

David A. Patterson John L. Hennessy

# Struttura e progetto dei calcolatori

Progettare con RISC-V

Edizione italiana a cura di Alberto Borghese



INFORMATICA ZANICHELLI

David A. Patterson John L. Hennessy

# Struttura e progetto dei calcolatori

Progettare con RISC-V

Edizione italiana a cura di Alberto Borgese

Biblioteca di Area  
Scientifico - Tecnologica

DI SIEI/A 1240  
UNIVERSITÀ



**Se vuoi accedere alle risorse online riservate**

1. Vai su [my.zanichelli.it](http://my.zanichelli.it)
2. Clicca su *Registrati*.
3. Scegli *Studente*.
4. Segui i passaggi richiesti per la registrazione.
5. Riceverai un'email: clicca sul link per completare la registrazione.
6. Cerca la tua chiave di attivazione stampata in verticale sul bollino argentato in questa pagina.
7. Inseriscila nella tua area personale su [my.zanichelli.it](http://my.zanichelli.it)

Se sei già registrato, per accedere ai contenuti riservati di altri volumi ti serve solo la relativa chiave di attivazione.

INFORMATICA ZANICHELLI

**Titolo originale: Computer and Organization Design RISC-V Edition, 1e** ISBN 978-0-12-812275-4  
Copyright © 2018 Elsevier Inc. All Rights Reserved.

Published by arrangement with Elsevier Inc., 230 Park Avenue South, New York, USA.

Morgan Kauffman is an imprint of Elsevier.

Questa traduzione di *Computer and Organization Design RISC-V Edition, 1e* di David Patterson e John Hennessy è pubblicata con l'autorizzazione di Elsevier Inc.

Questa traduzione è stata realizzata da Zanichelli editore ed è di sua esclusiva responsabilità.

Professionisti e ricercatori devono sempre basarsi sulla propria esperienza e sulle proprie conoscenze nel valutare e utilizzare qualsiasi informazione, metodo, composto chimico o esperimento qui descritto.

Nessuna responsabilità è assunta di fronte alla legge da Elsevier, dagli autori, redattori o collaboratori rispetto alla traduzione o in relazione a qualsiasi infortunio e/o danno arrecato a persone o proprietà a causa di prodotti difettosi, negligenza o altro, o in seguito all'uso o al funzionamento di qualsiasi metodo, prodotto, istruzione o idea descritti nel testo qui presente.

RISC-V and the RISC-V logo are registered trademarks managed by the RISC-V Foundation, used under permission of the RISC-V Foundation. All rights reserved.

This publication is independent of the RISC-V Foundation, which is not affiliated with the publisher and the RISC-V Foundation does not authorize, sponsor, endorse or otherwise approve this publication.

All material relating to the ARM technology has been reproduced with permission from ARM Limited, and should only be used for education purposes. All ARM-based models shown or referred to in the text must not be used, reproduced or distributed for commercial purposes, and in no event shall purchasing this textbook be construed as granting you or any third party, expressly or by implication, estoppel or otherwise, a license to use any other ARM technology or know how. Materials provided by ARM are copyright ©ARM Limited (or its affiliates).

**Traduzione:** Alberto Borghese

---

© 2019 Zanichelli editore S.p.A., via Irnerio 34, 40126 Bologna [82059]  
www.zanichelli.it

I diritti di elaborazione in qualsiasi forma o opera, di memorizzazione anche digitale su supporti di qualsiasi tipo (inclusi magnetici e ottici), di riproduzione e di adattamento totale o parziale con qualsiasi mezzo (compresi i microfilm e le copie fotostatiche), i diritti di noleggio, di prestito e di traduzione sono riservati per tutti i paesi. L'acquisto della presente copia dell'opera non implica il trasferimento dei suddetti diritti né li esaurisce.

Le fotocopie per uso personale (cioè privato e individuale, con esclusioni quindi di strumenti di uso collettivo) possono essere effettuate, nei limiti del 15% di ciascun volume, dietro pagamento alla S.I.A.E del compenso previsto dall'art. 68, commi 4 e 5, della legge 22 aprile 1941 n. 633. Tali fotocopie possono essere effettuate negli esercizi commerciali convenzionati S.I.A.E. o con altre modalità indicate da S.I.A.E. Per le riproduzioni ad uso non personale (ad esempio: professionale, economico, commerciale, strumenti di studio collettivi, come dispense e simili) l'editore potrà concedere a pagamento l'autorizzazione a riprodurre un numero di pagine non superiore al 15% delle pagine del presente volume.

Le richieste vanno inoltrate a:  
Centro Licenze e Autorizzazioni per le Riproduzioni Editoriali  
Corso di Porta Romana, n. 108  
20122 Milano  
e-mail autorizzazioni@clearedi.org e sito web www.clearedi.org.

L'autorizzazione non è concessa per un limitato numero di opere di carattere didattico riprodotte nell'elenco che si trova all'indirizzo http://zanichelli.it/fotocopie-opere-escluse

L'editore, per quanto di propria spettanza, considera rare le opere fuori del proprio catalogo editoriale. La loro fotocopia per i soli esemplari esistenti nelle biblioteche è consentita, oltre il limite del 15%, non essendo concorrentiale all'opera. Non possono considerarsi rare le opere di cui esiste, nel catalogo dell'editore, una successiva edizione, le opere presenti in cataloghi di altri editori o le opere antologiche. Nei contratti di cessione è esclusa, per biblioteche, istituti di istruzione, musei ed archivi, la facoltà di cui all'art. 71 - ter legge diritto d'autore.  
Per permessi di riproduzione, anche digitali, diversi dalle fotocopie rivolgersi a ufficiocontratti@zanichelli.it

---

**Realizzazione editoriale:** Epitesto, Milano

**Copertina:**

- **Progetto grafico:** Falcinelli & Co., Roma
- **Immagine di copertina:** © Vitaliy Golubtsov/iStockphoto

---

**Edizione italiana:** giugno 2019

**Ristampa: prima tiratura**

5    4    3    2    1              2019    2020    2021    2022    2023

Realizzare un libro è un'operazione complessa, che richiede numerosi controlli: sul testo, sulle immagini e sulle relazioni che si stabiliscono tra essi.

L'esperienza suggerisce che è praticamente impossibile pubblicare un libro privo di errori. Saremo quindi grati ai lettori che vorranno segnalarceli.

Per segnalazioni o suggerimenti relativi a questo libro scrivere al seguente indirizzo:

Zanichelli editore S.p.A.  
Via Irnerio 34  
40126 Bologna  
fax 051293322  
e-mail: linea\_universitaria@zanichelli.it  
sito web: www.zanichelli.it

Prima di effettuare una segnalazione è possibile verificare se questa sia già stata inviata in precedenza, identificando il libro interessato all'interno del nostro catalogo online per l'Università.

Per comunicazioni di tipo commerciale: universita@zanichelli.it

**Stampa:** Grafica Rago  
Via Lombardia 25, 40064 Tolara di Sotto, Ozzano Emilia (Bologna)  
per conto di Zanichelli editore S.p.A.  
Via Irnerio 34, 40126 Bologna

# Indice generale

Prefazione	VII	Misurare le prestazioni	28
		Prestazioni della CPU	29
		Misura delle prestazioni associate alle istruzioni	30
		Equazione classica di misura delle prestazioni	31
<b>1 Il calcolatore: astrazioni e tecnologia</b>	<b>1</b>	<b>1.7 Barriera dell'energia</b>	<b>35</b>
<b>1.1 Introduzione</b>	<b>1</b>	<b>1.8 Metamorfosi delle architetture: il passaggio dai sistemi uniprocesso ai sistemi multiprocesso</b>	<b>37</b>
Tipi di calcolatori e loro caratteristiche	3		
Benvenuti nell'era post-PC	4		
Che cosa si può imparare da questo libro	5		
<b>1.2 Otto grandi idee sull'architettura dei calcolatori</b>	<b>8</b>	<b>1.9 Un caso reale: la valutazione del Core i7 Intel</b>	<b>41</b>
Progettare tenendo conto della Legge di Moore	8	Benchmark SPEC per la CPU	41
Utilizzo delle astrazioni per semplificare il progetto	9	Benchmark SPEC sull'assorbimento di potenza	42
Rendere veloci le situazioni più comuni	9		
Prestazioni attraverso il parallelismo	9	<b>1.10 Errori e trabocchetti</b>	<b>43</b>
Prestazioni attraverso la pipeline	9		
Prestazioni attraverso la predizione	9	<b>1.11 Note conclusive</b>	<b>46</b>
Gerarchia delle memorie	10	Organizzazione del testo	47
Affidabilità e ridondanza	10	<b>1.12 Inquadramento storico e approfondimenti</b> 	<b>48</b>
<b>1.3 Che cosa c'è dietro un programma</b>	<b>10</b>		
Da un linguaggio ad alto livello al linguaggio dell'hardware	11	<b>1.13 Esercizi</b>	<b>48</b>
		Risposte alle domande di autovalutazione	52
<b>1.4 Componenti di un calcolatore</b>	<b>13</b>	<b>2 Le istruzioni: il linguaggio dei calcolatori</b>	<b>53</b>
Attraverso lo specchio	14	<b>2.1 Introduzione</b>	<b>53</b>
Touchscreen	15	<b>2.2 Operazioni svolte dall'hardware del calcolatore</b>	<b>54</b>
Dentro la scatola	16	<b>2.3 Operandi dell'hardware del calcolatore</b>	<b>58</b>
Un posto sicuro per i dati	19	Operandi allocati in memoria	60
Comunicare con gli altri calcolatori	20	Operandi immediati o costanti	63
<b>1.5 Tecnologie per la produzione di processori e memorie</b>	<b>21</b>	<b>2.4 Numeri con e senza segno</b>	<b>65</b>
<b>1.6 Prestazioni</b>	<b>24</b>	Riepilogo	70
Definizione delle prestazioni	25		

<b>2.5 Rappresentazione delle istruzioni nel calcolatore</b>	71	<b>3 L'aritmetica dei calcolatori</b>	151
Campi delle istruzioni RISC-V	73		
<b>2.6 Operazioni logiche</b>	79	<b>3.1 Introduzione</b>	151
<b>2.7 Istruzioni per prendere decisioni</b>	81	<b>3.2 Somme e sottrazioni</b>	151
Cicli	83	Riepilogo	154
Scorciatoie per il controllo dei confini di vettori e matrici	85	<b>3.3 Moltiplicazione</b>	155
Costrutto <i>case/switch</i>	85	Versione sequenziale dell'algoritmo della moltiplicazione e sua implementazione hardware	156
<b>2.8 Supporto hardware alle procedure</b>	86	Moltiplicazione di numeri dotati di segno	158
Utilizzo di più registri	88	Moltiplicazione veloce	159
Procedure annidate	90	Moltiplicazione nel RISC-V	159
Allocazione dello spazio nello stack per nuovi dati	92	Riepilogo	160
Allocazione dello spazio nello heap per nuovi dati	93	<b>3.4 Divisione</b>	160
<b>2.9 Comunicare con le persone</b>	95	Un algoritmo della divisione e l'hardware che lo implementa	161
Caratteri e stringhe in Java	98	Divisione di numeri dotati di segno	164
<b>2.10 Indirizzamento RISC-V di un campo immediato e di un indirizzo ampio</b>	100	Una divisione più veloce	165
Operandi immediati ampi	100	Divisione nel RISC-V	165
Indirizzamento nei salti	101	Riepilogo	165
Riassunto delle modalità di indirizzamento del RISC-V	104	<b>3.5 Numeri in virgola mobile</b>	168
Come decodificare il linguaggio macchina	104	Rappresentazione in virgola mobile	169
<b>2.11 Parallelismo e istruzioni: la sincronizzazione</b>	107	Eccezioni e Interrupt	170
<b>2.12 Tradurre e avviare un programma</b>	110	Standard IEE 754 per la virgola mobile	170
Compilatore	110	Addizione in virgola mobile	174
Assemblatore	111	Moltiplicazione in virgola mobile	178
Linker	113	Istruzioni in virgola mobile nel RISC-V	181
Loader	116	Accuratezza dell'aritmetica	188
Librerie a caricamento dinamico	116	Riepilogo	190
Come avviare un programma Java	118	<b>3.6 Parallelismo e aritmetica dei calcolatori: parallelismo a livello di parola</b>	192
<b>2.13 Un esempio riassuntivo in linguaggio C</b>	119	<b>3.7 Un caso reale: le estensioni SIMD per lo streaming e le estensioni avanzate dell'x86 per il calcolo vettoriale</b>	193
Procedura scambia	120	<b>3.8 Come andare più veloci: il parallelismo a livello di parola applicato alla moltiplicazione di matrici</b>	194
Procedura ordina	121	<b>3.9 Errori e trabocchetti</b>	197
<b>2.14 Confronto tra vettori e puntatori</b>	126	<b>3.10 Note conclusive</b>	200
Versione della procedura azzera che utilizza vettore e indice	127	<b>3.11 Inquadramento storico e approfondimenti</b> 	201
Versione della procedura azzera che utilizza i puntatori	128	<b>3.12 Esercizi</b>	202
Confronto tra le due versioni di azzera	129	Risposte alle domande di autovalutazione	205
<b>2.15 Approfondimento: compilazione del C e interpretazione di Java</b> 	130		
<b>2.16 Un caso reale: le istruzioni dell'architettura MIPS</b>	130		
<b>2.17 Un caso reale: le istruzioni dell'architettura x86</b>	131	<b>4 Il processore</b>	206
Evoluzione dell'Intel x86	131		
Registri e modalità di indirizzamento dell'x86	134	<b>4.1 Introduzione</b>	206
Operazioni su numeri interi dell'x86	134	Un'implementazione di base del RISC-V	207
Codifica delle istruzioni x86	138	<b>4.2 Convenzioni del progetto logico</b>	210
Conclusioni sull'x86	139	Metodologia di temporizzazione	211
<b>2.18 Un caso reale: le altre istruzioni dell'architettura RISC-V</b>	139	<b>4.3 Realizzazione di un'unità di elaborazione</b>	213
<b>2.19 Errori e trabocchetti</b>	141	Progettazione di un'unità di elaborazione unificata	218
<b>2.20 Note conclusive</b>	143	<b>4.4 Uno schema semplice di implementazione</b>	220
<b>2.21 Inquadramento storico e approfondimenti</b> 	144	Unità di controllo della ALU	221
<b>2.22 Esercizi</b>	146	Progettazione dell'unità di controllo principale	223
Risposte alle domande di autovalutazione	150	Funzionamento dell'unità di elaborazione	226
		Completamento dell'unità di controllo	228
		Perché non si utilizzano più implementazioni a singolo ciclo?	230

<b>4.5</b>	<b>Introduzione alla pipeline</b>	230	Un esempio di memoria cache: il processore FastMATH Intrinsity	342
	Progettazione dell'insieme di istruzioni per architetture dotate di pipeline	235	Riepilogo	344
	Hazard nelle pipeline	235		
	Hazard sul controllo	239		
	Riepilogo sulla pipeline	243		
<b>4.6</b>	<b>Unità di elaborazione con pipeline e unità di controllo associata</b>	244	<b>5.4</b> Come misurare e migliorare le prestazioni di una cache	345
	Rappresentazione grafica delle pipeline	253	Riduzione delle miss di una cache utilizzando un posizionamento più flessibile dei blocchi	348
	Unità di controllo della pipeline	257	Come trovare un blocco nella cache	353
<b>4.7</b>	<b>Hazard sui dati: propagazione o stallo</b>	261	Come scegliere il blocco da sostituire	354
	Hazard sui dati e stallo	268	Ridurre la penalità di miss utilizzando una cache multilivello	355
<b>4.8</b>	<b>Hazard sul controllo</b>	272	Ottimizzazione software mediante elaborazione a blocchi	357
	Ipotizzare che il salto condizionato non sia eseguito	273	Riepilogo	362
	Ridurre i ritardi associati ai salti condizionati	273		
	Predizione dinamica dei salti	275		
	Riepilogo sulla pipeline	278		
<b>4.9</b>	<b>Le eccezioni</b>	278	<b>5.5</b> Affidabilità delle gerarchie delle memorie	363
	Gestione delle eccezioni nelle architetture RISC-V	280	Definizione di malfunzionamento	363
	Eccezioni e loro gestione nella pipeline	281	Codice di Hamming per la correzione di errori singoli e identificazione di errori doppi (SEC/DED)	365
<b>4.10</b>	<b>Parallelismo a livello di istruzioni</b>	285		
	Concetto di speculazione	286	<b>5.6</b> Macchine virtuali	369
	Parallelizzazione statica dell'esecuzione	287	Requisiti del monitor di una macchina virtuale	370
	Processori dotati di parallelizzazione dinamica dell'esecuzione	291	Mancanza di supporto alle macchine virtuali da parte dell'architettura dell'insieme di istruzioni	371
	Efficienza energetica e pipeline avanzate	296	Protezione e architettura dell'insieme delle istruzioni	371
<b>4.11</b>	<b>Un caso reale: la pipeline del Cortex-A53 ARM e del Core i7 Intel</b>	297		
	Cortex-A53 ARM	297	<b>5.7</b> Memoria virtuale	372
	Core i7 920 di Intel	299	Come individuare la posizione di una pagina e come ritrovarla	376
	Prestazioni del Core i7 920 Intel	302	Page fault	377
<b>4.12</b>	<b>Come andare più veloci: parallelismo a livello di istruzioni e moltiplicazione di matrici</b>	304	Memoria virtuale per un insieme ampio di indirizzi virtuali	380
<b>4.13</b>	<b>Argomenti avanzati: un'introduzione alla progettazione digitale con un linguaggio di progettazione dell'hardware e un modello di pipeline, e approfondimenti sulla pipeline</b>	306	Che cosa succede in scrittura?	382
<b>4.14</b>	<b>Errori e trabocchetti</b>	307	Come rendere più veloce la traduzione degli indirizzi: il TLB	382
<b>4.15</b>	<b>Note conclusive</b>	308	TLB del processore FastMATH Intrinsity	384
<b>4.16</b>	<b>Inquadramento storico e approfondimenti</b>	309	Integrazione della memoria virtuale, dei TLB e delle cache	387
<b>4.17</b>	<b>Esercizi</b>	309	Mecanismi di protezione basati sulla memoria virtuale	388
	Risposte alle domande di autovalutazione	318	Gestione delle miss del TLB e dei page fault	390
		Riepilogo	393	
<b>5.8</b>	<b>Schema comune per le gerarchie delle memorie</b>	395		
	Domanda 1: dove può essere posizionato un blocco?	395		
	Domanda 2: come si individua un blocco?	396		
	Domanda 3: quale blocco deve essere sostituito in caso di miss della cache?	397		
	Domanda 4: come vengono gestite le scritture?	398		
	Le tre C: un modello intuitivo per comprendere il comportamento delle gerarchie delle memorie	399		
<b>5.9</b>	<b>Come utilizzare una macchina a stati finiti per controllare una cache semplificata</b>	401		
	Una cache semplificata	401		
	Macchine a stati finiti	402		
	FSM per il controllore semplificato della cache	404		
<b>5.10</b>	<b>Parallelismo e gerarchie delle memorie: coerenza delle cache</b>	405		
	Schemi di base per garantire la coerenza	407		
	Protocolli di snooping	407		
<b>5.11</b>	<b>Parallelismo e gerarchie delle memorie: i dischi RAID</b>	409		
<b>5.12</b>	<b>Argomenti avanzati: come implementare i controlleri delle cache</b>	409		

<b>5.13 Due casi reali: la gerarchia delle memorie del Cortex-A53 ARM e del Core i7 Intel</b>	410	<b>6.8 Introduzione alle topologie delle reti di calcolatori</b>	466
Prestazioni della gerarchia delle memorie del Cortex-A53 e del Core i7	412	Implementazione delle topologie di rete	468
<b>5.14 Un caso reale: il resto del sistema RISC-V e le istruzioni speciali</b>	414	<b>6.9 Come comunicare con il mondo esterno: le reti dei cluster</b>	469
<b>5.15 Come andare più veloci: blocchi di cache e moltiplicazione tra matrici</b>	415	<b>6.10 Benchmark per i multiprocessori</b>	470
<b>5.16 Errori e trabocchetti</b>	417	Modelli delle prestazioni	472
<b>5.17 Note conclusive</b>	422	Modello roofline	474
<b>5.18 Inquadramento storico e approfondimenti</b>	423	Confronto tra due generazioni di Opteron	475
<b>5.19 Esercizi</b>	423	<b>6.11 Un caso reale: il confronto mediante il modello roofline tra un Core i7 960 Intel e una GPU Tesla NVIDIA</b>	480
Risposte alle domande di autovalutazione	433	<b>6.12 Come andare più veloce: processori multipli e moltiplicazione di matrici</b>	484
<b>6 Processori paralleli: dai client al cloud</b>	434	<b>6.13 Errori e trabocchetti</b>	487
<b>6.1 Introduzione</b>	434	<b>6.14 Note conclusive</b>	489
<b>6.2 Le difficoltà nel creare programmi a esecuzione parallela</b>	436	<b>6.15 Inquadramento storico e approfondimenti</b>	492
<b>6.3 SISD, MIMD, SIMD, SPMD e processori vettoriali</b>	441	Bibliografia	492
SIMD negli x86: le estensioni multimediali	443	<b>6.16 Esercizi</b>	492
Architetture vettoriali	443	Risposte alle domande di autovalutazione	499
Confronto tra architetture vettoriali e scalari	445		
Processori vettoriali ed estensioni multimediali	446		
<b>6.4 Multithreading hardware</b>	448		
<b>6.5 I multicore e gli altri multiprocessori a memoria condivisa</b>	451		
<b>6.6 Introduzione alle unità di elaborazione grafica</b>	455	<b>Appendice A The Basics of Logic Design</b>	
Introduzione alle architetture GPU di NVIDIA	457	<b>Appendice B Mapping Control to Hardware</b>	
Strutture di memoria delle GPU NVIDIA	458	<b>Appendice C La grafica e il calcolo con la GPU</b>	
La prospettiva delle GPU	460	<b>Appendice D A Survey of RISC Architectures for Desktop, Server and Embedded Computers</b>	
<b>6.7 Cluster, calcolatori per centri di calcolo e altri multiprocessori a scambio di messaggi</b>	462		
Calcolatori per grandi centri di calcolo	463		

**Indice analitico**

501

**Manuale di riferimento RISC-V**

511

**Appendici**

- Appendice A The Basics of Logic Design**
- Appendice B Mapping Control to Hardware**
- Appendice C La grafica e il calcolo con la GPU**
- Appendice D A Survey of RISC Architectures for Desktop, Server and Embedded Computers**

# Prefazione

La cosa più bella che possiamo sperimentare è il mistero; esso è la fonte della vera arte e della vera scienza.

Albert Einstein, *What I Believe*, 1930

## Guida al libro

Crediamo che lo studio dell'informatica e dell'ingegneria informatica debba non solo riguardare i principi su cui è fondata l'elaborazione, ma anche riflettere lo stato delle conoscenze attuali in questi campi. Crediamo anche che i lettori, qualunque sia la branca dell'informatica nella quale lavorano, possano apprezzare i paradigmi secondo i quali sono organizzati i sistemi di elaborazione, che determinano le loro funzionalità e prestazioni, e ne decretano in ultima analisi il successo.

La tecnologia richiede oggi che i professionisti di tutte le branche dell'informatica conoscano sia il software sia l'hardware, la cui interazione a tutti i livelli è la chiave per capire i principi fondamentali dell'elaborazione. Inoltre, le idee che stanno alla base dell'organizzazione e della progettazione dei calcolatori valgono sia nell'ambito informatico sia in quello dell'ingegneria elettronica e sono le stesse sia che il vostro interesse principale sia il software sia che sia l'hardware. Per questo motivo, nel testo l'enfasi viene posta sulla relazione tra hardware e software e vengono approfonditi i concetti che stanno alla base dei calcolatori delle ultime generazioni.

Il passaggio recente dalle architetture uniprocessoare ai multiprocessori multicore ha confermato quanto sia corretta questa prospettiva, che abbiamo adottato fin dalla prima edizione di questo libro. Fino a poco tempo fa i programmati potevano fare affidamento sul lavoro dei progettisti delle architetture e dei compilatori e dei produttori dei chip, per rendere più veloci o più efficienti dal punto di vista energetico i propri programmi senza il bisogno di apportare alcuna modifica. Questa era è finita: affinché un programma possa essere eseguito più velocemente, deve diventare un programma parallelo. Anche se l'obiettivo di molti ricercatori è fare sì che i programmati non si accorgano della natura parallela dell'hardware per il quale scrivono i loro programmi, ci vorranno molti anni prima che ciò divenga effettivamente possibile. Crediamo

che nel prossimo decennio la maggior parte dei programmati dovrà capire a fondo il legame tra hardware e software perché i programmi vengano eseguiti in modo efficiente sui calcolatori paralleli.

Questo libro è rivolto principalmente a coloro che, pur avendo scarse conoscenze del linguaggio assembler e della logica digitale, vogliono capire i concetti di base dell'organizzazione degli elaboratori, e a quei lettori che sono interessati a capire il modo in cui si progetta un elaboratore, come funziona e perché si ottengono determinate prestazioni.

## L'altro testo degli stessi autori sulle architetture

Alcuni lettori conosceranno il testo degli stessi autori *Computer Architecture: A Quantitative Approach*, chiamato anche "Hennessy Patterson", mentre questo libro viene solitamente chiamato "Patterson Hennessy". Avevamo scritto quel libro con l'obiettivo di descrivere i principi su cui sono basate le architetture degli elaboratori utilizzando un robusto approccio ingegneristico per illustrare tutti i compromessi tra costi e prestazioni che si rendono necessari. In quel libro avevamo utilizzato un metodo basato su esempi e misure, effettuate su architetture disponibili sul mercato, per consentire al lettore di fare esperienza con la progettazione di sistemi realistici: l'obiettivo era dimostrare che le architetture degli elaboratori possono essere studiate utilizzando metodologie quantitative invece di un approccio descrittivo. Quel libro era stato pensato per i professionisti coscienziosi che volevano approfondire nei dettagli il funzionamento dei calcolatori.

La maggior parte dei lettori di questo libro non ha intenzione di diventare un progettista di calcolatori. Tuttavia, le prestazioni e l'efficienza energetica dei sistemi software prodotti nel prossimo futuro saranno enormemente influenzate da quanto bene i progettisti software avranno capito le strutture hardware di base che sono presenti in un calcolatore. Perciò, i progettisti dei compilatori e dei sistemi operativi, così come i programmati di database e della maggior parte delle applicazioni, hanno bisogno di conoscere bene i principi fondamentali in base ai quali funziona un calcolatore, presentati in questo libro. Analogamente, i progettisti hardware devono capire a fondo come i loro progetti andranno a influenzare le applicazioni software.

Ci siamo quindi resi conto che questo libro doveva essere molto di più di una semplice estensione del materiale contenuto nell'altro testo, per cui abbiamo riesaminato a fondo il contenuto e lo abbiamo modificato adattandolo ai lettori di questo libro. Il risultato ha avuto così tanto successo che abbiamo eliminato tutto il materiale introduttivo dalle versioni successive di *Computer Architecture* e, quindi, ora la sovrapposizione dei contenuti tra i due libri è minima.

## Perché il RISC-V per questa edizione?

La scelta dell'architettura dell'insieme di istruzioni è chiaramente critica per un libro di testo sulle architetture degli elaboratori. Non volevamo un insieme di istruzioni che richiedesse la descrizione di caratteristiche non fondamentali e barocche per chi si avvicina per la prima volta alle architetture, per quanto l'insieme di istruzioni sia popolare. Idealmente, il primo insieme di istruzioni dovrebbe fungere da modello, più o meno come il primo amore: sorprendentemente si ricorderanno entrambi con passione.

Dato che ci sono così tante possibili scelte oggigiorno, per la prima edizione del nostro testo *Computer Architecture: A Quantitative Approach*, abbiamo inventato il nostro personale insieme di istruzioni in stile RISC. Data la popolarità crescente, l'eleganza e la semplicità dell'insieme delle istruzioni MIPS siamo passati a utilizzare i processori MIPS per la prima edizione di questo libro di

testo e per le edizioni successive dell'altro libro. Il MIPS è stato di grande utilità per noi e i nostri lettori.

Sono trascorsi 20 anni da quando siamo passati al MIPS, e anche se miliardi di chip contenenti processori MIPS vengono ancora prodotti, questi si trovano tipicamente inseriti (*embedded*) in dispositivi dove l'insieme delle istruzioni è praticamente invisibile. E, quindi, da un po' di tempo è diventato difficile trovare un calcolatore reale sul quale i lettori possano scaricare un programma MIPS ed eseguirlo.

La buona notizia è che un insieme di istruzioni pubblico che aderisce da vicino ai principi RISC ha fatto il suo debutto e sta rapidamente guadagnando seguaci. Il RISC-V, che è stato sviluppato inizialmente all'Università di Berkeley, non solo ripulisce le stranezze dell'insieme delle istruzioni MIPS, ma offre un semplice, elegante e moderno esempio di quello a cui dovrebbe assomigliare un insieme di istruzioni nel 2017.

Inoltre, dato che non è un'architettura proprietaria, esistono simulatori, compilatori e debugger RISC-V "open source" facilmente reperibili e sono disponibili persino implementazioni RISC-V "open source" scritte nei linguaggi di descrizione dell'hardware. Inoltre, saranno presto disponibili delle piattaforme hardware a basso costo sulle quali si potranno eseguire programmi RISC-V. I lettori beneficeranno non solo dallo studio di queste architetture RISC-V, ma saranno anche in grado di modificarle percorrendo il processo di implementazione per capire l'impatto delle modifiche che propongono sulle prestazioni, sulla dimensione del chip e sull'energia assorbita. Questa è un'opportunità eccitante sia per l'industria dei calcolatori e sia per la didattica, e quindi quando è stato scritto questo libro, più di 40 società hanno aderito alla fondazione RISC-V. Questo elenco di sponsor comprende praticamente tutti i maggiori produttori tranne ARM e Intel, e comprende AMD, Google, Hewlett-Packard Enterprise, IBM, Microsoft, NVIDIA, Oracle e Qualcomm.

È per questi motivi che abbiamo scritto l'edizione RISC-V di questo libro, e stiamo per passare al RISC-V anche per il volume *Computer Architecture: A Quantitative Approach*.

Dato che il RISC-V offre praticamente con lo stesso insieme di istruzioni per indirizzi su 32 bit o su 64 bit, avremmo potuto passare all'insieme di istruzioni a 64 bit, ma abbiamo preferito mantenere la dimensione degli indirizzi a 32 bit. Tuttavia, il nostro editore ha intervistato i docenti che hanno adottato questo testo e hanno visto che il 75% di loro preferiva un indirizzo più ampio o erano indifferenti, per cui abbiamo allargato lo spazio di indirizzamento a 64 bit, che oggi può avere più senso di uno spazio di indirizzamento a 32 bit.

Le uniche modifiche dell'edizione RISC-V rispetto all'edizione MIPS sono quelle associate alla modifica dell'insieme delle istruzioni, che riguarda soprattutto il Capitolo 2, il Capitolo 3, la parte sulla memoria virtuale nel Capitolo 5, e i brevi esempi VMIPS del Capitolo 6. Nel Capitolo 4, siamo passati alle istruzioni RISC-V, abbiamo modificato diverse figure, e abbiamo aggiunto alcune sezioni di "Approfondimento", ma le modifiche sono state più semplici di quanto temevamo. Il Capitolo 1 e le altre Appendici sono rimaste praticamente invariate. L'estesa documentazione disponibile su web combinata con la complessità del RISC-V hanno reso difficile ottenere un'Appendice A (Gli assemblatori, i linker e il simulatore SPIM) come quella della quinta edizione MIPS. In compenso, nei Capitoli 2, 3 e 5 è riportato uno stuzzicante riassunto delle centinaia di istruzioni RISC-V che sono al di fuori delle istruzioni core RISC-V che abbiamo descritto in dettaglio nel resto del libro.

Si noti che non stiamo (ancora) dicendo che da qui in poi passeremo permanentemente al RISC-V. Ad esempio, oltre a questa nuova edizione RISC, si possono acquistare la versione ARMv8 e la versione MIPS. Una possibilità è che ci sarà richiesta per tutte e tre le edizioni di questi libri, oppure di una sola.

Decideremo quando sarà il momento. Per ora, attendiamo le vostre reazioni e i vostri feedback su questo sforzo.

## Le novità di questa edizione

Nella stesura di questa edizione di *Struttura e progetto dei calcolatori* abbiamo perseguito sei obiettivi principali:

1. dimostrare con esempi reali quanto sia importante comprendere il funzionamento dell'hardware;
2. evidenziare i temi principali di ogni argomento inserendo a bordo pagina le icone associate, che vengono introdotte nelle pagine iniziali;
3. proporre nuovi esempi per rispecchiare i cambiamenti occorsi nel passaggio dall'era dei PC all'era post-PC;
4. distribuire il materiale relativo all'I/O su tutto il libro invece di racchiuderlo in un unico capitolo;
5. aggiornare il contenuto tecnico per rispecchiare i cambiamenti nell'industria negli anni successivi alla pubblicazione della precedente edizione;
6. spostare online le Appendici e gli altri capitoli invece di includerli in un CD, per contenere il prezzo e rendere questa edizione disponibile anche in forma elettronica.

Prima di descrivere più in dettaglio questi obiettivi, analizziamo capitolo per capitolo l'approccio adottato nella descrizione sia dell'hardware sia del software (*vedi la tabella alla fine della Prefazione*). I Capitoli 1, 4, 5 e 6 coprono entrambe le aree, indipendentemente dall'argomento trattato.

Il **Capitolo 1** spiega perché sia diventato importante il consumo di energia e perché questo abbia causato il passaggio dai microprocessori a processore singolo ai microprocessori multicore. In questo capitolo vengono anche introdotte le otto grandi idee sulla progettazione delle architetture degli elaboratori.

Il **Capitolo 2** è principalmente un capitolo introduttivo all'hardware, ma contiene anche una descrizione dei compilatori e dei linguaggi di programmazione a oggetti, materiale fondamentale per i lettori più interessati agli aspetti software.

Il **Capitolo 3** è dedicato all'aritmetica in virgola mobile e all'unità di elaborazione dati. I lettori possono tralasciare le parti di questo capitolo che non interessano o che contengono materiale introduttivo, ma non i paragrafi da 3.6 a 3.8. Questi descrivono una procedura che calcola il prodotto di due matrici, mostrando come il parallelismo a livello di parola consenta di migliorare le prestazioni di quattro volte.

Il **Capitolo 4** descrive i processori dotati di pipeline. I paragrafi 4.1, 4.5 e 4.10 contengono un riassunto degli argomenti trattati nel capitolo e nel paragrafo 4.12 vengono mostrate le modifiche alla procedura che esegue il prodotto di due matrici, che consentono un ulteriore aumento delle prestazioni. Anche i lettori interessati all'hardware troveranno del materiale fondamentale in questo capitolo; se le conoscenze sui circuiti logici non sono sufficienti, si può leggere l'**Appendice A** sulla progettazione dei circuiti logici, disponibile online, prima di avventurarsi nella lettura di questo capitolo.

Il **Capitolo 6**, su multicore, multiprocessori e cluster, contiene principalmente nuovo materiale e dovrebbe essere letto da tutti. Il capitolo è stato riorganizzato per rendere più naturale la successione degli argomenti e contiene una descrizione più approfondita di GPU, grandi centri di calcolo e delle interfacce hardware/software delle schede di rete, fondamentali per i cluster.

Il **primo obiettivo** di questa edizione è quello di utilizzare un esempio concreto per dimostrare quanto sia importante comprendere il funzionamento dell'hardware per ottenere buone prestazioni ed elevata efficienza energetica.

Come già accennato, inizieremo descrivendo il parallelismo a livello di parola nel Capitolo 3, che consente di ottenere un miglioramento delle prestazioni di un fattore 4 sulla moltiplicazione di due matrici. Miglioreremo ulteriormente le prestazioni nel Capitolo 4, attraverso l'espansione dei cicli, dimostrando così l'importanza del parallelismo a livello di istruzioni. Nel Capitolo 5 raddoppieremo ancora le prestazioni, ottimizzando l'utilizzo della cache mediante l'accesso a blocchi. Infine, nel Capitolo 6 mostreremo come ottenere un miglioramento di 14 volte utilizzando 16 processori e il parallelismo a livello di thread. Per ottenere tutte queste ottimizzazioni aggiungeremo solamente 24 linee di codice C alla procedura iniziale.

Il **secondo obiettivo** è aiutare i lettori a distinguere le idee fondamentali identificando all'inizio otto grandi idee nella progettazione delle architetture che vengono richiamate nel resto del volume. Abbiamo inserito le relative icone a bordo pagina ed evidenziato nel testo le parole corrispondenti per ricordare ai lettori questi otto argomenti. Questo libro contiene quasi 100 citazioni. Tutti i capitoli contengono almeno sette esempi di grandi idee, e ciascuna idea viene citata almeno cinque volte. Le prestazioni attraverso il parallelismo, la pipeline e la predizione sono tre delle idee più utilizzate, seguite da vicino dalla Legge di Moore. Il Capitolo 4, riservato al processore, è quello che contiene più esempi; ciò non deve sorprendere, poiché è il capitolo che ha probabilmente riscosso più successo presso i progettisti delle architetture digitali. La grande idea richiamata in tutti i capitoli è quella delle prestazioni attraverso il parallelismo, ben allineata all'enfasi recente sul parallelismo.

Il **terzo obiettivo** è evidenziare il passaggio generazionale nel campo dell'elaborazione dalla generazione dei PC a quella post-PC, illustrato da esempi e commenti: il Capitolo 1 analizza in dettaglio un calcolatore tablet invece di un PC e il Capitolo 6 descrive l'infrastruttura di calcolo di un cloud. In questo libro trattiamo anche l'ARM, che è l'insieme di istruzioni utilizzato nei dispositivi mobili personali dell'era post-PC, esattamente come l'insieme di istruzioni x86 ha dominato l'era dei PC e (fino a oggi) domina il mondo del cloud computing.

Il **quarto obiettivo** consiste nel distribuire il materiale relativo all'I/O su tutto il libro invece di concentrarlo in un unico capitolo, come abbiamo fatto per il parallelismo nell'edizione precedente: potete trovare il materiale sull'I/O nei paragrafi 1.4, 4.9, 5.2, 5.5, 5.11 e 6.9. Pensiamo che i lettori (e i docenti) leggeranno più volentieri il materiale sull'I/O se non è contenuto in un capitolo a sé.

Il mondo degli elaboratori è un mondo in evoluzione molto veloce e quindi, come succede sempre per le nuove edizioni, uno degli obiettivi più importanti è utilizzare contenuto tecnico aggiornato. Relativamente al **quinto obiettivo**, come esempio di architetture in questo libro abbiamo utilizzato il Cortex-A53 ARM e il Core i7 Intel, esempi tipici dei calcolatori dell'era post-PC. Altri contenuti importanti un'introduzione sulla GPU che spiega anche la sua terminologia, una trattazione più approfondita dei calcolatori dei grandi centri di calcolo che sono alla base del cloud e una descrizione dettagliata delle schede Ethernet a 10 Gigabyte.

Per il **sesto obiettivo**, contenere il prezzo e mantenere il libro di dimensioni ragionevoli e compatibile con una versione elettronica, abbiamo reso disponibili online il materiale complementare e le Appendici invece di inserirli in un CD allegato, come avveniva nelle edizioni precedenti.

Infine, abbiamo aggiornato tutti gli esercizi riportati nel libro.

Abbiamo conservato le caratteristiche delle edizioni precedenti rivelatesi più utili: abbiamo mantenuto la definizione delle parole chiave a margine del testo la prima volta che compaiono, le sezioni *Capire le prestazioni dei programmi* (dedicate alle prestazioni e a come migliorarle), le sezioni *Interfaccia hardware/software* (sui compromessi da adottare a livello di questa interfaccia), le sezioni

*Quadro d'insieme* (che riepilogano i concetti principali espressi nel testo) e le sezioni *Autovalutazione*, che aiutano il lettore a valutare la comprensione degli argomenti trattati (con le relative risposte esatte alla fine di ogni capitolo). Anche questa edizione contiene nell'ultima pagina una scheda tecnica riasuntiva del RISC-V. Il contenuto della scheda è stato aggiornato e costituisce un riferimento immediato per coloro che scrivono programmi nell'assembler del RISC-V.

## Le risorse multimediali

[online.universita.zanichelli.it/patterson5e](http://online.universita.zanichelli.it/patterson5e)

A questo indirizzo sono disponibili le risorse multimediali di complemento al libro. Per accedere alle risorse protette è necessario registrarsi su [my.zanichelli.it](http://my.zanichelli.it) inserendo la chiave di attivazione personale contenuta nel libro.

### Libro con ebook

Chi acquista il libro può scaricare gratuitamente l'ebook, seguendo le istruzioni presenti nel sito. L'ebook si legge con l'applicazione BooktabZ, che si scarica gratis da App Store (sistemi operativi Apple) o da Google Play (sistemi operativi Android).

## Considerazioni conclusive

Se leggerete il successivo paragrafo di ringraziamenti, vi renderete conto che abbiamo corretto moltissimi errori. E quando un libro passa attraverso diverse ristampe, si ha la possibilità di correggere un numero ancora maggiore di errori. Se dovete trovare altri errori, per cortesia, contattate direttamente l'editore attraverso la posta elettronica o la posta ordinaria, utilizzando gli indirizzi riportati nella pagina del copyright.

Questa edizione è la seconda dopo l'interruzione della lunga collaborazione tra Hennessy e Patterson, iniziata nel 1989. Gli impegni richiesti per dirigere una delle più importanti università del mondo hanno tolto a Hennessy il tempo necessario per lavorare alle nuove edizioni: il suo coautore, Patterson, si è sentito ancora una volta come un acrobata che si esibisce senza rete. Per questo motivo, le persone elencate nel paragrafo dei ringraziamenti e i colleghi di Berkeley hanno avuto un ruolo ancora maggiore nel dare forma al contenuto di questo libro. Ciò nonostante, uno solo è l'autore responsabile del nuovo materiale che vi apprestate a leggere.

## Ringraziamenti per questa edizione

Siamo stati davvero fortunati a ricevere diversi contributi da molti lettori, redattori e ricercatori per ciascuna edizione di questo libro. Ciascuno di loro ha contribuito a rendere migliore il libro.

Siamo grati per l'assistenza di Khaled Benkrid e ai suoi colleghi di ARM Ltd. che hanno revisionato attentamente il materiale relativo agli ARM e hanno fornito interessanti suggerimenti.

Il Capitolo 6 è stato rivisto a fondo e abbiamo analizzato separatamente le idee e i contenuti; infine ho apportato le modifiche in base ai suggerimenti di ogni revisore. Vorrei ringraziare Christos Kozyrakis dell'Università di Stanford per il suggerimento di utilizzare l'interfaccia di rete dei cluster per illustrare l'interfaccia hardware/software dell'I/O e per i suggerimenti su come organizzare il capitolo; Mario Flagsilk dell'Università di Stanford, che ha fornito dettagli, diagrammi e misure delle prestazioni del NIC NetFPGA; e i seguenti ricercatori per i loro suggerimenti su come migliorare il capitolo: David Kaeli della

Northeastern University, Partha Ranganathan degli HP Labs, David Wood dell'Università del Wisconsin e i miei colleghi di Berkeley Siamak Faridani, Shoaib Kamil, Yunsup Lee, Zhangxi Tan e Andrew Waterman.

