processo di vedere la memoria al di fuori della porzione a lui dedicata. La cosa è particolarmente grave nel caso di macchine virtuali e sistemi in cloud dove i dati di un'azienda potrebbero essere spiati da altri processi in esecuzione sulla stessa macchina fisica.

In sostanza:

- si generano accessi a locazioni proibite (che si vogliono spiare).
- l'esecuzione fuori ordine permette la lettura di queste locazioni prima ancora della verifica dei permessi (che avviene concorrentemente).
- una volta verificato che i permessi non consentono l'accesso si genera un'eccezione e non si rende disponibile il contenuto della memoria.
- nel frattempo però la lettura ha provocato effetti collaterali nella cache che non sono annullati.
- interrogando la cache in modo opportuno (se una riga è presente la risposta è molto più veloce) si può risalire al contenuto della cella di RAM.

Dettagli nel paper: <https://meltdownattack.com/meltdown.pdf>

# Architetture Pentium a confronto

Diverse **versioni** (identificate da un “**core**” con nome diverso, es: **Willamette, Northwood, Prescott, Extreme Edition** per **Pentium 4**) sono state prodotte da Intel nell’ambito di ciascuna famiglia Pentium. Le differenze tra le versioni sono talvolta anche molto significative. Nella tabella sono indicati i parametri della versione con core più potente nella rispettiva famiglia. Tutte le versioni hanno:

- **architettura interna** (registri principali) a **32 bit**.
- **bus dati** verso l’esterno di **64 bit** e **bus indirizzi** di **36 bit** (64 GB indirizzabili).
- supporto **IA-32**.

	<b>Freq. Core e Front Side Bus</b>	<b>Cache</b>	<b>Peculiarità architettura</b>	<b>Istruzioni speciali</b>
<b>Pentium (1993)</b>	Core: 233 MHz FSB: 66 MHz	L1: 32 KB	2 Pipeline Superscalare	MMX
<b>Pentium II (1997)</b>	Core: 450 MHz FSB: 100 MHz	L1: 32 KB L2: 512 KB	2 Pipeline Superscalare	MMX
<b>Pentium III (1999)</b>	Core: 1,4 GHz FSB: 100 MHz	L1: 32 KB L2: 512 KB	2 Pipeline Superscalare	MMX, SSE
<b>Pentium 4 (2000)</b>	Core: 3.8 GHz FSB: 4×266 MHz (Quad Pumped)	L1: 32 KB L2: 512 KB L3: 2 MB	<b>NetBurst:</b> 2 Pipeline 20 stadi 5 unità esecuzione Migliore predizione <b>Hyper Threading*</b> <b>NX-bit**</b>	MMX, SSE, SSE2, SSE3, EM64T***

\* **Hyper Threading** il processore è visto dal sistema operativo come due processori “virtuali”, sui quali allocare processi diversi evitando sovra-allocazione di risorse.

\*\* **NX-bit** (No eXecute) protezione contro l'esecuzione di codice malevolo: le sezioni di memoria contrassegnate con l'NX bit sono dedicate al deposito di soli dati, e le istruzioni non dovrebbero risiedervi.

\*\*\* **EM64T** (Extended Memory 64 Technology) implementazione Intel di x64 (ora denominato Intel 64).