Un particolare ringraziamento va a Rimas Avizensis dell'Università di Berkeley, che ha sviluppato le diverse versioni della procedura di moltiplicazione di matrici e fornito la misura delle prestazioni. Avendo lavorato con suo padre quando ero studente di dottorato a UCLA, è stato particolarmente piacevole lavorare con Rimas a Berkeley.

Vorrei ringraziare anche il mio collaboratore di lunga data Randy Katz dell'Università di Berkeley, che mi ha aiutato a sviluppare il concetto delle grandi idee nelle architetture degli elaboratori all'interno della preparazione del corso della laurea di primo livello che abbiamo tenuto assieme.

Desidero ringraziare David Kirk, John Nickolls e i loro colleghi di NVIDIA (Micheal Garland, John Montrym, Doug Voorhies, Lars Nyland, Erik Lontholm, Laulius Micikevicius, Massimiliano Fatica, Stuart Oberman e Vasily Volkov) per avere scritto l'Appendice C, che tratta approfonditamente le GPU. Vorrei rinnovare il mio apprezzamento a Jim Larus, recentemente nominato direttore della Facoltà di Computer and Communication Science di EPFL, per il suo desiderio di contribuire con la sua esperienza sulla programmazione in linguaggio assembler e per avere messo a disposizione dei lettori il simulatore MIPS da lui sviluppato e costantemente aggiornato.

Sono anche molto grato a Zachary Kurmas dell'Università Statale della Grand Valley, che ha aggiornato gli esercizi e ne ha creati di nuovi, a partire dagli esercizi originali prodotti da Perry Alexander (Università del Kansas), Jason Bakos (Università della South Carolina), Javier Bruguera (Università di Santiago di Compostela), Matthew Farrens (Università di California, Davis), David Kaeli (Northeastern University), Nicole Kaiyan (Università di Adelaide), John Oliver (Cal Poly, San Luis Obispo), Milos Prvulovic (Georgia Tech), Jichuan Chang (Google), Jacob Leverich (Stanford), Kevin Lim (Hewlett-Packard) e Partha Ranganathan (Google).

Un ulteriore ringraziamento va a Peter Ashenden per avere sviluppato nuove slide per le lezioni.

Sono particolarmente grato ai diversi docenti che hanno risposto ai questionari fatti circolare dall'editore, hanno commentato la nostra proposta editoriale e hanno partecipato ai gruppi di discussione per analizzare e fornire suggerimenti per questa edizione. Questi docenti sono: Bruce Barton (Suffolk County Community College), Jeff Braun (Montana Tech), Ed Gehring (North Carolina State), Michael Goldweber (Xavier University), Ed Harcourt (St. Lawrence University), Mark Hill (University of Wisconsin, Madison), Patrick Homer (University of Arizona), Norm Jouppi (HP Labs), Dave Kaeli (Northeastern University), Christos Kozyrakis (Stanford University), Jae C. Oh (Syracuse University), Lu Peng (LSU), Milos Prvulovic (Georgia Tech), Partha Ranganathan (HP Labs), David Wood (University of Wisconsin), Craig Zilles (University of Illinois at Urbana-Champaign). Surveys and Reviews: Mahmoud Abou-Nasr (Wayne State University), Perry Alexander (The University of Kansas), Behnam Arad (Sacramento State University), Hakan Aydin (George Mason University), Hussein Badr (State University of New York at Stony Brook), Mac Baker (Virginia Military Institute), Ron Barnes (George Mason University), Douglas Blough (Georgia Institute of Technology), Kevin Bolding (Seattle Pacific University), Miodrag Bolic (University of Ottawa), John Bonomo (Westminster College), Jeff Braun (Montana Tech), Tom Briggs (Shippensburg University), Mike Bright (Grove City College), Scott Burgess (Humboldt State University), Fazli Can (Bilkent University), Warren R. Carithers (Rochester Institute of Technology), Bruce Carlton (Mesa Community College), Nicholas Carter (University of Illinois at Urbana-Champaign), Anthony Cocchi (The City University

of New York), Don Cooley (Utah State University), Gene Cooperman (Northeastern University), Robert D. Cupper (Allegheny College), Amy Csizmar Dalal (Carleton College), Daniel Dalle (Université de Sherbrooke), Edward W. Davis (North Carolina State University), Nathaniel J. Davis (Air Force Institute of Technology), Molisa Derk (Oklahoma City University), Andrea Di Blas (Stanford University), Derek Eager (University of Saskatchewan), Ata Elahi (Southern Connecticut State University), Ernest Ferguson (Northwest Missouri State University), Rhonda Kay Gaede (The University of Alabama), Etienne M. Gagnon (L'Université du Québec à Montréal), Costa Gerousis (Christopher Newport University), Paul Gillard (Memorial University of Newfoundland), Michael Goldweber (Xavier University), Georgia Grant (College of San Mateo), Paul V. Gratz (Texas A&M University), Merrill Hall (The Master's College), Tyson Hall (Southern Adventist University), Ed Harcourt (St. Lawrence University), Justin E. Harlow (University of South Florida), Paul F. Hemler (Hampden-Sydney College), Jayantha Herath (St. Cloud State University), Martin Herbordt (Boston University), Steve J. Hodges (Cabrillo College), Kenneth Hopkinson (Cornell University), Bill Hsu (San Francisco State University), Dalton Hunkins (St. Bonaventure University), Baback Izadi (State University of New York-New Paltz), Reza Jafari, Robert W. Johnson (Colorado Technical University), Bharat Joshi (University of North Carolina, Charlotte), Nagarajan Kandasamy (Drexel University), Rajiv Kapadia, Ryan Kastner (University of California, Santa Barbara), E.J. Kim (Texas A&M University), Jihong Kim (Seoul National University), Jim Kirk (Union University), Geoffrey S. Knauth (Lycoming College), Manish M. Kochhal (Wayne State), Suzan Koknar-Tezel (Saint Joseph's University), Angkul Kongmunvattana (Columbus State University), April Kontostathis (Ursinus College), Christos Kozyrakis (Stanford University), Danny Krizanc (Wesleyan University), Ashok Kumar, S. Kumar (The University of Texas), Zachary Kurmas (Grand Valley State University), Adrian Lauf (University of Louisville), Robert N. Lea (University of Houston), Alvin Lebeck (Duke University), Baoxin Li (Arizona State University), Li Liao (University of Delaware), Gary Livingston (University of Massachusetts), Michael Lyle, Douglas W. Lynn (Oregon Institute of Technology), Yashwant K Malaiya (Colorado State University), Stephen Mann (University of Waterloo), Bill Mark (University of Texas at Austin), Ananda Mondal (Claflin University), Alvin Moser (Seattle University), Walid Najjar (University of California, Riverside), Vijaykrishnan Narayanan (Penn State University), Danial J. Neebel (Loras College), Victor Nelson (Auburn University), John Nestor (Lafayette College), Jae C. Oh (Syracuse University), Joe Oldham (Centre College), Timour Paltashev, James Parkerson (University of Arkansas), Shaunak Pawagi (SUNY at Stony Brook), Steve Pearce, Ted Pedersen (University of Minnesota), Lu Peng (Louisiana State University), Gregory D. Peterson (The University of Tennessee), William Pierce (Hood College), Milos Prvulovic (Georgia Tech), Partha Ranganathan (HP Labs), Dejan Raskovic (University of Alaska, Fairbanks) Brad Richards (University of Puget Sound), Roman Rozanov, Louis Rubinfeld (Villanova University), Md Abdus Salam (Southern University), Augustine Samba (Kent State University), Robert Schaefer (Daniel Webster College), Carolyn J. C. Schable (Colorado State University), Keith Schubert (CSU San Bernardino), William L. Schultz, Kelly Shaw (University of Richmond), Shahram Shirani (McMaster University), Scott Sigman (Drury University), Shai Simonson (Stonehill College), Bruce Smith, David Smith, Jeff W. Smith (University of Georgia, Athens), Mark Smotherman (Clemson University), Philip Snyder (Johns Hopkins University), Alex Sprintson (Texas A&M), Timothy D. Stanley (Brigham Young University), Dean Stevens (Morningside College), Nozar Tabrizi (Kettering University), Yuval Tamir (UCLA), Alexander Taubin (Boston University), Will Thacker (Winthrop University), Mithuna Thottethodi (Purdue University),

Manghui Tu (Southern Utah University), Dean Tullsen (UC San Diego), Steve VanderLeest (Calvin College), Christopher Vickery (Queens College of CUNY), Rama Viswanathan (Beloit College), Ken Vollmar (Missouri State University), Guoping Wang (Indiana-Purdue University), Patricia Wenner (Bucknell University), Kent Wilken (University of California, Davis), David Wolfe (Gustavus Adolphus College), David Wood (University of Wisconsin, Madison), Ki Hwan Yum (University of Texas, San Antonio), Mohamed Zahran (City College of New York), Amr Zaky (Santa Clara University), Gerald D. Zarnett (Ryerson University), Nian Zhang (South Dakota School of Mines & Technology), Jiling Zhong (Troy University), Huiyang Zhou (North Carolina State University), Weiyu Zhu (Illinois Wesleyan University).

Desidero ringraziare in particolar modo Mark Smotherman per le successive revisioni della parte tecnica e per avere fornito spunti che hanno migliorato significativamente la qualità di questa edizione.

Desideriamo ringraziare la grande famiglia di Morgan Kaufmann per avere accettato di pubblicare questa nuova edizione del libro sotto la guida capace di Katy Birtcher, Steve Merken e Nate McFadden: non sarei riuscito a terminare questo libro senza il loro aiuto. Vogliamo ringraziare anche Lisa Jones che ha gestito il processo di produzione del libro e Victoria Pearson Esser che ha disegnato la copertina dell'edizione americana, che congiunge in modo intelligente il contenuto dell'era post-PC con quello della prima edizione.

Infine ho un enorme debito di riconoscenza con Yunsup Lee e Andrew Waterman per aver affrontato questa conversione al RISC-V nel tempo libero mentre stavano fondando la loro start-up. E anche con Erich Love che ha realizzato la versione RISC-V degli esercizi mentre stava terminando il dottorato. Siamo tutti impazienti di vedere cosa succederà con il RISC-V all'interno dell'Università e al di fuori di essa.

Il contributo delle quasi 150 persone elencate qui ci ha consentito di realizzare questa edizione, che spero risulti l'edizione migliore. Buona lettura!

David A. Patterson

Capitolo o Appendice	Paragrafi	Focalizzazione sul software	Focalizzazione sull'hardware
1. Il calcolatore: astrazioni e tecnologia	Da 1.1 a 1.11 1.12 (Storia)		
2. Le istruzioni: il linguaggio dei calcolatori	Da 2.1 a 2.14 2.15 (Compilatori e Java) Da 2.16 a 2.20 2.21 (Storia)		
D. RISC Instruction-Set Architectures	Da D.1 a D.17		
3. L'aritmetica dei calcolatori	Da 3.1 a 3.5 Da 3.6 a 3.8 (Parallelismo a livello di parola) Da 3.9 a 3.10 (Errori e trabocchetti) 3.11 (Storia)		
A. The Basics of Logic Design	Da A.1 ad A.13 4.1 (Panoramica) 4.2 (Convenzioni logiche) Da 4.3 a 4.4 (Una semplice implementazione) 4.5 (Introduzione alla pipeline) 4.6 (Unità di elaborazione con pipeline) Da 4.7 a 4.9 (Hazard ed eccezioni) Da 4.10 a 4.12 (Parallelismo, esempi reali) 4.13 (Unità di controllo della pipeline in Verilog) Da 4.14 a 4.15 (Errori e trabocchetti) 4.16 (Storia)		
4. Il processore	Da B.1 a B.6 Da 5.1 a 5.10 5.11 (Parallelismo e gerarchie delle memorie: i dischi RAID) 5.12 (Controllori delle cache in Verilog) Da 5.13 a 5.17 5.18 (Storia)		
5. Grande e veloce: la gerarchia delle memorie	Da 6.1 a 6.8 6.9 (Le reti) Da 6.10 a 6.14 6.15 (Storia)		
C. La grafica e il calcolo con la GPU	Da C.1 a C.13		

# 1

## Il calcolatore: astrazioni e tecnologia

*La civiltà progredisce estendendo il numero delle operazioni complesse che possiamo effettuare senza dover pensare ad esse.*

Alfred North Whitehead, *An Introduction to Mathematics*, 1911

### 1.1 | Introduzione

Benvenuti alla lettura di questo libro! È per noi un piacere mostrarvi l'eccitante mondo dei calcolatori. Questo, infatti, è tutt'altro che un mondo arido e scoraggiante, dove il progresso procede con estrema lentezza e le nuove idee finiscono per atrofizzarsi perché vengono trascurate. Tutt'altro! I calcolatori sono invece il prodotto principale di una tecnologia straordinariamente vitale, quella dell'informazione, che contribuisce per circa il 10% al prodotto interno lordo degli Stati Uniti d'America, la cui economia stessa è diventata in parte dipendente da quel rapido miglioramento della tecnologia dell'informazione promesso dalla Legge di Moore.

In questo particolare settore industriale, l'innovazione viene introdotta con una frequenza estremamente elevata. Negli ultimi 30 anni sono stati proposti diversi nuovi calcolatori la cui introduzione sembrava dovesse rivoluzionare l'intera industria degli elaboratori; invece queste rivoluzioni hanno avuto vita breve, ma solo perché altri calcolatori, più performanti, hanno sostituito rapidamente i calcolatori più lenti.

Questa corsa all'innovazione ha portato a progressi mai visti prima, già a partire dalla nascita del primo calcolatore elettronico, alla fine degli anni '40 del secolo scorso. Se l'industria dei trasporti avesse tenuto il passo di quella dei calcolatori, oggi si potrebbe andare da New York a Londra in circa un secondo spendendo solo qualche centesimo di dollaro. Fermadoci un momento a riflettere su come ciò modificherebbe la nostra società – si potrebbe

abitare a Tahiti, lavorare a San Francisco e recarsi a Mosca la sera per vedere un balletto al Bolshoi – è abbastanza facile apprezzare le opportunità che verrebbero offerte.

I calcolatori hanno dato vita alla terza rivoluzione nella nostra società, quella dell'informazione, che segue la rivoluzione agraria e quella industriale. La conseguente moltiplicazione delle potenzialità e delle capacità intellettive del genere umano ha avuto un impatto profondo sulla nostra vita quotidiana e ha cambiato il modo in cui viene generata nuova conoscenza. Esiste oggi un nuovo filone della ricerca scientifica in cui scienziati esperti di calcolo automatico collaborano con quelli che si occupano degli aspetti teorici e sperimentali al fine di esplorare nuove frontiere dell'astronomia, della biologia, della chimica, della fisica ecc.

La rivoluzione introdotta dai calcolatori è ancora in corso e ogni volta che il costo di elaborazione diminuisce di un fattore 10, i possibili utilizzi dei calcolatori si moltiplicano: applicazioni che un tempo non erano economicamente convenienti diventano improvvisamente possibili. Fino a pochi anni fa, le seguenti applicazioni erano considerate “fantascienza informatica”:

- *Calcolatori nelle automobili*: fino a quando i progressi nel campo dei microprocessori non hanno prodotto, agli inizi degli anni '80, un abbattimento dei prezzi e un aumento delle prestazioni, l'idea di un controllo computerizzato delle automobili era semplicemente ridicola. Oggi gli elaboratori consentono di ridurre l'inquinamento, di migliorare l'efficienza del carburante controllando l'iniezione e la combustione e di incrementare la sicurezza. Gli elaboratori avvisano il conducente quando qualcosa è presente nell'angolo cieco durante un parcheggio o quando un veicolo si trova nella corsia adiacente quando si inizia il cambio della corsia di marcia, e controllano l'apertura degli air bag per proteggere gli occupanti della vettura in caso di incidente.
- *Telefoni cellulari*: chi avrebbe potuto immaginare che i progressi nelle architetture di elaborazione avrebbero fatto sì che più di metà degli abitanti della Terra possedesse un telefono cellulare, consentendo alle persone di comunicare tra loro praticamente da qualsiasi parte del globo?
- *Mappatura del genoma umano*: il costo dei sistemi di calcolo necessari a mappare e analizzare le sequenze del DNA umano era dell'ordine delle centinaia di milioni di dollari. È inverosimile che qualcuno potesse lanciare un simile progetto se il costo dei calcolatori fosse stato 10 o 100 volte superiore a quello attuale, come lo era 15 o 25 anni fa. Inoltre, il costo della mappatura continua a diminuire, al punto che potrete presto essere in grado di acquisire il vostro genoma, e ciò vi consentirà di personalizzare le vostre cure mediche.
- *World Wide Web*: il World Wide Web, che non esisteva quando nacque la prima edizione di questo libro, ha trasformato la nostra società. Per molte persone, il web ha sostituito le biblioteche tradizionali e i quotidiani.
- *Motori di ricerca*: con l'aumento delle dimensioni del web, e del suo valore, trovare informazioni rilevanti è diventato sempre più importante. Oggigiorno, molte persone basano sui motori di ricerca una buona parte delle loro attività, al punto che non potrebbero più farne a meno.

**Augmented reality (realità aumentata)**: l'arricchimento della percezione sensoriale umana, visiva in particolare, mediante informazioni pertinenti convogliate elettronicamente.

È evidente come i progressi di questa tecnologia riguardino praticamente tutti gli aspetti della nostra società. L'evoluzione degli elaboratori ha consentito la creazione di programmi meravigliosamente utili, il che spiega come mai i calcolatori siano diventati onnipresenti. Le applicazioni immaginate oggi dalla fantascienza suggeriscono le applicazioni vincenti del futuro: per esempio sono già in fase di sviluppo avanzato occhiali per l'**augmented reality** (realità aumentata), modelli di società senza denaro contante e automobili che si guidano da sole.

## Tipi di calcolatori e loro caratteristiche

Nonostante calcolatori molto diversi tra loro condividano la stessa tecnologia hardware (parr. 1.4 e 1.5), da quelli usati negli elettrodomestici più avanzati ai telefoni cellulari, fino ai supercomputer più performanti, nella maggior parte dei casi le soluzioni utilizzate non sono identiche. Infatti, queste applicazioni sono caratterizzate da requisiti di progetto differenti che implicano un diverso utilizzo dell'hardware. A grandi linee, i calcolatori vengono raggruppati in tre classi ben distinte.

I **personal computer** (PC – calcolatori personali) rappresentano il tipo di calcolatore più conosciuto, che molti lettori di questo libro avranno già ampiamente utilizzato. I personal computer offrono buone prestazioni a un singolo utente mantenendo il costo limitato; inoltre, vengono solitamente utilizzati per eseguire software scritto da terze parti. L'evoluzione di molte tecnologie legate ai sistemi di elaborazione è guidata proprio dai personal computer, che hanno appena 35 anni di vita!

I **server** sono la forma moderna di quelli che un tempo erano calcolatori di dimensioni decisamente maggiori, e, di norma, ad essi si accede solo attraverso la rete. I server sono orientati all'elaborazione di carichi di lavoro di grosse dimensioni, rappresentati sia da singole applicazioni complesse, quali tipicamente quelle scientifiche o ingegneristiche, sia dalla gestione di tante piccole applicazioni, come nel caso di un grande server per il web. Queste applicazioni sono spesso basate su software proveniente da terze parti (per es. un database oppure un sistema di simulazione), ma sono il più delle volte modificate e adattate per svolgere una particolare funzione. I server sono realizzati con le stesse tecnologie base dei personal computer, ma offrono una maggiore potenza di calcolo, una maggiore velocità di input/output e una maggiore capacità della memoria. In generale, anche i progettisti dei server danno più importanza all'affidabilità, poiché un blocco del funzionamento di uno di questi sistemi è di solito più problematico del blocco di un PC, utilizzato da un singolo utente.

I server coprono il più ampio spettro per quanto riguarda i costi e le prestazioni. I server di fascia inferiore possono essere costruiti con poco più di un comune PC senza schermo e tastiera, e costare un migliaio di dollari. Questi server di fascia bassa sono tipicamente usati per il salvataggio dei dati, per piccole applicazioni commerciali, o come semplici web server. All'estremo opposto ci sono i **supercomputer**, che attualmente sono formati da diverse decine di migliaia di processori, con memoria principale di diversi **terabyte**. Il costo di un supercomputer va dalle decine alle centinaia di milioni di dollari. I supercomputer sono tipicamente usati nel calcolo intensivo di tipo scientifico e ingegneristico, come per esempio per le previsioni meteorologiche, nelle esplorazioni petrolifere, nella determinazione delle strutture delle proteine e in altri problemi di grandi dimensioni. Nonostante i supercomputer consentano la massima capacità di calcolo, essi costituiscono solo una frazione relativamente piccola dei server esistenti e, in termini di fatturato, del mercato complessivo dei calcolatori.

I **calcolatori embedded** (cioè dedicati) sono i più numerosi e coprono un ampio spettro di applicazioni e prestazioni. Essi comprendono tutti i microprocessori che potete trovare nella vostra automobile, i calcolatori presenti nei televisori e la rete di processori che controlla i moderni aeroplani o le navi commerciali. I sistemi di calcolo di tipo embedded sono progettati per eseguire una singola applicazione o un insieme di applicazioni correlate tra loro; queste applicazioni sono di norma integrate con l'hardware e si presentano all'utente come un sistema monolitico. Nonostante l'enorme numero di calcolatori di tipo embedded che usano quotidianamente, molti utenti non si renderanno mai conto che stanno in realtà usando un vero e proprio computer!

**Personal Computer (PC):** un calcolatore progettato per essere utilizzato da un unico utente alla volta; solitamente comprende anche un monitor grafico, una tastiera e un mouse.

**Server:** un calcolatore progettato per l'esecuzione di programmi di grosse dimensioni e per servire più utenti spesso simultaneamente. Tipicamente a questi calcolatori si accede solo via rete.

**Supercomputer:** calcolatori con il costo più elevato e le prestazioni migliori; sono tipicamente configurati come server e possono costare decine o centinaia di milioni di dollari.

**Terabyte (TB):** misura di capacità pari a  $1\,099\,511\,627\,776$  ( $2^{40}$ ) byte; gli sviluppatori di sistemi di comunicazione e di memoria di massa hanno iniziato a utilizzare questo termine per indicare la cifra di  $1\,000\,000\,000\,000$  ( $10^{12}$ ) byte. Per evitare confusione, utilizziamo qui il termine **Tebibyte (TiB)** per indicare  $2^{40}$  byte e il termine **Terabyte (TB)** per indicare  $10^{12}$  byte. L'elenco completo delle unità di misura decimali e binarie e del loro valore è riportato in Figura 1.1.

**Calcolatore embedded:** un calcolatore posto all'interno di un altro dispositivo e usato esclusivamente per eseguire una predeterminata applicazione o un insieme di programmi.

Termine decimale	Abbreviazione	Valore	Termine binario	Abbreviazione	Valore	% maggiore
Kilobyte	KB	$10^3$	Kibibyte	KiB	$2^{10}$	2%
Megabyte	MB	$10^6$	Mebibyte	MiB	$2^{20}$	5%
Gigabyte	GB	$10^9$	Gibibyte	GiB	$2^{30}$	7%
Terabyte	TB	$10^{12}$	Tebibyte	TiB	$2^{40}$	10%
Petabyte	PB	$10^{15}$	Pebibyte	PiB	$2^{50}$	13%
Exabyte	EB	$10^{18}$	Exbibyte	EiB	$2^{60}$	15%
Zettabyte	ZB	$10^{21}$	Zebibyte	ZiB	$2^{70}$	18%
Yottabyte	YB	$10^{24}$	Yobibyte	YiB	$2^{80}$	21%

**Figura 1.1** L'ambiguità nell'utilizzo della notazione  $2^x$  o  $10^x$  è stata risolta affiancando un'unità di misura binaria all'unità di misura decimale corrispondente per tutti i termini di uso comune. Nell'ultima colonna viene riportato, in percentuale, quanto l'unità di misura binaria è più grande dell'unità di misura decimale corrispondente. La differenza percentuale aumenta procedendo verso il basso. I prefissi valgono sia per i bit che per i byte, per cui 1 Gigabit rappresenta  $10^9$  bit, mentre 1 Gibibit rappresenta  $2^{30}$  bit.

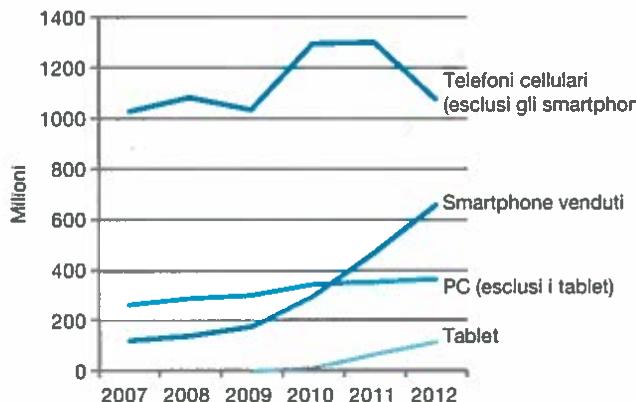
Le applicazioni di tipo embedded richiedono spesso prestazioni limitate con vincoli stringenti sul costo e sulla potenza assorbita dal dispositivo. Consideriamo, per esempio, un lettore digitale di musica: il processore deve solo essere abbastanza veloce da gestire le funzioni di lettura e riproduzione dei brani musicali; inoltre, in fase di progettazione si curerà molto attentamente la minimizzazione del costo e della potenza assorbita. Nonostante il loro basso costo, i calcolatori embedded sono spesso quelli che meno tollerano l'insorgere di malfunzionamenti, dato che le conseguenze dei guasti possono andare da una semplice seccatura (se il televisore si rompe) a veri disastri (come nel caso di un guasto al calcolatore che governa un aeroplano o una nave mercantile). Nelle applicazioni embedded di massa, come nel caso di un elettrodomestico con controllo digitale, l'affidabilità si ottiene innanzitutto attraverso la semplicità; l'enfasi è posta sull'implementazione di una funzione nella maniera migliore possibile. In sistemi embedded di grandi dimensioni, per assicurare una maggiore garanzia di corretto funzionamento, si applicano spesso tecniche che sfruttano la ridondanza, tipicamente utilizzate nei sistemi server. Anche se questo libro è focalizzato sui calcolatori ad ampio spettro di applicazioni, molti dei concetti che verranno affrontati possono venire applicati ai calcolatori embedded, direttamente o con qualche piccola modifica.

**Approfondimento.** Gli approfondimenti sono brevi sezioni presenti in tutto il testo che forniscono maggiori dettagli su argomenti che possono essere di particolare interesse. I lettori non interessati possono saltarli, perché i paragrafi successivi non faranno mai riferimento al contenuto di questi approfondimenti.

Molti processori embedded vengono progettati utilizzando i cosiddetti *processor cores* (processori nucleo), particolari versioni di processore che possono essere descritte in uno dei linguaggi di descrizione dell'hardware come il Verilog o il VHDL (vedi Cap. 4). Questi particolari nuclei permettono al progettista di integrare il processore con l'altro hardware richiesto dall'applicazione con lo scopo di realizzare tutto il sistema su un singolo chip.

## Benvenuti nell'era post-PC

Il progresso continuo della tecnologia ha prodotto dei ricambi generazionali dell'hardware dei calcolatori che hanno scosso l'intera industria informatica fin nelle fondamenta. In questo libro documentiamo uno di questi ricambi



**Figura 1.2** Il numero di tablet e di smartphone prodotti in un anno confrontato con il numero di PC e di telefoni cellulari tradizionali, illustra bene cosa sia l'era post-PC. Gli smartphone costituiscono la componente in crescita dell'industria dei telefoni cellulari e hanno superato il numero dei PC nel 2011. I tablet sono la categoria a più alto tasso di crescita: sono quasi raddoppiati dal 2011 al 2012. Negli ultimi anni la produzione dei PC e dei telefoni cellulari tradizionali si è mantenuta costante o leggermente declinante.

generazionali che si sta rivelando tanto importante quanto lo è stata l'introduzione del PC 30 anni fa: i **dispositivi mobili (PMD, Personal Mobile Devices)** stanno sostituendo i PC. I PMD sono dispositivi alimentati da una batteria, con connessione **wireless** (senza fili) a Internet e un costo tipicamente nell'ordine delle centinaia di dollari. Sui PMD gli utenti possono scaricare dalla rete un software ("apps") ed eseguirlo come nei PC, ma a differenza dei PC, i PMD non possiedono una tastiera e un mouse, ma ricevono l'input attraverso un touchscreen o, più recentemente, attraverso la voce. I PMD di oggi sono tipicamente gli smartphone (letteralmente, telefoni intelligenti) e i calcolatori tablet, ma in un prossimo futuro potranno comprendere, per esempio, gli occhialini elettronici. La Figura 1.2 mostra la rapida crescita dei tablet e degli smartphone in confronto a quella dei PC e dei telefoni cellulari tradizionali.

I server tradizionali sono stati sostituiti dal **cloud computing**, che è basato su enormi centri di calcolo noti come *Warehouse Scale Computers* (WSC). Società quali Amazon e Google hanno costruito WSC contenenti centinaia di migliaia di server, una parte dei quali viene poi noleggiata a società terze le quali possono così offrire ai PMD servizi software senza che esse debbano costruirsi il proprio WSC. I **software come servizio (SaaS, Software as a Service)** forniti attraverso il cloud stanno sicuramente rivoluzionando l'industria del software così come i PMD e i WSC stanno rivoluzionando l'industria dell'hardware. Oggi giorno gli sviluppatori software vedono una parte del loro software eseguita su un PMD e una parte sul cloud.

## Che cosa si può imparare da questo libro

I programmatore di maggior successo hanno sempre tenuto in considerazione le prestazioni dei loro programmi, perché fornire velocemente i risultati all'utente è un fattore critico per la creazione di software di successo. Negli anni '60 e '70, uno dei principali vincoli che limitava le prestazioni era costituito dalla dimensione della memoria degli elaboratori. I programmatore, di conseguenza, seguivano spesso questa logica: bisogna minimizzare lo spazio di memoria per rendere veloci i programmi. Nell'ultimo decennio, invece, i miglioramenti sia nel campo della progettazione dei calcolatori sia in quello delle tecnologie delle memorie hanno sensibilmente ridotto l'importanza di limitare le dimensioni della memoria in molte applicazioni tranne che in quelle rivolte ai sistemi embedded.

**Dispositivi mobili (PMD):** sono piccole dispositivi wireless che vengono collegati a Internet. Sono alimentati da un accumulatore e il software viene installato scaricando delle applicazioni, dette *app*, dalla rete. Tipici esempi sono gli smartphone e i tablet.

**Wireless:** indica una comunicazione radio tra un dispositivo elettronico e un altro o Internet senza uso di cavi elettrici.

**Cloud computing:** si riferisce a un insieme di molti server che forniscono servizi attraverso Internet; alcuni provider noleggiano dinamicamente un numero variabile di questi server a terze parti che forniscono servizi agli utenti.

**Software come servizio (SaaS):** sono il software e i dati forniti come servizio attraverso Internet, solitamente attraverso un piccolo programma, per esempio un browser, che viene eseguito come client su un dispositivo dell'utente. Non viene richiesto di installare tutto il codice binario sul dispositivo ed eseguirlo localmente. Esempi di servizi software sono la ricerca sul web e i social network.

I programmati interessati alle prestazioni cercano ora di comprendere le problematiche che hanno sostituito il semplice modello di memoria utilizzato negli anni '60: il parallelismo dei processori e la gerarchia delle memorie. Dimostriamo quanto sia importante comprendere a fondo questi concetti nei Capitoli 3 e 6 dove viene indicato come migliorare le prestazioni di un programma C di un fattore 200. Inoltre, come illustrato nel paragrafo 1.7, i programmati di oggi devono preoccuparsi del consumo di energia dei loro programmi, sia che vengano eseguiti su un PMD che sul cloud, e questo richiede di conoscere quello che sta sotto il software. I programmati che vogliono realizzare un software competitivo dovranno quindi migliorare la loro conoscenza dell'organizzazione degli elaboratori.

Siamo quindi onorati di avere l'opportunità di illustrare il contenuto di quella macchina rivoluzionaria che è il calcolatore, mostrando il software che sta sotto i programmi utente e l'hardware che sta all'interno dell'elaboratore. Al termine della lettura di questo libro, sarete in grado di rispondere alle seguenti domande:

- Come vengono tradotti nel linguaggio proprio del calcolatore i programmi scritti in un linguaggio ad alto livello, quali per esempio C o Java, e come fa il calcolatore a eseguirli? Questi concetti costituiscono la base per capire quali aspetti dell'hardware e del software hanno un impatto sulle prestazioni dei programmi.
- Qual è l'interfaccia tra hardware e software e come può il software fare in modo che l'hardware esegua le funzioni richieste? Questi concetti sono di vitale importanza per capire come sia possibile scrivere molti tipi di programmi.
- Cosa determina le prestazioni di un programma e come può un programmatore migliorarle? Come vedremo, le prestazioni dipendono dal programma originale, dalla traduzione del programma da un linguaggio ad alto livello nel linguaggio macchina e dall'efficienza con cui l'hardware esegue il programma.
- Quali tecniche può utilizzare il progettista di un calcolatore per migliorarne le prestazioni? Questo libro introdurrà solo i concetti di base della moderna progettazione di un'architettura di elaborazione, ma i lettori più interessati potranno trovare molto più materiale su questo argomento nella versione più recente del nostro libro *Computer Architecture: A Quantitative Approach*.
- Quali tecniche possono essere utilizzate dai progettisti hardware per migliorare l'efficienza energetica? Che cosa può fare un programmatore per migliorare o peggiorare l'efficienza energetica?
- Quali sono i motivi che hanno determinato recentemente il passaggio dall'elaborazione sequenziale a quella parallela e quali sono le conseguenze? In questo libro scoprirete le motivazioni e i meccanismi hardware che gestiscono oggi il parallelismo. Inoltre il Capitolo 6 è quasi interamente dedicato alla nuova generazione di microprocessori, i **processori multicore (multiprocessori)**.
- Quali sono state le migliori idee partorite dai progettisti hardware a partire dal primo calcolatore commerciale prodotto nel 1951 e che sono alla base dei calcolatori moderni?

Se non si capiscono le risposte a queste domande, qualsiasi tentativo di migliorare le prestazioni di un programma su un calcolatore moderno, oppure di valutare quali caratteristiche rendano un calcolatore migliore di un altro per una particolare applicazione, diventerà un complicato processo di tentativi alla cieca, invece di una procedura scientifica guidata da intuito e analisi.

Questo primo capitolo è fondamentale per capire il resto del libro: introduce le idee e le definizioni di base, mette nella giusta prospettiva le principali componenti hardware e software, mostra come valutare le prestazioni e il consumo di energia, introduce i circuiti integrati, cioè la tecnologia che spinge la rivoluzione digitale, e spiega il passaggio ai processori multicore.

**Processore multicore (multiprocessore):** un microprocessore che contiene più processori, detti *core*, in un unico circuito integrato.

In questo e nei capitoli successivi, troverete probabilmente molte parole nuove oppure parole che potete aver sentito ma del cui significato non siete sicuri: niente paura! È vero, ci sono molti termini speciali utilizzati per descrivere i calcolatori moderni ma, come in tutti gli ambiti scientifici, la terminologia in realtà aiuta perché consente di descrivere con precisione le funzionalità. Inoltre i progettisti di calcolatori (autori inclusi) *amano usare gli acronimi*, che sono facili da capire una volta che si conosce il loro significato. Per aiutare il lettore a ricordare e riconoscere i vari termini, il significato di ogni acronimo verrà riportato a lato della pagina la prima volta che l'acronimo apparirà nel testo. In breve tempo il lettore sarà in grado di utilizzare facilmente gli acronimi, e potrà impressionare gli amici con la sua capacità di usare correttamente termini quali BIOS, CPU, DIMM, DRAM, PCIe, SATA ecc.

Per sottolineare ulteriormente come software e hardware possano influenzare le prestazioni di un elaboratore, troverete in diversi punti del libro sezioni speciali intitolate *Capire le prestazioni dei programmi* che riassumono gli elementi determinanti per le prestazioni dei programmi. La prima di queste sezioni è riportata qui di seguito.

Le prestazioni di un programma dipendono dalla combinazione di diversi fattori: efficienza degli algoritmi implementati nel programma, efficienza del software utilizzato per creare e tradurre il programma in linguaggio macchina, ed efficienza con cui il calcolatore esegue le istruzioni, che possono comprendere anche operazioni di input/output (I/O). Questi concetti sono riassunti nella tabella che segue.

**Acronimo:** una parola creata prendendo le lettere iniziali di una sequenza di parole. Per esempio: **RAM** è l'acronimo di *Random Access Memory* e **CPU** è l'acronimo di *Central Processing Unit*.

## Capire le prestazioni dei programmi

Componente hardware o software	Come la componente influenza le prestazioni	Dove viene trattato questo argomento
Algoritmi	Determina il numero di linee del codice ad alto livello e il numero di operazioni di I/O da eseguire	Altri libri!
Linguaggi di programmazione, compilatori e architetture	Determinano il numero di istruzioni in linguaggio macchina per ogni istruzione ad alto livello	Capitoli 2 e 3
Processore e sistema di memoria	Determinano quanto velocemente possono essere eseguite le istruzioni ad alto livello	Capitoli 4, 5 e 6
Sistema di I/O (hardware e sistema operativo)	Determina quanto velocemente possono essere eseguite le operazioni di I/O	Capitoli 4, 5 e 6

Per mostrare l'impatto delle idee descritte in questo libro, come accennato poco fa, mostreremo nei diversi capitoli come possano essere progressivamente migliorate le prestazioni di un programma C che moltiplica una matrice con un vettore fino a ottenere un miglioramento finale di un fattore 200! Ciascun miglioramento è basato sullo sfruttamento ottimale di un particolare componente dell'hardware del microprocessore sottostante.

- Nell'ambito del *parallelismo a livello di dati*, nel Capitolo 3, utilizzeremo le *funzioni intrinseche del C* per implementare il *parallelismo a livello di parola*, aumentando le prestazioni di un fattore 3,8.
- Nell'ambito del *parallelismo a livello di istruzioni*, nel Capitolo 4, utilizzeremo l'*espansione dei cicli per sfruttare l'esecuzione di pacchetti di istruzioni multiple* e l'*esecuzione hardware fuori-ordine* per aumentare le prestazioni di un ulteriore fattore 2,3.

- Nell'ambito dell'*ottimizzazione della gerarchia delle memorie*, nel Capitolo 5, utilizzeremo il *blocco della cache* per aumentare le prestazioni su matrici di grandi dimensioni, con un miglioramento di un ulteriore fattore da 2 a 2,5.
- Nell'ambito del *parallelismo a livello di thread*, nel Capitolo 6, utilizzeremo *cicli for paralleli implementati in OpenMP per sfruttare l'hardware multicore*, aumentando le prestazioni di un altro fattore da 4 a 14.

## Autovalutazione

Le sezioni denominate *Autovalutazione* sono pensate per aiutare il lettore a valutare se abbia compreso i principali concetti introdotti nel capitolo e per mostrare a fondo le implicazioni. Alcune domande di queste sezioni sono molto semplici mentre altre hanno lo scopo di aprire una discussione con altre persone. Potete trovare le risposte alle singole domande alla fine del capitolo. Le domande di autovalutazione sono inserite al termine dei paragrafi, in modo da poterle evitare se siete sicuri di aver appreso quanto avete letto.

1. Il numero di calcolatori embedded venduti ogni anno supera largamente quello dei PC e persino quello dei calcolatori post-PC. Siete in grado di confermare o confutare questa affermazione basandovi sulla vostra esperienza? Provate a contare il numero di processori embedded presenti nella vostra abitazione e confrontatelo con il numero di calcolatori tradizionali. Qual è il risultato di questo confronto?
2. Come accennato precedentemente, sia il software sia l'hardware influenzano le prestazioni di un programma. Provate a pensare ad alcuni esempi pratici nei quali ognuno dei seguenti elementi può rappresentare un "collo di bottiglia" per le prestazioni:
  - algoritmo implementato;
  - linguaggio di programmazione o compilatore;
  - sistema operativo;
  - processore;
  - sistema di I/O e periferiche.

## 1.2 | Otto grandi idee sull'architettura dei calcolatori

Presentiamo ora otto grandi idee dei progettisti dei calcolatori degli ultimi 60 anni. Queste idee sono state così innovative da durare a lungo dopo la loro implementazione nelle prime architetture: i progettisti più giovani hanno ripreso queste idee nella progettazione delle architetture più recenti. Queste idee sono dei principi ricorrenti che compariranno spesso in questo capitolo e nei capitoli successivi; per identificarle introduciamo qui i simboli che le rappresentano e i termini ad esse associati. Utilizzeremo questi simboli per individuare i quasi 100 paragrafi di questo libro che descrivono le caratteristiche basate su queste idee.

### Progettare tenendo conto della Legge di Moore

Una delle costanti dei progettisti di calcolatori è la rapida evoluzione, descritta principalmente dalla **Legge di Moore**. Essa stabilisce che le risorse messe a disposizione dai circuiti integrati vengano duplicate ogni 18-24 mesi. La Legge di Moore deriva da una predizione sull'aumento della capacità dei circuiti integrati fatta nel 1965 da Gordon Moore, uno dei fondatori di Intel. Dato che la progettazione di un calcolatore può richiedere anni, le risorse messe a disposizione da un chip possono facilmente raddoppiare o quadruplicare prima della realizzazione finale del calcolatore. Come nel tiro al piattello, i progettisti dei



calcolatori devono indovinare dove sarà la tecnologia quando sarà terminata la fase di progettazione e non progettare in base alla tecnologia attuale. Utilizzeremo una freccia verso l'alto e verso destra come simbolo che rappresenta la Legge di Moore e l'evoluzione rapida che essa rappresenta.

## Utilizzo delle astrazioni per semplificare il progetto

Sia i progettisti dei calcolatori sia i programmati hanno dovuto inventare delle tecniche che li rendessero più produttivi, per evitare che il tempo riservato alla progettazione crescesse enormemente con il crescere secondo la Legge di Moore delle risorse rese disponibili. Una di queste tecniche è l'**astrazione**, utilizzata per rappresentare il progetto sia hardware sia software a diversi livelli di definizione: i dettagli vengono nascosti ai livelli più bassi per offrire un modello più semplice ai livelli più alti. Utilizzeremo un disegno astratto come simbolo per questa seconda grande idea.



ASTRAZIONE

## Rendere veloci le situazioni più comuni

**Rendere veloci le situazioni più comuni** tende a fare aumentare le prestazioni più dell'ottimizzazione delle funzionalità utilizzate raramente. Ironicamente, le situazioni più comuni sono spesso più semplici di quelle rare e quindi di solito sono più semplici da migliorare. Questo consiglio dettato dal buon senso prevede che sappiate quali sono le situazioni più comuni, cosa che è possibile solo attraverso un'attenta sperimentazione e analisi (par. 1.6). Utilizzeremo il simbolo di un'auto sportiva per rappresentare quest'idea, dato che i viaggi più frequenti prevedono uno o due passeggeri, ed è sicuramente più facile realizzare una macchina sportiva veloce che un minivan veloce!