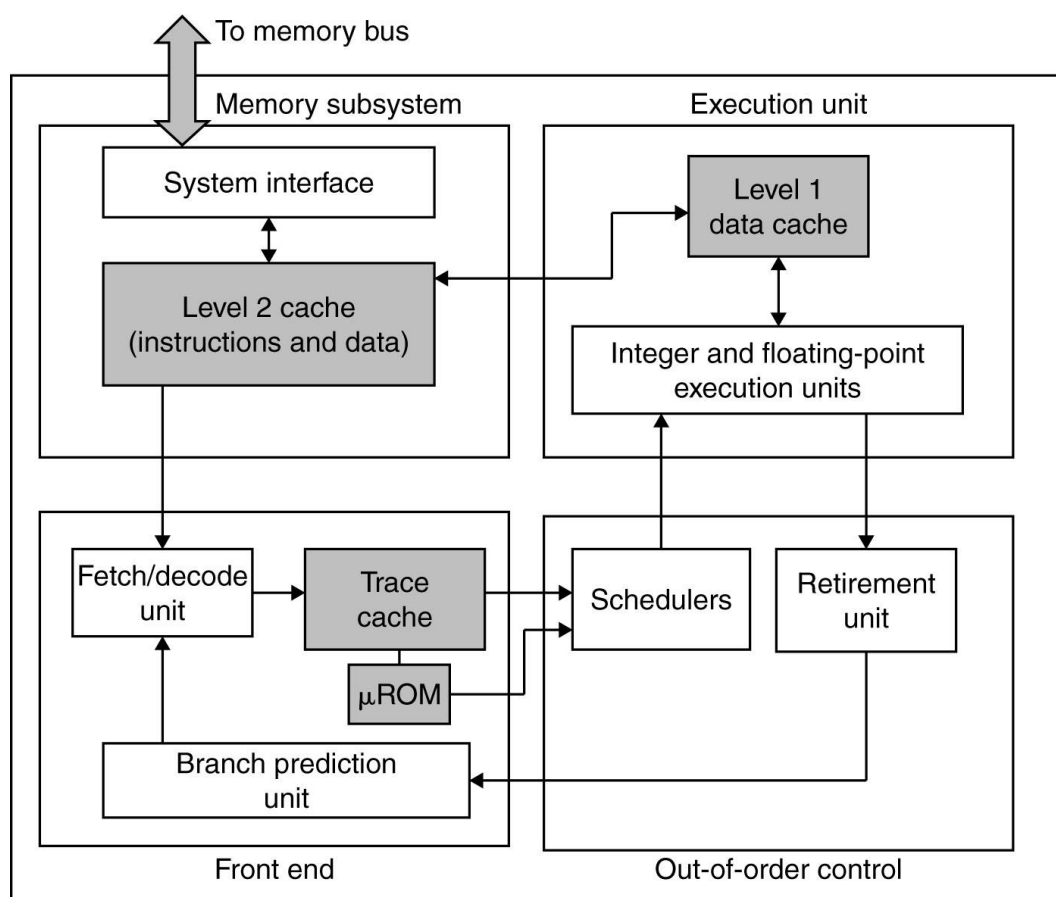
# Architettura Pentium 4

Pentium 4 è una CPU **Intel** completa a **32 bit**, (che appare una architettura **CISC** tradizionale, per problemi di compatibilità) e ha lo **stesso ISA (IA-32)** dei modelli precedenti compresi gli stessi registri istruzioni e standard floating point IEEE 754.

- supporta operazioni aritmetiche su **interi 32 bit** e **floating-point 64 bit**.
- i registri **visibili** sono **8**.
- la lunghezza delle **istruzioni (codifica)** è variabile da **1 a 17 byte**.

Al suo interno però contiene un **nucleo RISC moderno** che fa largo uso di pipelining e architettura superscalare.

Diverse revisioni architetturali: quella qui considerata, denominata **NetBurst**, rappresenta un forte cambiamento rispetto al passato, la precedente (denominata **P6**) era utilizzata sin dal Pentium Pro.



# Architettura Pentium 4 (2)

È costituito da quattro parti principali: il sottosistema di memoria, il front end, il controllo dell'esecuzione fuori sequenza e le unità esecutive.

## Sottosistema di memoria

- include cache unificata (dati e istruzioni insieme) di secondo livello L2 (la dimensione varia da 256KB a 2MB a seconda dei modelli). Le linee di cache hanno dimensione pari a 128 byte e in caso di fallimento nel reperire una informazione (miss) partono due trasferimenti, di 64 byte ciascuno, dalla memoria principale (RAM esterna alla CPU).

## Front end

- ha il compito di prelevare le istruzioni dalla cache L2 e di decodificarle nell'ordine del programma. La decodifica scompone ciascuna istruzione (CISC) in una sequenza di micro-operazioni (tipo RISC) successivamente eseguite dal cuore RISC del Pentium 4. Per la decodifica delle istruzioni più complesse si usa una micro-ROM che memorizza le equivalenze: Istruzione – Micro-operazioni. Le micro-operazioni vengono salvate nella cache delle tracce che è una cache di istruzioni di primo livello (L1), questo consente di evitare la decodifica per istruzioni successive identiche se queste si trovano già in cache. Fa parte di questo livello anche la predizione di salto.

## Controllo dell'esecuzione fuori sequenza

- Le micro-operazioni sono trasferite dalla cache delle tracce allo schedulatore che può mandarle in esecuzione anche fuori ordine (per ottimizzare le prestazioni). L'unità di ritiro garantirà poi che i risultati prodotti all'esterno siano equivalenti rispetto a una esecuzione in ordine.