SITUAZIONI COMUNI

## Prestazioni attraverso il parallelismo

Sin dagli albori, i progettisti hanno realizzato calcolatori che ottengono prestazioni più elevate eseguendo le operazioni in parallelo. Vedremo molti esempi di parallelismo in questo libro. Utilizzeremo il simbolo di un aereo a reazione con più motori per rappresentare le prestazioni attraverso il parallelismo.



PARALLELISMO

## Prestazioni attraverso la pipeline

Una particolare forma di parallelismo è così diffusa nelle architetture da meritarsi un nome proprio: **pipeline**. Per esempio, una squadra di uomini, prima di chiamare i vigili del fuoco, può cercare di spegnere il fuoco con i secchi, come avviene tipicamente nei film western quando il cattivo appicca il fuoco: gli abitanti del paese formano una catena umana per trasportare l'acqua dalla sorgente al fuoco. Infatti, in questo modo si riesce a trasportare l'acqua dalla sorgente al fuoco più velocemente che correndo avanti e indietro. Il simbolo della pipeline è una sequenza di condotti, dove ciascun condotto rappresenta uno stadio della pipeline.



PIPELINE

## Prestazioni attraverso la predizione

Seguendo il detto secondo cui “è meglio chiedere perdono che chiedere permesso”, un'altra grande idea è la **predizione**. In alcuni casi, è in genere più veloce tirare a indovinare e iniziare a lavorare di conseguenza, che aspettare di sapere con certezza. Questo è vero quando il meccanismo per recuperare una predizione sbagliata non è troppo costoso e la predizione è sufficientemente accurata. Utilizzeremo la sfera di cristallo come simbolo della predizione.



PREDIZIONE



## Gerarchia delle memorie

I programmatori vogliono che la memoria sia di grandi dimensioni, veloce e poco costosa, dato che la sua velocità spesso determina le prestazioni, la sua capacità limita la dimensione dei problemi che possono essere risolti e il suo costo rappresenta la voce di costo maggiore di un'architettura. I progettisti hanno scoperto che possono soddisfare queste esigenze contrapposte con la **gerarchia delle memorie**, nella quale la memoria più veloce, piccola e con il maggior costo per bit si trova in cima alla gerarchia e la memoria più lenta, più grande e con il minore costo per bit alla base. Come vedremo nel Capitolo 5, le memorie cache forniscono al programmatore l'illusione che la memoria principale sia veloce quasi quanto la memoria in cima alla gerarchia e sia economica e grande quasi quanto quella alla base della gerarchia. Utilizzeremo come simbolo della gerarchia delle memorie un triangolo a strati. La forma indica la velocità, il costo e la dimensione: più vicino è lo strato di memoria alla cima, maggiore sarà la sua velocità e il suo costo; più larga è la base, maggiore sarà la capacità della memoria.



## Affidabilità e ridondanza

I calcolatori non devono solo essere veloci, devono anche essere affidabili. Dato che tutti i dispositivi elettronici sono soggetti a guasti, per renderli affidabili dobbiamo introdurre dei componenti ridondanti che possano essere attivati quando si verifica un guasto e per aiutare a identificare i guasti stessi. Utilizzeremo come simbolo di affidabilità e ridondanza un bilico, dato che la doppia coppia di ruote degli assi posteriori consente al camion di continuare a viaggiare anche quando una ruota si sgonfia (anche se immaginiamo che il camionista si recherà subito presso un gommista per riparare la ruota, restaurando così la ridondanza!).

## 1.3 Che cosa c'è dietro un programma

*A Parigi la gente semplicemente si stupiva quando le parlavo in francese; non riuscii mai a fare comprendere a quegli idioti la loro lingua.*

Mark Twain, *Gli innocenti all'estero*, 1869.



ASTRAZIONE

**Software di sistema:** software che fornisce servizi di comune utilità. Comprende i sistemi operativi, i compilatori e gli assemblatori.

**Sistema operativo:** programma di supervisione che gestisce le risorse di un calcolatore a vantaggio dei programmi eseguiti su quel calcolatore.

Una tipica applicazione, come un programma per la scrittura di testi o un grande database, è costituita da milioni di linee di codice e si serve di sofisticate librerie software che implementano funzioni complesse di supporto all'applicazione stessa. Come vedremo più avanti, il calcolatore può solo eseguire istruzioni di basso livello estremamente semplici. Passare da un'applicazione complessa (scritta in un linguaggio vicino a quello dell'essere umano) alle semplici istruzioni che possono essere comprese dal calcolatore è un processo che coinvolge diversi strati di software che interpretano e traducono le operazioni definite ad alto livello nelle semplici istruzioni comprensibili al calcolatore. Questo è un esempio di **astrazione**.

Questi strati di software sono organizzati principalmente in maniera gerarchica come mostrato in Figura 1.3. Nel cerchio più esterno compaiono le applicazioni, mentre i diversi componenti del **software di sistema** sono posizionati nel cerchio intermedio tra l'hardware e le applicazioni software.

Il software di sistema ha diversi componenti, ma due sono quelli essenziali per tutti i calcolatori moderni: il sistema operativo e il compilatore.

Il **sistema operativo** permette di interfacciare i programmi utente con l'hardware del calcolatore, fornendo un gran numero di servizi e funzioni di supervisione. Alcune tra le funzioni più importanti del sistema operativo sono:

- gestire le operazioni base di input/output;
- allocare spazio nella memoria principale e nei dispositivi di memoria di massa;
- consentire a più applicazioni di utilizzare simultaneamente e in modo sicuro lo stesso calcolatore.

Linux, iOS e Windows sono tre esempi di sistemi operativi oggi largamente utilizzati.

I **compilatori** eseguono un'altra funzione vitale: la traduzione di un programma scritto in un linguaggio ad alto livello, come per esempio C, C++, Java o Visual Basic, in istruzioni eseguibili dall'hardware. Data la complessità dei moderni linguaggi di programmazione e la semplicità delle istruzioni che possono essere eseguite dall'hardware, la traduzione del codice scritto ad alto livello in istruzioni hardware è una funzione complessa. Introduciamo ora brevemente il processo di traduzione, che sarà trattato in modo più approfondito nel Capitolo 2.



## Da un linguaggio ad alto livello al linguaggio dell'hardware

Per parlare con una macchina elettronica è necessario inviare segnali elettrici; i segnali che un calcolatore può comprendere più facilmente sono *on* e *off* (acceso e spento), quindi l'alfabeto della macchina consta di due soli simboli. Proprio come le 26 lettere che costituiscono l'alfabeto inglese non pongono alcun limite a quante parole si possano scrivere, il fatto che l'alfabeto dei calcolatori sia costituito soltanto da due simboli non pone un limite a quello che può fare un calcolatore. Questi due simboli sono rappresentati dai numeri 0 e 1; il linguaggio macchina può quindi essere visto come una composizione di numeri in base 2, detti anche *numeri binari*, dove ogni "lettera" è una **cifra binaria** (*binary digit*), detta **bit**. I calcolatori obbediscono ai nostri comandi, chiamati istruzioni. Un'**istruzione** è semplicemente una stringa, o insieme, di bit comprensibile dal calcolatore e, perciò, può essere vista come un numero. Per esempio, la sequenza di bit:

```
1001010100101110
```

può ordinare a un particolare calcolatore di sommare due numeri. Il motivo per cui si utilizzano i numeri per rappresentare sia le istruzioni sia i dati verrà spiegato nel Capitolo 2. Non vogliamo anticiparne il contenuto, qui osserviamo solamente che ciò costituisce una delle idee fondamentali attorno alle quali sono costruiti i calcolatori.

I primi programmatore comunicavano con il calcolatore direttamente tramite numeri binari, ma era così macchinoso che in breve tempo furono inventate nuove notazioni assai più vicine al modo di pensare dell'essere umano. All'inizio tali notazioni venivano tradotte manualmente nel codice binario corrispondente, ma questo processo era ancora estenuante.

I pionieri dell'informatica allora inventarono del software in grado di tradurre una notazione simbolica in codice binario, usando il calcolatore stesso per programmare il calcolatore. Il primo di questi programmi fu chiamato **assembler (assemblatore)**: esso traduceva la versione simbolica delle istruzioni nella corrispondente forma binaria. Per esempio, un programmatore potrebbe scrivere:

```
add A, B
```

e l'assemblatore tradurrebbe questa notazione in:

```
1001010100101110
```

Questa istruzione dice al calcolatore di sommare i due numeri A e B. Il nome coniato per questo tipo di linguaggio simbolico, usato ancora oggi, è **linguaggio assembler (assembly language)**, mentre il linguaggio binario, compreso dal calcolatore, prende il nome di **linguaggio macchina**.

Sebbene il linguaggio assembler costituisca un notevole progresso, è ancora lontano dal costituire una notazione che uno scienziato vorrebbe utilizzare per

**Figura 1.3** Schema semplificato della relazione gerarchica tra hardware e software rappresentati come cerchi concentrici: l'hardware è rappresentato dal cerchio più interno e il software applicativo da quello più esterno. In applicazioni complesse ci sono spesso più livelli di applicazioni software. Per esempio, un database può essere eseguito al di sopra del software di sistema, che ha lanciato un'applicazione, a sua volta eseguita sopra il database.

**Compilatore:** programma che traduce le istruzioni scritte in linguaggio ad alto livello in istruzioni assembler.

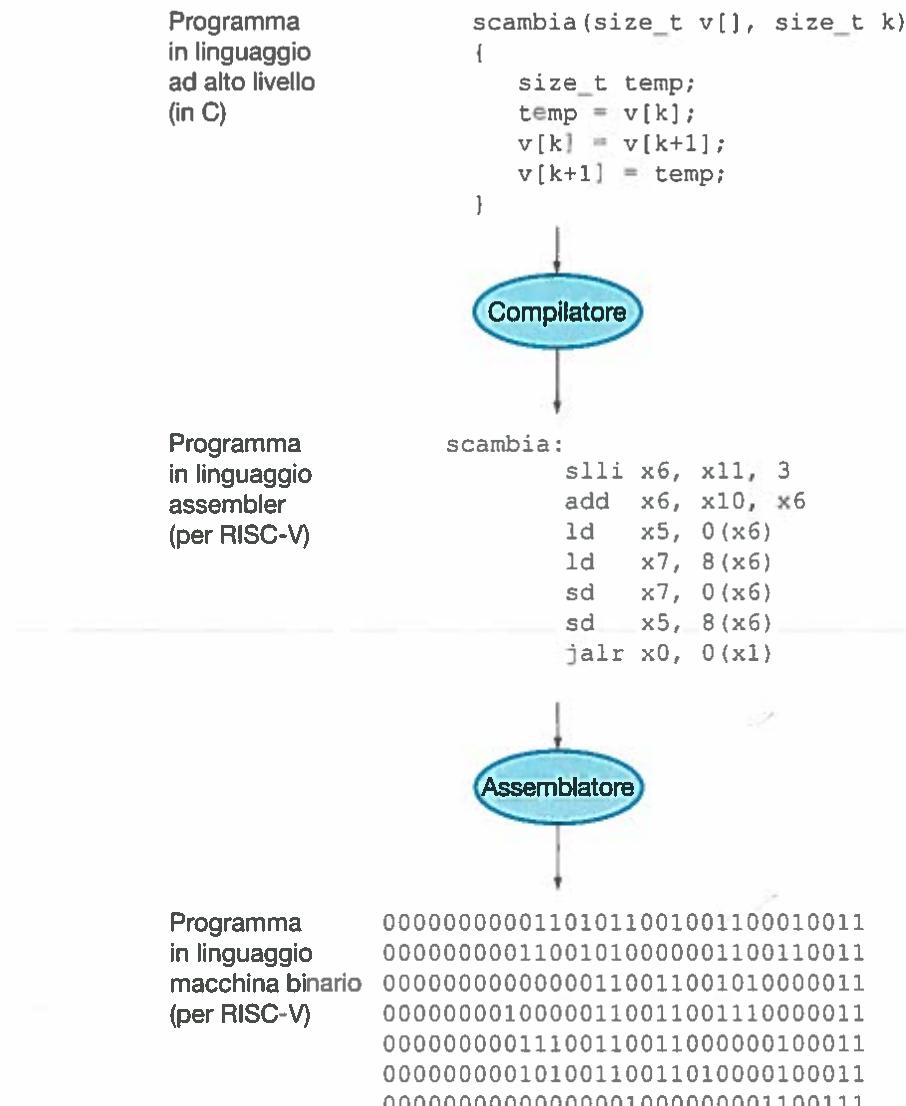
**Cifra binaria:** detta anche **bit**, è uno dei due numeri previsti dalla base 2 (0 o 1). È il componente elementare dell'informazione digitale.

**Istruzione:** un comando che l'hardware del computer comprende ed esegue.

**Assembler (assemblatore):** un programma che converte una versione simbolica delle istruzioni in una versione binaria comprensibile alla macchina.

**Linguaggio assembler:** rappresentazione simbolica delle istruzioni in linguaggio macchina.

**Linguaggio macchina:** rappresentazione binaria delle istruzioni comprensibile dall'hardware di un calcolatore.



**Figura 1.4** Programma in linguaggio C, compilato in linguaggio assembler e quindi tradotto (assemblato) in linguaggio macchina. Anche se la traduzione da linguaggio ad alto livello a linguaggio macchina viene qui suddivisa in due fasi, alcuni compilatori eliminano lo stadio intermedio producendo direttamente il codice in linguaggio macchina binario. Questi linguaggi e questo programma verranno esaminati più nel dettaglio nel Capitolo 2.

**Linguaggio di programmazione ad alto livello:** un linguaggio portabile come C, C++, Java o Visual Basic, composto da parole ed espressioni algebriche che possono essere tradotte da un compilatore in linguaggio assembler.



ASTRAZIONE

simulare il flusso di un fluido o che un amministratore potrebbe impiegare per la gestione dei propri registri contabili. Il linguaggio assembler richiede la scrittura di una linea per ogni istruzione che il calcolatore eseguirà, forzando così il programmatore a pensare come la macchina.

Avere compreso che si poteva scrivere un programma per tradurre in linguaggio macchina istruzioni scritte in un linguaggio più potente è stato uno dei passi più importanti compiuti all'inizio dell'era informatica. I programmatori di oggi devono non solo la loro produttività ma anche la loro salute mentale alla creazione di **linguaggi di programmazione ad alto livello** e di compilatori che traducono i programmi scritti in tali linguaggi nelle istruzioni del linguaggio macchina. La Figura 1.4 mostra la relazione tra i programmi scritti nei tre diversi linguaggi. Questo è un altro esempio della potenza dell'astrazione.

Un compilatore consente al programmatore di scrivere un'espressione come la seguente in un linguaggio ad alto livello:

A + B

che il compilatore tradurrebbe nel seguente codice assembler:

add A, B

Come mostrato nella figura precedente, l'assemblatore tradurrebbe infine questa istruzione nell'istruzione binaria che dice al calcolatore di sommare i due numeri A e B.

I linguaggi di programmazione ad alto livello offrono diversi e importanti vantaggi. Innanzitutto permettono al programmatore di ragionare in un linguaggio più naturale, basato su parole della lingua inglese e su notazioni algebriche, con il risultato di produrre programmi assai più simili a un testo che alle criptiche tabelle di simboli del linguaggio macchina (Figura 1.4). Inoltre permettono di progettare linguaggi specifici per l'uso al quale sono destinati, come è accaduto con il Fortran, progettato per i calcoli scientifici, o il Cobol, orientato al trattamento dei dati finanziari, o il Lisp, concepito per la manipolazione simbolica ecc. Sono nati anche linguaggi specifici per particolari settori e per gruppi di utenti ancora più ristretti, come la piccola comunità dei ricercatori interessati alla simulazione dei fluidi.

Il secondo vantaggio dei linguaggi ad alto livello è l'incremento di produttività dei programmatori. Uno dei pochi argomenti sullo sviluppo del software sul quale esiste un consenso generalizzato è la convinzione che occorre meno tempo per scrivere un programma utilizzando un linguaggio che richiede meno linee di codice per esprimere un dato concetto. La brevità è un evidente vantaggio dei linguaggi di programmazione ad alto livello rispetto al linguaggio assembler.

L'ultimo beneficio consiste nel fatto che un programma scritto in un linguaggio ad alto livello non dipende dal calcolatore sul quale viene sviluppato, poiché i compilatori e gli assemblatori possono tradurre il codice scritto in quel linguaggio nelle istruzioni in linguaggio macchina di un qualsiasi calcolatore. Questi tre vantaggi sono talmente importanti che oggi solo una minima parte dei programmi viene scritta in linguaggio assembler.

## 1.4 Componenti di un calcolatore

Dopo aver dato un'occhiata a quello che c'è dietro i vostri programmi e aver scoperto il software sottostante, apriamo la scatola del calcolatore per imparare qualcosa riguardo all'hardware. L'hardware di qualsiasi calcolatore svolge le stesse funzioni di base: acquisire dati dall'esterno (input), fornire dati all'esterno (output), elaborare i dati e memorizzarne i risultati. Lo scopo principale di questo libro è mostrarvi come queste funzioni sono realizzate e i capitoli seguenti tratteranno le diverse parti del calcolatore che svolgono questi quattro compiti.

Quando ci troveremo ad affrontare un argomento importante, in cui i concetti espressi sono fondamentali per la comprensione del calcolatore, lo enfatizzeremo attraverso un quadro d'insieme. In questo libro sono presenti una dozzina di sezioni *Quadro d'insieme*; il primo, qui di seguito, descrive i cinque componenti di un calcolatore che consentono di realizzare le funzioni di input, output, elaborazione e memorizzazione dei dati.

Due componenti chiave dei calcolatori sono i **dispositivi di input**, come per esempio il microfono, e i **dispositivi di output**, come per esempio gli altoparlanti. Come suggerito dal nome, un dispositivo di input fornisce i dati al

**Dispositivo di input:** meccanismo mediante il quale un calcolatore riceve i dati dall'esterno. Esempi di questi dispositivi sono la tastiera e il mouse.

**Dispositivo di output:** meccanismo che invia il risultato dell'elaborazione all'utente o a un altro calcolatore.

calcolatore mentre un dispositivo di output invia all'utente i risultati dell'elaborazione. Alcuni dispositivi, come le reti wireless, hanno funzioni sia di input sia di output.

I dispositivi di input/output (I/O) verranno descritti con maggior dettaglio nei Capitoli 5 e 6, mentre ora ci limiteremo a una panoramica introduttiva dell'hardware del calcolatore partendo dai dispositivi esterni di I/O.

## QUADRO D'INSIEME

I cinque componenti classici di un calcolatore sono input, output, memoria, unità di elaborazione dati (*datapath*) e unità di controllo. Gli ultimi due componenti vengono considerati, a volte, un unico componente, chiamato processore. La Figura 1.5 mostra l'organizzazione tipica di un calcolatore; questa è indipendente dalla tecnologia hardware: è infatti possibile attribuire a una di queste cinque categorie ogni parte di ogni calcolatore passato e presente. ■

### Attraverso lo specchio<sup>1</sup>

**Schermo a cristalli liquidi (LCD):** una tecnologia con cui sono costruiti gli schermi dei calcolatori che adotta un sottile strato di polimeri liquidi che può essere utilizzato per trasmettere o bloccare la luce a seconda che si applichi o meno una carica elettrica.

**Schermo a matrice attiva:** schermo a cristalli liquidi che utilizza un transistor per controllare la trasmissione della luce in ogni singolo pixel.

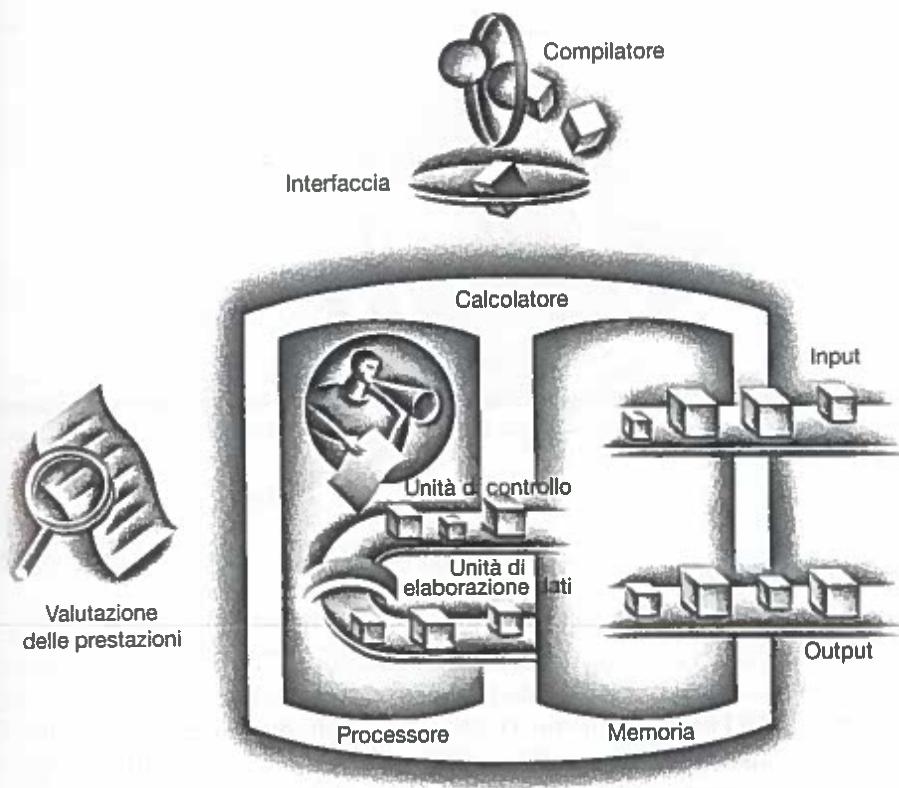
**Pixel:** il più piccolo elemento di un'immagine digitale. Gli schermi contengono dalle centinaia di migliaia ai milioni di pixel, organizzati in una matrice.

Il dispositivo di I/O più affascinante è probabilmente lo schermo grafico. La maggior parte dei dispositivi mobili personali utilizza uno **schermo a cristalli liquidi (LCD, Liquid Crystal Display)** per avere uno schermo sottile a basso consumo energetico. Lo schermo LCD non genera luce, ma controlla la trasmissione della luce attraverso di esso. Un tipico pannello LCD è composto da molecole a forma di bastoncino immerse in un liquido, nel quale assumono una forma a elica. Queste molecole deviano la luce che colpisce il pannello, emessa da una sorgente posta dietro lo schermo o, più raramente, riflessa; invece, quando si applica una corrente, le molecole a bastoncino si raddrizzano e non deviano più la luce. Dato che il materiale di cui sono costituiti i cristalli liquidi è inserito fra due schermi polarizzati a 90° l'uno rispetto all'altro, la luce – se non è deviata da uno dei due schermi – non può passare attraverso lo schermo. Oggi quasi tutti gli LCD sono basati su **matrici attive**, che hanno un piccolo interruttore a transistor associato a ogni punto, o pixel, che controlla con precisione la corrente e consente così di ottenere immagini più nitide. A ogni punto è associata una maschera rossa-verde-blu che determina l'intensità delle tre componenti additive del colore del punto nell'immagine visualizzata; in un LCD a colori a matrice attiva ci sono tre interruttori a transistor per ogni pixel.

Un'immagine è costituita da una matrice di elementi, detti **pixel**, che possono essere rappresentati come una matrice di bit, detta *bitmap*. A seconda delle dimensioni e della risoluzione dello schermo, la bitmap di un tablet può variare da  $1024 \times 768$  pixel a  $2048 \times 1536$  pixel. Uno schermo a colori utilizza tipicamente 8 bit per ciascuno dei tre componenti (rosso, verde e blu) per un totale di 24 bit per pixel, consentendo quindi la visualizzazione di milioni di colori diversi.

In un calcolatore, il supporto hardware alla grafica viene fornito principalmente dal *frame buffer* (buffer di immagine), detto anche *raster refresh buffer*,

<sup>1</sup> Il testo inglese gioca sulla definizione di *looking glass* = “(schermo di) vetro (del monitor) per guardare” ma anche “specchio”; il titolo *Attraverso lo specchio* è lo stesso del secondo dei “Libri di Alice” di Lewis Carrol (che, non dimentichiamolo, in quanto grande logico matematico, può essere considerato un precursore degli informatici).



**Figura 1.5** Tipica organizzazione di un calcolatore con i cinque componenti classici. Il processore riceve le istruzioni e i dati dalla memoria; l'unità di input scrive i dati nella memoria e l'unità di output legge i dati dalla memoria. L'unità di controllo invia i segnali che determinano le operazioni che devono essere svolte dall'unità di elaborazione dati, dalla memoria e dalle unità di input e output.

nel quale viene memorizzata la bitmap, associata all'immagine da visualizzare a schermo; per visualizzare l'immagine, quindi, l'elettronica di controllo dello schermo legge dal frame buffer il vettore di bit associato a ogni pixel e lo invia al terminale alla frequenza di aggiornamento stabilita. La Figura 1.6 mostra un frame buffer semplificato che opera su immagini di 4 bit per pixel.

Lo scopo della bitmap è quello di rappresentare fedelmente ciò che compare sullo schermo. Le difficoltà nella realizzazione dei sistemi grafici nascono dal fatto che l'occhio umano è molto sensibile nel rilevare cambiamenti anche minimi sullo schermo.

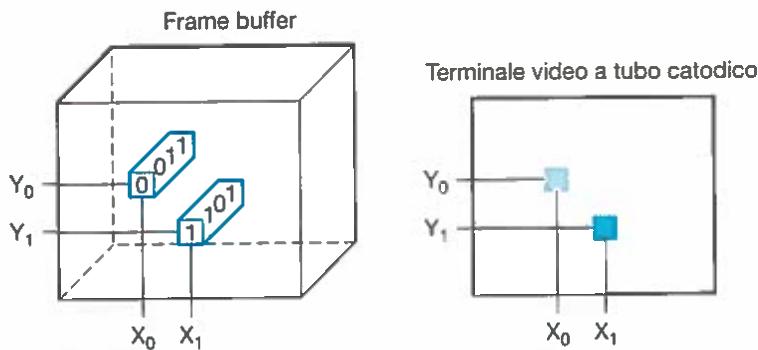
## Touchscreen

Anche i PC utilizzano gli LCD ma i tablet e gli smartphone dell'era post-PC hanno sostituito la tastiera e il mouse con schermi sensibili al tatto. Questi hanno il grande vantaggio di consentire all'utente di indicare direttamente quello a cui è interessato invece di farlo attraverso il mouse.

Anche se ci sono diversi modi di realizzare un touchscreen, molti tablet utilizzano oggi la tecnologia capacitiva. Se noi, che siamo buoni conduttori dell'elettricità, tocchiamo un isolante, come il vetro, ricoperto di materiale conduttivo trasparente, il suo campo elettrico superficiale viene distorto, provocando una variazione della capacità. Questa tecnologia consente di toccare lo schermo in più punti simultaneamente permettendo quindi l'interazione mediante gesti e la realizzazione di interfacce utente particolarmente attraenti.

*Attraverso lo schermo del calcolatore ho fatto atterrare un aereo sul ponte di una portaerei, ho osservato una particella nucleare colpire una buca di potenziale, ho volato in un razzo quasi alla velocità della luce e ho visto un calcolatore svelare i suoi meccanismi più interni.*

*Ivan Sutherland, il "padre" della computer graphics, citato in "Computer Software for Graphics", Scientific American, 1984*



**Figura 1.6** Ciascuna coordinata nel frame buffer a sinistra determina il livello di grigio nella posizione corrispondente dello schermo a tubo catodico (CRT) posto a destra. Il pixel  $(X_0, Y_0)$  contiene la sequenza 0011, che corrisponde sullo schermo a un colore più chiaro della sequenza 1101 associata al pixel  $(X_1, Y_1)$ .

## Dentro la scatola

La Figura 1.7 mostra il contenuto di un tablet, un iPad2 di Apple. Dei cinque componenti classici che costituiscono un calcolatore, l'I/O è quello preponderante in questo calcolatore che è dedicato soprattutto alla lettura. L'elenco dei dispositivi di I/O comprende: un LCD multi-touch, due videocamere, una di fronte e una sul retro, un microfono, un connettore jack per le cuffie, altoparlanti, un accelerometro, un giroscopio e la connessione Wi-Fi e Bluetooth alla rete. L'unità di elaborazione dati, l'unità di controllo e la memoria hanno una dimensione ridotta.

I rettangolini di Figura 1.8 sono i dispositivi che guidano il nostro sviluppo tecnologico e sono chiamati **circuiti integrati** o **chip**. L'integrato A5 al centro di Figura 1.8 contiene due processori ARM con una frequenza di clock di 1 GHz. Il **processore** è la parte attiva del calcolatore, quella che esegue fedelmente le istruzioni di un programma: è in grado di effettuare la somma di numeri, compiere test su di essi, inviare segnali per attivare i dispositivi di I/O e così via. A volte il processore viene indicato con l'acronimo **CPU** che sta per **central processing unit** (*unità centrale di elaborazione*).

Entrando ancora più in dettaglio nell'hardware di un calcolatore, la Figura 1.9 rivela alcuni particolari di un microprocessore. Il processore comprende due componenti principali: l'**unità di elaborazione dati (datapath)** e l'**unità di controllo**, rispettivamente il braccio e la mente del processore. Il datapath provvede a eseguire le operazioni aritmetiche mentre l'unità di controllo indica al datapath, alla memoria e ai dispositivi di I/O che cosa fare a seconda di quanto specificato nelle istruzioni di un programma. Il datapath e l'unità di controllo di un calcolatore progettato per ottenere prestazioni elevate saranno descritti nel Capitolo 4.

L'integrato A5 di Figura 1.8 contiene anche 2 chip di memoria, ciascuno con una capacità di 2 Gibibit, e fornisce quindi una capacità totale di 512 MiB. La **memoria** è il luogo dove vengono tenuti i programmi in esecuzione assieme ai loro dati ed è costituita da chip di DRAM. DRAM sta per **memoria dinamica ad accesso casuale**; più DRAM possono essere utilizzate assieme per memorizzare istruzioni e dati di un programma. A differenza delle memorie ad accesso sequenziale, quali i nastri magnetici, le memorie DRAM, e le RAM in genere, sono memorie il cui accesso richiede lo stesso tempo indipendentemente dalla particolare area di memoria cui si accede.

Analizzando più da vicino i singoli componenti dell'hardware, possiamo scoprire ulteriori dettagli sul calcolatore (Figura 1.9). All'interno del proces-

**Circuito integrato:** detto anche **chip**, è un dispositivo costituito da decine di milioni di transistor.

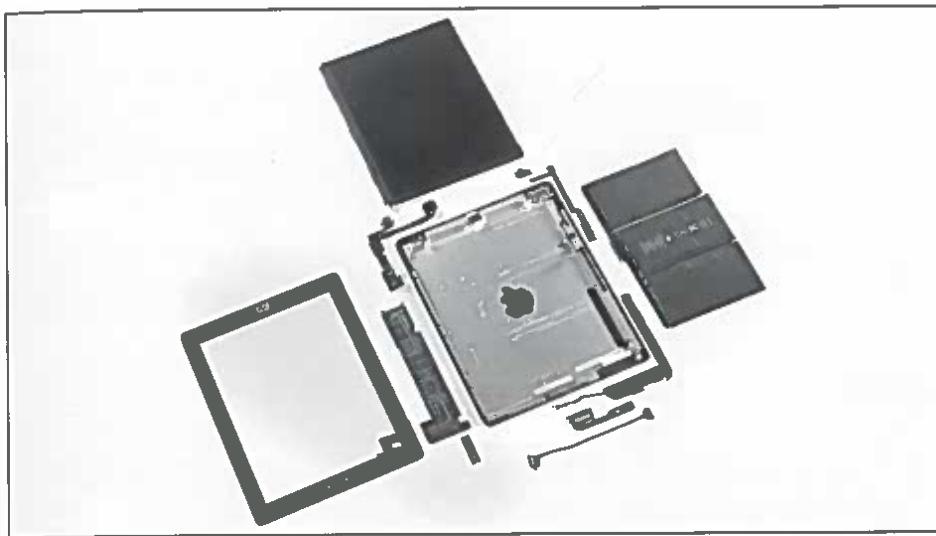
**Processore (CPU):** è la parte attiva del calcolatore: contiene le unità di controllo e di elaborazione dei dati. È in grado di effettuare la somma di numeri, test su di essi, inviare segnali per attivare i dispositivi di I/O ecc.

**Datapath:** detta anche **unità di elaborazione dati**, è il componente della CPU che provvede a eseguire le operazioni aritmetico-logiche sui dati.

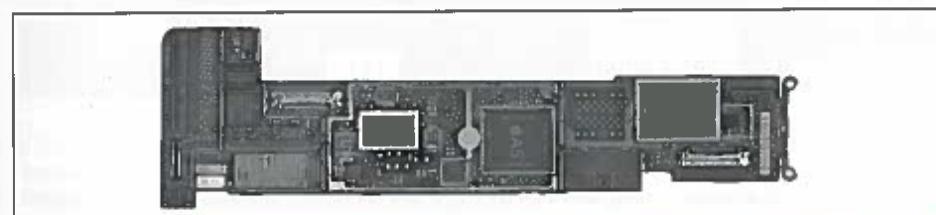
**Unità di controllo:** è il componente del processore che invia i comandi all'unità di elaborazione dati, alla memoria e ai dispositivi di I/O secondo le istruzioni del programma.

**Memoria:** luogo in cui vengono tenuti i programmi quando vengono eseguiti, assieme ai dati richiesti dalla loro esecuzione.

**DRAM (memoria dinamica ad accesso casuale):** memoria costruita con circuiti integrati che permette l'accesso a ogni locazione nello stesso tempo. Il tempo di accesso è dell'ordine dei 50 nanosecondi e il costo per Gigabyte era compreso tra i 5 e i 10 dollari americani nel 2012.



**Figura 1.7 Componenti di un iPad2 A1395 di Apple.** Il retro di metallo dell'iPad (con il logo rovesciato della Apple nel mezzo) si trova al centro. In alto potete vedere lo schermo LCD multi-touch; sulla destra trovate la batteria a polimeri da 3,8 V e 25 W/h, che consiste in tre moduli agli ioni di litio e consente un'autonomia di 10 ore; sulla sinistra potete vedere la cornice che collega lo schermo al corpo dell'iPad. Attorno al retro dell'iPad, potete vedere dei minuscoli componenti: questi sono i componenti che immaginiamo quando pensiamo a un calcolatore; hanno spesso una forma a L per potere essere alloggiati all'interno del calcolatore di fianco alla batteria. In Figura 1.8 viene mostrato uno zoom della schedina a forma di L che si trova nell'angolo in basso a sinistra dell'iPad e che contiene il processore e la memoria. Il piccolo rettangolo che vedete in Figura 1.7 sotto questa schedina contiene il chip che provvede alla comunicazione wireless: Wi-Fi, Bluetooth o radio in modulazione di frequenza (FM) e viene collocato in un piccolo alloggiamento posto nell'angolo a sinistra in basso della schedina del processore. Vicino allo spigolo superiore sinistro del corpo dell'iPad potete vedere un altro componente a forma di L: è la videocamera frontale con i suoi accessori che sono il jack per le cuffie e il microfono, mentre vicino allo spigolo superiore destro potete vedere la schedina che contiene il controllo del volume e il pulsante di rotazione dello schermo con un giroscopio e un accelerometro che, combinati, consentono di riconoscere i 6 gradi di libertà del movimento. Il piccolo rettangolo a fianco è la videocamera posteriore, mentre vicini allo spigolo in basso a destra potete vedere gli altoparlanti. Il cavo in basso connette la schedina della CPU con la schedina che controlla la videocamera e il volume. La schedina posta tra il cavo e gli altoparlanti contiene il controllore del touchscreen capacitivo (Fonte: iFixit, [www.ifixit.com](http://www.ifixit.com)).



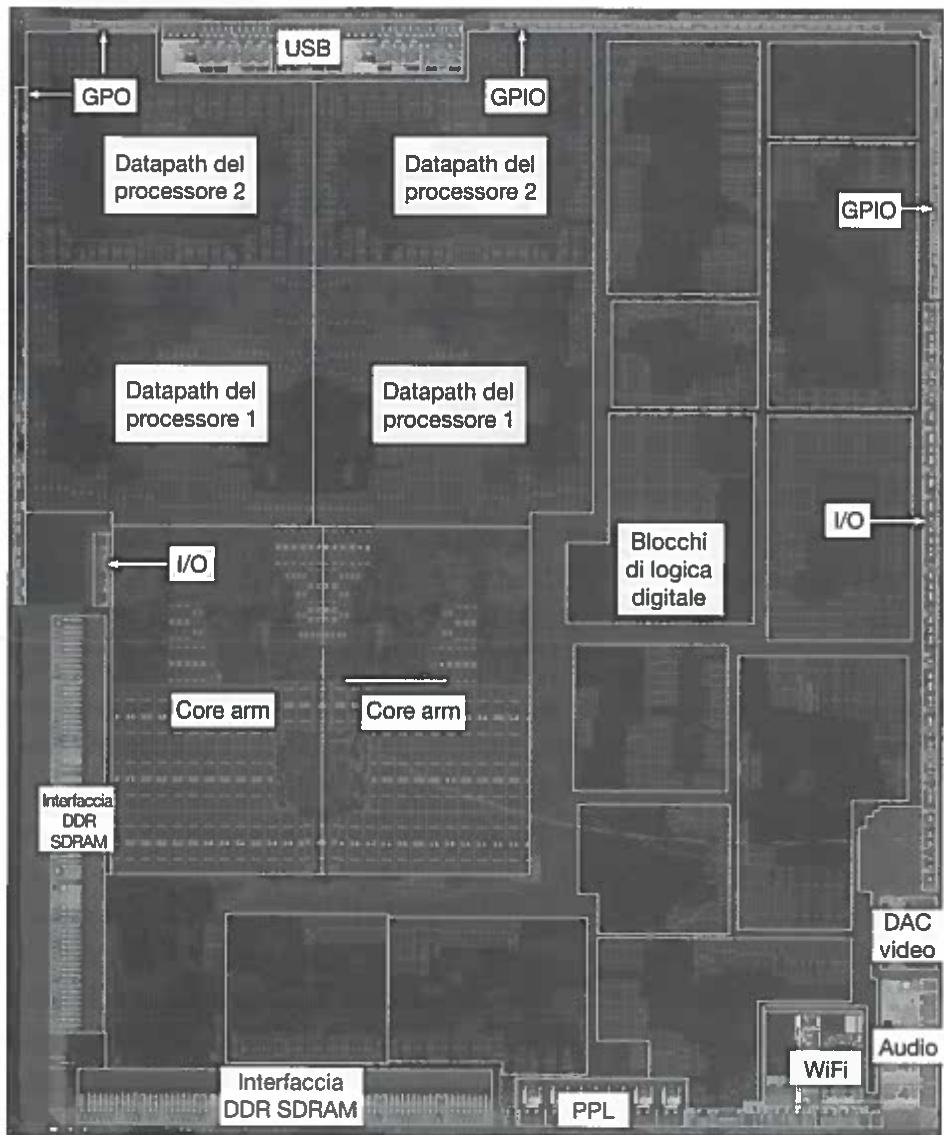
**Figura 1.8 Schedina del processore dell'iPad2 di Figura 1.7.** La fotografia mostra cinque circuiti integrati. Il grande circuito integrato al centro è un chip A5 di Apple, che contiene un processore ARM a due core, che lavorano alla frequenza di clock di 1 GHz, e 512 MB di memoria principale. La Figura 1.9 mostra una fotografia del processore contenuto nel chip A5. Il chip di dimensioni simili a sinistra contiene una memoria flash, non volatile, da 32 GB. È presente dello spazio vuoto tra i due chip, all'interno del quale si può inserire un secondo chip di memoria flash per duplicare la capacità di memorizzazione dell'iPad. I chip sulla destra dell'A5 contengono il controllore della tensione e i chip di I/O (Fonte: iFixit, [www.ifixit.com](http://www.ifixit.com)).

sore c'è un secondo tipo di memoria, la **memoria cache**, che consiste in una memoria piccola ma veloce che funge da "tampone" nei confronti della memoria DRAM (una definizione non tecnica di *memoria cache* è: un posto sicuro dove nascondere le cose). La cache è costruita usando una tecnologia di memoria differente detta **SRAM**, **memoria statica ad accesso casuale**. Le memorie di tipo SRAM sono più veloci, ma meno dense e quindi più costose, delle memorie di tipo DRAM. SRAM e DRAM sono due strati della **gerarchia delle memorie**.

**Memoria cache:** una memoria piccola e veloce che agisce come un buffer per una memoria più grande e più lenta.

**SRAM (RAM statica):** memoria costruita con circuiti integrati, più veloce ma meno densa delle memorie DRAM.





**Architettura dell'insieme di istruzioni (ISA):** detta anche semplicemente **architettura**. Interfaccia astratta tra l'hardware e il livello più basso del software di un calcolatore. Comprende tutte le informazioni necessarie per scrivere un programma in linguaggio macchina funzionante in modo corretto, comprese le istruzioni, i registri, gli accessi a memoria, l'I/O ecc.



ASTRAZIONE

Come abbiamo appena visto, una delle grandi idee per migliorare i calcolatori è l'astrazione e una delle astrazioni più importanti è l'interfaccia tra l'hardware e il software di più basso livello. Per la sua importanza, è stato associato un nome specifico a questa interfaccia: **architettura dell'insieme di istruzioni** del calcolatore (**ISA**, *Instruction Set Architecture*), o più semplicemente **architettura** del calcolatore. L'architettura dell'insieme di istruzioni comprende tutto ciò che i programmatore devono sapere per scrivere un programma in linguaggio macchina correttamente funzionante e comprende quindi le istruzioni, i dispositivi di I/O ecc. Tipicamente, un sistema operativo incapsula i dettagli sull'I/O, sull'allocazione della memoria e su altre informazioni di basso livello, evitando quindi al programmatore di dover gestire questi dettagli. La combinazione dell'insieme di base delle istruzioni e dell'interfaccia del sistema operativo fornita ai programmatore per scrivere le applicazioni

**Figura 1.9** Circuito integrato del processore contenuto nel chip dell'A5. Il chip ha dimensioni di 12,1 × 10,1 mm ed è stato inizialmente prodotto con la tecnologia a 45 nm (par. 1.5). Contiene due processori, o core, ARM che sono alloggiati in mezzo a sinistra del chip e una GPU (unità di elaborazione grafica) PowerVR contenente quattro unità di elaborazione dati che potete vedere in alto a sinistra. A sinistra in basso potete vedere le interfacce verso la memoria principale (DRAM) (Fonte: Chipworks, [www.chipworks.com](http://www.chipworks.com)).

è chiamata **interfaccia binaria delle applicazioni (ABI, Application Binary Interface)**.

Un'architettura dell'insieme di istruzioni permette ai progettisti di descrivere le funzionalità di un calcolatore in maniera completamente indipendente dall'hardware che le implementa. Per esempio, si può parlare delle funzioni di un orologio digitale (misurare il tempo, visualizzare l'ora, programmare la sveglia) indipendentemente dalla tecnologia hardware con cui l'orologio viene realizzato (schermo a cristalli liquidi, schermo a LED, tasti in plastica). Tutti i progettisti di calcolatori distinguono fra architettura e **implementazione** di un'architettura allo stesso modo: un'implementazione di un'architettura è un hardware che realizza la descrizione astratta dell'architettura. Questi concetti ci portano verso un altro *Quadro d'insieme*.

## QUADRO D'INSIEME

Sia l'hardware sia il software consistono in diversi livelli gerarchici basati su diversi livelli di astrazione, ognuno dei quali nasconde alcuni dettagli al livello superiore. Questo principio di astrazione consente sia ai progettisti hardware sia ai progettisti software di affrontare la complessità della progettazione dei calcolatori. Un'interfaccia chiave tra i due livelli di astrazione è l'**architettura dell'insieme di istruzioni (ISA, Instruction Set Architecture)**, ossia l'interfaccia tra l'hardware e il software di basso livello. Quest'interfaccia astratta permette diverse *implementazioni* hardware, caratterizzate da diversi costi e prestazioni, per eseguire lo stesso software. ■

### Un posto sicuro per i dati

Finora abbiamo visto come inserire dati, eseguire calcoli su di essi e visualizzare il risultato. Tuttavia, se dovesse interrompersi l'alimentazione del calcolatore, tutto andrebbe perso perché la memoria interna del calcolatore è una **memoria volatile**. Questo significa che quando viene a mancare l'alimentazione, la memoria dimentica tutto. Al contrario, quando si spegne il lettore DVD, il disco DVD non dimentica mai il film che vi è stato registrato sopra; questo è un esempio di **memoria non volatile**.

Per distinguere tra la memoria utilizzata per memorizzare i programmi mentre questi vengono eseguiti e la memoria non volatile utilizzata per salvare i programmi tra un'esecuzione e l'altra, si utilizza il termine **memoria principale** (o **memoria primaria**) per la prima e **memoria di massa** (o **memoria secondaria**) per la seconda. La memoria di massa costituisce il secondo strato della gerarchia delle memorie. Le DRAM hanno dominato il settore delle memorie principali sin dal 1975, mentre i **dischi magnetici** (o **disco rigido – hard disk**) hanno dominato quello delle memorie di massa ben da prima. Per la loro dimensione e la loro forma, i dispositivi mobili utilizzano le **memorie flash**, memorie non volatili a semiconduttore, al posto dei dischi rigidi. La Figura 1.8 mostra un chip che contiene la memoria flash di un iPad2. Sebbene più lente delle DRAM, le memorie flash costano molto meno e possiedono anch'esse la caratteristica della non volatilità. Anche se hanno un costo per bit maggiore dei dischi, le memorie flash sono più compatte, sono disponibili con capacità inferiori, sono più robuste e consumano meno corrente dei dischi. A differenza dei dischi e delle DRAM, il contenuto di un bit di una memoria flash tende a svanire dopo essere stato riscritto da 100 000 a 1 000 000 di volte. Per evitare ciò, il file system deve tenere traccia del numero di scritture e adottare una strategia

**Interfaccia binaria delle applicazioni (ABI):** è il sottoinsieme di istruzioni dedicate all'utente più l'interfaccia del sistema operativo fornita al programmatore per scrivere le applicazioni. Definisce uno standard per la portabilità del codice binario tra calcolatori differenti.

**Implementazione:** hardware che rispetta le specifiche imposte dalla descrizione astratta di un'architettura.

**Memoria volatile:** memoria in grado di mantenere i dati solamente se è alimentata. Un tipico esempio di memoria volatile sono le DRAM.

**Memoria non volatile:** memoria che conserva i dati anche quando viene a mancare l'alimentazione; viene utilizzata per conservare i dati fra un'esecuzione e l'altra. I dischi DVD costituiscono una memoria non volatile.

**Memoria principale:** detta anche **memoria primaria**, viene utilizzata per contenere i programmi durante la loro esecuzione. Nei calcolatori moderni è tipicamente costituita da DRAM.

**Memoria di massa:** detta anche **memoria secondaria**, è una memoria non volatile utilizzata per conservare i programmi e i dati fra un'esecuzione e l'altra. È tipicamente una memoria flash nei PMD mentre è costituita da dischi magnetici nei server.

**Disco magnetico (o disco rigido – hard disk):** memoria secondaria non volatile costituita da piatti rotanti ricoperti da materiale magnetico in grado di memorizzare informazioni binarie. Dato che i dischi ruotano meccanicamente, il tempo di accesso è compreso tra i 5 e i 20 millisecondi. Il costo per Gigabyte era compreso tra 0,05 e 0,10 dollari americani nel 2012.



**Memoria flash:** una memoria non volatile a semiconduttori. È più lenta e meno costosa delle memorie DRAM, ma più costosa e più veloce dei dischi rigidi. Il tempo di accesso è compreso tra i 5 e i 50 microsecondi e il costo per Gigabyte era compreso tra 0,75 e 1 dollaro americano nel 2012.

per evitare che i dati memorizzati svaniscano, per esempio trasferendo i dati che vengono scritti più di frequente. I dischi e le memorie flash verranno descritti in maggiore dettaglio nel Capitolo 5.

## Comunicare con gli altri calcolatori

Abbiamo visto come sia possibile leggere, elaborare, mostrare e memorizzare dati; rimane ancora da esaminare un ultimo elemento importante che fa parte dei moderni elaboratori: il collegamento alla rete. Come il processore mostrato in Figura 1.5 è collegato alla memoria e agli altri dispositivi di I/O, la rete collega tra loro i calcolatori consentendo agli utenti di estendere la capacità di elaborazione attraverso la comunicazione. Le reti sono divenute talmente diffuse da costituire la spina dorsale dei moderni sistemi di elaborazione; una nuova macchina priva della possibilità di interfacciarsi alla rete sarebbe improponibile. I vantaggi derivanti dalla presenza di un collegamento in rete fra calcolatori sono diversi:

- *comunicazione*: le informazioni vengono scambiate fra i calcolatori a velocità elevata;
- *condivisione delle risorse*: i vari calcolatori collegati alla rete possono condividere gli stessi dispositivi di I/O, invece di possederne ciascuno una copia privata;
- *accesso non locale*: collegando fra loro elaboratori posti a grande distanza, gli utenti non devono più essere fisicamente vicini al calcolatore che stanno utilizzando.

**Rete locale (LAN)**: rete progettata per trasportare dati in un'area geograficamente limitata, tipicamente all'interno di un singolo edificio.

**Rete geografica (WAN)**: rete che si può estendere per centinaia di chilometri e può arrivare a coprire un intero continente.

Le distanze tra i nodi e le prestazioni di una rete sono variabili; il costo della comunicazione aumenta proporzionalmente sia alla velocità della trasmissione sia alla distanza del ricevitore. Il tipo di rete forse più diffusa è *Ethernet*, che prevede una distanza massima tra due nodi di un chilometro e permette di trasferire dati fino a una velocità massima di 40 Gigabit al secondo. Questi elementi rendono Ethernet adatta a connettere calcolatori posti sullo stesso piano di un edificio: Ethernet è quindi un esempio di quella che viene generalmente chiamata una **rete locale (LAN, Local Area Network)**. Le reti locali possono essere interconnesse tra loro con interruptori di rete, che forniscono sia il corretto instradamento della trasmissione sia la sicurezza nella comunicazione stessa. Le **reti geografiche (WAN, Wide Area Network)**, invece, attraversano i continenti e sono l'ossatura di Internet, che supporta il web. Esse sono tipicamente realizzate con fibre ottiche e sono noleggiate dalle società di telecomunicazioni.

Le reti hanno cambiato la faccia dell'informatica negli ultimi 30 anni sia per la diffusione di programmi, sia per l'aumento esponenziale delle prestazioni. Negli anni '70 solo pochi individui avevano accesso a una casella di posta elettronica, Internet e il web non esistevano, e spedire fisicamente nastri magnetici era il modo principale per trasferire tra due luoghi grandi quantità di dati. Le reti locali erano quasi inesistenti e c'erano solo pochissime reti di tipo geografico che avevano modesta capacità e accesso limitato.

Con il miglioramento della tecnologia, le reti divennero molto più economiche e con capacità significativamente più elevata. Per esempio, la prima tecnologia standardizzata per reti locali, sviluppata circa 30 anni fa, era una versione di Ethernet con capacità massima, chiamata **banda**, di 10 milioni di bit al secondo, tipicamente condivisa da qualche decina se non centinaia di calcolatori. Oggi la tecnologia delle reti locali offre capacità che vanno da 100 milioni di bit al secondo a 40 Gigabit al secondo, di solito condivisi da pochi calcolatori. La tecnologia ottica di comunicazione ha permesso una crescita simile della banda delle reti geografiche, che è passata da qualche centinaia di

kilobit ai Gigabit al secondo, e nel numero di calcolatori connessi, passato da qualche centinaio fino ai milioni di calcolatori oggi connessi alla rete mondiale. Il rapido incremento della copertura territoriale delle reti, unito al rapido incremento della banda di trasmissione, ha reso centrale la tecnologia delle reti per la rivoluzione dell'informatica negli ultimi 30 anni.

Nell'ultimo decennio un'altra innovazione nell'ambito delle reti sta cambiando il modo in cui i calcolatori comunicano tra loro. La tecnologia wireless (radio) si è rapidamente diffusa portandoci nell'era post-PC. La possibilità di realizzare una radio con la stessa tecnologia a basso costo utilizzata per produrre i semiconduttori CMOS, utilizzati per memorie e microprocessori, ha permesso di ridurre significativamente il prezzo e ciò ha portato all'esplosione della diffusione dei dispositivi wireless. La tecnologia wireless attualmente disponibile, chiamata 802.11 che è il nome del suo standard IEEE, permette frequenze di trasmissione da 1 a poco meno di 100 milioni di bit al secondo. La tecnologia wireless è notevolmente diversa da quella delle reti basate su cavi, dato che tutti gli utenti presenti in una certa area condividono le stesse onde radio.

### Autovalutazione

Le memorie di tipo DRAM e flash a semiconduttore e quelle a disco differiscono in maniera significativa. Descrivere le differenze fondamentali per ognuna delle seguenti caratteristiche: volatilità, tempo di accesso medio e costo considerando come riferimento le DRAM.

## 1.5 | Tecnologie per la produzione di processori e memorie

I processori e le memorie sono progrediti con incredibile velocità, poiché i progettisti hanno da molto tempo adottato ogni novità della tecnologia elettronica nel tentativo di vincere la sfida di realizzare il calcolatore migliore. La Figura 1.10 riporta le tecnologie che sono state via via utilizzate negli anni, con accanto una stima delle prestazioni relative per unità di costo. Dato che la tecnologia determina quello che i calcolatori potranno fare e quanto velocemente potranno evolvere, tutti i professionisti dell'informatica dovrebbero avere una certa familiarità con le nozioni fondamentali dei circuiti integrati.

Un **transistor** è semplicemente un interruttore acceso/spento, controllato da un segnale elettrico. Il *circuito integrato* (IC) ha permesso di inserire da

**Transistor:** interruttore di tipo on/off controllato da un segnale elettrico.

Anno	Tecnologia utilizzata nei calcolatori	Prestazioni relative per unità di costo
1951	Tubo a vuoto (valvola)	1
1965	Transistor	35
1975	Circuito integrato	900
1995	Circuito integrato a grandissima scala di integrazione	2 400 000
2013	Circuito integrato a scala di integrazione estremamente grande	250 000 000 000

**Figura 1.10** Prestazioni relative per unità di costo delle tecnologie via via utilizzate dai calcolatori nel tempo. *Fonte:* Computer Museum, Boston. Il valore del 2013 è stato estrapolato dagli autori (par. 1.12).

**Circuito integrato a grandissima scala di integrazione (VLSI):** dispositivo che contiene da centinaia di migliaia a milioni di transistor.

**Silicio:** un elemento naturale che è semiconduttore.

**Semiconduttore:** una sostanza che non conduce bene l'elettricità.

**Lingotto di silicio cristallino:** una barra composta da un unico cristallo di silicio di diametro variabile tra 15 e 30 cm e una lunghezza compresa tra 30 e 60 cm.

**Wafer:** fetta di un lingotto di silicio, di spessore non superiore a 2,5 mm, utilizzata per costruire i circuiti integrati.

qualche decina a qualche centinaio di transistor in un singolo frammento di silicio, detto chip. Quando Gordon Moore predisse il raddoppio continuo delle risorse, predisse anche il continuo raddoppio del numero di transistor per chip. Per descrivere l'enorme incremento nel numero di transistor che si possono inserire sulla stessa piastrina di silicio, passato da alcune centinaia di migliaia ad alcuni milioni, si utilizza l'espressione *a grande scala di integrazione*, da cui il termine **circuito integrato a grandissima scala di integrazione (VLSI, Very Large-Scale Integrated Circuit)**.

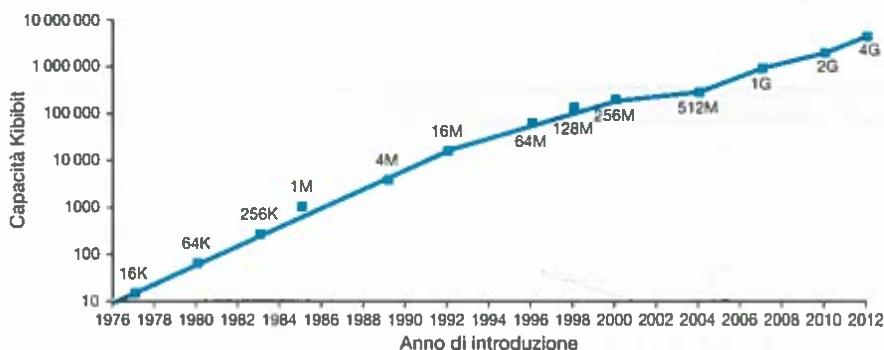
La crescita della scala di integrazione è rimasta notevolmente stabile negli anni. La Figura 1.11 mostra la crescita della capacità delle DRAM dal 1977. Per anni i produttori hanno puntualmente quadruplicato la capacità dei loro dispositivi ogni 3 anni, arrivando a ottenere un incremento finale superiore a un fattore 16 000!

Per capire come vengono fabbricati i circuiti integrati, partiamo dal principio. La produzione di un chip inizia dal **silicio**, una sostanza che si trova nella sabbia; dato che il silicio non conduce bene l'elettricità, è detto **semiconduttore**. Con uno speciale processo chimico è possibile aggiungere al silicio dei materiali che permettono a minuscole aree di trasformarsi in:

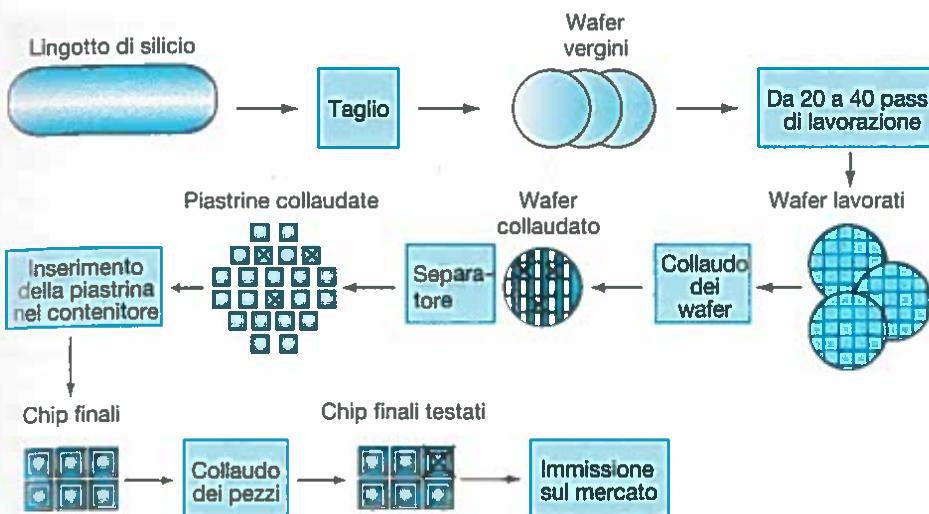
- eccellenti conduttori elettrici (utilizzando microscopici fili di rame o alluminio);
- eccellenti isolanti elettrici (paragonabili alle guaine in plastica o al vetro);
- aree che in particolari condizioni possono condurre *oppure* non condurre (analogamente a un interruttore).

I transistor rientrano nell'ultima categoria. Quindi, un circuito VLSI non è altro che una combinazione di milioni di elementi conduttori, isolanti e interruttori inseriti all'interno di un unico piccolo blocchetto.

Il processo produttivo dei circuiti integrati è fondamentale per determinare il costo dei chip ed è quindi importante per i progettisti dei calcolatori. La Figura 1.12 mostra schematicamente questo processo, che inizia con un **lingotto di silicio cristallino**, simile a una gigantesca salsiccia. Attualmente i lingotti hanno un diametro di 15-30 cm e sono lunghi 30-60 cm. Il lingotto viene tagliato in dischi sottili, detti **wafer**, spessi non più di 2,5 mm. I wafer attraversano una serie di passaggi di lavorazione, nel corso dei quali vengono depositate, secondo un disegno prestabilito, sostanze chimiche adatte a creare i transistor, gli elementi conduttori e gli elementi isolanti di cui abbiamo parlato prima; questi passi di lavorazione prendono il nome di **mascheratura**. I circuiti integrati attuali contengono solo uno strato di transistor, ma possono avere da due a otto strati di metallo conduttore, separati da strati di isolante.



**Figura 1.11** Crescita della capacità di un chip di DRAM in funzione del tempo. L'asse delle ordinate è misurato in Kibibit ( $2^{10}$  bit). L'industria ha quadruplicato la capacità delle DRAM circa ogni tre anni, con un incremento del 60% annuo, per 20 anni. Negli ultimi anni il tasso di crescita ha rallentato e la capacità raddoppia ogni due, tre anni.



**Figura 1.12 Processo produttivo di un chip.** Dopo essere stati tagliati dal lingotto di silicio, i wafer vergini subiscono dai 20 ai 40 processi di mascheratura sino a ottenere i circuiti desiderati (Figura 1.13). I wafer vengono quindi collaudati con un tester di wafer per stilare una mappa delle parti correttamente funzionanti. Vengono poi ritagliati in chip (Figura 1.9). Il wafer ha prodotto 20 chip, di cui 17 hanno passato il collaudo (una X indica che il chip è risultato difettoso). Il rendimento è in questo caso di 17/20, pari all'85%. I chip buoni vengono poi inseriti nei loro contenitori (*package*) e ulteriormente collaudati prima della spedizione ai clienti. In questo collaudo finale è stata trovata una parte danneggiata durante il confezionamento del chip.

Una singola microscopica falla nel wafer o una sbavatura in una sola delle decine di maschere utilizzate possono far sì che una particolare area del wafer risulti mal funzionante; questi **difetti** rendono virtualmente impossibile produrre un wafer perfetto. Per gestire queste imperfezioni si utilizzano diverse strategie, la più semplice delle quali consiste nel porre molti componenti indipendenti sullo stesso wafer; questo viene poi tagliato separando i diversi componenti, chiamati **piastrine** (*dies*), più informalmente note come **chip**. La Figura 1.13 mostra la fotografia di un wafer contenente i microprocessori prima che vengano separati in diverse piastrine, mentre la Figura 1.9 mostrava la fotografia della piastrina di un singolo microprocessore e i suoi componenti principali.

Il processo di separazione delle piastrine permette di scartare solo quei chip che contengono difetti, anziché l'intero wafer; l'entità di questo scarto è misurata dal cosiddetto **rendimento** del processo, definito come la percentuale di chip funzionanti rispetto al numero totale di chip costruiti su un singolo wafer.

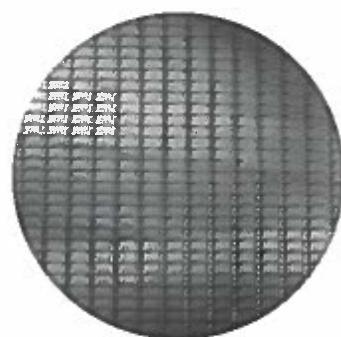
Il costo di un circuito integrato cresce rapidamente con le dimensioni della piastrina; ciò è dovuto sia al più basso rendimento che al minore numero di piastrine che si ricavano dal singolo wafer. Per ridurre i costi, una piastrina di grandi dimensioni viene spesso "rimpicciolita" facendo ricorso al processo di produzione della generazione successiva, in grado di creare transistor e connessioni di dimensioni inferiori. Ciò aumenta la resa e il numero di piastrine per wafer. La tecnologia disponibile nel 2012 era tipicamente a 32 nanometri (nm); ciò significa essenzialmente che il più piccolo componente disponibile sul chip era di 32 nm.

Una volta determinate le piastrine funzionanti, queste vengono connesse ai piedini di input/output del loro contenitore, utilizzando un processo denominato *bonding* (saldatura). I chip così inscatolati vengono collaudati un'ultima volta, per controllare che non si siano introdotti difetti nel processo di confezionamento, e immessi sul mercato.

**Difetto:** imperfezione microscopica in un wafer o che si è formata nei passi di lavorazione dei circuiti in esso contenuti. Può provocare il malfunzionamento del chip contenente il difetto.

**Piastrina:** sezione rettangolare ritagliata da un wafer, più informalmente denominata **chip**.

**Rendimento:** percentuale di chip funzionanti rispetto al numero totale di chip costruiti su un singolo wafer.



**Figura 1.13 Wafer di 300 mm (12 pollici)** di diametro contenente diversi processori Core Intel i7 (per gentile concessione di Intel). Il numero di piastrine per wafer ottenute con un rendimento del 100% sarebbe 280, ciascuna di dimensioni 20,7 per 10,5 mm. Le decine di chip che si trovano sul bordo del wafer sono inutilizzabili; questi vengono creati lo stesso solamente perché rendono molto più semplici le operazioni di mascheratura. Questo chip utilizza la tecnologia di 32 nanometri, che significa che i componenti più piccoli hanno una dimensione di circa 32 nm. In realtà questi sono solitamente ancora più piccoli delle dimensioni nominali, dato che le dimensioni dei transistor vengono arrotondate per eccesso alle dimensioni nominali di produzione.

**Approfondimento.** Il costo di un circuito integrato si può esprimere con le seguenti tre semplici equazioni:

$$\text{Costo per piastrina} = \frac{\text{Costo per wafer}}{\text{Piastrine per wafer} \times \text{rendimento}}$$

$$\text{Piastrine per wafer} \approx \frac{\text{Superficie del wafer}}{\text{Superficie della piastrina}}$$

$$\text{Rendimento} = \frac{1}{(1 + (\text{difetti per area} \times \text{area della piastrina}/2))^2}$$

La prima equazione è immediata da derivare. La seconda è un'approssimazione, dato che non viene sottratta la superficie vicina alla circonferenza del wafer che, essendo di forma circolare, non può contenere esattamente piastrine rettangolari (Figura 1.13). La terza equazione è basata sull'osservazione empirica del rendimento nelle fabbriche di circuiti integrati, e l'esponente dipende dal numero dei passi produttivi critici.

Perciò, a seconda della percentuale di difetti e delle dimensioni della piastrina e del wafer, i costi variano generalmente in modo non lineare con la superficie della piastrina.

### Autovalutazione

Un fattore chiave per determinare il costo di un circuito integrato è il suo volume di produzione. Quale dei seguenti motivi può essere la ragione per cui un chip prodotto in grandi volumi costa di meno?

1. Con grandi volumi, il processo di produzione può essere messo a punto su quel particolare dispositivo, aumentandone il rendimento.
2. Viene richiesto meno lavoro per progettare un componente con grandi volumi di produzione che un componente prodotto in piccole quantità.
3. Le maschere utilizzate per realizzare i chip sono costose e quindi il costo di un chip diminuisce per grandi volumi di produzione.
4. I costi di sviluppo e ingegnerizzazione sono alti e praticamente indipendenti dai volumi di produzione; quindi il costo di sviluppo di una singola piastrina è più basso per chip prodotti in grandi volumi.
5. I chip prodotti in grandi volumi di solito richiedono piastrine di dimensioni più piccole e quindi un singolo wafer offre un rendimento maggiore.

## 1.6 | Prestazioni

Valutare le prestazioni di un calcolatore può rivelarsi molto difficile. Le dimensioni e la complessità dei moderni sistemi software, unite all'ampia varietà di tecniche di ottimizzazione implementate dai progettisti hardware per migliorare le prestazioni, hanno reso la valutazione delle prestazioni ancora più difficile.

Quando vi trovate a decidere quale calcolatore acquistare, le prestazioni sono sicuramente uno degli elementi più importanti. Misurarle in modo accurato e compararle con quelle degli altri calcolatori è fondamentale per la scelta della macchina da acquistare e lo è quindi anche per i progettisti e, da un altro punto di vista, per i venditori. Questi ultimi cercano di far apparire i loro calcolatori come la migliore scelta possibile, anche se non sempre soddisfano le esigenze di chi li acquista; in alcuni casi le caratteristiche dei calcolatori che il venditore mette in risalto non hanno alcun legame con l'uso che ne deve fare l'acquirente.

Tipo di aeroplano	Capacità di trasporto (numero di passeggeri)	Autonomia di volo (km)	Velocità di crociera (km/h)	Portata (passeggeri × km/h)
Boeing 777	375	7400	980	367 500
Boeing 747	470	6640	980	460 600
BAC/Sud Concorde	132	6400	2160	285 120
Douglas DC-8-50	146	13 950	870	127 020

**Figura 1.14 Capacità, autonomia di volo e velocità di alcuni aeroplani commerciali.** L'ultima colonna mostra l'efficienza con la quale l'aeroplano trasporta i passeggeri, cioè la capacità moltiplicata per la velocità di crociera (non vengono considerati i tempi di decollo e atterraggio).

Perciò, quando si sceglie un calcolatore, è importante capire come si misurano le prestazioni e quali siano i limiti intrinseci di queste misure.

Questo paragrafo descrive differenti modi di misurare le prestazioni; vedremo, quindi, le diverse metriche utilizzate per valutare le prestazioni dal punto di vista sia dell'utente del calcolatore sia del progettista. Analizzeremo anche come queste metriche siano correlate e presenteremo le equazioni classiche utilizzate per misurare le prestazioni, che useremo in tutto il libro.

## Definizione delle prestazioni

Che cosa si intende quando si dice che un calcolatore ha prestazioni superiori a un altro? Benché la risposta possa sembrare semplice, l'analogia con il trasporto aereo illustrerà quanto questa domanda sia sottile. In Figura 1.14 sono riportati alcuni modelli di aeroplani per il trasporto di passeggeri, unitamente alla loro velocità di crociera, autonomia di volo e capacità di trasporto. Se vogliamo identificare quale di questi aeroplani abbia le prestazioni migliori, occorre anzitutto definire che cosa si intende per prestazioni. Infatti, a seconda della metrica utilizzata, otteniamo risultati diversi: scopriamo che l'aereo con la velocità di crociera più elevata è il Concorde (dismesso nel 2003), quello con maggiore autonomia di volo è il DC-8 e quello in grado di trasportare il maggior numero di passeggeri è il 747.

Supponiamo di definire le prestazioni in termini di velocità: questa metrica lascia spazio a due possibili definizioni. Si può, per esempio, definire come aereo più veloce quello che ha la velocità di crociera più elevata e che è in grado di trasportare un singolo passeggero da un luogo a un altro nel minor tempo possibile; se però l'obiettivo fosse trasportare 450 passeggeri, allora chiaramente il 747 sarebbe l'aereo più veloce, come mostrato dall'ultima colonna della tabella. Analogamente è possibile definire le prestazioni di un calcolatore in molti modi diversi. Se facessimo eseguire un programma su due calcolatori desktop, potremmo affermare che il più veloce è quello che termina prima il programma. Se dirigessimo un centro di calcolo che ha alcuni server che eseguono attività richieste da molti utenti, potremmo invece considerare come calcolatore più veloce quello che in una giornata porta a termine il maggior numero di attività. L'utente singolo di un calcolatore è interessato a ridurre il **tempo di esecuzione**, cioè il tempo fra l'inizio e il completamento di un task, detto anche **tempo di risposta**. I gestori dei centri di calcolo invece sono spesso interessati ad aumentare il **throughput** o la **larghezza di banda** (*bandwidth*), cioè il numero di task eseguiti nell'unità di tempo. Perciò, nella maggior parte dei casi, avremo bisogno di metriche e di applicazioni diverse per valutare le prestazioni dei PMD, più orientati a minimizzare il tempo di esecuzione, e dei server, più orientati a massimizzare il throughput.

**Tempo di risposta:** detto anche **tempo di esecuzione**, è il tempo totale richiesto da un calcolatore per completare un task; esso comprende gli accessi a disco, gli accessi a memoria, le attività di I/O, il tempo richiesto dal sistema operativo, il tempo d'esecuzione della CPU ecc.

**Throughput:** detto anche **bandwidth** (larghezza di banda), è un'altra misura delle prestazioni e rappresenta il numero di programmi completati per unità di tempo.

## Throughput e tempo di esecuzione

### ESEMPIO

Le seguenti modifiche apportate a un calcolatore aumentano il throughput, diminuiscono il tempo di esecuzione o entrambi?

1. Sostituire il processore del calcolatore con una versione più veloce.
2. Aggiungere processori addizionali a un sistema che utilizza più processori per task separati, per esempio per la ricerca sul web.

### SOLUZIONE

Diminuire il tempo di esecuzione di un sistema significa quasi sempre aumentare il throughput. Quindi nel caso 1 migliorano sia il tempo di esecuzione sia il throughput. Nel caso 2 nessun task viene eseguito più velocemente e di conseguenza aumenta soltanto il throughput.

Tuttavia, nel secondo caso, se le richieste di utilizzo del calcolatore fossero vicine al suo throughput, il sistema potrebbe doverle accodare. In questo caso, migliorare il throughput potrebbe migliorare anche il tempo di esecuzione, poiché ridurrebbe il tempo di attesa dei task in coda. Quindi, in molti calcolatori reali, modificare il tempo di esecuzione spesso influenza il throughput e viceversa.

Esaminando le prestazioni dei calcolatori, in questi primi capitoli ci occuperemo soprattutto del tempo di esecuzione. Per massimizzare le prestazioni vogliamo quindi minimizzare il tempo di esecuzione, o tempo di risposta, richiesto da un dato task. Possiamo quindi mettere in relazione le prestazioni con il tempo di esecuzione; per un generico calcolatore X varrà la relazione:

$$\text{Prestazioni}_X = \frac{1}{\text{Tempo di esecuzione}_X}$$

Questo significa che se le prestazioni del calcolatore X sono maggiori di quelle del calcolatore Y, avremo:

$$\begin{aligned} \text{Prestazioni}_X &> \text{Prestazioni}_Y \\ \frac{1}{\text{Tempo di esecuzione}_X} &> \frac{1}{\text{Tempo di esecuzione}_Y} \\ \text{Tempo di esecuzione}_Y &> \text{Tempo di esecuzione}_X \end{aligned}$$

Cioè, se X è più veloce di Y, il tempo di esecuzione su Y è maggiore di quello su X. Per valutare il progetto di un calcolatore, si usa spesso confrontare quantitativamente le prestazioni di due diversi calcolatori. Useremo la frase "X è  $n$  volte più veloce di Y" – o, equivalentemente, "X è veloce quanto  $n$  volte Y" – per indicare che:

$$\frac{\text{Prestazioni}_X}{\text{Prestazioni}_Y} = n$$

Se X è  $n$  volte più veloce di Y, allora il tempo di esecuzione su Y è  $n$  volte maggiore di quello su X:

$$\frac{\text{Prestazioni}_X}{\text{Prestazioni}_Y} = \frac{\text{Tempo di esecuzione}_Y}{\text{Tempo di esecuzione}_X} = n$$

## Prestazioni relative

Se il calcolatore A esegue un programma in 10 secondi e il calcolatore B esegue lo stesso programma in 15 secondi, quanto è più veloce A rispetto a B?

Sappiamo che A è  $n$  volte più veloce di B se:

$$\frac{\text{Prestazioni}_A}{\text{Prestazioni}_B} = \frac{\text{Tempo di esecuzione}_B}{\text{Tempo di esecuzione}_A} = n$$

Quindi il rapporto tra le prestazioni dei due calcolatori sarà:

$$\frac{15}{10} = 1,5$$

cioè A è 1,5 volte più veloce di B.

### ESEMPIO

### SOLUZIONE

Nell'esempio appena riportato, possiamo dire che il calcolatore B è 1,5 volte più lento del calcolatore A, poiché:

$$\frac{\text{Prestazioni}_A}{\text{Prestazioni}_B} = 1,5$$

significa che:

$$\frac{\text{Prestazioni}_A}{1,5} = \text{Prestazioni}_B$$

Per semplicità cercheremo di utilizzare sempre il termine *più veloce* nei confronti quantitativi fra calcolatori. Poiché le prestazioni e il tempo di esecuzione sono uno il reciproco dell'altro, aumentare le prestazioni richiede la diminuzione del tempo di esecuzione. Per evitare di fare confusione con i termini *aumentare* e *diminuire*, verranno utilizzate le frasi “migliorare le prestazioni” o “migliorare il tempo di esecuzione”, intendendo rispettivamente “aumentare le prestazioni” e “diminuire il tempo di esecuzione”.

Applicazioni differenti sono sensibili a differenti aspetti delle prestazioni di un calcolatore. Molte applicazioni, soprattutto quelle eseguite dai server, dipendono in larga misura anche dalle prestazioni del sistema di I/O, che, a loro volta, dipendono sia dall'hardware sia dal software. Il tempo assoluto è la misura che ci interessa. In alcuni ambienti applicativi, l'utente può essere più interessato al throughput, al tempo di esecuzione, o a una combinazione di entrambi (per es., il throughput massimo combinato con il tempo di esecuzione peggiore). Per migliorare le prestazioni di un programma occorre innanzitutto identificare in modo chiaro quale metrica sia la più adatta. Solo successivamente è lecito iniziare a cercare i colli di bottiglia nell'esecuzione dei programmi; allo scopo, durante l'esecuzione del programma vengono prese le misure con la metrica identificata e queste vengono poi analizzate per cercare i punti responsabili del calo di prestazioni. Nei capitoli seguenti descriveremo come ricercare i colli di bottiglia e migliorare le prestazioni delle varie parti del sistema.

## Capire le prestazioni dei programmi

## Misurare le prestazioni

La grandezza utilizzata per misurare le prestazioni di un calcolatore è il tempo: il calcolatore che esegue un certo lavoro nel tempo minore è il più veloce; per esempio, il *tempo di esecuzione* di un programma è misurato in secondi. Tuttavia il tempo può essere definito in diversi modi, a seconda di ciò che interessa misurare: la definizione più immediata è il *tempo assoluto*, detto anche *tempo di esecuzione* o *tempo di risposta* o anche *tempo relativo*. Esso è il tempo totale richiesto per completare un task, compreso il tempo di accesso al disco e alla memoria, di input/output (I/O), e il tempo richiesto dal sistema operativo, in sostanza il tempo di tutte le attività richieste per completare il task.

Tuttavia, i calcolatori lavorano spesso in condivisione e può accadere che un processore stia lavorando su più programmi contemporaneamente; in questo caso, il sistema può cercare di massimizzare il throughput anziché cercare di minimizzare il tempo di esecuzione di un singolo programma. Di conseguenza, si distingue spesso tra tempo assoluto di esecuzione di un programma e tempo durante il quale il processore ha effettivamente lavorato su quel programma.

Il **tempo di esecuzione della CPU**, o più semplicemente **tempo di CPU**, tiene conto di questa distinzione: è il tempo effettivamente speso dalla CPU nella computazione richiesta dal programma e non comprende il tempo speso per le operazioni di I/O o nell'esecuzione di altri programmi (ricordate comunque che il tempo di esecuzione percepito dall'utente sarà il tempo assoluto e non il tempo di CPU). Il tempo di CPU può essere ulteriormente suddiviso nel tempo speso dalla CPU per l'effettiva esecuzione del programma, denominato **tempo di CPU utente**, e nel tempo speso dalla CPU per eseguire le funzioni del sistema operativo richieste per l'esecuzione del programma, denominato **tempo di CPU di sistema**. A volte può essere difficile distinguere con precisione il tempo di CPU relativo all'utente e quello relativo al sistema operativo, dal momento che è spesso arduo assegnare una certa funzione del sistema operativo a un certo programma utente piuttosto che a un altro, e che esistono differenze di funzionamento tra i vari sistemi operativi.

Per mantenere una coerenza interna, manterremo la distinzione tra le prestazioni misurate con il tempo assoluto e quelle misurate mediante il tempo di esecuzione della CPU. Inoltre, utilizzeremo il termine *prestazioni di sistema* per fare riferimento al tempo assoluto misurato su un sistema scarico e il termine *prestazioni della CPU* per indicare il tempo di CPU utente. In questo capitolo ci concentreremo sulle prestazioni della CPU, benché gli argomenti relativi alla valutazione delle prestazioni possano essere applicati sia al tempo assoluto sia al tempo di CPU.

Benché agli utenti dei calcolatori interessi principalmente il parametro tempo, quando si esamina un calcolatore nei dettagli è meglio ragionare in termini di prestazioni utilizzando altre metriche; in particolare, i progettisti di calcolatori preferiscono utilizzare misure che esprimano la velocità con cui l'hardware è in grado di eseguire certe operazioni elementari.

Quasi tutti i calcolatori sono costruiti utilizzando un segnale che sincronizza le varie funzioni implementate nell'hardware; questo segnale è periodico nel tempo e i relativi intervalli di tempo sono denominati **cicli di clock** (o **colpi, colpi di clock, periodi di clock, clock, cicli**). I progettisti utilizzano il termine **periodo di clock** per indicare il tempo necessario per completare un intero ciclo di clock (per es. 250 picosecondi o 250 ps), o il termine *frequenza di clock* (per es. 4 Gigahertz o 4 GHz), che è l'inverso del periodo di clock. Nel prossimo paragrafo formalizzeremo la relazione tra i cicli di clock misurati dal progettista hardware e il tempo assoluto, misurato in secondi, dell'utente di un calcolatore.

**Tempo di esecuzione della CPU:** detto anche **tempo di CPU**, è il tempo che la CPU utilizza effettivamente nella computazione richiesta da un certo task.

**Tempo di CPU utente:** tempo effettivamente speso dalla CPU nella computazione richiesta da un programma.

**Tempo di CPU di sistema:** tempo speso dalla CPU per eseguire le funzioni del sistema operativo richieste per l'esecuzione di un programma.

**Ciclo di clock:** detto anche **colpo, colpo di clock, periodo di clock, ciclo**, è la durata di un periodo del clock; solitamente viene misurato il periodo del clock del processore. Il clock è un dispositivo che genera un segnale a frequenza costante.

**Periodo di clock:** durata di un ciclo di clock.

## Autovalutazione

1. Un'applicazione che utilizza il cloud e un PMD è limitata dalle prestazioni della rete. Stabilire se le seguenti modifiche migliorino il solo throughput, il throughput e il tempo di esecuzione, oppure nessuno dei due.
  - a. Un canale di rete supplementare viene inserito tra il PMD e il cloud, aumentando così il throughput totale della rete e riducendo i ritardi di accesso alla rete (dato che sono ora disponibili due canali di rete).
  - b. Il software di rete viene migliorato, riducendo così i ritardi di comunicazione sulla rete; non viene invece migliorato il throughput.
  - c. Viene aggiunta altra memoria al calcolatore.
2. Le prestazioni del calcolatore C sono 4 volte migliori di quelle del calcolatore B, che esegue una certa applicazione in 28 secondi. Quanto impiega il calcolatore C per eseguire quell'applicazione?

## Prestazioni della CPU

Gli utenti e i progettisti utilizzano di solito metriche di tipo diverso per valutare le prestazioni. Se si riuscissero a mettere in relazione le diverse metriche, sarebbe possibile determinare l'effetto di una modifica dell'architettura sulle prestazioni che vengono percepite dai diversi utenti. Dal momento che ora ci stiamo occupando delle prestazioni della CPU, possiamo affermare che la misura fondamentale è il tempo di esecuzione della CPU (o tempo di CPU). Le metriche più elementari (numero di cicli di clock e periodo del clock) sono legate al tempo di CPU tramite una semplice equazione:

$$\frac{\text{Tempo di CPU relativo}}{\text{a un programma}} = \frac{\text{Cicli di clock della CPU}}{\text{relativi al programma}} \times \text{Periodo di clock}$$

Poiché la frequenza di clock è il reciproco del suo periodo, è possibile utilizzare in alternativa la seguente equazione:

$$\frac{\text{Tempo di CPU relativo}}{\text{a un programma}} = \frac{\text{Cicli di clock della CPU relativi a un programma}}{\text{Frequenza di clock}}$$

Questa ultima equazione chiarisce che un progettista hardware può migliorare le prestazioni riducendo il numero di cicli di clock necessari per eseguire un programma oppure la durata del ciclo di clock stesso. Come vedremo nei prossimi capitoli, i progettisti spesso devono trovare un compromesso tra il numero di cicli di clock e la durata del singolo ciclo. Molte delle tecniche che riducono il numero di cicli di clock, infatti, tendono ad aumentare la durata del periodo del clock.

## Migliorare le prestazioni

Il nostro programma preferito viene eseguito in 10 secondi dal calcolatore A, che è dotato di un clock a 2 GHz. Stiamo cercando di aiutare un progettista a costruire un calcolatore B in grado di eseguire lo stesso programma in 6 secondi. Il progettista ha concluso che è possibile aumentare in modo significativo la frequenza di clock; questa modifica

### ESEMPIO

(continua)

(continua)

avrà però un'influenza su tutto il progetto della CPU, facendo sì che il calcolatore B richieda un numero di cicli di clock maggiore di un fattore 1,2 rispetto al calcolatore A per eseguire il programma. Dovendo dare un consiglio al progettista, quale sarà la frequenza di clock che permette di raggiungere l'obiettivo?

**SOLUZIONE**

Innanzitutto determiniamo il numero di cicli di clock necessari all'esecuzione del programma su A:

$$\text{Tempo di CPU}_A = \frac{\text{Cicli di clock CPU}_A}{\text{Frequenza di clock}_A}$$

$$10 \text{ secondi} = \frac{\text{Cicli di clock CPU}_A}{2 \times 10^9 \frac{\text{cicli}}{\text{secondo}}}$$

$$\text{Cicli di clock CPU}_A = 10 \text{ secondi} \times 2 \times 10^9 \frac{\text{cicli}}{\text{secondo}} = 20 \times 10^9 \text{ cicli}$$

Il tempo di CPU per B si può calcolare tramite la seguente equazione:

$$\text{Tempo di CPU}_B = \frac{1,2 \times \text{Cicli di clock CPU}_A}{\text{Frequenza di clock}_B}$$

$$6 \text{ secondi} = \frac{1,2 \times 20 \times 10^9 \text{ cicli}}{\text{Frequenza di clock}_B}$$

$$\text{Frequenza di clock}_B = \frac{1,2 \times 20 \times 10^9 \text{ cicli}}{6 \text{ secondi}} = \frac{0,2 \times 20 \times 10^9 \text{ cicli}}{\text{secondi}} = \frac{4 \times 10^9 \text{ cicli}}{\text{secondi}} = 4 \text{ GHz}$$

Quindi il calcolatore B dovrà lavorare a una frequenza di clock doppia rispetto ad A per eseguire il programma in 6 secondi.