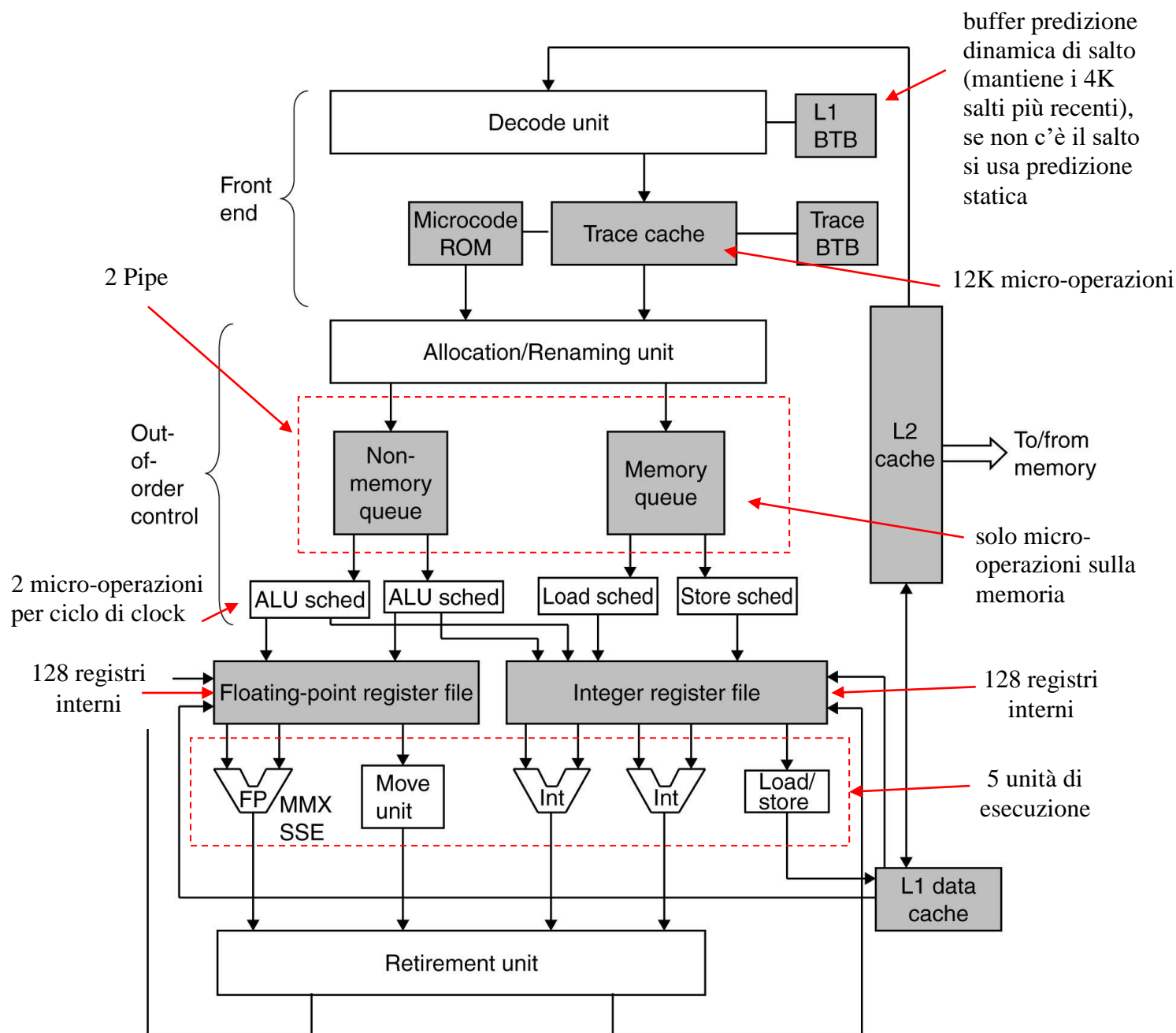
## Unità esecutive

- Le unità di esecuzione eseguono le operazioni tra interi, in virgola mobile o operazioni speciali (es. MMX). Gli operandi vengono reperiti tramite la cache dati L1.



# Architettura Pentium 4 (3)

La figura illustra l'architettura mettendo in evidenza il **pipelining** (20 stadi) e le componenti **superscalari** (5 unità funzionali a doppia frequenza di clock).



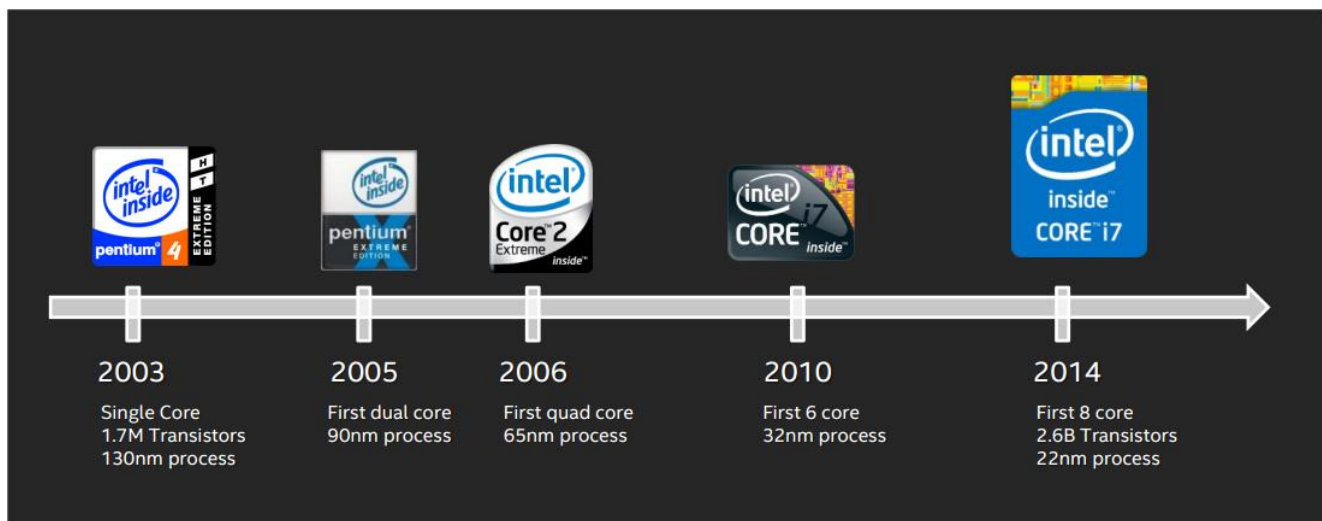


# Intel multi-core (1)

Non riuscendo più ad aumentare significativamente le prestazioni con l'aumento di frequenza, Intel ha deciso di **puntare al parallelismo** proponendo come successori del Pentium IV dapprima **Pentium D** (2005), poi a partire dal 2006: **Core 2**, **Core i3**, **Core i5**, **Core i7** e, più recentemente, **Core i9**.

Una **CPU multi core** (fino a 18 nei modelli sopracitati) unisce più processori indipendenti, le rispettive Cache e i cache controller in un singolo package (rendendo possibile tra l'altro un'implementazione "fisica" dell'Hyper Threading): consente di aumentare la potenza di calcolo senza aumentare la frequenza di lavoro, a tutto vantaggio del calore dissipato (**che aveva oramai superato i 100W per CPU single core**).

Con l'avvento dei processori della serie **Core 2** per la prima volta sia il mercato dei sistemi **portatili** che dei sistemi **desktop**, si basano su un **unico processore**.