### Misura delle prestazioni associate alle istruzioni

Le equazioni precedenti non contenevano alcun riferimento al numero di istruzioni presenti nel programma. Tuttavia, dato che il compilatore genera delle istruzioni da eseguire e il calcolatore esegue le singole istruzioni per eseguire il programma, il tempo di esecuzione dovrà necessariamente dipendere dal numero delle istruzioni che compongono il programma. È possibile esprimere il tempo di esecuzione totale come il risultato del prodotto del numero di istruzioni da eseguire per il tempo medio di esecuzione di ciascuna istruzione. Di conseguenza, il numero di cicli di clock necessari per l'esecuzione di un programma si può scrivere come:

$$\text{Cicli di clock della CPU} = \frac{\text{Numero di istruzioni}}{\text{del programma}} \times \frac{\text{Numero medio di cicli}}{\text{di clock per istruzione}}$$

**Cicli di clock per istruzione (CPI):** numero medio di cicli di clock per istruzione di un programma o di un frammento di programma.

Il termine **cicli di clock per istruzione**, calcolato come media del numero di cicli di clock che le diverse istruzioni richiedono per essere completate, è spesso abbreviato con la sigla **CPI**. Dato che istruzioni diverse possono richiedere un tempo di esecuzione differente in funzione del compito che svolgono, il CPI assume il valore medio calcolato su tutte le istruzioni dal programma. Il CPI è una quantità utile a confrontare due calcolatori diversi che condividono la stessa architettura dell'insieme di istruzioni, dal momento che il numero di istruzioni contenute in un programma sarà, in questo caso, lo stesso.

## Come utilizzare l'equazione delle prestazioni

Siano date due implementazioni diverse della stessa architettura dell'insieme di istruzioni. Il calcolatore A ha un ciclo di clock di 250 ps e un CPI pari a 2,0 per un determinato programma; il calcolatore B, invece, ha un ciclo di clock pari a 500 ps e un CPI pari a 1,2 misurato sullo stesso programma. Quale dei due calcolatori è più veloce nell'esecuzione del programma, e di quanto?

Sappiamo che i due calcolatori eseguono lo stesso numero di istruzioni elementari quando eseguono il programma; chiamiamo  $I$  il numero di istruzioni. Per prima cosa occorre calcolare il numero di cicli di clock richiesti per l'esecuzione da parte di ciascun calcolatore:

$$\text{Cicli di clock della CPU}_A = I \times 2,0$$

$$\text{Cicli di clock della CPU}_B = I \times 1,2$$

Si può ora calcolare il tempo di CPU per ciascun calcolatore:

$$\begin{aligned}\text{Tempo di CPU}_A &= \text{Cicli di clock della CPU}_A \times \text{Periodo di clock} \\ &= I \times 2,0 \times 250 \text{ ps} = 500 \times I \text{ ps}\end{aligned}$$

Analogamente, per B si ha:

$$\text{Tempo di CPU}_B = I \times 1,2 \times 500 \text{ ps} = 600 \times I \text{ ps}$$

È chiaro che il calcolatore A è più veloce di B, di un fattore dato dal rapporto tra i tempi di esecuzione:

$$\frac{\text{Prestazioni della CPU}_A}{\text{Prestazioni della CPU}_B} = \frac{\text{Tempo di esecuzione}_B}{\text{Tempo di esecuzione}_A} = \frac{600 \times I \text{ ps}}{500 \times I \text{ ps}} = 1,2$$

Si può dunque concludere che il calcolatore A è 1,2 volte più veloce di B nell'esecuzione di questo programma.

### ESEMPIO

### SOLUZIONE

## Equazione classica di misura delle prestazioni

Possiamo quindi scrivere questa equazione fondamentale in funzione del **numero di istruzioni** (eseguite da un programma), del CPI e del periodo di clock:

$$\text{Tempo di CPU} = \text{Numero di istruzioni} \times \text{CPI} \times \text{Periodo di clock}$$

Oppure, poiché la frequenza di clock è l'inverso del periodo di clock:

$$\text{Tempo di CPU} = \frac{\text{Numero di istruzioni} \times \text{CPI}}{\text{Frequenza di clock}}$$

Le formule precedenti sono particolarmente utili in quanto evidenziano i tre fattori principali che influenzano le prestazioni. Utilizzeremo queste formule per confrontare tra loro due architetture o per valutare un'alternativa progettuale conoscendone l'impatto sui tre parametri.

**Numero di istruzioni:** il numero delle istruzioni eseguite da un programma.

## ESEMPIO

## Confronto di frammenti di codice

Un progettista di compilatori deve decidere quale tra due sequenze di codice implementare per un certo calcolatore. I progettisti dell'hardware gli hanno fornito queste informazioni:

	CPI per ciascun tipo di istruzione		
	A	B	C
CPI	1	2	3

Per una certa istruzione in linguaggio ad alto livello, il progettista sta considerando due sequenze di codice in linguaggio macchina che richiedono un numero diverso di istruzioni dei tre tipi:

Sequenza di istruzioni	Numero di istruzioni in linguaggio macchina per ciascun tipo di istruzione		
	A	B	C
Sequenza 1	2	1	2
Sequenza 2	4	1	1

Quale sequenza richiede l'esecuzione di un maggior numero di istruzioni? Quale verrà eseguita più velocemente? Qual è il CPI delle due sequenze?

## SOLUZIONE

Nella sequenza 1 vengono eseguite:  $2 + 1 + 2 = 5$  istruzioni, mentre la sequenza 2 richiede  $4 + 1 + 1 = 6$  istruzioni; perciò la sequenza 1 richiede un numero minore di istruzioni.

Per calcolare il numero totale di cicli di clock per ciascuna sequenza, possiamo utilizzare l'equazione che misura il numero di cicli di clock della CPU in funzione del numero di istruzioni e del CPI:

$$\text{Cicli di clock della CPU} = \sum_{i=1}^n (\text{CPI}_i \times C_i)$$

Da questa equazione possiamo ricavare il numero totale dei cicli di clock:

$$\begin{aligned}\text{Cicli di clock della CPU}_1 &= (2 \times 1) + (1 \times 2) + (2 \times 3) = 2 + 2 + 6 = 10 \text{ cicli} \\ \text{Cicli di clock della CPU}_2 &= (4 \times 1) + (1 \times 2) + (1 \times 3) = 4 + 2 + 3 = 9 \text{ cicli}\end{aligned}$$

Pertanto possiamo concludere che la sequenza 2 è più veloce, anche se richiede l'esecuzione di un'istruzione in più. Dato che tale sequenza richiede meno cicli di clock ma ha un'istruzione in più, deve avere un CPI più basso. Si può calcolare il valore del CPI come:

$$\text{CPI} = \frac{\text{Cicli di clock della CPU}}{\text{Numero di istruzioni}}$$

$$\text{CPI}_1 = \frac{\text{Cicli di clock della CPU}_1}{\text{Numero di istruzioni}_1} = \frac{10}{5} = 2,0$$

$$\text{CPI}_2 = \frac{\text{Cicli di clock della CPU}_2}{\text{Numero di istruzioni}_2} = \frac{9}{6} = 1,5$$

Componente delle prestazioni	Unità di misura
Tempo di esecuzione della CPU per un dato programma	Secondi per programma
Numero di istruzioni	Istruzioni eseguite per singolo programma
Cicli di clock per istruzione (CPI)	Numero medio di cicli di clock per istruzione
Durata del ciclo di clock	Secondi per ciclo di clock

Figura 1.15 I componenti di base delle prestazioni e come vengono misurati.

## QUADRO D'INSIEME

La Figura 1.15 riporta le misure di base che vengono utilizzate per i calcolatori a livelli diversi e che cosa viene misurato in ciascun caso. Potete notare come i diversi fattori si possano combinare tra loro per formare il tempo di esecuzione di un programma, che viene misurato in secondi:

$$\text{Tempo} = \frac{\text{Istruzioni}}{\text{Programma}} \times \frac{\text{Cicli di clock}}{\text{Istruzioni}} \times \frac{\text{Secondi}}{\text{Ciclo di clock}}$$

Occorre tenere sempre presente che il tempo è l'unica misura completa e affidabile per valutare le prestazioni di un calcolatore. Per esempio, modificare l'insieme delle istruzioni per diminuire il numero di istruzioni può portare a un'architettura con periodo di clock più lungo o a un CPI più elevato, che vanifica la riduzione del numero di istruzioni. Analogamente, poiché il CPI dipende dal tipo di istruzione eseguita, il codice che esegue il minor numero di istruzioni potrebbe non essere il più veloce. ■

Come si possono determinare i valori dei fattori che compaiono nell'equazione per il calcolo delle prestazioni? Il tempo di esecuzione della CPU può essere misurato facendo eseguire il programma, mentre la durata del ciclo di clock viene normalmente riportata nella documentazione tecnica del calcolatore; il numero di istruzioni e il CPI sono invece più difficili da ottenere. Se la frequenza di clock e il tempo di esecuzione della CPU fossero noti, sarebbe sufficiente conoscere il CPI o il numero di istruzioni per conoscere anche l'altro termine.

Il numero di istruzioni può essere determinato attraverso opportuni strumenti software di profiling che osservano l'esecuzione del programma oppure attraverso un simulatore dell'architettura. Alternativamente si potrebbero utilizzare i contatori hardware, contenuti nella maggior parte dei processori, per misurare il numero di istruzioni eseguite, il CPI medio e, spesso, l'origine stessa della perdita di prestazioni. Essendo il numero di istruzioni dipendente dall'architettura ma indipendente dalla particolare implementazione, possiamo misurarlo anche senza conoscere nei dettagli l'implementazione delle istruzioni. Il CPI, viceversa, dipende da un ampio spettro di dettagli progettuali del calcolatore, tra cui la modalità di realizzazione del sottosistema di memoria e la struttura del processore (come vedremo nei Capitoli 4 e 5); inoltre, dipende anche dalla combinazione dei tipi di istruzioni eseguite dall'applicazione specifica. Il CPI dunque varia da applicazione ad applicazione, come pure tra diverse implementazioni dello stesso insieme di istruzioni.

**Composizione delle istruzioni (instruction mix):** una misura della frequenza dinamica delle istruzioni in uno o più programmi.

L'esempio precedente evidenzia il pericolo che si corre nell'utilizzare uno solo dei fattori, per esempio il numero di istruzioni, quando si valutano le prestazioni. Nel confrontare due calcolatori è necessario tenere conto di tutte e tre le componenti, che, combinate tra loro, forniscono il tempo di esecuzione. Se alcuni di questi fattori sono identici, come lo era la frequenza di clock nell'esempio precedente, le prestazioni dei due calcolatori si possono determinare confrontando solamente gli altri fattori. Dato che il CPI dipende dalla **composizione delle istruzioni (instruction mix)**, occorre considerare sia il numero di istruzioni sia il CPI, anche se il clock ha la stessa frequenza. Alcuni esercizi riportati alla fine di questo capitolo vi permetteranno di capire meglio come una serie di miglioramenti del calcolatore e del compilatore influisca sulla frequenza di clock, sul CPI e sul numero di istruzioni. Nel paragrafo 1.10 , riportiamo una misura delle prestazioni che, pur essendo di uso comune, non incorpora tutti e tre questi fattori e può quindi essere fuorviante.

## Capire le prestazioni dei programmi

Le prestazioni di un programma dipendono dall'algoritmo, dal linguaggio, dal compilatore, dall'architettura e dall'hardware del calcolatore. La seguente tabella riassume come questi componenti influenzino i tre fattori dell'equazione delle prestazioni.

Componente hardware o software	Che cosa influenza?	Come?
Algoritmo	Numero di istruzioni, eventualmente il CPI	L'algoritmo determina il numero di istruzioni del programma sorgente e quindi il numero di istruzioni in linguaggio macchina che vengono eseguite dal processore. L'algoritmo può anche influenzare il CPI, favorendo l'utilizzo di istruzioni più o meno veloci. Per esempio, se l'algoritmo utilizza più operazioni di divisione, tenderà ad avere un CPI più elevato
Linguaggio di programmazione	Numero di istruzioni, CPI	Il linguaggio di programmazione influenza certamente il numero di istruzioni, dal momento che i costrutti del linguaggio ad alto livello sono tradotti in istruzioni in linguaggio macchina e queste determinano il numero di istruzioni eseguite. Il linguaggio ad alto livello può anche influenzare il CPI a seconda delle sue caratteristiche; per esempio, un linguaggio con un esteso supporto per i dati astratti (per es. Java) richiederà chiamate indirette a funzione, che sono caratterizzate da un CPI più alto
Compilatore	Numero di istruzioni, CPI	L'efficienza del compilatore influenza sia il numero di istruzioni sia il numero medio di cicli per istruzione, dal momento che il compilatore traduce le istruzioni dal linguaggio sorgente ad alto livello nelle istruzioni in linguaggio macchina. Il ruolo del compilatore può essere molto complesso e può influenzare il CPI in maniera complessa
Architettura dell'insieme delle istruzioni	Numero di istruzioni, frequenza di clock, CPI	L'architettura dell'insieme di istruzioni influenza tutti e tre i fattori delle prestazioni della CPU, dato che influenza le istruzioni richieste da una data funzione, il costo in numero di cicli di ogni istruzione e la frequenza del processore

**Approfondimento.** Potreste aspettarvi che il minimo CPI sia 1,0; invece, come vedremo nel Capitolo 4, alcuni processori leggono ed eseguono più istruzioni in un unico ciclo di clock. A seguito di ciò, alcuni progettisti invertono il CPI e misurano l'*IPC*, o *numero di istruzioni per ciclo di clock*. Se un processore esegue in media 2 istruzioni per ciclo di clock, esso ha un IPC di 2 e quindi un CPI di 0,5.

**Approfondimento.** Anche se la frequenza di clock è tradizionalmente fissa, per risparmiare energia o per migliorare temporaneamente le prestazioni, i moderni processori possono variare la loro frequenza di clock. Dobbiamo perciò considerare la frequenza di clock media per valutare un programma. Per esempio, il Core i7 di Intel può aumentare la sua frequenza di clock fino al 10% fino a quando il chip non inizia a scaldare troppo. Questa modalità di funzionamento viene chiamata *Turbo Mode* da Intel.

### Autovalutazione

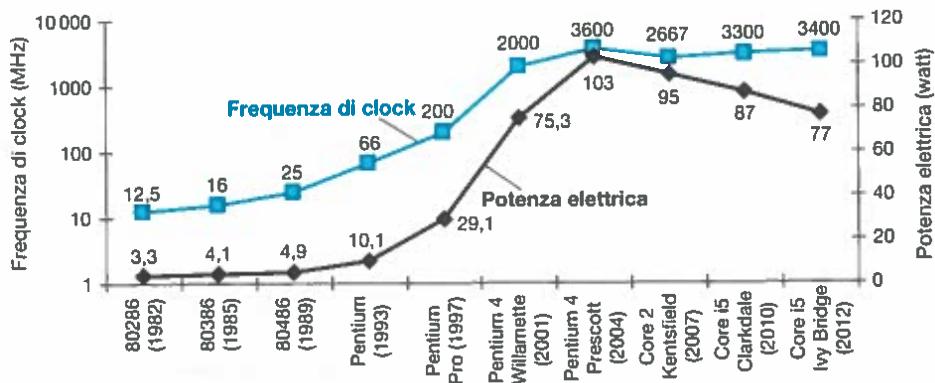
La durata dell'esecuzione di una data applicazione scritta in Java è di 15 secondi sul processore di un calcolatore desktop. Un nuovo compilatore Java richiede solo il 60% delle istruzioni del vecchio compilatore. Sfortunatamente il CPI aumenta di 1,1. Quanto più velocemente verrà eseguita l'applicazione utilizzando il nuovo compilatore?

- a.  $\frac{15 \times 0,6}{1,1} = 8,2 \text{ s}$
- b.  $15 \times 0,6 \times 1,1 = 9,9 \text{ s}$
- c.  $\frac{1,5 \times 1,1}{0,6} = 27,5 \text{ s}$

## 1.7 Barriera dell'energia

La Figura 1.16 mostra l'aumento della frequenza di clock e della potenza elettrica richiesta dalle otto generazioni di microprocessori Intel prodotte negli ultimi 30 anni. Entrambe sono cresciute rapidamente per anni, ma recentemente hanno smesso di aumentare. Il motivo per cui sono cresciute in maniera proporzionale è che sono correlate, e il motivo per cui hanno smesso di aumentare è che è stata raggiunta la massima potenza dissipabile dai sistemi di raffreddamento montati nei microprocessori.

Anche se la potenza elettrica fornisce un limite a quanto possiamo raffreddare, nell'era post-PC la risorsa veramente preziosa è l'energia. La durata della batteria può risultare la carta vincente di un dispositivo mobile e i progettisti



**Figura 1.16** Frequenza di clock e potenza assorbita dalle otto generazioni di processori Intel x86 prodotti negli ultimi 25 anni. Il Pentium 4 ha segnato un incremento marcato nella frequenza di clock e nella quantità di potenza assorbita; tuttavia, il miglioramento delle prestazioni è stato meno accentuato. I problemi di dissipazione termica della versione Prescott del Pentium 4 hanno portato all'abbandono dello sviluppo ulteriore del Pentium 4 a favore di un approccio multicore: con il microprocessore Core 2 si è tornati a una pipeline più semplice, con una frequenza di clock inferiore, e sono stati montati più processori, detti core, sullo stesso chip. La pipeline del Core i5 segue gli stessi principi.

dei centri di elaborazione dati cercano di ridurre i costi dell'alimentazione e del raffreddamento di centinaia di migliaia di server, dato che i costi per questi numeri sono alti. Così come misurare il tempo in secondi è una misura delle prestazioni di un programma più prudente dei MIPS (par. 1.10), il consumo di energia misurato in joule è migliore dell'assorbimento istantaneo di potenza misurato in watt cioè in joule/secondo.

La tecnologia dominante per la fabbricazione dei circuiti integrati è la tecnologia CMOS (*Complementary Metal Oxid Semiconductors* – semiconduttori a ossido metallico complementari). I CMOS assorbono energia elettrica soprattutto durante la fase di commutazione; questa energia viene chiamata *energia dinamica*, cioè energia assorbita durante la commutazione. L'energia dinamica dipende dal carico capacitivo dei transistor e dalla tensione applicata, secondo la formula:

$$\text{Energia} \propto \text{Carico capacitivo} \times \text{Tensione}^2$$

Questa equazione misura l'energia di un impulso durante la commutazione  $0 \rightarrow 1 \rightarrow 0$  o  $1 \rightarrow 0 \rightarrow 1$ . L'energia richiesta da una singola transizione è perciò:

$$\text{Energia} \propto 1/2 \times \text{Carico capacitivo} \times \text{Tensione}^2$$

La potenza richiesta da ciascun transistor è quindi il prodotto dell'energia richiesta per una transizione per la frequenza delle transizioni:

$$\text{Potenza} \propto 1/2 \times \text{Carico capacitivo} \times \text{Tensione}^2 \times \text{Frequenza di commutazione}$$

La frequenza di commutazione è funzione della frequenza di clock. Il carico capacitivo dei transistor è funzione sia del numero di transistor connessi allo stesso output, detto *fanout*, sia della tecnologia con cui sono costruiti, che determina a sua volta anche la capacità delle connessioni e dei transistor stessi.

Guardando la Figura 1.16 si può notare che la frequenza di clock è aumentata di 1000 volte mentre l'assorbimento di potenza è aumentato solamente di un fattore 30. Come è stato possibile? L'assorbimento si può ridurre diminuendo la tensione di alimentazione, cosa che è puntualmente avvenuta ad ogni nuova generazione della tecnologia, e la potenza è una funzione della tensione al

### Potenza relativa

#### ESEMPIO

Supponiamo di volere sviluppare un nuovo processore, più semplice, con l'85% del carico capacitivo dei vecchi processori più complessi. Inoltre, ipotizziamo che la tensione sia regolabile e che possa essere ridotta del 15% rispetto a un vecchio processore B, il che consentirebbe una riduzione del 15% della frequenza di clock. Quale sarà l'impatto sulla potenza assorbita?

#### SOLUZIONE

$$\frac{\text{Potenza}_{\text{nuova}}}{\text{Potenza}_{\text{vecchia}}} = \frac{(\text{Carico capacitivo} \times 0,85) \times (\text{Tensione} \times 0,85)^2 \times (\text{Frequenza commutazione} \times 0,85)}{\text{Carico capacitivo} \times \text{Tensione}^2 \times \text{Frequenza commutazione}}$$

da cui segue che il rapporto tra le potenze assorbite dai due calcolatori è:

$$0,85^4 = 0,52$$

Quindi, il nuovo processore richiederebbe circa metà della potenza del vecchio processore.

quadrato. Perciò, dato che in media la tensione di alimentazione è stata ridotta del 15% in ogni nuova generazione, in 20 anni la tensione è scesa da 5 V a 1 V; questo spiega come mai la potenza sia cresciuta solo di 30 volte.

Il problema ad oggi è che diminuendo ulteriormente la tensione di alimentazione i transistor tenderebbero a disperdere troppa corrente, come un rubinetto che non può mai essere completamente chiuso. Già il 40% della potenza assorbita da un transistor è dovuto alla dispersione di corrente; se i transistor avessero perdite ancora maggiori, i processori diventerebbero ingestibili.

Per cercare di risolvere il problema dell'assorbimento di potenza, i progettisti hanno provato a inserire dispositivi più grandi per la dissipazione del calore e meccanismi di controllo per spegnere le parti dei circuiti che non vengono utilizzate in un dato ciclo di clock. Molte altre tecniche potrebbero essere espilate per raffreddare i chip e aumentare la potenza che può essere assorbita, per esempio fino a 300 watt; queste tecniche, però, sono troppo costose non solo per i PC e i PMD, ma anche per i server.

L'assorbimento di potenza si è quindi rivelato una barriera invalicabile, e i progettisti hanno cercato nuove strade per migliorare i calcolatori, sviluppando un nuovo modo di progettare i microprocessori, diverso da quello utilizzato nei primi 30 anni.

**Approfondimento.** Anche se un transistor CMOS assorbe energia principalmente durante la fase di commutazione, in condizioni statiche una certa quantità di energia viene comunque dissipata a causa delle correnti di dispersione che esistono anche quando il transistor è spento; queste perdite sono responsabili del 40% dell'energia totale assorbita da un server. Perciò, aumentando il numero di transistor, aumenta l'energia assorbita anche se i transistor sono sempre spenti. Diversi accorgimenti nella progettazione dei chip e ulteriori innovazioni tecnologiche sono stati introdotti per limitare la dispersione, ma rimane difficile ridurre ulteriormente la tensione di alimentazione.

**Approfondimento.** L'alimentazione elettrica è una sfida per i circuiti integrati per due motivi. In primo luogo, l'alimentazione deve essere distribuita ai vari chip: i moderni microprocessori utilizzano centinaia di pin solamente per l'alimentazione e la terra! Analogamente, livelli multipli di interconnessioni tra i chip vengono utilizzati solamente per distribuire l'alimentazione e la terra a porzioni del chip. In secondo luogo, la potenza elettrica viene dissipata come calore e il calore deve essere rimosso. I chip dei server possono consumare più di 100 watt e raffreddare il chip e i sistemi attorno rappresenta una delle spese maggiori nei calcolatori dei centri di calcolo (vedi Cap. 6).

## 1.8 Metamorfosi delle architetture: il passaggio dai sistemi uniprocessoare ai sistemi multiprocessoore

Il limite all'assorbimento di potenza ha obbligato i progettisti a cambiare radicalmente il modo di progettare i microprocessori. La Figura 1.17 mostra il miglioramento, negli anni, del tempo di esecuzione di un programma eseguito su processore desktop: dal 2002 il miglioramento è sceso da un fattore 1,5 per anno a un fattore 1,2 per anno.

Invece di continuare a diminuire il tempo di esecuzione del singolo programma eseguito su un processore singolo, dal 2006 tutti i produttori di desktop e server hanno iniziato a distribuire microprocessori contenenti più processori sul singolo chip; in questi microprocessori il miglioramento si verifica spesso più nel throughput che nel tempo di esecuzione. Per eliminare ogni possibile confusione tra i termini processore e microprocessore, i produttori hanno iniziato a chiamare i processori *core*, e quindi questi microprocessori di nuova

*Fino a oggi, la maggior parte del software era come uno spartito scritto per un solista; con la nuova generazione di chip, incominciamo a fare esperienza con i duetti, i quartetti e con altre piccole formazioni musicali; ma scrivere uno spartito per una grande orchestra e coro è ben altra cosa.*

Brian Hayes, *Computing in Parallel Universe*, 2007.

generazione vengono chiamati *multicore* (multiprocessori): un microprocessore *quadcore* è perciò un chip che contiene al suo interno quattro core (ovvero quattro processori).

In passato, i programmati potevano confidare nelle innovazioni dell'hardware, delle architetture e dei compilatori per raddoppiare le prestazioni dei loro programmi ogni 18 mesi senza dover modificare una sola linea di codice. Oggi, invece, per ottenere significativi miglioramenti nel tempo di esecuzione dei loro programmi, devono riscrivere il codice per poter sfruttare al meglio i diversi core. Per di più, i programmati devono continuare ad aggiornare il proprio codice per migliorarne le prestazioni ogni volta che viene aumentato il numero di core, in modo tale che i loro programmi possano essere eseguiti sempre più velocemente via via che vengono introdotti nuovi microprocessori.

Per sottolineare come i sistemi software e hardware lavorino in stretta sinergia, abbiamo inserito delle sezioni speciali denominate *Interfaccia hardware/software* all'interno del testo. In queste sezioni vengono riassunti i concetti fondamentali che stanno alla base dell'interfaccia tra hardware e software. La prima riguarda il parallelismo.



## Interfaccia hardware/software



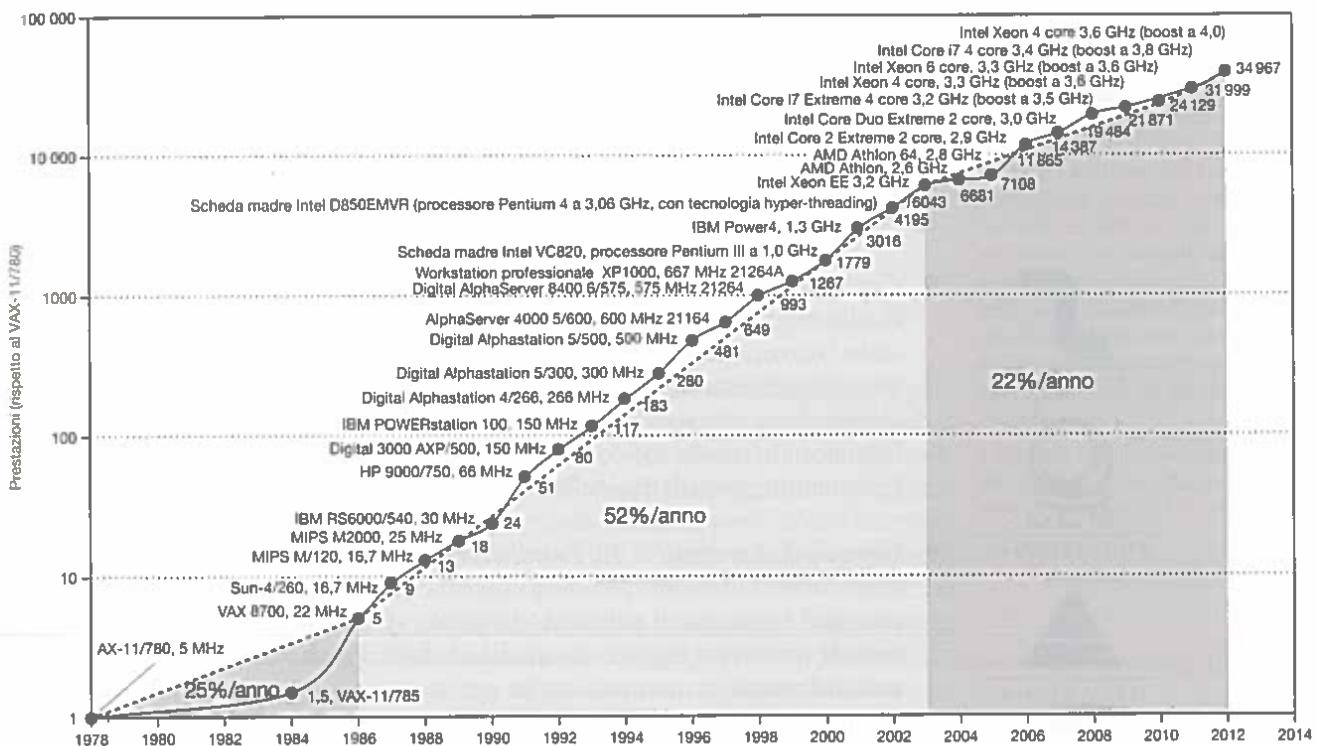
**Il parallelismo** è stato sempre un elemento critico per le prestazioni dei calcolatori, ma spesso è nascosto. Nel Capitolo 4 illustreremo la **pipeline**, una tecnica elegante per eseguire i programmi più velocemente sovrapponendo parzialmente l'esecuzione delle diverse istruzioni. La pipeline è un esempio di *parallelismo a livello di istruzioni*; in questo approccio la natura parallela dell'hardware viene nascosta in modo tale che il programmatore e il compilatore possano ragionare come se l'hardware eseguisse le istruzioni sequenzialmente.

Obbligare i programmati a tenere conto del parallelismo dell'hardware e a riscrivere i loro programmi per potere essere eseguiti in parallelo non è stata ritenuta la strada giusta nello sviluppo delle architetture: società che avevano incentrato il loro business sulla parallelizzazione del codice, infatti, non hanno avuto successo (vedi par. 6.15). Dalla prospettiva storica, è quindi sorprendente che l'intera industria dell'IT (*Information Technology* – tecnologia dell'informazione) abbia scommesso il proprio futuro sul fatto che i programmati si convertiranno alla programmazione parallela.

Perché è così difficile per un programmatore scrivere un programma esplicitamente parallelo? Il primo motivo è che la programmazione parallela è, per definizione, ad alte prestazioni, caratteristica che aumenta la difficoltà della programmazione stessa. Non solo il programma deve essere corretto, risolvere un problema importante e fornire un'interfaccia efficace alle persone che lo utilizzeranno, ma deve anche essere veloce. Altrimenti, se la velocità di esecuzione non è un parametro fondamentale, tanto vale scrivere un classico programma sequenziale.

Il secondo motivo è che veloce per l'hardware parallelo significa che il programmatore deve suddividere la sua applicazione in modo tale che ciascun core abbia all'incirca lo stesso lavoro da compiere nello stesso intervallo di tempo, e che il tempo aggiuntivo richiesto per la pianificazione (*scheduling*) delle attività e il coordinamento non comprometta tutto il beneficio potenziale del parallelismo.

Come esempio immaginiamo di dover scrivere un articolo per un quotidiano. Se riunissimo otto giornalisti, potenzialmente, l'articolo potrebbe essere terminato otto volte più velocemente. Per ottenere questo aumento di velocità occorre però che l'articolo sia scomposto in modo tale che tutti i giornalisti possano lavorare contemporaneamente. Perciò, abbiamo bisogno di pianificare



**Figura 1.17** Crescita nelle prestazioni dei processori a partire dalla metà degli anni '80. Questo grafico riporta le prestazioni relative al VAX11/780, misurate attraverso i benchmark SPECint (par. 1.10). Prima della metà degli anni '80, l'aumento delle prestazioni dei processori era dovuto principalmente alla tecnologia, ed era dell'ordine del 25% all'anno. L'aumento delle prestazioni nel periodo seguente è stato di circa il 52% all'anno, grazie a nuove idee nella progettazione delle architetture e nell'organizzazione dei calcolatori. Questo ha portato a un aumento delle prestazioni che nel 2002 era di sette volte l'aumento che si sarebbe verificato con un aumento di prestazioni del 25% all'anno. A partire dall'anno 2002 il limite sulla potenza assorbita, il parallelismo implicito delle istruzioni e la latenza della memoria hanno rallentato l'aumento delle prestazioni delle architetture monoprocessoressi a un 22% all'anno.

(*schedule*) l'attività di ciascuno di loro. Se qualcosa andasse storto e anche un solo giornalista terminasse la sua attività dopo gli altri sette, i benefici dell'avere otto giornalisti a disposizione svanirebbero. Occorre perciò *bilanciare il carico* di lavoro degli otto giornalisti per ottenere l'aumento di velocità massimo. Un altro problema si verificherebbe se i giornalisti dovessero passare troppo tempo a parlarsi tra loro per decidere come scrivere le loro parti. Non si riuscirebbe a raggiungere l'obiettivo neppure se una parte dell'articolo, per esempio la conclusione, non potesse essere scritta fino a quando le altre parti non fossero state terminate. Perciò, molta attenzione deve essere riposta nella *riduzione del tempo di comunicazione e di sincronizzazione*. Sia per gli otto giornalisti sia per la programmazione parallela, la sfida è rappresentata dallo *scheduling*, dal bilanciamento del carico e dal tempo richiesto per la sincronizzazione e la comunicazione tra gli attori. Come potete indovinare, la sfida diventa sempre più difficile all'aumentare del numero di giornalisti che scrivono l'articolo o dei core da programmare in parallelo.

Per meglio descrivere la metamorfosi delle architetture, i prossimi cinque capitoli contengono ciascuno un paragrafo che descrive le conseguenze della rivoluzione del parallelismo sugli argomenti trattati nel capitolo:

- **Capitolo 2, Paragrafo 2.11. Parallelismo e istruzioni: la sincronizzazione.** In alcuni istanti, attività parallele e indipendenti devono essere coordinate, per esempio quando è richiesto che forniscano un segnale al termine della loro attività. Questo capitolo illustra le istruzioni utilizzate dai processori multicore per sincronizzare i task.



- *Capitolo 3, Paragrafo 3.6. Parallelismo e aritmetica dei calcolatori: parallelismo a livello di parola.* Forse la forma di parallelismo più semplice da costruire è quella che consente di effettuare calcoli in parallelo su diversi elementi, come nel caso della moltiplicazione di due vettori tra loro. Il parallelismo a livello di parola sfrutta le risorse date dalla Legge di Moore per fornire delle unità di elaborazione che possano operare su più operandi simultaneamente.
- *Capitolo 4, Paragrafo 4.10. Parallelismo a livello di istruzioni.* Vista la difficoltà nel programmare esplicitamente in modo parallelo, negli anni '90 è stato fatto un grandissimo sforzo per rendere evidente il parallelismo implicito, inizialmente attraverso la pipeline. In questo capitolo vengono descritte alcune di queste tecniche aggressive, come la lettura ed esecuzione di più istruzioni in modo concorrente, la predizione del risultato delle decisioni e l'esecuzione speculativa delle istruzioni mediante la predizione.
- *Capitolo 5, Paragrafo 5.10. Parallelismo e gerarchie delle memorie: coerenza delle cache.* Un modo per diminuire il costo della comunicazione è fare sì che tutti i processori utilizzino lo stesso spazio di memoria, permettendo loro di scrivere o leggere un qualsiasi dato. Poiché tutti gli attuali processori utilizzano le memorie cache per avere una copia temporanea dei dati in una memoria più veloce e più vicina al processore, si può immaginare che la programmazione parallela sarebbe ancora più difficile se le memorie cache contenessero valori dei dati inconsistenti tra i vari processori. Questo capitolo descrive i meccanismi utilizzati per mantenere la coerenza dei dati tra le varie memorie cache.
- *Capitolo 5, Paragrafo 5.11. Parallelismo e gerarchie delle memorie: i dischi RAID.* Questo paragrafo descrive come si possano utilizzare più dischi assieme per ottenere un throughput molto maggiore; questa è stata l'ispirazione originale dei dischi RAID (*Redundant Arrays of Inexpensive Disks*, pile ridondante di dischi economici). La popolarità raggiunta dai RAID è però dipesa dalla loro maggiore affidabilità, ottenuta aggiungendo un limitato numero di dischi ridondanti. Nel paragrafo vengono spiegate le differenze nelle prestazioni, nei costi e nell'affidabilità dei vari livelli RAID.

Oltre a questi paragrafi, troverete un intero capitolo sull'elaborazione parallela. Il Capitolo 6 offre maggiori dettagli sulle sue sfide: presenta i due approcci opposti alla comunicazione degli indirizzi condivisi e al passaggio esplicito dei messaggi. Viene anche descritto un modello ridotto di parallelismo, più facile da programmare, viene spiegata la difficoltà nel confrontare i processori paralleli, viene introdotto un nuovo semplice modello per valutare le prestazioni dei processori multicore e, infine, vengono descritti e valutati quattro modelli di processori multicore utilizzando questo modello.

Come visto in precedenza, nei Capitoli 3 e 6 verrà utilizzato l'esempio della moltiplicazione di una matrice per un vettore per mostrare come il parallelismo possa aumentare significativamente le prestazioni.

Nell'Appendice C viene descritto un componente hardware, presente nei calcolatori desktop, che sta acquistando sempre maggiore popolarità: la GPU (*Graphics Processing Unit* – unità di elaborazione grafica). Le GPU sono state inventate per accelerare la grafica ma stanno diventando delle vere e proprie piattaforme di programmazione. Come si può immaginare vista l'attuale tendenza nello sviluppo delle architetture, le GPU dipendono fortemente dal parallelismo. L'Appendice C descrive la GPU dell'NVIDIA e mostra le parti principali del suo ambiente di sviluppo.

## 1.9 | Un caso reale: la valutazione del Core i7 Intel

Ogni capitolo contiene un paragrafo dedicato a un “caso reale”, il cui scopo è mettere in relazione i concetti descritti nel libro con un calcolatore che potreste utilizzare tutti i giorni. Questi paragrafi prendono in esame le tecnologie su cui sono basati i calcolatori moderni. Nel primo esaminiamo come vengono prodotti i circuiti integrati e in che modo si misurano le prestazioni e la potenza, utilizzando come esempio il Core i7 di Intel.

### Benchmark SPEC per la CPU

Un utente che tutti i giorni utilizzi gli stessi programmi è il candidato ideale per valutare le prestazioni di un nuovo calcolatore. L’insieme dei programmi che esegue rappresenta il **carico di lavoro** (*workload*) e per valutare due elaboratori l’utente dovrebbe semplicemente confrontare il tempo di esecuzione del carico di lavoro sulle due macchine. La maggior parte degli utenti, però, non si trova in questa situazione e deve affidarsi ad altri metodi per la misura delle prestazioni di una macchina, sperando che il metodo prescelto misuri effettivamente le prestazioni del nuovo calcolatore sul carico di lavoro di interesse. Solitamente si sceglie questa alternativa, valutando il calcolatore attraverso un insieme di **benchmark**, ossia programmi campione, appositamente scelti per misurare le prestazioni: i benchmark sono pensati per fornire un carico di lavoro che possa essere significativo al fine di stimare le prestazioni sui carichi di lavoro tipici. Come abbiamo già avuto modo di vedere, per rendere veloce la situazione più comune, occorre conoscere bene quale essa sia. Per questo motivo i benchmark giocano un ruolo importante nell’architettura dei calcolatori.

Lo SPEC (System Performance Evaluation Cooperative – Cooperativa per la Valutazione delle Prestazioni dei Sistemi) rappresenta uno sforzo congiunto, sovvenzionato e supportato da un certo numero di industrie produttrici di calcolatori, per creare un insieme standard di benchmark per i moderni calcolatori. I primi benchmark SPEC, nati nel 1989, erano focalizzati sulla valutazione delle prestazioni dei processori e sono ora indicati con la sigla SPEC89. Sono seguite cinque generazioni di benchmark SPEC di cui l’ultima, SPEC CPU2006, consiste in 12 programmi di benchmark per il calcolo intero (CINT2006) e in 17 programmi di benchmark per il calcolo in virgola mobile (CFP2006). I benchmark per il calcolo intero spaziano da parti di un compilatore C a un programma per giocare a scacchi e a un simulatore di calcolo quantistico. I benchmark per il calcolo in virgola mobile comprendono programmi di calcolo su griglia per la modellazione a elementi finiti, programmi che implementano metodi particellari per la dinamica molecolare e programmi di algebra lineare per la dinamica dei fluidi che implementano il calcolo con matrici sparse.

I programmi di benchmark SPEC per il calcolo intero e il loro tempo di esecuzione su un Core i7 sono riportati in Figura 1.18. Vengono anche riportati i diversi fattori che determinano il tempo di esecuzione: numero di istruzioni, CPI e periodo di clock. Si noti come il CPI vari di un fattore 5 tra i vari programmi.

Per semplificare la valutazione di un calcolatore, il consorzio SPEC ha deciso di riportare anche un singolo valore che riassumesse i risultati ottenuti nei 12 programmi di benchmark. Per ottenere tale valore, il tempo di esecuzione viene dapprima normalizzato, dividendolo per il tempo di esecuzione misurato su un calcolatore di riferimento; questa operazione fornisce una misura, chiamata **SPECratio** (rapporto tra misure SPEC), che ha il vantaggio di assumere un valore tanto maggiore quanto più veloce è l’esecuzione (lo SPECratio è l’inverso del tempo di esecuzione). Per ottenere una valutazione globale secondo CINT2006 o CFP2006, viene considerata la media geometrica degli SPECratio misurati sui diversi programmi di benchmark.

*Pensavo che [i calcolatori] sarebbero stati un’idea universalmente applicabile, come lo è un libro. Ma non pensavo che si sarebbe sviluppata così velocemente, perché non prevedevo che saremmo stati in grado di integrare su un chip tanti pezzi quanti in effetti è stato possibile. Il transistor comparve inatteso. Avvenne tutto molto più rapidamente di quanto noi ci attendessimo.*

J. Presper Eckert, coinventore dell’ENIAC, in un discorso del 1991.



### SITUAZIONI COMUNI

**Carico di lavoro:** un insieme di programmi eseguiti su un calcolatore; può essere costituito dall’insieme dei programmi utilizzati solitamente dall’utente oppure costruito a partire da programmi reali. Un tipico carico di lavoro specifica sia i programmi che la frequenza relativa con cui vengono eseguiti.

**Benchmark:** programma selezionato per comparare le prestazioni dei calcolatori.

Descrizione	Nome	Numero di istruzioni $\times 10^9$	CPI	Periodo di clock (secondi $\times 10^{-9}$ )	Tempo di esecuzione (secondi)	Tempo di riferimento (secondi)	SPECratio
Elaborazione mediante interpretazione di stringhe	perl	2252	0,60	0,376	508	9770	19,2
Compressione mediante <i>block-sorting</i>	bzip2	2390	0,70	0,376	629	9650	15,4
Compilatore GNU C	gcc	794	1,20	0,376	358	8050	22,5
Ottimizzazione combinatoria	mcf	221	2,66	0,376	221	9120	41,2
Gioco del go (AI, Artificial Intelligence)	go	1274	1,10	0,376	527	10 490	19,9
Ricerca di sequenze genetiche	hmmer	2616	0,60	0,376	590	9330	15,8
Gioco degli scacchi (AI)	sjeng	1948	0,80	0,376	586	12 100	20,7
Simulazione di calcolo quantistico	libquantum	659	0,44	0,376	109	20 720	190,0
Compressione video	h264avc	3793	0,50	0,376	713	22 130	31,0
Libreria di simulazione di eventi discreti	omnetpp	367	2,10	0,376	290	6250	21,5
Giochi/ricerca di percorsi	astar	1250	1,00	0,376	470	7020	14,9
Parsing XML	xalancbmk	1045	0,70	0,376	275	6900	25,1
Media geometrica	-	-	-	-	-	-	25,7

**Figura 1.18** Programmi benchmark che costituiscono lo SPECINT2006 e loro tempo di esecuzione sul Core i7 di Intel, modello 920. Come si è visto a pagina 31, l’equazione che calcola il tempo di esecuzione dipende da tre fattori: il numero di istruzioni (misurate qui in miliardi), il numero di cicli per istruzione (CPI) e il periodo di clock, misurato in nanosecondi. Lo SPECratio (rapporto tra misure SPEC) è semplicemente il rapporto tra il tempo di esecuzione di riferimento, fornito da SPEC, e il tempo di esecuzione misurato. L’unico numero riportato in corrispondenza della voce “Media geometrica” è la media geometrica degli SPECratio.

**Approfondimento.** Quando si vogliono confrontare due calcolatori attraverso lo SPECratio, si utilizza la media geometrica in modo da ottenere lo stesso risultato indipendentemente dal calcolatore utilizzato per normalizzare i risultati. Se facessemo la media dei tempi di esecuzione normalizzati utilizzando la media aritmetica, il risultato varierebbe a seconda del calcolatore utilizzato come riferimento.

La formula della media geometrica è:

$$\sqrt[n]{\prod_{i=1}^n \text{Rapporto tempo di esecuzione}_i}$$

dove *rapporto del tempo di esecuzione* rappresenta il rapporto tra il tempo di esecuzione dell’*i*-esimo programma di benchmark sul calcolatore da valutare e quello sul calcolatore di riferimento, *n* rappresenta il numero di programmi di benchmark e:

$$\prod_{i=1}^n a_i \text{ indica il prodotto } a_1 \times a_2 \times \dots \times a_n$$

### Benchmark SPEC sull’assorbimento di potenza

Data l’importanza che l’assorbimento di energia e di potenza ricopre nei calcolatori moderni, SPEC ha creato dei benchmark specifici per misurare la potenza. Questi benchmark misurano l’assorbimento di potenza di server in condizioni reali con diversi livelli di carico di lavoro, suddivisi in intervalli con ampiezza del 10%. In Figura 1.19 vengono riportate le misure relative a un server che utilizza il processore Xeon di Intel, simile al Core i7.

% Carico di lavoro	Prestazioni (ssj_ops)	Potenza media (watt)
100%	865 618	258
90%	786 688	242
80%	698 051	224
70%	607 826	204
60%	521 391	185
50%	436 757	170
40%	345 919	157
30%	262 071	146
20%	176 061	135
10%	86 784	121
0%	0	80
Somma totale	4 787 166	1922
$\Sigma \text{ssj\_ops} / \Sigma \text{potenza} =$		2490

**Figura 1.19** Pogramma ssj2008 del banchmark SPECpower eseguito su uno Xeon X5650 di Intel, 2,66 GHz, doppio connettore, con 16 GB di DRAM e un disco SSD da 100 GB.

Il primo benchmark sull'assorbimento di potenza era presente nello SPECJBB2005, proposto per valutare le applicazioni Java per il mondo degli affari, e impegnava il processore, la cache e la memoria principale, nonché la Java virtual machine (macchina virtuale Java), il compilatore, il programma di pulizia della memoria (*garbage collector*) e parti del sistema operativo. Le prestazioni vengono misurate attraverso il throughput e l'unità di misura utilizzata è il numero di transazioni al secondo. Anche qui, per semplificare la valutazione, SPEC propone un singolo numero, chiamato “ssj\_ops totali per watt”, che viene calcolato come segue:

$$\text{ssj\_ops totali per watt} = \left( \sum_{i=0}^{10} \text{ssj\_ops}_i \right) / \left( \sum_{i=0}^{10} \text{potenza}_i \right)$$

dove  $\text{ssj\_ops}_i$  rappresenta le prestazioni per ciascuno degli intervalli con ampiezza del 10% del carico di lavoro e  $\text{potenza}_i$  rappresenta la potenza assorbita per ciascun livello di carico.

## 1.10 Errori e trabocchetti

L'obiettivo di questo paragrafo è quello di confutare alcune credenze sbagliate che si possono incontrare quando si parla di architetture degli elaboratori. Ogni capitolo di questo libro conterrà un paragrafo di questo tipo, e ogni volta che illustreremo un errore che nasce da una credenza sbagliata cercheremo di fornire un controcuento. Illustreremo anche alcuni *trabocchetti*, ossia errori in cui è facile incappare: questi sono spesso originati dalla generalizzazione di principi che sono veri solo in un contesto limitato. Lo scopo di questi paragrafi è quello di aiutare il lettore a evitare di cadere in questi errori nella progettazione e nella valutazione dei calcolatori. Errori e trabocchetti che riguardano il rapporto costo/prestazioni hanno tratto in inganno molti progettisti, compresi

*La scienza deve iniziare dai miti e dalla discussione critica dei miti.*

Karl Popper, *The Philosophy of Science*, 1957

gli autori di questo libro, perciò questo paragrafo non soffre della mancanza di esempi significativi. Iniziamo da un trabocchetto che ha fatto inciampare molti progettisti e rivela una caratteristica importante della progettazione dei calcolatori.

*Trabocchetto: ci si aspetta che il miglioramento di uno dei componenti di un calcolatore produca un aumento delle prestazioni proporzionale alla dimensione del miglioramento.*



**Legge di Amdahl:** una regola che afferma che il miglioramento delle prestazioni reso possibile da una data modifica è limitato dalla quantità tempo in cui quella modifica è effettivamente sfruttata. È la versione quantitativa della legge dei rendimenti decrescenti.

La grande idea di rendere veloce la situazione più comune ha un corollario che tormenta i progettisti dell'hardware e del software e che ci ricorda che l'impatto del miglioramento di una funzionalità dipende dal tempo per il quale quella funzionalità verrà utilizzata.

Lo illustriamo con un semplice esempio. Si supponga che un programma venga eseguito in 100 secondi su un dato calcolatore e che 80 di questi siano impiegati in operazioni di moltiplicazione. Di quanto bisogna migliorare la velocità di esecuzione delle moltiplicazioni se si vuole che il programma venga eseguito 5 volte più velocemente?

Il tempo di esecuzione del programma dopo il miglioramento è dato dalla seguente semplice equazione, nota come **Legge di Amdahl:**

$$\text{Tempo di esecuzione dopo il miglioramento} = \frac{\text{Tempo di esecuzione influenzato dal miglioramento}}{\text{Miglioramento}} + \text{Tempo di esecuzione non influenzato dal miglioramento}$$

Nel nostro caso si ottiene:

$$\text{Tempo di esecuzione dopo il miglioramento} = \frac{80 \text{ secondi}}{n} + (100 - 80 \text{ secondi})$$

Dato che vogliamo che l'esecuzione diventi cinque volte più veloce, il tempo di esecuzione dovrà scendere a 20 secondi, e quindi:

$$\begin{aligned} 20 \text{ secondi} &= \frac{80 \text{ secondi}}{n} + 20 \text{ secondi} \\ 0 &= \frac{80 \text{ secondi}}{n} \end{aligned}$$

Ciò significa che *non esiste* un miglioramento nella velocità di esecuzione delle moltiplicazioni,  $n$ , che permetta di ottenere un aumento delle prestazioni di un fattore 5 se le moltiplicazioni rappresentano solo l'80% del carico di lavoro del processore. Il possibile incremento totale delle prestazioni indotto dal miglioramento di una caratteristica è limitato dal tempo di utilizzo di quella caratteristica. Questo concetto è comunemente noto come la legge dei rendimenti decrescenti (*diminishing returns*) e si applica a molti casi della vita reale.

Possiamo utilizzare la Legge di Amdahl per stimare il miglioramento delle prestazioni quando conosciamo il tempo in cui viene utilizzata una certa funzionalità e il suo potenziale incremento di velocità. La Legge di Amdahl, insieme all'equazione delle prestazioni della CPU, è uno strumento facile per valutare potenziali miglioramenti; essa verrà analizzata più nel dettaglio negli esercizi. Una delle linee guida nella progettazione dell'hardware è descritta da un corollario della Legge di Amdahl: rendere veloce la situazione più comune. Essa invita a tenere conto del fatto che in molti casi la frequenza con cui accade un evento può essere molto più elevata di un'altra.

La Legge di Amdahl viene utilizzata anche per individuare il numero limite pratico di processori paralleli; esamineremo questo problema nel paragrafo *Errori e trabocchetti* del Capitolo 6.

**Errore:** i calcolatori che vengono utilizzati poco dissipano poca potenza.

L'efficienza energetica è importante anche quando un calcolatore viene utilizzato poco, dato che il carico di lavoro dei server è variabile: l'utilizzo dei server dei centri di calcolo di Google, per esempio, varia tra il 10 e il 50% per la maggior parte del tempo e solamente in meno dell'1% del tempo raggiunge il 100%. Anche dopo cinque anni di esperienza nell'utilizzo dei benchmark SPECpower, i calcolatori meglio ottimizzati del 2012 utilizzavano il 33% della potenza di picco a fronte di un 10% di carico di lavoro. Sistemi non configurati attraverso i benchmark SPECpower assorbono sicuramente maggiore potenza.

Dato che viene assorbita una frazione significativa della potenza di picco, nonostante il carico di lavoro dei server sia variabile, Luiz Barroso e Urs Hözle [2007] hanno proposto di riprogettare l'hardware in modo che "il calcolo sia proporzionale all'energia". Se i server del futuro utilizzeranno, per esempio, il 10% della potenza di picco per eseguire il 10% del carico di lavoro, si potrebbe ridurre il consumo energetico dei centri di elaborazione dati: questi diventerebbero società virtuose in un periodo di aumento dell'attenzione sulle emissioni di CO<sub>2</sub>.

**Errore:** progettare un'architettura che massimizzi le prestazioni e un'architettura che massimizzi l'efficienza energetica sono obiettivi non collegati.

Dato che l'energia esprime l'assorbimento di potenza nel tempo, ottimizzazioni hardware e software che riducono il tempo di esecuzione consentono un risparmio globale di energia anche se i componenti che sono stati ottimizzati possono richiedere una quantità maggiore di energia. Il motivo è che tutto il resto del calcolatore consuma energia mentre un programma è in esecuzione, per cui anche se i componenti ottimizzati consumano un po' più di energia, il minore tempo di esecuzione fa sì che tutti gli altri componenti consumino globalmente meno energia; il bilancio energetico complessivo è quindi positivo.

**Trabocchetto:** utilizzare solo una parte dell'equazione delle prestazioni come metrica per valutare le prestazioni.

Abbiamo già visto l'errore che si commette se si misurano le prestazioni utilizzando solamente la frequenza di clock, il numero di istruzioni oppure il CPI. Un altro errore comune è utilizzare solamente due dei tre fattori coinvolti per confrontare le prestazioni. Benché l'utilizzo di due fattori possa essere valido in un contesto circoscritto, spesso viene fatto a sproposito; quasi tutte le alternative al tempo proposte come metrica per misurare le prestazioni hanno portato ad affermazioni errate, risultati distorti o interpretazioni scorrette.

Un'alternativa alle misure basate sul tempo per valutare le prestazioni è rappresentata dai **MIPS** (*Million Instructions Per Second – milioni di istruzioni al secondo*). Dato un programma, i MIPS sono definiti semplicemente come:

$$\text{MIPS} = \frac{\text{Numero di istruzioni}}{\text{Tempo di esecuzione} \times 10^6}$$

**MIPS (milioni di istruzioni al secondo):** una misura della velocità di esecuzione di un programma basata sul numero (milioni) di istruzioni. I MIPS vengono calcolati come numero totale delle istruzioni diviso per il prodotto del tempo di esecuzione e 10<sup>6</sup>.

Dal momento che i MIPS sono una misura della velocità di esecuzione delle istruzioni, rappresentano il reciproco del tempo di esecuzione: le macchine più veloci hanno un valore più elevato di MIPS. L'aspetto positivo dei MIPS è che rappresentano un concetto facile da comprendere: le macchine più veloci avranno un valore di MIPS più elevato, e ciò risulta intuitivo.

Tuttavia ci sono tre problemi nell'utilizzo dei MIPS come misura per confrontare i calcolatori. In primo luogo, essi misurano la frequenza di esecuzione delle istruzioni ma non tengono conto delle caratteristiche funzionali delle istruzioni: non si possono confrontare due calcolatori con insiemi di istruzioni differenti utilizzando i MIPS, dal momento che il numero di istruzioni risulterebbe sicuramente diverso. In secondo luogo, nello stesso calcolatore i MIPS variano a seconda del programma, quindi un calcolatore non ha un'unica valutazione in MIPS valida per tutti i programmi. Per esempio, possiamo ottenere una relazione tra MIPS, frequenza di clock e CPI sostituendo l'espressione del tempo di esecuzione ricavata dall'equazione fondamentale delle prestazioni:

$$\text{MIPS} = \frac{\text{Numero di istruzioni}}{\frac{\text{Numero di istruzioni} \times \text{CPI}}{\text{Frequenza di clock}} \times 10^6} = \frac{\text{Frequenza di clock}}{\text{CPI} \times 10^6}$$

Dato che il CPI varia di un fattore 5 per le misure SPEC CPU2006 sul Core i7 di Intel (Figura 1.18), i MIPS varieranno anch'essi di un fattore 5 tra i vari programmi di benchmark. Inoltre, se un programma esegue più istruzioni, ciascuna delle quali è più veloce, avremo che i MIPS possono variare anche indipendentemente dalle prestazioni!

## Autovalutazione

Considerate le seguenti misure di prestazioni di un programma:

Misura	Calcolatore A	Calcolatore B
Numero di istruzioni	10 miliardi	8 miliardi
Frequenza di clock	4 GHz	4 GHz
CPI	1,0	1,1

- Quale calcolatore ha il più alto numero di MIPS?
- Quale calcolatore è più veloce?

Mentre ... l'ENIAC è dotato di 18 000 valvole e pesa 30 tonnellate, i calcolatori del futuro potranno avere solo 1000 valvole e pesare solamente una tonnellata e mezza.  
Popular Mechanics, marzo 1949.

## 1.11 Note conclusive

Sebbene sia difficile prevedere esattamente il costo e le prestazioni dei calcolatori futuri, si può scommettere tranquillamente che saranno molto migliori degli attuali. Per giocare una parte attiva in questo progresso, i progettisti di calcolatori e i programmati dovranno considerare un insieme di problematiche più vasto.

I progettisti sia dell'hardware sia del software costruiscono i sistemi di elaborazione per strati gerarchici, detti livelli di astrazione, in cui ciascuno strato nasconde dei dettagli al livello superiore. Questa grande idea, l'**astrazione**, è fondamentale per comprendere i calcolatori odierni, ma non significa che i progettisti debbano conoscere solo la tecnologia utilizzata nello strato di loro



ASTRAZIONE

competenza. Forse l'esempio più importante di astrazione è l'interfaccia tra l'hardware e il software di basso livello, detta *architettura dell'insieme di istruzioni*. Se l'architettura dell'insieme di istruzioni rimane invariata, uno stesso programma sarà eseguibile su diverse implementazioni dell'architettura, che presumibilmente avranno costi e prestazioni diverse. D'altro canto, l'architettura potrebbe costituire un ostacolo all'introduzione di innovazioni che porterebbero alla modifica di questa interfaccia.

C'è un metodo affidabile per determinare e misurare le prestazioni utilizzando come metrica il tempo di esecuzione di programmi reali. Questo è legato ad altre grandezze importanti attraverso la seguente equazione:

$$\frac{\text{Secondi}}{\text{Programma}} = \frac{\text{Istruzioni}}{\text{Programma}} \times \frac{\text{Cicli di clock}}{\text{Istruzione}} \times \frac{\text{Secondi}}{\text{Ciclo di clock}}$$

Utilizzeremo più e più volte questa equazione e i suoi tre termini. Ricordatevi che i singoli fattori non determinano le prestazioni: solamente il loro prodotto costituisce una misura affidabile delle prestazioni ed è uguale al tempo di esecuzione.

## QUADRO D'INSIEME

Il tempo di esecuzione è l'unica misura valida e inconfutabile delle prestazioni. Molte altre metriche sono state proposte e sono state inizialmente considerate interessanti. A volte queste metriche risultavano sbagliate fin dall'inizio perché non riflettevano il tempo di esecuzione; altre volte una metrica valida in un contesto circoscritto veniva estesa e utilizzata al di fuori di quel contesto oppure estesa senza le specifiche necessarie per utilizzarla validamente. ■

La tecnologia fondamentale dell'hardware dei calcolatori moderni è basata sul silicio. Di pari importanza rispetto alla comprensione della tecnologia dei circuiti integrati è la comprensione del tasso di evoluzione tecnologica previsto dalla **Legge di Moore**. Mentre il silicio continua a sostenere una rapida evoluzione dell'hardware, nuove idee sull'organizzazione dei calcolatori hanno migliorato il rapporto prezzo/prestazioni. Le due idee fondamentali sono lo sfruttamento del parallelismo nei programmi, tipicamente ottenuto oggi grazie all'utilizzo di processori multicore, e lo sfruttamento della località nell'accesso alla **gerarchia delle memorie**, raggiunto soprattutto attraverso le memorie cache.

Il consumo di energia ha sostituito le dimensioni del chip come elemento critico nella progettazione dei microprocessori. Aumentare le prestazioni senza aumentare l'assorbimento di potenza ha obbligato i produttori di hardware a passare ai processori multicore, richiedendo quindi che i produttori di software passassero al **parallelismo**.

I diversi calcolatori sono stati sempre valutati non solo in termini di costi e prestazioni, ma anche di altri importanti fattori quali il consumo di energia, l'affidabilità, il costo e la scalabilità. Anche se questo capitolo era principalmente focalizzato sul costo, sulle prestazioni e sull'energia, un progetto vincente permetterà il giusto equilibrio tra i tre fattori per il mercato a cui si rivolge.

## Organizzazione del testo

Alla base di queste astrazioni ci sono i cinque componenti classici di un calcolatore: unità di elaborazione dati (o datapath), unità di controllo, memoria,



LEGGE DI MOORE



GERARCHIA



PARALLELISMO

input e output (Figura 1.5). Questi cinque componenti servono anche come quadro di riferimento per gli altri capitoli di questo libro:

- *unità di elaborazione dati* (o datapath): Capitoli 3, 4, 6 e Appendice C 
- *unità di controllo*: Capitoli 4, 6 e Appendice C 
- *memoria*: Capitolo 5
- *input*: Capitoli 5 e 6
- *output*: Capitoli 5 e 6

Il Capitolo 4 descrive come i processori sfruttano il parallelismo implicito, mentre nel Capitolo 6 vengono descritti i processori multicore, caratterizzati da un parallelismo esplicito, che sono il nocciolo della rivoluzione del parallelismo; un processore grafico ad alto grado di parallelismo viene descritto nell'Appendice C . Il Capitolo 5 descrive come le gerarchie delle memorie sfruttano la località. Il Capitolo 2 descrive gli insiemi di istruzioni – cioè l'interfaccia tra i compilatori e il calcolatore – ed enfatizza il ruolo dei compilatori e dei linguaggi di programmazione nello sfruttare le caratteristiche dell'insieme delle istruzioni. Il Capitolo 3 mostra come i calcolatori gestiscono i numeri e realizzano le operazioni aritmetiche. L'Appendice A  descrive le tecniche di progettazione dei circuiti logici. Potete trovare tutte le Appendici online sul sito web del libro:

[online.universita.zanichelli.it/patterson5e](http://online.universita.zanichelli.it/patterson5e)

## 1.12 Inquadramento storico e approfondimenti

*Una branca della scienza in piena attività è come un immenso formicaio: l'individuo quasi svanisce nell'ammasso di menti che turbinano, trasmettendo le informazioni da una parte all'altra, diffondendole alla velocità della luce.*

Lewis Thomas, "Natural Science", in *The Lives of a Cell*, 1974.

Per ogni capitolo, troverete disponibile online sul sito web del libro un paragrafo dedicato all'inquadramento storico. In questi paragrafi illustriamo lo sviluppo di un'idea anche attraverso i calcolatori che si sono succeduti negli anni, oppure descriviamo progetti particolarmente importanti; inoltre, vengono forniti riferimenti bibliografici utili per ulteriori approfondimenti.

L'inquadramento storico relativo a questo capitolo presenta lo stato dell'arte di alcune delle idee chiave presentate. Lo scopo è illustrare la storia del progresso tecnologico attraverso la storia degli scienziati e degli ingegneri che lo hanno determinato e collocare le loro conquiste nel contesto storico. Comprendendo il passato, si può capire meglio quali aspetti guideranno lo sviluppo dell'informatica nel futuro. Al termine di ogni paragrafo di inquadramento storico, vengono riportati suggerimenti per ulteriori approfondimenti; tutti i suggerimenti dei vari capitoli sono raggruppati nella sezione *Approfondimenti*, disponibile online .

## 1.13 | Esercizi

Il tempo relativo stimato per completare un esercizio viene mostrato tra parentesi quadre a fianco del numero dell'esercizio: in generale, un esercizio valutato [10] richiederà, per essere svolto, il doppio del tempo di un esercizio valutato [5]. Le parti del libro che dovreste leggere prima di svolgere un esercizio sono riportate tra parentesi uncinate; per esempio, <1.4> significa che dovreste avere letto il paragrafo 1.4, "Componenti di un calcolatore", prima di cercare di risolvere l'esercizio.

**1.1** [2] <1.1> Elencare e descrivere altri quattro tipi di dispositivi che utilizzano calcolatori, oltre agli smartphone utilizzati da un miliardo di persone.

**1.2** [5] <1.2> Le otto grandi idee nella progettazione delle architetture sono simili a soluzioni proposte in altri campi. Associare le otto grandi idee nelle architetture: "Legge di Moore", "Rendere veloci le situazioni più comuni", "Prestazioni attraverso il parallelismo", "Prestazioni attraverso la pipeline", "Prestazioni at-

traverso la predizione”, “Gerarchie delle memorie”, “Affidabilità attraverso la ridondanza”, alle seguenti idee applicate in altri campi:

- linee di produzione nelle fabbriche di automobili;
- cavi per i ponti sospesi;
- sistemi di navigazione per gli aerei e le navi che incorporano informazioni sui venti;
- ascensori veloci per gli edifici;
- banco di prenotazione dei libri in una biblioteca;
- aumentare l'area del gate di un transistor CMOS per diminuire il suo tempo di commutazione;
- introdurre per il lancio degli aerei delle catapulte azionate dall'elettromagnetismo (alimentate dalla corrente elettrica) al posto di quelle attuali, alimentate dal vapore in pressione; questo viene reso possibile dall'aumento della potenza generata dai reattori costruiti con la nuova tecnologia;
- costruire automobili che si guidano da sole il cui sistema di controllo è basato in parte sui sensori che vengono già installati sui veicoli, quali i sistemi per il controllo dello spostamento di corsia e i sistemi intelligenti di controllo automatico della velocità.

**1.3** [2] <1.3> Descrivere i passi necessari per trasformare un programma scritto in linguaggio ad alto livello, per esempio in C, in una rappresentazione che possa essere eseguita direttamente dal processore di un calcolatore.

**1.4** [2] < 1.4> Un terminale video a colori utilizza 8 bit per pixel per ciascuno dei tre colori primari (rosso, verde e blu) e ha una risoluzione di  $1280 \times 1024$  pixel.

- Qual è la dimensione (in byte) del frame buffer associato?
- Quanto tempo occorre per trasmettere un frame attraverso una rete da 100 Mbit/s?

**1.5** [4] <1.6> Si considerino tre diversi processori P1, P2 e P3 che eseguono lo stesso insieme di istruzioni. P1 ha una frequenza di clock di 3 GHz e un CPI di 1,5, P2 ha una frequenza di clock di 2,5 GHz e un CPI di 1,0 e P3 ha una frequenza di clock di 4,0 GHz e un CPI di 2,2.

- Quale processore ha le prestazioni migliori espresse in numero di istruzioni al secondo?
- Determinare il numero di cicli di clock utilizzati e di istruzioni eseguite da ciascun processore, supponendo che tutti eseguano un programma in 10 secondi.
- Si vorrebbe ridurre il tempo di esecuzione del 30%, ma per raggiungere questo obiettivo il CPI aumenterebbe del 20%. Quale sarebbe la frequenza di clock che consentirebbe questa riduzione del tempo di esecuzione?

**1.6** [20] <1.6> Si considerino due differenti implementazioni dello stesso insieme di istruzioni che si possono suddividere in quattro classi: A, B, C e D, a seconda del loro CPI. L'implementazione P1 ha una frequenza di clock di 2,5 GHz e un CPI rispettivamente di 1, 2, 3 e 3, mentre l'implementazione P2 ha una frequenza di clock di 3 GHz e un CPI rispettivamente di 2, 2, 2 e 2.

Si consideri un programma costituito da  $1,0E6$  istruzioni così suddivise: 10% di classe A, 20% di classe B, 50% di classe C e 20% di classe D. Quale delle due è più veloce: P1 o P2?

- Qual è il CPI totale per ciascuna delle due implementazioni?
- Determinare il numero di cicli di clock richiesti dalle due implementazioni per eseguire il programma.

**1.7** [15] <1.6> I compilatori possono avere un profondo impatto sulle prestazioni delle applicazioni per un certo processore. Si supponga che, per un certo programma, il compilatore A produca un numero di istruzioni pari a  $1,0E9$  con un tempo di esecuzione di 1,1 s, mentre il compilatore B produca un numero di istruzioni pari a  $1,2E9$  con un tempo di esecuzione di 1,5 s.

- Determinare il CPI medio per ciascuno dei due codici compilati, sapendo che il processore ha un periodo di clock di 1 ns.
- Si supponga che il programma compilato venga eseguito su due processori diversi. Se il tempo di esecuzione sui due processori risulta uguale, di quanto dovrà essere più veloce il clock del processore che esegue il codice compilato da A rispetto al clock del processore che esegue il codice compilato da B?
- Viene sviluppato un nuovo compilatore che crea un programma costituito solamente da 600 milioni di istruzioni con un CPI medio di 1,1. Quale sarà lo speedup ottenuto con il nuovo compilatore rispetto al tempo di esecuzione dei programmi compilati con i compilatori A e B eseguiti sul calcolatore originario?

**1.8** Il processore Pentium 4 Prescott, lanciato nel 2004, aveva una frequenza di clock di 3,6 GHz e una tensione di alimentazione di 1,25 V. Si supponga che il suo assorbimento fosse, in media, di 10 W di potenza statica e di 90 W di potenza dinamica.

Il calcolatore Core i5 Ivy Bridge, lanciato nel 2012, aveva una frequenza di clock di 3,4 GHz e una tensione di alimentazione di 0,9 V. Si supponga che il suo assorbimento sia, in media, di 30 W di potenza statica e di 40 W di potenza dinamica.

**1.8.1 [5] <1.7>** Determinare il carico capacitivo medio per ciascuno dei due processori.

**1.8.2 [5] <1.7>** Determinare la percentuale di assorbimento di potenza statica rispetto all'assorbimento totale e il rapporto tra assorbimento statico e dinamico per ciascuno dei due processori.

**1.8.3 [15] <1.7>** Di quanto occorre ridurre la tensione di alimentazione per mantenere lo stesso livello di correnti di dispersione se si vuole diminuire l'assorbimento totale di potenza del 10%?

**1.9** Si supponga che un processore abbia un CPI rispettivamente di 1, 1,2 e 1,5 per le istruzioni aritmetiche, di lettura/scrittura e di salto condizionato. Si supponga anche che un programma richieda l'esecuzione di 2,56E9 istruzioni aritmetiche, 1,28E9 istruzioni di lettura/scrittura e 256 milioni di istruzioni di salto condizionato, su un certo processore, e che il processore abbia una frequenza di clock di 2 GHz.

Si supponga che il numero di istruzioni aritmetiche e di lettura/scrittura del programma parallelizzato su più core si riduca di un fattore  $0,7 \times p$ , dove  $p$  è il numero di core, e che il numero di istruzioni di salto condizionato rimanga invariato.

**1.9.1 [5] <1.7>** Determinare il tempo di esecuzione totale del programma su un processore a 1, 2, 4 e 8 core e calcolare lo speedup ottenuto con i processori a 2, 4 e 8 core rispetto al processore a core singolo.

**1.9.2 [10] <1.6, 1.8>** Quale sarebbe l'impatto del raddoppio del CPI delle istruzioni aritmetiche sul tempo di esecuzione del programma eseguito su un processore a 1, 2, 4 e 8 core?

**1.9.3 [10] <1.6, 1.8>** Di quanto si dovrebbe ridurre il CPI delle istruzioni di lettura/scrittura perché le prestazioni del processore a singolo core siano uguali a quelle del processore con 4 core sul quale vengono eseguite le istruzioni con il CPI originario?

**1.10** Si supponga che un wafer del diametro di 15 cm abbia un costo di 12, contenga 84 chip e presenti 0,020 difetti/cm<sup>2</sup>, e che un wafer del diametro di 20 cm abbia un costo di 15, contenga 100 chip e presenti 0,031 difetti/cm<sup>2</sup>.

**1.10.1 [10] <1.5>** Determinare la resa di entrambi i wafer.

**1.10.2 [5] <1.5>** Determinare il costo di ciascun chip per entrambi i wafer.

**1.10.3 [5] <1.5>** Determinare il numero di chip che si riescono a produrre e la resa quando il numero di

chip per wafer viene aumentato del 10% e il numero di difetti per unità di area aumenta del 15%.

**1.10.4 [5] <1.5>** Si supponga che il processo di fabbricazione migliori la resa da 0,92 a 0,95. Determinare il numero di difetti per unità di area per ciascuno dei due processi per un chip di 200 mm<sup>2</sup>.

**1.11** I risultati misurati sul programma di benchmark bzip2 dello SPEC CPU2006, eseguito su un processore AMD Barcelona, riportano un numero di istruzioni pari a 2,389E12, un tempo di esecuzione di 750 s e un tempo di riferimento pari a 9650 s.

**1.11.1 [5] <1.6, 1.9>** Determinare il CPI sapendo che il periodo del clock è di 0,333 ns.

**1.11.2 [5] <1.9>** Determinare lo SPECratio.

**1.11.3 [5] <1.6, 1.9>** Determinare l'aumento del tempo di CPU se il numero di istruzioni del programma di benchmark viene aumentato del 10% senza alterare il CPI.

**1.11.4 [5] <1.6, 1.9>** Determinare l'aumento del tempo di CPU se il numero di istruzioni del programma di benchmark viene aumentato del 10% e il CPI viene aumentato del 5%.

**1.11.5 [5] <1.6, 1.9>** Determinare la variazione dello SPECratio conseguente alle modifiche.

**1.11.6 [10] <1.6>** Si supponga di sviluppare una nuova versione del processore AMD Barcelona con una frequenza di clock di 4 GHz, aggiungendo alcune istruzioni all'insieme originario. Con il nuovo insieme di istruzioni il numero di istruzioni costituenti il programma di benchmark viene diminuito del 15%. Il tempo di esecuzione scende a 700 s e il nuovo SPECratio diventa di 13,7. Determinare il CPI.

**1.11.7 [10] <1.6>** Il CPI è maggiore di quello ottenuto nell'esercizio 1.11.1 dato che la frequenza di clock è passata da 3 GHz a 4 GHz. Determinare se l'aumento del CPI è proporzionale all'aumento della frequenza di clock; in caso contrario, fornire una spiegazione.

**1.11.8 [5] <1.6>** Di quanto è stato ridotto il tempo di CPU?

**1.11.9 [10] <1.6>** Si consideri un secondo programma del benchmark SPEC CPU2006: libquantum. Si supponga che per questo programma il tempo di esecuzione sia di 960 s, il CPI sia di 1,61 e la frequenza di clock sia sempre di 3 GHz. Determinare il numero di istruzioni se il tempo di esecuzione viene ridotto di un altro 10%, senza modificare il CPI, e la frequenza di clock aumentata a 4 GHz.

**1.11.10 [10] <1.6>** Determinare la frequenza di clock richiesta per ottenere un'ulteriore riduzione del 10% del tempo di CPU senza modificare il CPI e il numero di istruzioni.

**1.11.11 [10] <1.6>** Determinare la frequenza di clock se il CPI viene ridotto del 15% e il tempo di CPU del 20% senza modificare il numero di istruzioni.

**1.12** Il paragrafo 1.10 descrive l'errore in cui si può incorrere quando viene utilizzata solamente una parte dell'equazione delle prestazioni per misurare le prestazioni stesse. Per illustrare ciò, si considerino i seguenti processori: P1 ha una frequenza di clock di 4 GHz, un CPI medio di 0,9 e richiede l'esecuzione di 5,0E9 istruzioni, mentre P2 ha una frequenza di clock di 3 GHz, un CPI medio di 0,75 e richiede l'esecuzione di 1,0E9 istruzioni.

**1.12.1 [5] <1.6, 1.10>** Uno degli errori più comuni è quello di pensare che il calcolatore con il clock a frequenza maggiore abbia anche le prestazioni maggiori. Verificare se ciò è vero per i calcolatori P1 e P2.

**1.12.2 [10] <1.6, 1.10>** Un altro errore è quello di pensare che il processore che esegue il numero di istruzioni maggiore impieghi un tempo di CPU maggiore. Supponendo che il processore P1 esegua una sequenza di 1,0E9 istruzioni e che il CPI dei due processori, P1 e P2, non venga modificato, determinare il numero di istruzioni che P2 può eseguire nell'intervallo di tempo in cui P1 esegue 1,0E9 istruzioni.

**1.12.3 [10] <1.6, 1.10>** Un errore comune è quello di utilizzare i MIPS (milioni di istruzioni per secondo) per confrontare le prestazioni di due processori e concludere che il processore con il valore di MIPS più elevato sia anche quello più performante. Verificare se ciò è vero per i processori P1 e P2.

**1.12.4 [10] <1.10>** Un'altra unità di misura delle prestazioni di uso comune è il MFLOPS (milioni di istruzioni in virgola mobile per secondo, *Millions of Floating point Operations Per Second*), definito come:

$$\text{MFLOPS} = \text{N}^{\circ} \text{ di operazioni FP} / (\text{tempo di esecuzione} \times 1\text{E}6)$$

Tuttavia questa unità di misura ha gli stessi problemi del MIPS. Si supponga che il 40% delle istruzioni eseguite da P1 e P2 sia costituito da istruzioni *floating-point* e si determinino i MFLOPS dei processori.

**1.13** Un altro trabocchetto riportato nel paragrafo 1.10 è quello per cui ci si aspetta un miglioramento delle prestazioni complessive modificando solamente

una parte del calcolatore. Questo non è sempre vero. Si consideri un calcolatore che esegue un programma che richiede 250 s, 70 s dei quali sono spesi per l'esecuzione di istruzioni FP, 85 s per l'esecuzione di istruzioni di lettura/scrittura e 40 s per l'esecuzione di istruzioni di salto condizionato.

**1.13.1 [5] <1.10>** Di quanto verrebbe ridotto il tempo totale di esecuzione se il tempo di esecuzione delle operazioni FP venisse ridotto del 20%?

**1.13.2 [5] <1.10>** Di quanto verrebbe ridotto il tempo di esecuzione delle operazioni su interi (INT) se il tempo totale di esecuzione venisse ridotto del 20%?

**1.13.3 [5] <1.10>** Si può ridurre il tempo totale del 20% riducendo solamente il tempo di esecuzione dei salti condizionati?

**1.14** Si supponga che un programma richieda l'esecuzione di  $50 \times 10^6$  istruzioni FP,  $100 \times 10^6$  istruzioni INT,  $80 \times 10^6$  istruzioni di lettura/scrittura e  $16 \times 10^6$  istruzioni di salto condizionato. Il CPI per ciascun tipo di istruzione è rispettivamente di 1, 1, 4 e 2. Si supponga che il processore abbia una frequenza di clock di 2 GHz.

**1.14.1 [10] <1.10>** Di quanto occorre aumentare il CPI delle istruzioni FP se si vuole che il programma venga eseguito al doppio della velocità?

**1.14.2 [10] <1.10>** Di quanto occorre aumentare il CPI delle istruzioni di lettura/scrittura se si vuole che il programma venga eseguito al doppio della velocità?

**1.14.3 [5] <1.10>** Di quanto aumenterà la velocità di esecuzione del programma se si riduce il CPI delle istruzioni INT e FP del 40% e quello delle istruzioni di lettura/scrittura e salto condizionato del 30%?

**1.15 [5] <1.8>** Quando un programma viene adattato per essere eseguito sui diversi processori di un'architettura multiprocessore, il tempo di esecuzione su ciascun processore comprende il tempo di calcolo, il tempo aggiuntivo richiesto per il blocco (*lock*) delle sezioni critiche e/o per inviare i dati da un processore all'altro. Si supponga che un programma richieda un tempo di esecuzione,  $t = 100$  s su un processore. Quando lo stesso programma viene eseguito su  $p$  processori, ciascun processore richiede un tempo pari a  $t/p$  secondi, più un tempo aggiuntivo di 4 s indipendente dal numero di processori. Calcolare il tempo di esecuzione per il processore in architetture a 2, 4, 8, 16, 32, 64 e 128 processori. Riportare lo speedup rispetto all'architettura monoprocessoresso per ciascuna di queste architetture e il rapporto tra lo speedup ottenuto e quello ideale, ottenuto quando non si ha tempo aggiuntivo.

## Risposte alle domande di autovalutazione

**Paragrafo 1.1, pagina 8** – Domande generali: sono accettabili molte risposte.

**Paragrafo 1.4, pagina 21** – *Memoria DRAM*: volatile, con tempo di accesso ridotto compreso tra i 50 e i 70 nanosecondi, e un costo per GB di 5-10 dollari. *Memoria su disco*: non volatile, tempi di accesso da 100 000 a 400 000 volte maggiori di quelli della DRAM e un costo per GB 100 volte inferiore a quello delle DRAM. *Memoria flash*: non volatili, tempi di accesso da 100 a 1000 volte maggiori delle DRAM e costi per GB da 7 a 10 volte inferiori delle DRAM.

**Paragrafo 1.5, pagina 24** – Le risposte 1, 3 e 4 sono motivi validi. La risposta 5 può essere vera in generale perché con grandi volumi sarebbe una buona decisione prevedere investimenti aggiuntivi per ridurre la dimensione delle piastrine, diciamo del 10%, ma ciò non è automaticamente vero.

**Paragrafo 1.6, pagina 29** – a. Entrambe. b. Latenza. c. Nessuna delle due. 7 secondi.

**Paragrafo 1.6, pagina 35** – b.

**Paragrafo 1.10, pagina 46** – a. Il calcolatore A ha il valore più alto di MIPS. b. Il calcolatore B è più veloce.

# 2

## Le istruzioni: il linguaggio dei calcolatori

*Parlo spagnolo con Dio, italiano con le donne, francese con gli uomini e tedesco col mio cavallo.*

*Carlo V, imperatore del Sacro Romano Impero, 1500-1558*

### 2.1 Introduzione

Per impartire comandi all'hardware di un calcolatore occorre saper parlare il suo linguaggio. Le parole del linguaggio del calcolatore sono dette istruzioni e l'intero vocabolario viene chiamato **insieme delle istruzioni** (*instruction set*). In questo capitolo verrà descritto l'insieme delle istruzioni di un calcolatore reale, sia nella forma scritta dall'uomo sia in quella interpretata dal calcolatore. Le istruzioni verranno presentate "dall'alto al basso": partendo dalla notazione tipica di un linguaggio di programmazione ad alto livello, anche se limitato, le trasformeremo passo dopo passo fino a rivelare il vero linguaggio di un calcolatore reale.

Nel Capitolo 3 il nostro viaggio continuerà verso il basso per scoprire come è fatto l'hardware utilizzato per i calcoli e come vengono rappresentati i numeri in virgola mobile.

Si potrebbe pensare che i linguaggi dei calcolatori siano molti e variegati come quelli degli uomini. In realtà sono abbastanza simili tra loro: sono più affini a dialetti regionali derivanti da un'unica radice linguistica che a vere e proprie lingue differenti. Perciò, una volta appreso un linguaggio, è molto semplice imparare anche gli altri.

L'insieme di istruzioni prescelto è il RISC-V, che è stato inizialmente sviluppato all'Università di Berkeley a partire dal 2010.

Per mostrare come sia facile imparare anche altri insiemi di istruzioni, illustreremo anche altri due insiemi di istruzioni molto diffusi.

**Insieme delle istruzioni:** il vocabolario dei comandi compresi da una data architettura.

- Il MIPS è un esempio elegante di insieme di istruzioni progettato a partire dal 1980. Sotto molti aspetti, il RISC-V è molto simile.
- L'Intel x86 è nato nel 1970, ma alimenta ancora oggi sia i PC che i calcolatori dei cloud dell'era post-PC.

Questa somiglianza tra insiemi di istruzioni è dovuta al fatto che tutti i calcolatori sono costruiti a partire da tecnologie hardware basate sugli stessi principi fondamentali e al fatto che esistono alcune operazioni di base che ogni calcolatore deve poter eseguire. Inoltre, i progettisti dei calcolatori hanno un obiettivo comune: trovare un linguaggio che renda semplice costruire l'hardware e realizzare i compilatori e, allo stesso tempo, consenta di massimizzare le prestazioni e minimizzare costi ed energia assorbita. Quest'obiettivo è rimasto valido nel tempo; la seguente citazione fu scritta ben prima che voi poteste acquistare un calcolatore ed è vera oggi come lo era nel 1947:

*È facile verificare tramite la logica formale che esistono certi [insiemi di istruzioni] che in linea di principio si prestano bene a controllare l'hardware per eseguire una qualsiasi sequenza di operazioni [...] Le considerazioni davvero decisive, dal punto di vista attuale, per la scelta di un insieme di istruzioni sono di natura molto più pratica: la semplicità dei dispositivi richiesti [dall'insieme delle istruzioni], la chiarezza della sua applicazione alla soluzione di problemi realmente importanti e la velocità con cui questi problemi vengono gestiti.*

Burks, Goldstine e von Neuman, 1947

La “semplicità dei dispositivi” viene considerata un parametro di grande valore per i calcolatori moderni come lo era per quelli degli anni '50. L'obiettivo del presente capitolo è mostrare un insieme di istruzioni progettato seguendo queste indicazioni: vedremo come viene rappresentato nell'hardware e qual è la relazione tra i linguaggi di programmazione ad alto livello e questo linguaggio più primitivo. Gli esempi che mostreremo nel corso del capitolo sono realizzati utilizzando il C come linguaggio di programmazione ad alto livello; nel paragrafo 2.15 mostreremo gli stessi esempi utilizzando un linguaggio di programmazione orientato agli oggetti come Java.