*Per trarre pieno beneficio dalla tecnologia multi core i programmi devono essere pensati per un uso ottimale di multi-thread, in caso contrario essi impegneranno solo uno dei due core, lasciando l'altro pressoché inutilizzato.*

# Intel multi-core (2)

## Skylake-X (2017):

- fino a 18 core;
- supporto integrato DDR4-2666;
- 40 canali PCIe 3.0 integrati;
- Frequenza 2.6-4.3 GHz;
- Cache:
  - L1 (64 KB) per ogni core;
  - L2 (1MB) per ogni core;
  - L3 (1.375MB) per core;
- Istruzioni AVX512, SSE4.2, FMA3;
- 165 Watt;

In un'architettura multi-core moderna con istruzioni SIMD, il calcolo delle prestazioni teoriche è molto complesso:

CPU GHz \* number of cores \* vector ops (AVX) \* special instructions (FMA3)

che nel caso di questa CPU si istanzia con:

3.3 (Frequenza effettiva con tutti i core attivi) \* 18 \* 8 (AVX512) \* 2 (FMA3) = 950 GFLOPS

molto vicino alla prestazione di 977 GigaFlops misurata su Linpack nei primi test apparsi su Internet nel 2017.

## Sistemi multi-processore

Attenzione a non confondere sistemi **multi-core** (*più core nello stesso chip, che condividono il BUS*) e sistemi **multi-processore** (*chip separati*). La soluzione multi-processore è più costosa ma anche più potente.

Nella famiglia **Xeon**, erano da tempo previste soluzioni multi-processore:

- **Xeon DP** in ambito dual processor
- **Xeon MP** per sistemi a 4 o più vie

# Famiglie Intel x86

Famiglie di processori Intel con architettura x86 (aggiornato al 2017):

- Linea **Core**: CPU destinate a computer di fascia medio-alta e workstation. Le principali famiglie sono (Core i3, Core i5, Core i7, Core i9). Disponibili versioni fino a 18 core e 25MB di cache L3.



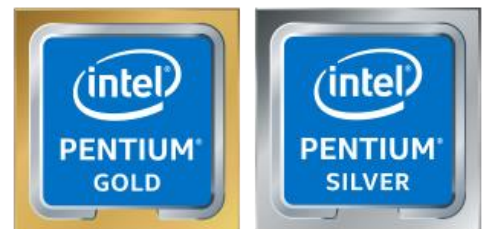
- Linea **Xeon**: CPU destinate a workstation e server di fascia alta, derivate dalla linea Core, ma con caratteristiche più avanzate, come il supporto per memoria ECC, numero di core più alto, memoria cache più ampia, supporto per sistemi multiprocessore. Disponibili versioni fino a 24 core e 60MB di cache L3 o 28 core e 39MB di cache L3.



- Linea **Atom**: CPU per dispositivi ultra-portatili (es. Netbook, schede per sistemi embedded), caratterizzate da piccole dimensioni e consumi molto contenuti.



- Linea **Pentium**: CPU destinate a computer di fascia medio-bassa, disponibili versioni fino a 4 core e 6MB di cache L2.



- Linea **Celeron**: CPU più economiche destinate a computer di fascia bassa, disponibili versioni fino a 4 core e 4MB di cache L2.



# Non solo Intel: le CPU AMD

**AMD** (**A**dvanced **M**icro **D**evelopments) è il secondo produttore mondiale di microprocessori (dopo Intel).

Fin dal 1975 AMD iniziò a produrre microprocessori **x86-compatibili** (partendo con una variante dell'Intel 8080 chiamata AMD8080) e instaurando una **competizione perenne** con Intel:

- un processore si dice **x86-compatibile** se in grado di supportare le istruzioni ISA dei processori x86 (i.e. correntemente **IA-32**). Tutti i programmi sviluppati per piattaforma IA-32 (i.e. piattaforma Windows) “girano” su microprocessori x86-compatibili.
- due processori **x86-compatibili**, **a seconda dell'architettura interna**, possono però avere prestazioni molto differenti ed eseguire lo stesso programma in tempi diversi (infatti anche quando operano alla stessa frequenza, il numero di cicli di clock per eseguire la stessa istruzione ISA può essere diverso).

Negli corso degli anni AMD **non si limita a produrre CPU gemelle** (e generalmente più economiche) rispetto a quelle di Intel, ma spesso **anticipa Intel** dal punto di vista tecnologico, costringendo quest'ultima a correre ai ripari.

CPU AMD di successo sono:

- la serie **AMD K6** (metà anni '90), la serie **Athlon** (1999)
- **Athlon 64** (2003) è il primo processore desktop con supporto 64 bit della famiglia x86. Il modello AMD64 viene presto “imitato” da Intel, che inserisce EM64T nel Pentium 4, per “ricucire lo strappo”.
- **Opteron** (2003, destinata alla fascia server, supporto 64 bit)
- **Phenom** (2007): quad core, **Phenom II** (2009): esa core
- **APU** (2011,...): **CPU + GPU** integrate → **PS4**, **XBOX ONE**.
- Le serie **Bulldozer**, **Jaguar**, **Puma** (2011-2015)
- Le serie **Zen** (2017, ...)

# CPU RISC non x86 (Server)

**SPARC** (Scalable Processor ARChitecture) è il nome di un'architettura "aperta" e non proprietaria per microprocessori RISC - big-endian (basata sul modello RISC originariamente sviluppato nell'università californiana di Berkley). L'architettura fu disegnata nel 1985 da Sun Microsystems; Workstation e server di Sun Microsystems e Fujitsu sono stati i maggiori fruitori di CPU Sparc. Nel corso degli anni l'architettura ha subito diverse revisioni, la modifica maggiore si è avuta con la versione 9 che ha introdotto la gestione dei dati a 64 bit. Al 2017, Oracle (che ha acquisito Sun nel 2010) è il principale fornitore di sistemi Sparc.

**Power** è un'architettura di microprocessori RISC (32 e 64 bit) creata da IBM. Nel 1991 l'alleanza Apple – IBM – Motorola (AIM) ne derivò un'implementazione di successo nota come PowerPC. Apple adottò PowerPC per il suo Machintosh; anche altri Sistemi Operativi (tra cui Windows NT) ben presto supportarono PowerPC, che a metà degli anni 90 era l'architettura più potente per personal computer. Sfortunatamente PowerPC non resse il confronto con la concorrenza (soprattutto per non compatibilità x86) e nel 2006 fu abbandonato anche da Apple. L'architettura Power è stata in seguito utilizzata come base per altri progetti non-x86 tra cui: Xenon (triple core dell'Xbox 360); Cell (cuore di Playstation 3). Al 2017, IBM propone sistemi con O.S. AIX basati su PowerPC.

**Alpha** è un'architettura di processori di tipo RISC (64 bit) sviluppata e prodotta dalla Digital Equipment Corporation (DEC) a partire dal 1992. Alpha, inizialmente supportato da diversi sistemi operativi (tra cui Windows NT, Unix e Linux), è stato poi supportato da HP (che nel frattempo aveva acquisito Compaq che a sua volta aveva acquisito DEC), per poi essere abbandonato pochi anni dopo il 2000. Alcune soluzioni architetturali innovative della piattaforma Alpha, che per alcuni anni le consentirono di primeggiare in quanto a prestazioni, sono state molto più tardi riprese per il progetto di nuove CPU a 64 bit (es. Itanium). Al 2017, HP propone sistemi con O.S. HP-UX basati su macchine Itanium.

In generale, al 2017, la quota di mercato di macchine server non-x86 è inferiore al 5% e non sembra destinata a crescere.



# IA-64

Caratteristiche peculiari di questi ISA (adottato da CPU **Itanium** e **Itanium II**) sono:

- **parallelismo a livello di istruzione** che è **esplicito** nelle istruzioni macchina piuttosto che determinato dal processore durante l'esecuzione. Questo tipo di parallelismo è noto come **EPIC** (Explicit Parallel Instruction Computing). In questo caso **è il compilatore (e non il processore) a dover capire quali istruzioni possono essere eseguite in parallelo** e a generare il codice macchina relativo, unendole in un'unica istruzione più grande, da cui il nome di Very Long Instruction Word (**VLIW**).
- Pertanto i processori EPIC non necessitano di circuiti complessi per l'esecuzione fuori sequenza.
- **sono previsti fino a 256 registri a 64 bit** (128 interi + 128 in virgola mobile nell'Itanium).
- **più pipeline parallele** (si prevedono implementazioni con 8 o più unità esecutive). Il pipelining può essere gestito a livello software.
- le istruzioni hanno un **formato a lunghezza fissa** di **41 bit**. Un **pacchetto** (bundle) contiene tre istruzioni. Il processore può caricare uno o più pacchetti per volta.
- l'utilizzo dei **predicati di salto** con la maggior parte delle istruzioni (concetto diverso dalla predizione di salto). Questo si realizza aggiungendo un **registro predicato** davanti alle istruzioni: durante l'esecuzione se il valore del predicato è FALSE l'istruzione non viene eseguita ma non si procede a nessun salto (dannoso per il pipelining) e si continua semplicemente con la successiva.
- **il caricamento speculativo** e l'esecuzione speculativa (lungo entrambi i rami di un salto). Sono previsti appositi formati (modificatori delle istruzioni) per la gestione efficiente dell'esecuzione speculativa.

# Cell: il cuore della PlayStation 3

**Cell** è il nome di una tipologia di processori multi-core sviluppati da **IBM** in cooperazione con **Sony** e **Toshiba**.