Imparando il modo in cui sono rappresentate le istruzioni, potrete anche scoprire il segreto dell'informatica: il **concetto di programma memorizzato**. Inoltre, potrete fare pratica di “lingua straniera” scrivendo programmi nel linguaggio del calcolatore, per farli poi eseguire al simulatore che troverete sul sito web del libro. Potrete anche osservare l'impatto dei linguaggi di programmazione e delle ottimizzazioni del compilatore sulle prestazioni. Il capitolo si conclude con uno sguardo all'evoluzione storica degli insiemi di istruzioni, con un accenno ad altri dialetti usati dai calcolatori. L'insieme di istruzioni RISC-V verrà presentato in maniera graduale, un sottoinsieme alla volta, spiegando ogni volta la motivazione e le strutture del calcolatore interessate. Questo approccio dall'alto al basso intreccia passo dopo passo la descrizione dei componenti con la spiegazione del loro funzionamento, rendendo il linguaggio assembler più “digeribile”. In Figura 2.1 viene anticipato l'insieme delle istruzioni presentate in questo capitolo.

### Concetto di programma

**memorizzato:** l'idea che le istruzioni e i dati di diverso tipo possano essere memorizzati come numeri; questa idea porta al calcolatore con programma memorizzato.

Devono certamente essere previste istruzioni per il calcolo delle operazioni aritmetiche fondamentali.

Burks, Goldstine e von Neumann, 1947

## 2.2 Operazioni svolte dall'hardware del calcolatore

Qualsiasi calcolatore deve saper eseguire le operazioni aritmetiche. La notazione del linguaggio assembler RISC-V:

`add a, b, c`

indica al calcolatore di sommare le due variabili `b` e `c` e di porre il risultato della somma nella variabile `a`.

Questa è una notazione rigida, nel senso che ciascuna istruzione aritmetica RISC-V esegue solo un'operazione e deve contenere esattamente tre variabili. Si supponga, per esempio, di volere memorizzare in  $x_0$  la somma delle variabili  $b$ ,  $c$ ,  $d$  ed  $e$  (in questo paragrafo il termine "variabile" sarà lasciato deliberatamente vago, rimandando al paragrafo successivo una spiegazione più precisa).

### Operandi RISC-V

Nome	Esempio	Commenti
32 registri	$x_0 \leftarrow x_1$	Accesso veloce ai dati. Nel RISC-V gli operandi devono essere contenuti nei registri per potere eseguire delle operazioni. Il registro $x_0$ contiene sempre il valore 0
$2^{31}$ parole di memoria	Memoria[0], Memoria[8], ... Memoria[18 446 744 073 709 551 608]	Alla memoria si accede solamente attraverso istruzioni di trasferimento dati. Il RISC-V utilizza l'indirizzamento al byte, perciò due variabili ampie due parole (double word) hanno indirizzi in memoria a distanza 8. La memoria consente di memorizzare strutture dati, vettori, o il contenuto dei registri

### Linguaggio assembler RISC-V

Tipo di istruzioni	Istruzioni	Esempio	Significato	Commenti
Aritmetiche	Somma	add x5, x6, x7	$x_5 = x_6 + x_7$	Operandi in tre registri
	Sottrazione	sub x5, x6, x7	$x_5 = x_6 - x_7$	Operandi in tre registri
	Somma immediata	addi x5, x6, 20	$x_5 = x_6 + 20$	Utilizzata per sommare delle costanti
Trasferimento dati	Lettura parola doppia	ld x5, 40(x6)	$x_5 = \text{Memoria}[x_6 + 40]$	Spostamento di una parola doppia da memoria a registro
	Memorizzazione parola doppia	sd x5, 40(x6)	$\text{Memoria}[x_6 + 40] = x_5$	Spostamento di una parola doppia da registro a memoria
	Lettura parola	lw x5, 40(x6)	$x_5 = \text{Memoria}[x_6 + 40]$	Spostamento di una parola da memoria a registro
	Lettura parola senza segno	lwu x5, 40(x6)	$x_5 = \text{Memoria}[x_6+40]$	Spostamento di una parola senza segno da memoria a registro
	Memorizzazione parola	sw x5, 40(x6)	$\text{Memoria}[x_6+40] = x_5$	Spostamento di una parola da registro a memoria
	Lettura mezza parola	lh x5, 40(x6)	$x_5 = \text{Memoria}[x_6+40]$	Spostamento di una mezza parola da memoria a registro
	Lettura mezza parola, senza segno	lhu x5, 40(x6)	$x_5 = \text{Memoria}[x_6+40]$	Spostamento di una mezza parola senza segno da memoria a registro
	Memorizzazione mezza parola	sh x5, 40(x6)	$\text{Memoria}[x_6+40] = x_5$	Spostamento di una mezza parola da registro a memoria
	Lettura byte	lb x5, 40(x6)	$x_5 = \text{Memoria}[x_6+40]$	Spostamento di un byte da memoria a registro
	Lettura byte, senza segno	lbu x5, 40(x6)	$x_5 = \text{Memoria}[x_6+40]$	Spostamento di un byte senza segno da memoria a registro
Operazioni atomiche	Memorizzazione byte	sb x5, 40(x6)	$\text{Memoria}[x_6+40] = x_5$	Spostamento di un byte da registro a memoria
	Lettura di una parola e blocco	lrd x5, (x6)	$x_5 = \text{Memoria}[x_6]$	Caricamento di una parola come prima fase di un'operazione atomica sulla memoria
	Memorizzazione condizionata di una parola	sc.d x7, x5, (x6)	$\text{Memoria}[x_6] = x_5; x_7 = 0/1$	Memorizzazione di una parola come seconda fase di un'operazione atomica sulla memoria
	Caricamento costante nella mezza parola superiore	lui x5, 0x12345	$x_5 = 0x12345000$	Caricamento di una costante su 20 bit nei 12 bit più significativi di una parola

## Risposte alle domande di autovalutazione

**Paragrafo 1.1, pagina 8** – Domande generali: sono accettabili molte risposte.

**Paragrafo 1.4, pagina 21** – *Memoria DRAM*: volatile, con tempo di accesso ridotto compreso tra i 50 e i 70 nanosecondi, e un costo per GB di 5-10 dollari. *Memoria su disco*: non volatile, tempi di accesso da 100 000 a 400 000 volte maggiori di quelli della DRAM e un costo per GB 100 volte inferiore a quello delle DRAM. *Memoria flash*: non volatili, tempi di accesso da 100 a 1000 volte maggiori delle DRAM e costi per GB da 7 a 10 volte inferiori delle DRAM.

**Paragrafo 1.5, pagina 24** – Le risposte 1, 3 e 4 sono motivi validi. La risposta 5 può essere vera in generale perché con grandi volumi sarebbe una buona decisione prevedere investimenti aggiuntivi per ridurre la dimensione delle piastrine, diciamo del 10%, ma ciò non è automaticamente vero.

**Paragrafo 1.6, pagina 29** – a. Entrambe. b. Latenza. c. Nessuna delle due. 7 secondi.

**Paragrafo 1.6, pagina 35** – b.

**Paragrafo 1.10, pagina 46** – a. Il calcolatore A ha il valore più alto di MIPS. b. Il calcolatore B è più veloce.

# 2

## Le istruzioni: il linguaggio dei calcolatori

*Parlo spagnolo con Dio, italiano con le donne, francese con gli uomini e tedesco col mio cavallo.*

*Carlo V, imperatore del Sacro Romano Impero, 1500-1558*

### 2.1 | Introduzione

Per impartire comandi all'hardware di un calcolatore occorre saper parlare il suo linguaggio. Le parole del linguaggio del calcolatore sono dette istruzioni e l'intero vocabolario viene chiamato **insieme delle istruzioni** (*instruction set*). In questo capitolo verrà descritto l'insieme delle istruzioni di un calcolatore reale, sia nella forma scritta dall'uomo sia in quella interpretata dal calcolatore. Le istruzioni verranno presentate "dall'alto al basso": partendo dalla notazione tipica di un linguaggio di programmazione ad alto livello, anche se limitato, le trasformeremo passo dopo passo fino a rivelare il vero linguaggio di un calcolatore reale.

Nel Capitolo 3 il nostro viaggio continuerà verso il basso per scoprire come è fatto l'hardware utilizzato per i calcoli e come vengono rappresentati i numeri in virgola mobile.

Si potrebbe pensare che i linguaggi dei calcolatori siano molti e variegati come quelli degli uomini. In realtà sono abbastanza simili tra loro: sono più affini a dialetti regionali derivanti da un'unica radice linguistica che a vere e proprie lingue differenti. Perciò, una volta appreso un linguaggio, è molto semplice imparare anche gli altri.

L'insieme di istruzioni prescelto è il RISC-V, che è stato inizialmente sviluppato all'Università di Berkeley a partire dal 2010.

Per mostrare come sia facile imparare anche altri insiemi di istruzioni, illustreremo anche altri due insiemi di istruzioni molto diffusi.

**Insieme delle istruzioni:** il vocabolario dei comandi compresi da una data architettura.

Tipo di istruzioni	Istruzioni	Esempio	Significato	Commenti
Logiche	And	and x5, x6, x7	x5 = x6 & x7	Operandi in tre registri; AND bit a bit
	Or inclusivo	or x5, x6, x8	x5 = x6   x8	Operandi in tre registri; OR bit a bit
	Or esclusivo (Xor)	xor x5, x6, x9	x5 = x6 ^ x9	Operandi in tre registri; XOR bit a bit
	And immediato	andi x5, x6, 20	x5 = x6 & 20	AND bit a bit tra un operando in registro e una costante
	Or immediato	ori x5, x6, 20	x5 = x6   20	OR bit a bit tra un operando in registro e una costante
	Xor immediato	xori x5, x6, 20	x5 = x6 ^ 20	XOR bit a bit tra un operando in registro e una costante
Scorrimento (shift)	Scorrimento logico a sinistra	sll x5, x6, x7	x5 = x6 << x7	Scorrimento a sinistra mediante registro
	Scorrimento logico a destra	srl x5, x6, x7	x5 = x6 >> x7	Scorrimento a destra mediante registro
	Scorrimento aritmetico a destra	sra x5, x6, x7	x5 = x6 >> x7	Scorrimento a destra aritmetico mediante registro
	Scorrimento logico a sinistra immediato	slli x5, x6, 3	x5 = x6 << 3	Scorrimento a sinistra mediante costante
	Scorrimento logico a destra immediato	srli x5, x6, 3	x5 = x6 >> 3	Scorrimento a destra mediante costante
	Scorrimento aritmetico a destra immediato	srai x5, x6, 3	x5 = x6 >> 3	Scorrimento a destra aritmetico mediante costante
Salti condizionati	Salta se uguale	beq x5, x6, 100	Se (x5 == x6) vai a PC+100	Test di uguaglianza; salto relativo al PC
	Salta se non è uguale	bne x5, x6, 100	Se (x5 != x6) vai a PC+100	Test di diseguaglianza; salto relativo al PC
	Salta se minore di	blt x5, x6, 100	Se (x5 < x6) vai a PC+100	Comparazione di minoranza; salto relativo al PC
	Salta se maggiore o uguale di	bge x5, x6, 100	Se (x5 >= x6) vai a PC+100	Comparazione di maggioranza o uguaglianza; salto relativo al PC
	Salta se minore di senza segno	bltu x5, x6, 100	Se (x5 < x6) vai a PC+100	Comparazione di minoranza senza segno; salto relativo al PC
	Salta se maggiore o uguale di senza segno	bgeu x5, x6, 100	Se (x5 >= x6) vai a PC+100	Comparazione di maggioranza o uguaglianza senza segno; salto relativo al PC
Salti incondizionati	Salta e collega	jal x1, 100	x1 = PC+4; vai a PC+100	Chiamata a procedura con indirizzamento relativo al PC
	Salta e collega mediante registro	jalr x1, 100(x5)	x1 = PC+4; vai a x5+100	Ritorno da procedura; chiamata indiretta

**Figura 2.1** Il linguaggio assembler del RISC-V esaminato in questo capitolo. Si possono trovare queste informazioni anche nella prima colonna della scheda tecnica riassuntiva del RISC-V in fondo al libro.

La seguente sequenza di istruzioni effettuerà la somma delle quattro variabili:

```
add a,b,c // la somma di b e c è posta in a
add a,a,d // la somma di b, c e d è ora posta in a
add a,a,e // la somma di b, c, d ed e è ora posta in a
```

Occorrono quindi tre istruzioni per sommare le quattro variabili.

Su ciascuna linea di codice, le parole a destra del simbolo // rappresentano i commenti rivolti a chi legge il programma e sono ignorate dal calcolatore. Si noti

che, a differenza di altri linguaggi di programmazione, in linguaggio Assembler ciascuna linea può contenere al massimo un'istruzione. Un'altra differenza rispetto al C è che i commenti terminano sempre con la fine della linea.

Il numero di operandi più naturale per un'operazione come la somma è pari a tre: i due numeri da sommare e il riferimento alla locazione in cui memorizzare il risultato. Il fatto di richiedere che tutte le istruzioni abbiano esattamente tre operandi, né uno di più né uno di meno, è conforme alla filosofia di mantenere l'hardware semplice: l'hardware richiesto per un numero variabile di operandi sarebbe più complesso di quello necessario per un numero fisso di operandi. Questa situazione illustra il primo dei tre principi fondamentali per la progettazione dell'hardware:

*Principio di progettazione numero 1: la semplicità favorisce la regolarità.*

Nei due esempi che seguono verrà mostrata la relazione tra programmi scritti in linguaggi ad alto livello e programmi scritti in questa notazione più primitiva.

### Compilazione di due assegnamenti in C e loro traduzione in istruzioni RISC-V

Questo frammento di programma C contiene cinque variabili a, b, c, d ed e. Dato che Java è nato dal C, questo esempio e i successivi saranno validi per entrambi questi linguaggi di programmazione ad alto livello:

```
a = b + c;
d = a - e;
```

La traduzione da C ad assembler RISC-V viene effettuata dal *compilatore*. Mostrare il codice RISC-V prodotto dal compilatore.

Un'istruzione RISC-V lavora su due operandi sorgente e inserisce il risultato in un operando destinazione. Il risultato della fase di compilazione delle due istruzioni mostrate sopra è quindi dato semplicemente da queste due istruzioni assembler RISC-V:

```
add a,b,c
sub d,a,e
```

ESEMPIO

SOLUZIONE

### Compilazione di un assegnamento complesso in C e sua traduzione in istruzioni RISC-V

Si consideri il seguente assegnamento complesso contenente le cinque variabili f, g, h, i, j:

$$f = (g + h) - (i + j);$$

Qual è il codice ottenuto in seguito alla compilazione?

Il compilatore deve spezzare questo assegnamento C in più istruzioni assembler, poiché si può effettuare una sola operazione per ogni istruzione RISC-V. La prima istruzione RISC-V calcolerà la somma di g e h. Dovendo memorizzare il risultato di questa operazione da qualche parte,

ESEMPIO

SOLUZIONE

(continua)

(continua)

il compilatore creerà una variabile temporanea chiamata t0:

```
add t0,g,h // la variabile temporanea t0 contiene g + h
```

Nonostante l'operazione successiva sia una sottrazione, dobbiamo, per motivi di precedenza, effettuare prima la somma tra i e j. Quindi, la seconda istruzione sarà composta dalla somma di i e j e verrà creata una nuova variabile temporanea t1 per memorizzare il risultato:

```
add t1,i,j // la variabile temporanea t1 contiene i + j
```

Alla fine, verrà effettuata la sottrazione tra la prima e la seconda somma e verrà scritto il risultato nella variabile f, completando così la traduzione del codice C:

```
sub f,t0,t1 // f assume il valore t0 - t1,  
// cioè (g + h) - (i + j)
```

## Autovalutazione

Per descrivere una data funzione, quale linguaggio di programmazione richiederà verosimilmente più linee di codice? Si ordinino i seguenti tre linguaggi utilizzando questa metrica.

1. Java
2. C
3. Linguaggio assembler RISC-V

**Approfondimento.** Per aumentarne la portabilità, Java fu originariamente concepito per essere “interpretato” da un software. L’insieme delle istruzioni di questo interprete fu chiamato *bytecode Java* (vedi par. 2.15), che è molto differente dall’insieme delle istruzioni del RISC-V. Per avvicinare le prestazioni a quelle di un equivalente programma C, i sistemi Java tipicamente compilano il bytecode Java, traducendolo nell’insieme delle istruzioni native del processore, per esempio in istruzioni RISC-V. Dato che i compilatori Java entrano in azione molto più tardi rispetto ai compilatori C, i compilatori Java vengono spesso chiamati compilatori *Just In Time (JIT)*. Il paragrafo 2.12 illustrerà che cosa si intende affermando che i compilatori di tipo JIT vengono utilizzati più tardi rispetto a quelli C nell'avviamento di un'applicazione; il paragrafo 2.13 mostrerà l'impatto sulle prestazioni della compilazione rispetto all'interpretazione di un programma Java.

## 2.3 Operandi dell'hardware del calcolatore

**Parola doppia:** un gruppo di bit di diversa numerosità a cui si accede naturalmente in un calcolatore; tipicamente è un gruppo di 64 bit. Nell'architettura RISC-V corrisponde alla dimensione dei registri. Nell'architettura MIPS corrisponde alla dimensione dei registri.

**Parola:** il numero di bit a cui si accede più naturalmente in un calcolatore; tipicamente è un gruppo di 32 bit.

A differenza dei programmi scritti nei linguaggi ad alto livello, gli operandi delle istruzioni aritmetiche del RISC-V devono obbedire ad alcune restrizioni: devono essere scelti tra un numero limitato di locazioni particolari, chiamate *registri*. I registri sono i mattoni che costituiscono il calcolatore, dal momento che rappresentano sia le primitive utilizzate nella progettazione dell'hardware sia gli elementi visibili al programmatore. La dimensione dei registri nell'architettura RISC-V è di 64 bit; gruppi di 64 bit sono così frequenti che, nelle architetture RISC-V, prendono il nome di **parola doppia (doubleword)**. Un altro gruppo di bit molto utilizzato nell'architettura RISC-V è la **parola (word)**, costituita da 32 bit.

Una delle differenze più importanti fra le variabili utilizzate nei linguaggi di programmazione e i registri è il numero limitato di questi ultimi, tipicamente 32 nei calcolatori della classe dei RISC-V (vedi par. 2.21 per una panoramica storica sul numero dei registri). In questo paragrafo, compiendo un altro passo verso la rappresentazione in linguaggio macchina delle istruzioni, introduciamo il vincolo addizionale per cui i tre operandi delle istruzioni aritmetiche RISC-V devono essere scelti tra i 32 registri a 64 bit.

La ragione della limitazione a 32 del numero di registri può essere trovata nel secondo dei tre principi fondamentali per la progettazione dell'hardware:

*Principio di progettazione numero 2: minori sono le dimensioni, maggiore è la velocità.*

Un numero di registri molto elevato potrebbe aumentare la durata del ciclo di clock semplicemente perché i segnali elettrici impiegherebbero un tempo maggiore a compiere il percorso assegnato, dovendo coprire una distanza maggiore.

Principi guida come “minorì sono le dimensioni, maggiore è la velocità” non sono validi in modo assoluto: un calcolatore con 31 registri potrebbe non essere più veloce di uno con 32. Tuttavia, tali osservazioni sono basate sulla realtà, tanto che i progettisti le prendono sul serio. Nel caso specifico, il progettista deve trovare un compromesso fra la richiesta da parte dei programmi di avere più registri a disposizione e il suo desiderio di mantenere veloce il clock. Un'altra ragione per non utilizzare più di 32 registri è il numero di bit che sarebbe necessario per indirizzarli all'interno dell'istruzione, come si vedrà nel paragrafo 2.5.

Il Capitolo 4 mostra il ruolo centrale dei registri nella costruzione dell'hardware; come si potrà vedere, un utilizzo efficiente dei registri è la chiave per ottenere buone prestazioni da un programma.

Anche se potremmo benissimo indicare ogni registro con un numero compreso tra 0 e 31, la convenzione RISC-V prevede l'utilizzo di una  $x$  seguita dal numero del registro, ad eccezione di alcuni registri che hanno un nome personalizzato e che introdurremo più avanti.

## Compilazione di un assegnamento in C utilizzando i registri

Associare i registri alle variabili di un programma è compito del compilatore. Si prenda come esempio l'assegnamento esaminato in precedenza:

$$f = (g + h) - (i + j);$$

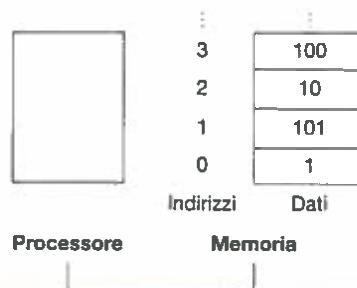
Le variabili  $f$ ,  $g$ ,  $h$ ,  $i$ ,  $j$  vengono assegnate ai registri  $x19$ ,  $x20$ ,  $x21$ ,  $x22$  e  $x23$ . Quale sarà il codice RISC-V ottenuto in seguito alla compilazione?

Il programma compilato è molto simile a quello dell'esempio precedente, ad eccezione del fatto che le variabili sono state sostituite dai nomi dei registri precedentemente definiti e che sono comparsi due registri temporanei,  $x5$  e  $x6$ , che corrispondono alle variabili temporanee:

```
add x5,x20,x21 // il registro temporaneo x5
                  // contiene il risultato g + h
add x6,x22,x23 // il registro temporaneo x6
                  // contiene il risultato i + j
sub x19,x5,x6   // f assume il valore x5 - x6,
                  // cioè (g + h) - (i + j)
```

### ESEMPIO

### SOLUZIONE



**Figura 2.2 Indirizzi di memoria e contenuto della memoria nelle posizioni corrispondenti.** Se i dati memorizzati fossero delle parole doppie, i loro indirizzi sarebbero errati, dato che il RISC-V utilizza l'indirizzamento della memoria al byte e ogni parola doppia contiene 8 byte. La Figura 2.3 mostra come vengano indirizzate in memoria parole consecutive.

### Operandi allocati in memoria

I linguaggi di programmazione possono utilizzare variabili semplici che, come nel caso appena esposto, sono composte da un singolo elemento, ma possono utilizzare anche strutture di dati più complesse, come i vettori e le strutture dati. Queste strutture complesse possono contenere un numero di elementi molto maggiore del numero di registri presenti nel calcolatore: come può un calcolatore rappresentare e accedere a strutture dati così grandi?

I cinque componenti del calcolatore, introdotti nel Capitolo 1, suggeriscono che, benché il processore possa contenere un limitato numero di dati nei registri, la memoria può contenere miliardi di dati. Di conseguenza le strutture dati (vettori e strutture) vengono allocate in memoria.

Come detto precedentemente, nel linguaggio assembler RISC-V le istruzioni aritmetiche richiedono che gli operandi siano memorizzati nei registri; l'assembler RISC-V deve quindi contenere delle **istruzioni di trasferimento dati** che trasferiscano dati fra la memoria e i registri. Per accedere a una parola singola o doppia in memoria, l'istruzione deve fornire l'**indirizzo** di memoria corrispondente. La memoria può essere vista come un grande vettore monodimensionale, con l'indirizzo che funge da indice e parte a contare da zero. Per esempio, in Figura 2.2, l'indirizzo del terzo dato è 2 e il valore di Memoria[2] è 10.

L'istruzione di trasferimento che sposta un dato dalla memoria a un registro viene chiamata tradizionalmente *load* (carica). Il formato dell'istruzione load prevede il nome dell'operazione seguito dal registro in cui deve essere trasferito il dato e da una coppia costante-registro utilizzata per accedere alla memoria. L'indirizzo del dato in memoria viene ottenuto dalla somma della costante e del contenuto del secondo registro. Il nome convenzionale dell'istruzione di load nel linguaggio RISC-V è **ld**, che sta per *load doubleword* (carica una parola doppia).

### Compilazione di un'istruzione di assegnamento quando uno degli operandi risiede in memoria

#### ESEMPIO

Si supponga che A sia un vettore di 100 parole e che il compilatore associa le variabili g e h ai registri **x20** e **x21**, come nel caso precedente. Si supponga inoltre che l'indirizzo di partenza del vettore, detto *indirizzo base*, sia memorizzato in **x22**. Si traduca in assembler il seguente assegnamento scritto in C:

**g = h + A[8];**

(continua)

Benché l'assegnamento in C sia costituito da una sola operazione, uno degli operandi si trova in memoria e quindi, per prima cosa, occorre trasferire il contenuto di A[8] in un registro. L'indirizzo di questo elemento del vettore è dato dalla somma dell'indirizzo base del vettore A, che si trova nel registro x22, e del numero che permette di selezionare l'elemento 8. Il dato letto dovrà essere memorizzato in un registro temporaneo e potrà essere utilizzato dall'istruzione successiva. Osservando la Figura 2.2, la prima istruzione in linguaggio assembler sarà:

```
ld x9, 8(x22) // il valore di A[8] viene caricato
                // nel registro temporaneo x9
```

(Nel prossimo esempio faremo una piccola modifica a questa istruzione, ma per ora utilizziamo questa versione semplificata.) L'istruzione successiva potrà utilizzare il valore memorizzato in x9 (che è uguale ad A[8]), dal momento che esso è ora contenuto in un registro. L'istruzione dovrà quindi sommare h (contenuto in x21) ad A[8] (contenuto in x9) e scrivere il risultato nel registro assegnato a g (x20):

```
add x20, x21, x9 // g = h + A[8]
```

Il registro utilizzato per costruire l'indirizzo (x22) è detto *registro base* e la costante in un'istruzione di trasferimento (8 in questo caso) si chiama *offset* (spiazzamento).

(continua)

**SOLUZIONE**

Oltre ad associare le variabili ai registri, il compilatore alloca strutture dati quali i vettori e le strutture dati più complesse in locazioni di memoria, e pone quindi il loro indirizzo di partenza corretto nelle istruzioni di trasferimento dati.

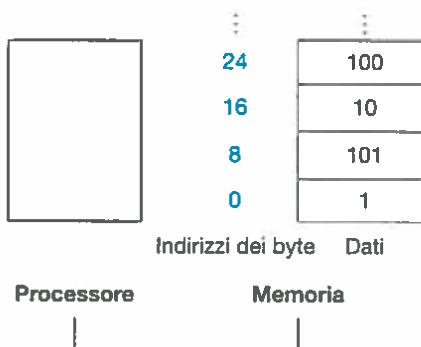
Dal momento che i byte sono utilizzati in molti programmi, virtualmente tutte le architetture indirizzano il singolo byte; di conseguenza, l'indirizzo di una parola doppia corrisponde all'indirizzo di uno degli 8 byte che la compongono e l'indirizzo di due parole doppie consecutive differisce di 8 unità. La Figura 2.3, per esempio, mostra i veri indirizzi RISC-V per le parole della Figura 2.2: l'indirizzo in byte della terza parola è ora 16.

I calcolatori con indirizzamento relativo al singolo byte si dividono fra quelli che utilizzano il byte più significativo, quello più a sinistra (*big end*), per specificare l'indirizzo della parola doppia, e quelli che utilizzano il byte meno significativo, quello più a destra (*little end*). Il RISC-V fa parte della seconda categoria detta *Little Endian*. Dato che l'ordinamento importa solo se si legge o scrive una parola byte per byte, a pochi interessa se la codifica sia big o little endian.

L'indirizzamento al byte influenza anche il calcolo degli indici dei vettori: per ottenere l'indirizzo in byte corretto per il frammento di codice precedente, l'offset da sommare al registro di base, x22, deve essere uguale a  $8 \times 8$ , cioè 64, in modo che l'istruzione load possa selezionare A[8] e non A[8/8]. (Per i trabocchetti che si possono generare *vedi* anche il paragrafo 2.19.)

L'istruzione complementare a quella di load è chiamata tradizionalmente *store* (memorizza) e trasferisce un dato da un registro alla memoria. Il formato di questa istruzione è simile a quello dell'istruzione load: il nome dell'istruzione seguito dal registro che contiene il dato da trasferire in memoria, dall'offset corrispondente all'elemento del vettore e dal registro base. Come nel caso precedente, l'indirizzo di memoria viene specificato in parte da una costante e in parte dal contenuto di un registro. Il nome convenzionale dell'istruzione store in linguaggio RISC-V è sd, che sta per *store doubleword* (memorizza una parola doppia).

## Interfaccia hardware/software



**Figura 2.3** Indirizzi di memoria effettivi nel RISC-V e contenuto delle locazioni corrispondenti. Gli indirizzi evidenziati sono quelli che cambiano rispetto alla Figura 2.2. Dato che il RISC-V indirizza il singolo byte, gli indirizzi delle parole doppie sono multipli di 8: ci sono 8 byte in una parola doppia.

**Vincolo di allineamento:** requisito per cui, in molte architetture, le parole devono iniziare sempre a indirizzi multipli di 4 e le parole doppie a indirizzi multipli di 8.

**Approfondimento.** In molte architetture, le parole devono iniziare sempre a indirizzi multipli di 4 e le parole doppie a indirizzi multipli di 8. Questo requisito si chiama **vincolo di allineamento** ed è comune a molte altre architetture (nel Capitolo 4 viene suggerito come l'allineamento produca una velocità maggiore nel trasferimento dei dati). I RISC-V e gli Intel x86 non hanno il vincolo di allineamento, mentre i MIPS sì.

## Interfaccia hardware/software

Dal momento che gli indirizzi delle istruzioni load e store sono numeri binari, si può capire perché gli indirizzi della memoria principale DRAM siano espressi in multipli di 2 invece che di 10. Vengono cioè espressi in Gibibyte ( $2^{30}$ ), o in Tebibyte ( $2^{40}$ ), e non in Gigabyte ( $10^9$ ) o Terabyte ( $10^{12}$ ). Vedi anche la Figura 1.1.

### Compilazione utilizzando load e store

#### ESEMPIO

Si supponga che la variabile `h` sia associata al registro `x21` e che il registro base del vettore `A` sia `x22`. Quale sarà il codice assembler RISC-V relativo al seguente assegnamento in linguaggio C?

`A[12] = h + A[8];`

#### SOLUZIONE

Benché l'assegnamento sia costituito in C da una sola operazione, adesso gli operandi che si trovano in memoria sono due e quindi occorrono ancora più istruzioni RISC-V che nel caso precedente. Le prime due istruzioni sono le stesse dell'esempio precedente, tranne che per l'utilizzo dell'offset corretto per realizzare l'indirizzamento al byte nell'istruzione di caricamento in registro di una parola doppia (ld) che seleziona `A[8]`. L'istruzione di somma add inserirà il risultato in `x9`:

```
ld x9, 64(x22) // il registro temporaneo x9
                  // assume il valore A[8]
add x9, x21, x9 // il registro temporaneo x9
                  // assume il valore h + A[8]
```

L'ultima istruzione memorizza la somma in `A[12]`, utilizzando 96 ( $8 \times 12$ ) come offset e `x22` come registro base:

```
sd x9, 96(x22) // memorizza h + A[8] in A[12]
```

Nelle architetture RISC-V le istruzioni load doubleword e store doubleword servono per trasferire una parola tra la memoria e i registri. Altre architetture utilizzano altre istruzioni, oltre a load e store, per trasferire dati; una di queste architetture è l'x86 Intel, che verrà descritta nel paragrafo 2.17.

Molti programmi contengono più variabili di quanti siano i registri del calcolatore. Di conseguenza il compilatore cerca di mantenere nei registri le variabili utilizzate più di frequente e mette le altre in memoria, utilizzando le istruzioni di trasferimento dati per spostare le variabili tra i registri e la memoria. Il processo di trasferimento in memoria delle variabili utilizzate meno di frequente (oppure di quelle che verranno utilizzate successivamente) si chiama *register spilling* (*versamento dei registri*).

Il già citato principio hardware che lega le dimensioni alla velocità suggerisce che l'accesso alla memoria sia più lento dell'accesso ai registri, dal momento che il numero dei registri è minore del numero di locazioni della memoria. Ed è effettivamente così: l'accesso ai dati risulta più veloce se i dati si trovano nei registri piuttosto che in memoria.

Inoltre i dati risultano più utili se si trovano nei registri. Le istruzioni aritmetiche del RISC-V nello stesso ciclo di clock possono leggere due registri, eseguire l'operazione e scrivere il risultato. Le istruzioni di trasferimento dati del RISC-V, invece, possono solamente leggere o scrivere un singolo operando, senza potere eseguire nessuna operazione su di esso.

I registri, quindi, richiedono un tempo di accesso minore della memoria e offrono un throughput considerevolmente più elevato, rendendo i dati in essi contenuti accessibili più velocemente e più semplici da utilizzare. Per ottenere le prestazioni migliori e ridurre il consumo di energia, un'architettura deve avere un sufficiente numero di registri e i compilatori devono farne uso in maniera efficiente.

## Interfaccia hardware/software

**Approfondimento.** Guardiamo in prospettiva l'energia e le prestazioni dei registri in funzione della memoria. Supponendo dati su 64 bit, i registri erano circa 200 volte più veloci (0,25 ns in confronto a 50 ns) e avevano un'efficienza energetica superiore di 10 000 volte (0,1 picoJoule in confronto a 1000 picoJoule) delle DRAM nel 2015. Questa grossa differenza ha portato allo sviluppo delle memorie cache, che riducono le penalità sulle prestazioni e sul consumo di energia quando si accede alla memoria (vedi Cap. 5).

### Operandi immediati o costanti

Spesso i programmi utilizzano all'interno di un'operazione una costante, per esempio per incrementare l'indice di un contatore in modo da puntare all'elemento successivo di un vettore. In più della metà delle operazioni aritmetiche RISC-V contenute nei programmi benchmark dello SPEC CPU2006, uno degli operandi è una costante.

Utilizzando solamente le istruzioni viste finora dovremmo caricare una costante dalla memoria per poterla utilizzare (la costante verrà scritta in memoria all'atto del caricamento del programma). Per esempio, per sommare 4 al registro x22 possiamo utilizzare le seguenti istruzioni:

```
ld x9, IndirizzoCostante4(x3) // x9 = costante 4
add x22, x22, x9           // x22 = x22 + x9
                           // (dove x9 == 4)
```

supponendo che x3 + IndirizzoCostante4 sia l'indirizzo di memoria della costante 4.

Un'alternativa che consente di evitare l'operazione di load è costituita da una versione delle istruzioni aritmetiche nella quale uno degli operandi è una costante. La versione dell'istruzione di somma in cui un operando è una costante è chiamata *add immediate* (somma immediata) o *addi*. Per sommare 4 al registro  $x_{22}$ , possiamo semplicemente scrivere:

```
addi x22, x22, 4      //  $x_{22} = x_{22} + 4$ 
```

Le operazioni su costanti sono molto frequenti (e l'istruzione *addi* è l'istruzione più utilizzata nella maggior parte dei programmi RISC-V). Inserendo le costanti all'interno delle istruzioni aritmetiche, le operazioni risultano molto più veloci e richiedono meno energia rispetto al caso in cui le costanti siano caricate dalla memoria.

La costante zero ha anche il ruolo di semplificare l'insieme delle istruzioni, consentendo di realizzare delle utili varianti. Per esempio, si può negare il contenuto di un registro utilizzando l'istruzione *sub* con zero come primo operando. Per questo motivo, i RISC-V dedicano il registro  $x_0$  a contenere il valore prefissato di zero. Giustificare l'introduzione di una costante attraverso la frequenza del suo utilizzo è un'ulteriore esempio della grande idea introdotta nel Capitolo 1 di rendere veloci le situazioni più comuni.



SITUAZIONI COMUNI



LEGGE DI MOORE

## Autovalutazione

Data l'importanza dei registri, qual è il tasso di crescita del numero dei registri su un singolo chip?

1. Molto veloce: il loro incremento segue la **Legge di Moore**, che predice il raddoppio del numero di transistor su un singolo chip ogni 18 mesi.
2. Molto lento: dato che i programmi sono tipicamente distribuiti ai clienti nella versione tradotta in linguaggio macchina, c'è una notevole inerzia nelle generazioni di nuove architetture di insiemi di istruzioni. Il numero di registri aumenta alla velocità con cui vengono resi disponibili nuovi insiemi di istruzioni.

**Approfondimento.** Nonostante i registri RISC-V considerati in questo libro abbiano un'ampiezza di 64 bit, gli architetti hanno pensato a diverse varianti dell'ISA. Oltre alla variante nota come RV64, la variante RV32 ha registri a 32 bit e il suo costo ridotto rende l'architettura RV32 più adatta ai processori a costo molto basso.

**Approfondimento.** Il sistema di indirizzamento dei RISC-V tramite offset e registro base è ottimo per accedere sia alle strutture dati sia ai vettori, dal momento che il registro può puntare all'inizio della struttura e l'offset può selezionare l'elemento desiderato. Troverete un esempio nel paragrafo 2.13.

**Approfondimento.** Nelle istruzioni per il trasferimento dati, i registri furono originariamente introdotti per memorizzare l'indice dell'elemento del vettore, mentre l'offset conteneva l'indirizzo di partenza del vettore stesso. Perciò il registro base è anche chiamato *registro indice*. Oggi le memorie sono molto più grandi e il modello software per l'allocazione dei dati è molto più sofisticato; l'indirizzo base di un vettore viene quindi solitamente memorizzato in un registro, dato che la dimensione del campo offset non sarebbe sufficiente a contenerlo, come si vedrà in seguito.

**Approfondimento.** La migrazione dai calcolatori da 32 bit a quelli da 64 bit ha lasciato ai programmati dei compilatori la scelta della dimensione dei tipi di dati in C. Chiaramente, i puntatori dovranno essere a 64 bit, ma che dimensione

dovranno avere gli interi? Inoltre, il C ha i tipi di dati `int`, `long int` e `long long int`. I problemi nascono quando si devono convertire i dati da un tipo all'altro e quando si verifica inaspettatamente un overflow in un programma in C non completamente ortodosso, cosa che sfortunatamente si verifica non raramente. La tabella seguente riporta due opzioni popolari:

Sistema operativo	Puntatori	<code>int</code>	<code>long int</code>	<code>long long int</code>
Windows Microsoft	64 bit	32 bit	32 bit	64 bit
Linux, la maggior parte di Unix	64 bit	32 bit	64 bit	64 bit

Anche se un compilatore ha a disposizione più scelte, di solito i compilatori associati a un dato sistema operativo operano la stessa scelta. Per fare un esempio semplice, in questo libro supporremo che i puntatori siano tutti su 64 bit e dichiareremo tutti gli interi C come `long long int` in modo tale che abbiano tutti la stessa dimensione. Seguiremo anche lo standard C99 e dichiareremo le variabili indici di vettori e matrici come `size_t`, che garantisce che abbiano la dimensione corretta indipendentemente dalla dimensione dei vettori e delle matrici. Tipicamente vengono dichiarate di dimensione uguale alle variabili `long int`.

## 2.4 Numeri con e senza segno

Per prima cosa ripassiamo velocemente la rappresentazione dei numeri nei calcolatori. Gli uomini imparano a contare in base 10, ma i numeri possono essere rappresentati in qualunque base; per esempio 123 in base 10 = 1111011 in base 2.

I numeri vengono rappresentati nell'hardware del calcolatore tramite una serie di segnali elettrici di livello alto o basso, e quindi la rappresentazione più naturale è quella in base 2 (come i numeri in base 10 sono chiamati numeri *decimali*, così i numeri in base 2 sono chiamati numeri *binari*).

Le cifre che costituiscono un numero binario costituiscono quindi le "unità atomiche" del calcolo, dato che tutte le informazioni in un calcolatore sono costituite da **cifre binarie** (*binary digits*) o *bit*. Le cifre binarie possono assumere due soli valori, che possono essere descritti in vario modo: alto o basso, acceso o spento, vero o falso, 1 o 0.

Qualunque sia la base, il valore della  $i$ -esima cifra  $d$  di un numero è pari a:

$$d \times \text{Base}^i$$

dove  $i$  assume un valore che parte da 0 e viene incrementato spostandosi da destra verso sinistra. Di conseguenza, un modo immediato per associare un numero ai bit all'interno di una parola è quello di utilizzare la potenza della base relativa a quel bit. Indicheremo i numeri decimali con il pedice *dec* (decimali) e i numeri binari con il pedice *due*. Per esempio,

$$1011_{\text{due}}$$

rappresenta:

$$\begin{aligned} & (1 \times 2^3) + (0 \times 2^2) + (1 \times 2^1) + (1 \times 2^0)_{\text{dec}} \\ &= (1 \times 8) + (0 \times 4) + (1 \times 2) + (1 \times 1)_{\text{dec}} \\ &= 8 + 0 + 2 + 1_{\text{dec}} \\ &= 11_{\text{dec}} \end{aligned}$$

**Cifra binaria o bit:** una delle due cifre binarie, 0 o 1, con cui vengono rappresentati i numeri in base 2, che sono il fondamento della rappresentazione dell'informazione.

Abbiamo assegnato un numero ai bit: 0, 1, 2, 3 ... analizzati da *destra verso sinistra* all'interno della parola doppia. Lo schema seguente illustra il valore assunto dai diversi bit all'interno di una parola doppia del RISC-V e la posizione che

assumerebbe il numero  $1011_{\text{due}}$  (sfortunatamente questa rappresentazione deve essere suddivisa su due righe per rientrare nella pagina del libro):

63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32			
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0			
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1

(ampiezza di 64 bit, suddiviso su due righe da 32 bit ciascuna)

**Bit meno significativo:** il bit più a destra in una parola RISC-V singola o doppia.

**Bit più significativo:** il bit più a sinistra in una parola RISC-V singola o doppia.

Dal momento che le parole vengono scritte a volte verticalmente e a volte orizzontalmente, i termini destra e sinistra possono risultare poco chiari. Di qui l'utilizzo dei termini **bit meno significativo** (*least significant bit*) per indicare il bit più a destra di una parola (il bit 0 nell'esempio precedente) e **bit più significativo** (*most significant bit*) per indicare il bit più a sinistra (il bit 63 per una parola doppia).