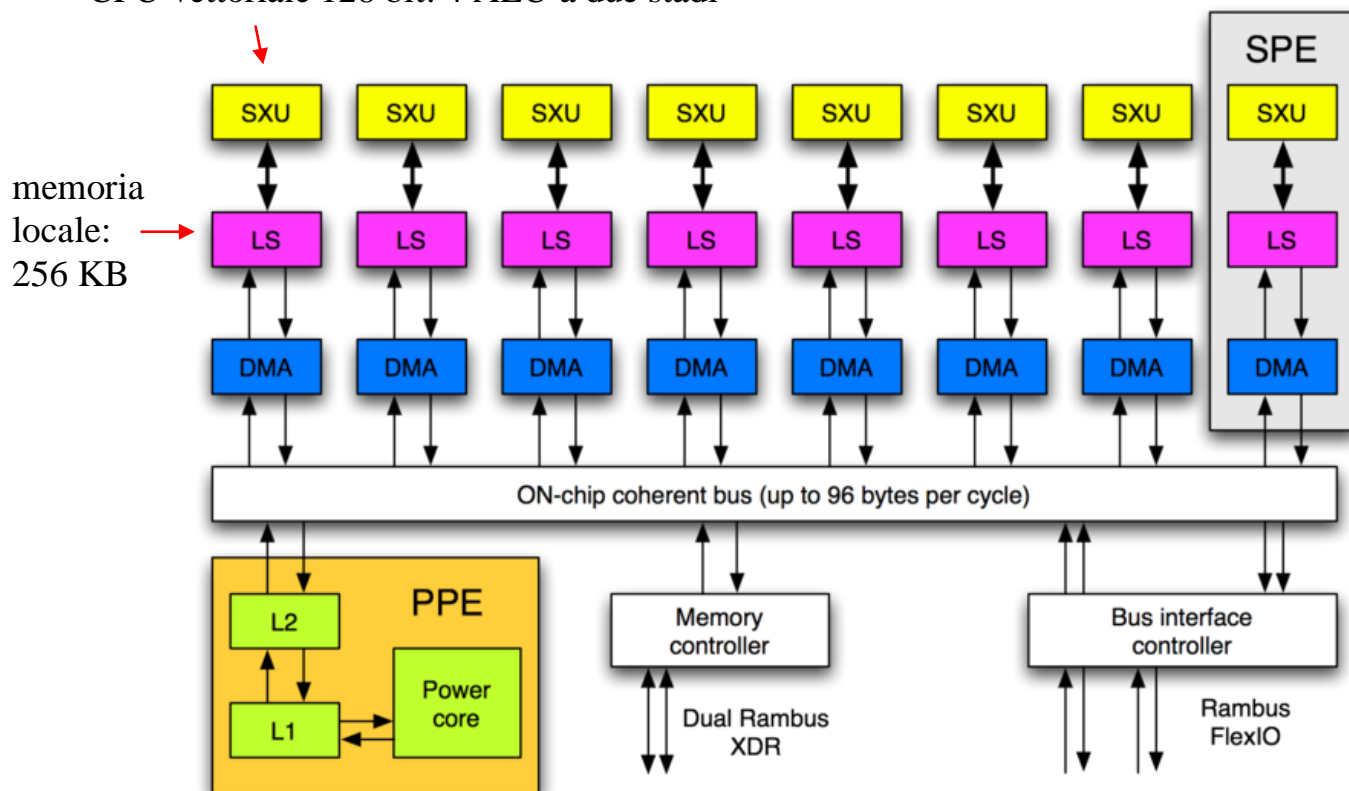
Cell è stato ideato per un utilizzo quasi universale, dalle applicazioni dedicate (**embedded**) fino al mercato dei **supercalcolatori**:

- **Sony** li utilizza per la console **PlayStation 3**
- **IBM** ha sviluppato **Roadrunner** (nel 2008 il più potente del mondo) utilizzando 13000 processori PowerCell (derivati da Cell): 1.7 PetaFlops.

Cell è un processore **multi-core** con **9 core**, dei quali uno (il **PPE**) è il coordinatore e **8** (gli **SPE**) eseguono parallelamente i compiti affidati.

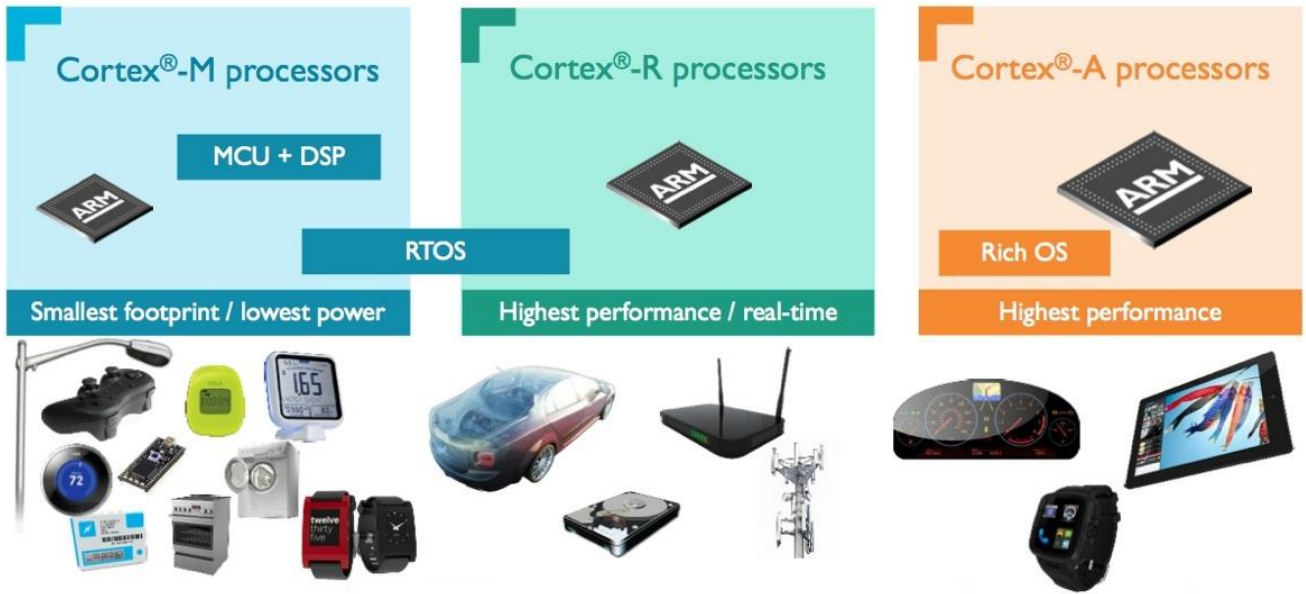
*A livello teorico Cell può arrivare sino a 204,8 GFLOPS lavorando in singola precisione, e 25 GFLOPS in precisione doppia, rispettivamente 64 e 8 volte quelle di un Pentium 4 con lo stesso clock.*

CPU vettoriale 128 bit: 4 ALU a due stadi



# CPU per sistemi embedded

**ARM®**



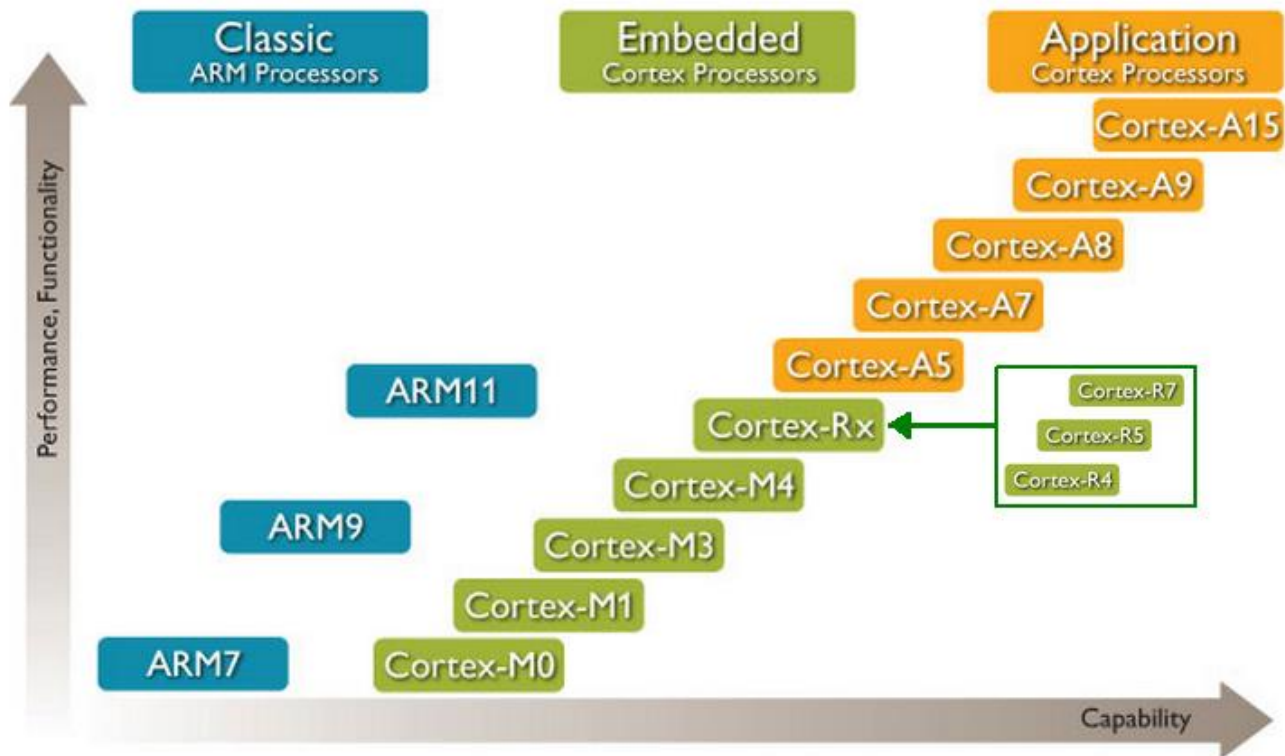
**ARM** (**A**dvanced **RISC** **M**achine) indica una famiglia di processori RISC a 32 bit sviluppata dall'azienda Inglese **ARM Holdings** che detiene la proprietà intellettuale di molti “core” di CPU embedded. ARM non produce direttamente hardware ma vende licenze di produzione dei suoi core a grandi produttori di CPU e **System-On-Chip (SOC)** tra cui: **FreeScale**, **STMicroelectronics**, **NXP**, **Texas Instruments**, **Samsung**, **Broadcom**, **Qualcomm**, **Nvidia**, ecc.

Si tratta di CPU RISC, di disegno architeturale semplice e pulito, caratterizzate da **basso consumo** (caratteristica fondamentale per dispositivi a batteria) e un buon compromesso: **prestazioni** ↔ **costo**

*Cortex A9 alla frequenza di 1 GHz consuma meno di 250 mW per core*  
*Un processore Intel Core-i7 multicore può consumare oltre 100 W*



# CPU per sistemi embedded (2)



## Core Classic più popolari:

Nati per sistemi embedded, ARM 9 e ARM 11 adottati da dispositivi mobili.

**ARM 7:** aritmetica intera, fino 20..60 MIPS

**ARM 9:** aritmetica intera, fino 100..200 MIPS

**ARM 11:** floating point, 300..1000 MIPS (2600 MIPS multicore). Modello a 700 MHz usato in **Raspberry Pi** (2012).

## Core Cortex:

Presentati nel 2005 come evoluzione futura ARM.

**Cortex-M** (M = **microcontroller**): naturale evoluzione di ARM 7 e ARM 9 per device embedded senza grandi necessità di performance (es. 100 MIPS).

**Cortex-R** (R = **real-time**): per sistemi embedded con requisiti prestazionali medio-elevati (es. 600 MIPS). Floating point opzionale.

# CPU per sistemi embedded (3)

**Cortex-A** (A = **applications**): per sistemi con requisiti prestazionali elevati (es. smartphone, tablet). Include **floating point**, **SIMD**.

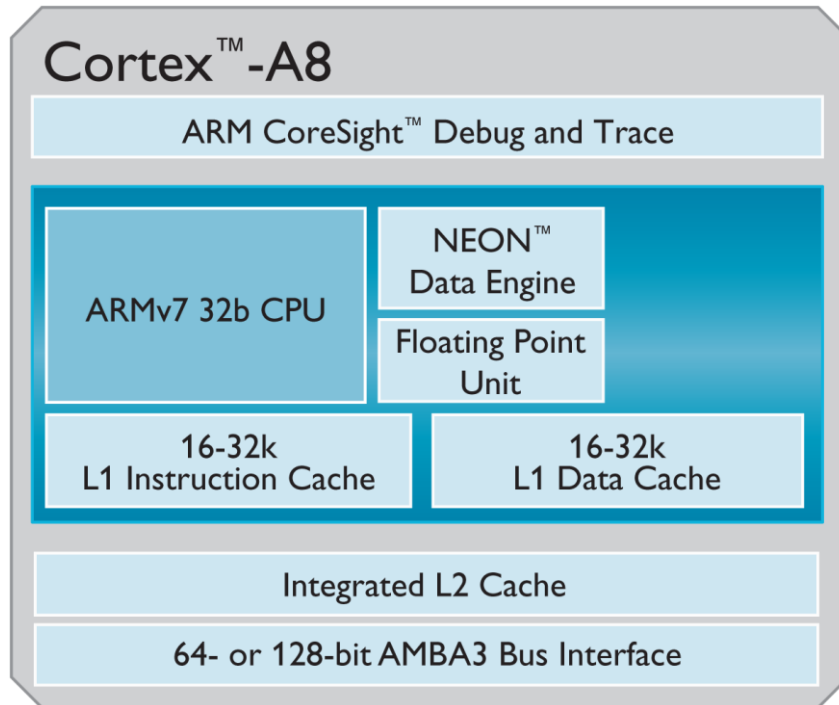
Disponibili versioni a **32** bit e **64** bit:

Announced (64-bit)		Announced (32-bit)	
Year	Core	Year	Core
2012	Cortex-A53	2005	Cortex-A8
2012	Cortex-A57	2007	Cortex-A9
2015	Cortex-A72	2009	Cortex-A5
2015	Cortex-A35	2010	Cortex-A15
2016	Cortex-A73	2011	Cortex-A7
2017	Cortex-A55	2013	Cortex-A12
2017	Cortex-A75	2014	Cortex-A17
		2016	Cortex-A32

- Il Cortex **A7** supporta fino a **4 core**, ciascuno con una pipeline parzialmente **superscalare** a **8 stadi** e unità **predizione salto** ma con **esecuzione in ordine**. 1.9 MIPS per MHz. Scelto per via del basso consumo, per il system-on-chip Broadcom per **Raspberry Pi 2**.
- Il Cortex **A15** supporta fino a **8 core**, ciascuno con **doppia pipeline** fino a **25 stadi**; supporta **esecuzione fuori ordine**; 3.5 MIPS per MHz. Una variante utilizzata nei system-on-chip di **iPhone5** (Apple A6).
- Il Cortex **A72** supporta fino a **4 core** per cluster (possibili più cluster). Ha una cache L2 fino a 4MB. Architettura **superscalare a tre vie**. Predizione di salto (sostanziale). Supporta **esecuzione fuori ordine** e **speculativa**. Rende disponibili 4.7 MIPS per MHz.
- Varianti dei recenti core “di punta” a 64 bit sono state usate nei SOC Exynos (di Samsung) e Snapdragon (di Qualcomm) usati in **Samsung 8**, e nel chip A11 Bionic di Apple usato in **iPhone X**.

# CPU per sistemi embedded (4)

## Cortex-A8



**Frequenza:** da 600 MHz a 1 GHz

**Consumo:** meno di 300 mW

**Prestazioni:** 2 MIPS per MHz (fino 2000 MIPS)

# CPU per sistemi embedded (5)

## trend corrente (2017) ...

Nel passato il mercato delle CPU per PC/Workstation e quello delle CPU per sistemi embedded sono stati **piuttosto separati**:

- **Intel** ha dominato il primo
- **ARM** ha dominato il secondo

I due mondi si stanno oggi **avvicinando**, ma entrambi i contendenti faticano a conquistare spazi nel terreno avversario:

- Intel è interessata al grosso business dei sistemi mobili:
  - **Intel Atom** sono CPU di basso consumo per dispositivi mobili.
  - **Android**, sebbene dominato da CPU ARM, supporta anche processori Intel (compatibilità applicazioni sul Play Store al 90%).
- ARM è interessata a estendere la propria influenza sul settore PC:
  - Nel 2011 l'amministratore delegato di ARM dichiarò che *nel 2015 ARM avrebbe conquistato il 50% delle quote di mercato nell'equipaggiamento dei PC portatili* (dichiarazione azzardata!)
  - **Window 8 RT** rilasciato nel 2012 per il mondo ARM, con scarso successo a causa della incompatibilità x86. Più recentemente (2017) Microsoft ha annunciato una versione di Windows denominata "**Windows 10 on ARM**" che sarebbe in grado di eseguire anche programmi x86 mediante emulazione software. Intel non sembra aver gradito la notizia e ha fatto notare che cercare di emulare l'ISA x86 senza la sua autorizzazione non è consentito, per via di alcuni brevetti...

Core e processori a basso consumo stanno diventando interessanti anche per applicazioni HPC e per la realizzazione di supercalcolatori dai consumi energetici più ridotti (**il top500 nel 2017 consuma 15 Mega Watt!**).

- **Tegra (Nvidia)** è un system-on-chip basato su Core ARM utilizzato, oltre che in dispositivi mobili, anche in server dedicati al calcolo scientifico.
- **Exynos 5 (Samsung)** utilizzato come base per la realizzazione di un supercalcolatore progettato per entrare nella Top500.