Poiché la generica parola doppia dei RISC-V è lunga 64 bit, possiamo rappresentare  $2^{64}$  diverse combinazioni di 64 bit ciascuna. È naturale utilizzare le combinazioni di questi 64 bit per rappresentare i numeri da 0 a  $2^{64} - 1$  (18 446 774 073 709 551 615<sub>dec</sub>):

00000000	00000000	00000000	00000000	00000000	00000000	00000000	00000000 <sub>due</sub>	= 0 <sub>dec</sub>
00000000	00000000	00000000	00000000	00000000	00000000	00000000	00000001 <sub>due</sub>	= 1 <sub>dec</sub>
00000000	00000000	00000000	00000000	00000000	00000000	00000000	00000010 <sub>due</sub>	= 2 <sub>dec</sub>
...	...	...	...	...	...	...	...	...
11111111	11111111	11111111	11111111	11111111	11111111	11111111	11111101 <sub>due</sub>	= 18 446 774 073 709 551 613 <sub>dec</sub>
11111111	11111111	11111111	11111111	11111111	11111111	11111111	11111110 <sub>due</sub>	= 18 446 774 073 709 551 614 <sub>dec</sub>
11111111	11111111	11111111	11111111	11111111	11111111	11111111	11111111 <sub>due</sub>	= 18 446 774 073 709 551 615 <sub>dec</sub>

Cioè, i numeri binari a 64 bit si possono rappresentare come somma di potenze di 2 (dove  $x_i$  sta per  $i$ -esimo bit di  $x$ ):

$$(x_{63} \times 2^{63}) + (x_{62} \times 2^{62}) + (x_{61} \times 2^{61}) + \dots + (x_1 \times 2^1) + (x_0 \times 2^0)$$

Per i motivi che vedremo fra poco, questi numeri sono chiamati numeri *unsigned* (senza segno).

## Interfaccia hardware/software

La base 2 non è naturale per gli esseri umani: avendo dieci dita, a noi risulta più naturale la base 10. Perché i calcolatori non utilizzano la numerazione in base 10? In realtà, i primi calcolatori commerciali utilizzavano l'aritmetica in base 10. Il problema è che un calcolatore utilizza sempre i segnali acceso e spento, per cui i numeri venivano comunque rappresentati da un insieme di cifre binarie. La codifica decimale si rivelò così inefficiente che i calcolatori delle generazioni successive ritornarono alla codifica binaria, convertendo i numeri nella rappresentazione in base 10 solamente per le operazioni, relativamente poco frequenti, di input/output.

Si tenga presente che le combinazioni di bit introdotte precedentemente sono semplicemente delle *rappresentazioni* dei numeri. In realtà i numeri hanno un numero infinito di cifre, delle quali quasi tutte assumono il valore 0 tranne quelle nelle posizioni più a destra. Normalmente non si considerano gli 0 iniziali.

È possibile progettare un circuito logico che sommi, sottragga, moltiplichi e divida queste combinazioni di bit. Se il risultato di queste operazioni non può essere rappresentato con i bit resi disponibili dall'hardware, si dice che

si è verificato un *overflow* (tracimazione). È compito del linguaggio ad alto livello, del sistema operativo e del programma determinare cosa fare quando si verifica un overflow.

I programmi dei calcolatori eseguono operazioni su numeri sia positivi sia negativi, ed è quindi necessaria una rappresentazione che permetta di distinguergli. La soluzione più ovvia consiste nell'aggiungere un'informazione di segno, che può essere convenientemente rappresentata da un singolo bit. Il nome di questa rappresentazione è *modulo e segno*.

Purtroppo la rappresentazione in modulo e segno ha diversi limiti. Anzitutto non è chiaro dove posizionare il bit di segno. A destra? A sinistra? Nei primi calcolatori vennero sperimentate ambedue le soluzioni. In secondo luogo, i sommatori che operano su numeri rappresentati mediante modulo e segno richiedono un passo aggiuntivo per determinare il segno, dal momento che questo non può essere determinato a priori. Infine, un bit di segno separato dagli altri significa che la rappresentazione in modulo e segno possiede sia uno 0 positivo che uno 0 negativo, cosa che può creare dei problemi ai programmatore meno attenti. La conseguenza di questi limiti è che la rappresentazione in modulo e segno fu molto presto abbandonata.

Mentre si cercavano alternative più convenienti, si presentò il problema di come si potesse rappresentare il risultato di un'operazione di sottrazione tra due numeri senza segno quando si sottraeva un numero grande da uno piccolo. Applicando le normali regole della sottrazione, il risultato è una stringa di 1 iniziali, ottenuti dai riporti generati dall'operazione.

Dal momento che non sembravano esserci alternative migliori, la rappresentazione adottata fu quella che rendeva l'hardware più semplice: prevedeva che una sequenza di 0 iniziali indicasse un valore positivo e una sequenza di 1 un numero negativo. Questa convenzione per la rappresentazione dei numeri binari dotati di segno prende il nome di rappresentazione in *complemento a 2*:

```

00000000 00000000 00000000 00000000 00000000due = 0dec
00000000 00000000 00000000 00000000 00000001due = 1dec
00000000 00000000 00000000 00000000 00000010due = 2dec
...
11111111 11111111 11111111 11111111 11111111 11111110due = 9 223 372 036 854 775 805dec
01111111 11111111 11111111 11111111 11111111 11111110due = 9 223 372 036 854 775 806dec
01111111 11111111 11111111 11111111 11111111 11111111due = 9 223 372 036 854 775 807dec
10000000 00000000 00000000 00000000 00000000 00000000due = -9 223 372 036 854 775 808dec
10000000 00000000 00000000 00000000 00000000 00000001due = -9 223 372 036 854 775 807dec
10000000 00000000 00000000 00000000 00000000 00000010due = -9 223 372 036 854 775 806dec
...
11111111 11111111 11111111 11111111 11111111 11111110due = -3dec
11111111 11111111 11111111 11111111 11111111 11111111due = -2dec
11111111 11111111 11111111 11111111 11111111 11111111due = -1dec

```

La metà positiva dei numeri, da 0 a 9 223 372 036 854 775 807<sub>dec</sub> ( $2^{63} - 1$ ), utilizza la stessa rappresentazione vista in precedenza. La combinazione successiva (1000 ... 0000<sub>due</sub>) rappresenta il numero negativo di valore assoluto maggiore, ossia -9 223 372 036 854 775 808<sub>dec</sub> ( $-2^{63}$ ). Subito dopo vi è una sequenza decrescente di numeri negativi, da -9 223 372 036 854 775 807<sub>dec</sub> (1000 ... 0001<sub>due</sub>) fino a -1<sub>dec</sub> (1111 ... 1111<sub>due</sub>).

La rappresentazione in complemento a 2 comprende un numero negativo, -9 223 372 036 854 775 808<sub>dec</sub>, che non ha un corrispondente numero positivo. Questo sbilanciamento richiede ai programmatore una certa attenzione, ma la rappresentazione in modulo e segno creava problemi sia ai programmatore sia ai progettisti hardware. Pertanto tutti i calcolatori di oggi utilizzano la rappresentazione binaria in complemento a 2 per i numeri dotati di segno.

La rappresentazione in complemento a due ha il vantaggio che tutti i numeri negativi hanno il bit più significativo uguale a 1. Perciò l'hardware deve controllare soltanto questo bit per verificare se il numero è positivo o negativo (un numero che ha inizio con lo zero viene considerato positivo). Questo bit è a volte chiamato *bit di segno*. Considerando il ruolo particolare assunto da questo bit, si possono rappresentare numeri positivi e negativi su 32 bit moltiplicando ciascun bit per una potenza di 2:

$$(x_{63} \times -2^{63}) + (x_{62} \times 2^{62}) + (x_{61} \times 2^{61}) + \dots + (x_1 \times 2^1) + (x_0 \times 2^0)$$

Il bit di segno è moltiplicato per  $-2^{63}$ , mentre tutti gli altri bit sono moltiplicati per il valore, positivo, delle rispettive basi.

### Conversione da binario a decimale

#### ESEMPIO

Qual è il numero decimale associato al seguente numero su 64 bit in complemento a 2?

11111111 11111111 11111111 11111111 11111111 11111111 11111111 111111100<sub>due</sub>

#### SOLUZIONE

Sostituendo il valore dei bit nella formula precedentemente introdotta si ha:

$$\begin{aligned} & (1 \times -2^{63}) + (1 \times 2^{62}) + (1 \times 2^{61}) + \dots + (1 \times 2^1) + (0 \times 2^0) + (0 \times 2^0) \\ & = -2^{63} + 2^{62} + 2^{61} + \dots + 2^1 + 0 + 0 \\ & = -9\,223\,372\,036\,854\,775\,808_{\text{dec}} + 9\,223\,372\,036\,854\,775\,804_{\text{dec}} = -4_{\text{dec}} \end{aligned}$$

Verrà introdotta a breve una scorciatoia che rende più semplice la conversione.

Come per il risultato di un'operazione su numeri senza segno, anche in un'operazione su numeri in complemento a 2 si può eccedere la capacità di rappresentazione del risultato da parte dell'hardware. La condizione di overflow si verifica quando il bit più a sinistra non è uguale agli infiniti bit presenti alla sinistra del numero stesso (ossia il bit di segno non è corretto): uno 0 a sinistra della combinazione di bit che rappresenta un numero negativo oppure un 1 a sinistra di quella che rappresenta un numero positivo.

## Interfaccia hardware/software

La distinzione tra numeri con e senza segno riguarda sia le istruzioni di trasferimento dalla memoria sia le istruzioni aritmetiche. Quando si effettua il caricamento di un numero dotato di segno, il bit di segno viene ripetutamente copiato fino a riempire il resto del registro. Questa operazione, chiamata *estensione del segno*, ha l'obiettivo di scrivere la rappresentazione corretta del numero nel registro. Quando invece si carica un numero senza segno, tutte le posizioni libere a sinistra del numero vengono riempite con degli 0, dato che il numero rappresentato è privo di segno.

Quando si carica una parola di memoria su 64 bit in un registro a 64 bit la questione diventa irrilevante, dato che l'operazione di caricamento diventa la stessa sia che il numero abbia segno sia che non ce l'abbia. Invece, per il caricamento dei byte, il RISC-V offre due istruzioni: *load byte unsigned* (`lbu`) tratta il byte come un numero privo di segno e perciò riempie con zeri gli altri bit alla sinistra del byte

(continua)

letto, mentre l'istruzione *load byte* (1b) considera il byte dotato di segno e copia il bit più significativo nelle posizioni del registro alla sinistra del byte letto. Dato che i programmi C utilizzano i byte quasi sempre per rappresentare caratteri alfanumerici e non per rappresentare numeri interi, viene utilizzata principalmente l'istruzione *lbu* per caricare i byte dalla memoria.

(continua)

A differenza dei numeri esaminati prima, gli indirizzi della memoria partono dal valore 0 e proseguono fino all'indirizzo più alto; in altre parole, gli indirizzi negativi non hanno senso. Per questo motivo, a volte, i programmi richiedono l'utilizzo di numeri che possono assumere valori positivi o negativi, mentre a volte richiedono l'utilizzo di numeri che possono essere solo positivi. Alcuni linguaggi di programmazione riflettono questa distinzione. Il C, per esempio, chiama i numeri del primo tipo *integers* (interi, dichiarati come *long long int* all'interno dei programmi), mentre chiama i numeri solo positivi *unsigned integers* (interi senza segno, dichiarati come *unsigned long long int* all'interno dei programmi). Inoltre, alcuni manuali di programmazione in C per evitare confusione suggeriscono di dichiarare i numeri dotati di segno come *signed long long int*.

## Interfaccia hardware/software

Esaminiamo alcune utili scorciatoie per lavorare con i numeri in complemento a 2. La prima scorciatoia è un modo veloce per cambiare il segno a un numero binario in complemento a 2. Basta semplicemente trasformare ogni 0 in 1 e ogni 1 in 0 e poi sommare 1 al risultato. Questa scorciatoia si basa sull'osservazione che la somma di un numero e del suo inverso dà sempre  $111\dots111_2$ , che rappresenta il numero  $-1$ . Dato che  $x + \bar{x} = -1$ , segue che  $x + \bar{x} + 1 = 0$ , ossia  $\bar{x} + 1 = -x$ . (Utilizziamo la notazione per indicare un numero binario ottenuto da  $\bar{x}$ , invertendo tutte le cifre binarie da 0 a 1 e viceversa.)

### Scorciatoia per cambiare di segno un numero

Si cambi di segno il numero  $2_{\text{dec}}$  e si verifichi poi il risultato cambiando di segno il numero  $-2_{\text{dec}}$ .

$$2_{\text{dec}} = 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000010_{\text{due}}$$

Invertendo questo numero tramite l'inversione dei bit e sommando poi 1, si ottiene:

$$\begin{array}{r} 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111101_{\text{due}} \\ + \qquad 1_{\text{due}} \\ \hline = 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111110_{\text{due}} \\ = -2_{\text{dec}} \end{array}$$

Procedendo nell'altra direzione, si parte da:

$$11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111110_{\text{due}}$$

che prima si inverte e poi si incrementa:

$$\begin{array}{r} 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000001_{\text{due}} \\ + \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad 1_{\text{due}} \\ \hline = 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 000000010_{\text{due}} \\ = 2_{\text{dec}} \end{array}$$

### ESEMPIO

### SOLUZIONE

La seconda scorciatoia serve per convertire un numero binario rappresentato su  $n$  bit in un numero rappresentato con più di  $n$  bit. La scorciatoia consiste nel prendere il bit più significativo (ossia il bit di segno) del numero su meno bit e nel replicarlo, in maniera da riempire gli altri bit presenti nel numero definito su un numero maggiore di bit.

Gli altri bit non contenenti il segno vengono semplicemente copiati nei bit più a destra della nuova parola doppia creata. Questa scorciatoia prende comunemente il nome di *estensione del segno*.

### Scorciatoia per l'estensione del segno

#### ESEMPIO

Si convertano i numeri  $2_{\text{dec}}$  e  $-2_{\text{dec}}$  codificati su 16 bit nei corrispondenti numeri binari codificati su 64 bit.

#### SOLUZIONE

La versione binaria su 16 bit del numero 2 è:

$$00000000\ 00000010_{\text{due}} = 2_{\text{dec}}$$

La conversione in un numero a 64 bit viene effettuata copiando 48 volte il valore del bit più significativo (0) nella metà sinistra della parola. La metà destra della parola riceve invece gli altri bit:

$$00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000000\ 00000010_{\text{due}} = 2_{\text{dec}}$$

Si esegua ora il cambiamento di segno del numero 2 su 16 bit, utilizzando la scorciatoia illustrata in precedenza. Quindi,

$$0000\ 0000\ 0000\ 0010_{\text{due}}$$

diventa:

$$\begin{array}{r} 1111\ 1111\ 1111\ 1101_{\text{due}} \\ + \qquad \qquad \qquad 1_{\text{due}} \\ \hline = 1111\ 1111\ 1111\ 1110_{\text{due}} \end{array}$$

Per creare la versione su 64 bit del numero negativo si deve copiare il bit di segno per 48 volte e metterlo nella parte sinistra della parola:

$$11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111110_{\text{due}} = -2_{\text{dec}}$$

Questo trucco funziona perché i numeri positivi in complemento a 2 hanno un numero infinito di 0 a sinistra, mentre quelli negativi un numero infinito di 1. Nella stringa di cifre binarie utilizzata per rappresentare un numero non vengono riportati i bit iniziali in modo tale da fare rientrare la lunghezza della cifra nell'ampiezza della parola dell'hardware; l'estensione del segno semplicemente ripristina alcuni di questi bit.

### Riepilogo

L'argomento principale trattato in questo paragrafo è la necessità di rappresentare nella parola di un calcolatore numeri interi sia positivi sia negativi. Benché in ogni decisione vi siano sempre dei pro e dei contro, la scelta unanime a partire dal 1965 è stata la rappresentazione in complemento a 2.

**Approfondimento.** Per quanto riguarda i numeri decimali dotati di segno, utilizzeremo il segno meno ( $-$ ) per rappresentare i numeri negativi, perché non ci sono limiti alla dimensione dei numeri. Con la notazione binaria ed esadecimale, invece, la codifica è su parole di un numero prefissato di bit (vedi Figura 2.4) e il segno viene espresso nella codifica stessa, e quindi i simboli  $+$  e  $-$  in queste notazioni normalmente non vengono utilizzati.

## Autovalutazione

Qual è il valore decimale di questo numero di 64 bit codificato in complemento a 2?

11111111 11111111 11111111 11111111 11111111 11111111 111111000<sub>due</sub>

1.  $-4_{dec}$
2.  $-8_{dec}$
3.  $-16_{dec}$
4. 18 446 744 073 709 551 609<sub>dec</sub>

**Approfondimento.** La rappresentazione in complemento a 2 deve il suo nome alla proprietà in base alla quale la somma senza segno di un numero di  $n$  bit e del suo complemento è pari a  $2^n$ , e quindi il complemento (o negazione) di un numero  $x$  in complemento a 2 è pari a  $2^n - x$ , ovvero il suo complemento a 2.

Una terza possibile rappresentazione dei numeri dotati di segno è denominata **complemento a 1**. Il numero negativo di un numero in complemento a 1 si ottiene invertendo ciascun bit, da 0 a 1 e da 1 a 0, e viene indicato con  $\bar{x}$ . Questo spiega il suo nome, dato che il complemento a 1 di  $x$  sarà  $2^n - x - 1$ . Questa rappresentazione, quando fu introdotta, era un tentativo di avere una codifica migliore della rappresentazione in modulo e segno, e diversi tra i primi calcolatori di tipo scientifico adottarono tale notazione. Essa è simile al complemento a 2, tranne che prevede anch'essa due rappresentazioni per lo 0: lo 0 positivo (00 ... 00<sub>due</sub>) e quello negativo (11 ... 11<sub>due</sub>). Il numero negativo di valore assoluto maggiore su 32 bit è 10 ... 000<sub>due</sub> e rappresenta il valore  $-2\ 147\ 483\ 647_{dec}$ , e quindi l'insieme dei numeri positivi e quello dei numeri negativi sono bilanciati. Tuttavia i sommatori in complemento a 1 richiedono un passo in più per eseguire la sottrazione di un numero, e per questo il complemento a 2 è diventata la scelta dominante.

L'ultima notazione prevede di rappresentare il numero più negativo con la sequenza 00 ... 000<sub>due</sub> e il numero più positivo con 11 ... 11<sub>due</sub>; lo 0 tipicamente assume il valore 10 ... 00<sub>due</sub>. Questa notazione è chiamata **notazione polarizzata**, dato che il numero viene polarizzato in modo tale che la somma del numero e della polarizzazione produca un numero non negativo.

**Complemento a 1:** codifica che rappresenta il minimo numero negativo come 10 ... 000<sub>due</sub> e il massimo numero positivo come 01 ... 11<sub>due</sub>. Gli insiemi dei numeri positivi e dei numeri negativi sono bilanciati ma hanno due codifiche per lo zero, lo zero positivo (00 ... 00<sub>due</sub>) e lo zero negativo (11 ... 11<sub>due</sub>). Complemento a 1 indica anche l'inversione dei bit di una stringa, da 0 a 1 e da 1 a 0.

**Notazione polarizzata:** notazione che prevede che il minimo valore negativo venga rappresentato da 00 ... 000<sub>due</sub> e il massimo valore positivo da 11 ... 11<sub>due</sub>, con lo 0 che tipicamente assume il valore 10 ... 00<sub>due</sub>; quindi, sommando un numero polarizzato con la polarizzazione si ottiene sempre un numero non negativo.

## 2.5 | Rappresentazione delle istruzioni nel calcolatore

Siamo pronti ora a esaminare la differenza tra il modo in cui gli esseri umani danno istruzioni a un calcolatore e quello in cui il calcolatore vede queste istruzioni.

Anche le istruzioni sono memorizzate nel calcolatore come una sequenza di segnali elettrici alto/basso e possono quindi essere viste come numeri. In effetti, ciascuna parte di un'istruzione può essere vista come un numero e l'istruzione vera e propria si può ottenere disponendo questi numeri uno di fianco all'altro. I 32 registri dei RISC-V sono semplicemente individuati dal loro numero progressivo da 0 a 31.

## Traduzione di un'istruzione assembler RISC-V nell'istruzione macchina corrispondente

### ESEMPIO

Facciamo il prossimo passo nella definizione del linguaggio RISC-V mediante un esempio. Mostreremo nel vero linguaggio RISC-V la versione dell'istruzione rappresentata simbolicamente dall'espressione:

`add x9, x20, x21`

dapprima come sequenza di numeri decimali e poi come sequenza di numeri binari.

La rappresentazione decimale è:

0	21	20	0	9	51
---	----	----	---	---	----

Ciascuna delle parti che costituiscono un'istruzione è chiamato *campo* (*field*). La combinazione del primo, del quarto e del sesto campo (che in questo caso contengono rispettivamente 0, 0 e 51) indica al processore RISC-V che questa istruzione esegue una somma. Il secondo campo indica il numero del registro che contiene il primo operando (21 = x21), il terzo campo il secondo operando della somma (20 = x20) e il quinto campo contiene il numero del registro in cui viene memorizzata la somma (9 = x9). Questa istruzione dunque somma il contenuto dei registri x21 e x20 e pone il risultato in x9.

Questa stessa istruzione può anche essere rappresentata come sequenza di campi contenenti numeri binari invece che decimali:

0000000	10101	10100	000	01001	0110011
7 bit	5 bit	5 bit	3 bit	5 bit	7 bit

**Formato dell'istruzione:** un modo di rappresentare un'istruzione attraverso un insieme di campi di numeri binari.

**Linguaggio macchina:** rappresentazione binaria utilizzata per la comunicazione all'interno dei calcolatori.

**Numeri esadecimali:** numeri in base 16.

Questa scomposizione in campi di un'istruzione è chiamata **formato dell'istruzione**. Come si può notare contando il numero di bit, l'istruzione RISC-V richiede esattamente 32 bit, una parola, o una mezza parola doppia. In accordo con il principio di progetto per cui la semplicità favorisce la regolarità, tutte le istruzioni RISC-V sono lunghe 32 bit.

Per distinguere la versione numerica delle istruzioni da quella in linguaggio assembler, la prima viene chiamata **linguaggio macchina**; una sequenza di istruzioni in linguaggio macchina viene definita *codice macchina*.

Da quanto abbiamo visto, sembrerebbe che si debbano leggere e scrivere lunghe e noiose stringhe di numeri binari; per evitare ciò si ricorre all'uso di una base di numerazione più alta di quella binaria, che possa però essere convertita facilmente in binario. Dato che quasi tutti i calcolatori hanno una dimensione dei dati multipla di 4, i **numeri esadecimali** (in base 16) sono molto diffusi. Inoltre, essendo la base 16 una potenza di 2, si può facilmente convertire un numero binario in esadecimale raggruppando le cifre binarie a gruppi di 4 e sostituendo a ciascun gruppo una singola cifra esadecimale, e viceversa. La Figura 2.4 mostra la conversione da esadecimale a binario (e viceversa) dei primi numeri interi.

Poiché utilizzeremo spesso numeri in basi diverse, per evitare di fare confusione i decimali saranno indicati con il pedice *dec*, i binari con il pedice *due* e gli esadecimali con il pedice *esa*. (I numeri senza pedici vanno considerati in base dieci.) A questo proposito, sia C che Java utilizzano la notazione `0xnnnn` per identificare un numero esadecimale.

Esadecimale	Binario	Esadecimale	Binario	Esadecimale	Binario	Esadecimale	Binario
0 <sub>esa</sub>	0000 <sub>due</sub>	4 <sub>esa</sub>	0100 <sub>due</sub>	8 <sub>esa</sub>	1000 <sub>due</sub>	c <sub>esa</sub>	1100 <sub>due</sub>
1 <sub>esa</sub>	0001 <sub>due</sub>	5 <sub>esa</sub>	0101 <sub>due</sub>	9 <sub>esa</sub>	1001 <sub>due</sub>	d <sub>esa</sub>	1101 <sub>due</sub>
2 <sub>esa</sub>	0010 <sub>due</sub>	6 <sub>esa</sub>	0110 <sub>due</sub>	a <sub>esa</sub>	1010 <sub>due</sub>	e <sub>esa</sub>	1110 <sub>due</sub>
3 <sub>esa</sub>	0011 <sub>due</sub>	7 <sub>esa</sub>	0111 <sub>due</sub>	b <sub>esa</sub>	1011 <sub>due</sub>	f <sub>esa</sub>	1111 <sub>due</sub>

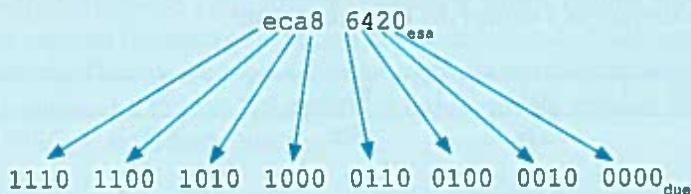
**Figura 2.4** Tabella di conversione tra binario ed esadecimale. È sufficiente sostituire una cifra esadecimale con le corrispondenti quattro cifre binarie e viceversa. Se la lunghezza del numero binario non è un multiplo di 4, si procede da destra a sinistra.

## Conversione da binario a esadecimale e viceversa

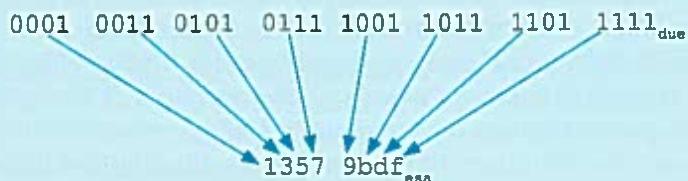
Convertire i seguenti numeri, uno esadecimale su 8 cifre e l'altro binario su 32 bit, rispettivamente in base 2 e 16:

eca8 6420<sub>esa</sub>  
0001 0011 0101 0111 1001 1011 1101 1111<sub>due</sub>

Utilizzando la Figura 2.4, la risposta è una semplice trasformazione diretta dei simboli:



La conversione da base 2 a base 16 sarà:



## Campi delle istruzioni RISC-V

Per rendere più facile la trattazione, ai diversi campi delle istruzioni RISC-V viene associato un nome:

funz7	rs2	rs1	funz3	rd	codop
7 bit	5 bit	5 bit	3 bit	5 bit	7 bit

Il significato del nome dato a ciascun campo è il seguente:

- **codop**: operazione base dell'istruzione, tradizionalmente chiamata **codice operativo** o **codop** (da *operating code*);
- **rd**: registro destinazione: riceve il risultato dell'operazione;
- **funz3**: un codice operativo aggiuntivo;
- **rs1**: registro contenente il primo operando sorgente;
- **rs2**: registro contenente il secondo operando sorgente;
- **funz7**: un codice operativo aggiuntivo.

### ESEMPIO

### SOLUZIONE

**Codice operativo (codop):**  
il campo di un'istruzione che specifica l'operazione e il formato dell'istruzione stessa.

Può nascere un problema quando un'istruzione richiede campi di dimensioni maggiori rispetto a quelle sopra specificate. Per esempio, l'istruzione load richiede di specificare due registri e una costante: se la costante venisse inserita in uno dei campi da 5 bit, non potrebbe superare il valore di  $2^5 - 1$ , cioè 31. Questa costante viene utilizzata di solito per selezionare un elemento specifico all'interno di vettori o altri tipi di strutture dati e spesso deve assumere valori molto maggiori di 31: un campo da 5 bit risulterebbe troppo piccolo per essere utile.

Nasce dunque un conflitto fra il desiderio di mantenere la stessa lunghezza per tutte le istruzioni e quello di avere un unico formato. Ciò conduce all'ultimo principio della progettazione dell'hardware:

*Principio di progettazione numero 3: un buon progetto richiede buoni compromessi.*

Il compromesso scelto dai progettisti del RISC-V è quello di mantenere uguale la lunghezza di tutte le istruzioni, predisponendo formati differenti per tipi di istruzioni diverse. Per esempio, il formato descritto precedentemente è chiamato di *tipo R* (R sta per registro). Un secondo tipo di formato è detto di *tipo I* (I sta per immediato) e viene utilizzato dalle operazioni aritmetiche nelle quali un operando è una costante, compreso `addi`, e dalle istruzioni di trasferimento dati con la memoria. I campi del formato I sono:

immediato	rs1	funz3	rd	codop
12 bit	5 bit	3 bit	5 bit	7 bit

Il campo immediato a 12 bit viene interpretato come un valore in complemento a due, per cui rappresenta gli interi compresi tra  $-2^{11}$  e  $+2^{11} - 1$ . Quando il formato di tipo I viene utilizzato per le istruzioni di trasferimento dati, il campo immediato rappresenta l'offset in byte, per cui le istruzioni di caricamento di una parola doppia possono indirizzare tutte le doppie parole in un intervallo di  $\pm 2^{11}$ , cioè 2048 byte ( $\pm 2^8$  cioè 256 doppie parole) rispetto all'indirizzo contenuto nel registro di base rd. Osserviamo che avere più di 32 registri sarebbe difficile con questo formato di istruzioni, poiché occorrerebbe almeno un bit in più per ognuno dei campi rd e rs1, e sarebbe quindi più difficile far stare tutto in una sola parola.

Si consideri l'istruzione di load esaminata a pagina 62:

```
ld x9, 64(x22) // il registro temporaneo x9
                  // assume il valore di A[8]
```

In questo caso il numero 22 (che rappresenta il registro x22) viene posto nel campo rs1, il numero 64 nel campo immediato e il numero 9 (che rappresenta x9) nel campo rd. Occorre definire un formato anche per l'istruzione di trasferimento in memoria di una parola doppia, `sd`, che ha bisogno di due registri sorgenti (uno per l'indirizzo base e l'altro per il dato da memorizzare) e un valore immediato per l'offset dell'indirizzo. I campi del formato di tipo S sono:

immediato[11:5]	rs2	rs1	funz3	Immediato[4:0]	codop
7 bit	5 bit	5 bit	3 bit	5 bit	7 bit

Il valore immediato su 12 bit nel formato di tipo S è suddiviso in due campi, che forniscono rispettivamente i 5 bit meno significativi e i 7 bit più signifi-

Istruzione	Formato	funz7	rs2	rs1	funz3	rd	codop
add	R	0000000	reg	reg	000	reg	0110011
sub (sottrazione)	R	0100000	reg	reg	000	reg	0110011
Istruzione	Formato	immediato		rs1	funz3	rd	codop
addi (addizione immediata)	I	costante		reg	000	reg	0010011
ld (caricamento di parola doppia)	I	indirizzo		reg	011	reg	0000011
Istruzione	Formato	immediato	rs2	rs1	funz3	immediato	codop
sd (memorizzazione di parola doppia)	S	indirizzo	reg	reg	011	indirizzo	0100011

**Figura 2.5 Codifica delle istruzioni RISC-V.** In questa tabella “reg” indica il numero di un registro ed è compreso tra 0 e 31, e “indirizzo” è un indirizzo su 12 bit o una costante. I campi funz3 e funz7 agiscono come codici operativi aggiuntivi.

tivi. I progettisti del RISC-V hanno scelto questa soluzione perché consente di mantenere i campi rs1 e rs2 nella stessa posizione in tutti i formati di istruzioni. Avere dei formati delle diverse istruzioni più simili possibili tra loro riduce la complessità. Analogamente, il codice operativo e il campo funct3 hanno la stessa dimensione e si trovano nella stessa posizione in tutte le istruzioni.

Si possono distinguere i due formati in base al valore contenuto nel codice operativo: a ciascun formato è assegnato un insieme di valori del campo codop, in modo tale che l’hardware sappia esattamente come deve trattare i rimanenti bit dell’istruzione. La Figura 2.5 mostra il contenuto di ciascun campo delle istruzioni RISC-V introdotte finora.

## Traduzione in linguaggio macchina delle istruzioni assembler del RISC-V

È ora possibile mostrare un esempio completo, partendo da ciò che scrive il programmatore fino ad arrivare a ciò che viene eseguito dal calcolatore. Se x10 contiene l’indirizzo base del vettore A e x21 corrisponde a h, l’istruzione di assegnamento:

A[30] = h + A[30] + 1;

viene compilata in:

```
ld  x9, 240(x10) // il registro temporaneo x9
    // assume il valore A[30]
add x9, x21, x9 // il registro temporaneo x9
    // assume il valore h + A[30]
addi x9, x9, 1   // il registro temporaneo x9
    // assume il valore h + A[30] + 1
sd  x9, 240(x10) // memorizza h + A[30] + 1 nuovamente
    // in A[30]
```

Qual è il codice macchina RISC-V corrispondente a queste quattro istruzioni?

Per comodità conviene inizialmente rappresentare le istruzioni in linguaggio macchina nel formato decimale. Dalla Figura 2.5 possiamo determinare queste quattro istruzioni in linguaggio macchina:

### ESEMPIO

### SOLUZIONE

(continua)

(continua)

immediato	rs1	funz3	rd	codop
240	10	3	9	3
funz7	rs2	rs1	funz3	rd
0	9	21	0	9
immediato	rs1	funz3	rd	codop
1	9	0	9	19
immediato[11:5]	rs2	rs1	funz3	immediato[4:0]
7	9	10	3	16
				35

L'istruzione **ld** è identificata dal numero 3 (Figura 2.5), contenuto nel campo codop e dal numero 3 contenuto nel campo funz3. Il registro base 10 è specificato nel campo rs1 e il registro destinazione 9 è specificato nel campo rd. L'offset per selezionare A[30] ( $240 = 300 \times 8$ ) si trova nel campo immediato.

La successiva istruzione di somma, **add**, è identificata dal numero 51 nel campo codop, dal numero 0 nel campo funz3 e dal numero 0 nel campo funz7. I tre registri degli operandi (9, 21 e 9) sono contenuti rispettivamente nei campi rd, rs1 e rs2.

L'istruzione seguente, **addi**, è identificata dal numero 19 contenuto nel campo codop e dal numero 0 contenuto nel campo funz3. I registri degli operandi (9 e 9) sono specificati nei campi rd e rs1, e la costante, l'addendo 1, si trova nel campo immediato.

L'istruzione **sd** è identificata dal numero 35 contenuto nel campo codop e dal numero 3 contenuto nel campo funz3. I registri degli operandi (9 e 10) sono specificati rispettivamente nei campi rs2 e rs1. L'offset dell'indirizzo è suddiviso tra i due campi immediato. Dato che il campo immediato con i bit più significativi contiene i bit dal 5 all'11, possiamo decomporre l'offset 240 dividendolo per  $2^5$ . La parte con i bit più significativi conterrà il quoziente della divisione, cioè 7, mentre la parte con i bit meno significativi conterrà il resto cioè 16.

Dato che  $240_{dec} = 0000\ 1111\ 0000_{due}$ , la rappresentazione binaria equivalente delle quattro istruzioni, definite mediante numeri decimali, è:

immediato	rs1	funz3	Rd	codop
000011110000	01010	011	01001	0000011
funz7	rs2	rs1	funz3	rd
0000000	01001	10101	000	01001
immediato	rs1	funz3	rd	codop
000000000001	01001	000	01001	0010011
immediato[11:5]	rs2	rs1	funz3	Immediato[4:0]
0000111	01001	01010	011	10000
				0100011

**Approfondimento.** I programmatore assembler in linguaggio RISC-V non sono obbligati a utilizzare addi quando lavorano con le costanti. Il programmatore può semplicemente scrivere add, e l'assemblatore genera il codice operativo e il formato adeguato all'istruzione a seconda che gli operandi siano tutti contenuti nei registri (tipo R) oppure uno di essi sia una costante (tipo I).

Quando introduciamo il linguaggio assembler e lo confrontiamo con il linguaggio macchina utilizziamo esplicitamente i diversi nomi delle varianti delle istruzioni RISC-V con i loro codici operativi e formati, perché riteniamo che sia più chiaro per il lettore.

**Approfondimento.** Anche se il RISC-V contiene sia le istruzioni di add che di sub, non contiene l'istruzione di subi associata alla addi. Questo perché il numero contenuto nel campo immediato può essere rappresentato come intero in complemento a due, per cui addi può essere utilizzata per sottrarre una costante.

Il desiderio di avere tutte le istruzioni della stessa dimensione è in conflitto con il desiderio di avere un numero di registri grande a piacere. L'aumento del numero di registri richiede l'aggiunta di almeno un bit all'interno del formato dell'istruzione per tutti i campi che contengono un numero di registro. Dati questi vincoli e il principio di progettazione per il quale più piccolo implica più veloce, la maggior parte delle istruzioni, oggi, ha un numero di registri variabile tra 16 e 32.

## Interfaccia hardware/software

La Figura 2.6 riassume le parti di linguaggio macchina RISC-V introdotte in questo paragrafo. Come si vedrà nel Capitolo 4, la somiglianza della rappresentazione binaria di istruzioni correlate consente di semplificare la progettazione dell'hardware. Queste somiglianze sono un altro esempio della regolarità dell'architettura RISC-V.

### QUADRO D'INSIEME

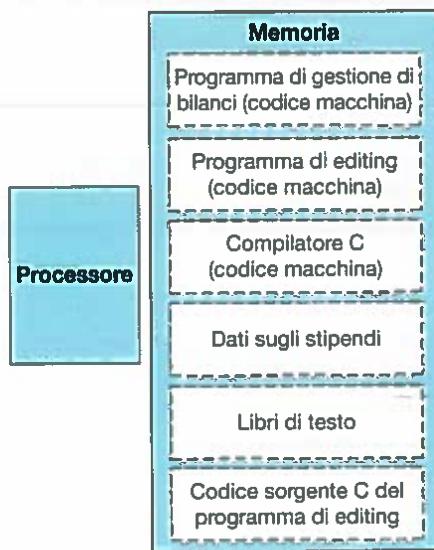
Oggi i calcolatori sono costruiti secondo due principi chiave:

1. Le istruzioni vengono rappresentate come numeri.
2. I programmi sono contenuti in memoria e possono essere letti e scritti, proprio come dati.

Questi principi esprimono il concetto di programma memorizzato (*stored program*): un'invenzione fondamentale per l'informatica. La Figura 2.7 mostra la portata di questo concetto; nel caso specifico, la memoria contiene il codice sorgente di un programma per l'elaborazione di testi, il corrispondente codice macchina compilato, il testo su cui il programma compilato sta lavorando e anche il compilatore che ha generato il codice macchina. Una conseguenza del fatto che le istruzioni sono rappresentate come numeri è che i programmi possono essere trasmessi attraverso file di numeri binari. L'implicazione commerciale è che un calcolatore può ereditare del codice già pronto, purché le istruzioni che costituiscono il codice siano compatibili con il suo insieme di istruzioni. Questa "compatibilità binaria" ha spesso portato l'industria a sviluppare software per un numero ristretto di architetture di insiemi di istruzioni. ■

Istruzione (R)	funz7	rs2	rs1	funz3	rd	codop	Esempio
add	0000000	00011	00010	000	00001	0110011	add x1, x2, x3
sub (sottrazione)	0100000	00011	00010	000	00001	0110011	sub x1, x2, x3
Istruzione (I)	immediato		rs1	funz3	rd	codop	Esempio
addi (addizione immediata)	001111101000		00010	000	00001	0010011	addi x1,x2,1000
ld (caricamento di parola doppia)	001111101000		00010	011	00001	0000011	ld x1, 1000 (x2)
Istruzione (S)	Immediato	rs2	rs1	funz3	immediato	codop	Esempio
sd (memorizzazione di parola doppia)	0011111	00001	00010	011	01000	0100011	sd x1, 1000 (x2)

**Figura 2.6** L'architettura RISC-V descritta fino al paragrafo 2.5. I tre formati delle istruzioni RISC-V visti finora sono R, I e S. Il formato R ha due registri sorgenti per gli operandi e un registro destinazione. Il formato I sostituisce uno dei due registri sorgente con la parte meno significativa del campo *immediato* su 12 bit. Il formato S ha due operandi sorgente e un campo *immediato* su 12 bit, ma non ha un registro destinazione. Il campo *immediato* delle istruzioni di tipo S viene suddiviso in due parti: i bit 11-5 sono contenuti nel campo *immediato* più a sinistra e i bit 4-0 nel secondo campo da destra.



**Figura 2.7** Concetto di programma memorizzato. La memorizzazione dei programmi consente a un calcolatore che gestisce bilanci di trasformarsi improvvisamente in un calcolatore che serve a un autore per scrivere un libro. Il cambiamento avviene semplicemente caricando in memoria programmi e dati e poi specificando al calcolatore l'indirizzo di memoria da cui fare partire l'esecuzione. Trattare i programmi e i dati allo stesso modo semplifica sia l'hardware della memoria sia il software dei sistemi operativi. Nel caso specifico, la stessa tecnologia di memoria utilizzata per i dati si può applicare anche ai programmi; programmi quali i compilatori, per esempio, possono tradurre il codice scritto in una notazione molto più comprensibile all'uomo in un codice che il calcolatore è in grado di capire.

## Autovalutazione

Quale istruzione RISC-V è rappresentata sotto? Selezionare una delle quattro risposte riportate.

funz7	rs2	rs1	funz3	rd	codop
32	9	10	0	11	51

1. add x9, x10, x11
2. add x11, x9, x10
3. sub x11, x10, x9
4. sub x11, x9, x10

## 2.6 Operazioni logiche

Sebbene i primi calcolatori effettuassero operazioni su parole intere, divenne presto evidente che era utile riuscire a operare anche su gruppi di bit o addirittura sui singoli bit di una parola. Esaminare i caratteri all'interno di una parola, dove ciascun carattere è memorizzato su 8 bit, è uno degli esempi di questo tipo di operazioni (vedi par. 2.9). Il risultato fu che vennero aggiunte ai linguaggi di programmazione e alle architetture degli insiemi di istruzioni operazioni che permettevano di semplificare, per esempio, l'inserimento e la lettura di gruppi di bit all'interno di una parola. Queste istruzioni sono chiamate *operazioni logiche*. In Figura 2.8 vengono riportate le *operazioni logiche* in C, Java e nel RISC-V.

La prima tipologia di queste operazioni è detta *shift* (scorrimento) e consiste nello spostare tutti i bit di una parola a sinistra o a destra, riempiendo i bit vuoti con degli zeri. Per esempio, supponiamo che il registro x19 contenga:

00000000 00000000 00000000 00000000 00000000 00000000 00001001<sub>bin</sub> = 9<sub>dec</sub>

ed eseguiamo l'istruzione di shift a sinistra di 4; il numero ottenuto sarà il seguente:

00000000 00000000 00000000 00000000 00000000 00000000 10010000<sub>bin</sub> = 144<sub>dec</sub>

L'operazione duale dello shift a sinistra è lo shift a destra. Il nome delle due istruzioni RISC-V di shift è *shift a sinistra logico immediato* (slli, *shift left logical immediate*) e *shift a destra logico immediato* (srli, *shift right logical immediate*). L'istruzione seguente permette di eseguire l'operazione riportata in precedenza. Si suppone che l'istruzione legga l'operando dal registro x19 e metta il risultato nel registro x11:

```
slli x11, x19, 4 // reg x11 = reg x19 << 4 bit
```

Queste istruzioni di shift utilizzano il formato di tipo I. Dato che non serve fare scorrere i bit di un registro a 64 bit per più di 63 posizioni, solamente i 6 bit meno significativi del campo immediato del formato di tipo I vengono effettivamente utilizzati. I rimanenti 6 bit vengono utilizzati come un campo aggiuntivo di codice operativo.

funz6	immediato	rs1	funz3	rd	codop
0	4	19	1	11	19

L'istruzione slli prevede 19 nel campo codop, 11 nel campo rd, 1 nel campo funz3, 19 nel campo rs1, 4 nel campo immmediato e 0 nel campo funz6.

*"Al contrario", continuò Tweedledee, "se fosse così, potrebbe essere; e se così fosse, sarebbe; ma poiché non è, non è. Questa è logica".*

Lewis Carrol, *Alice nel paese delle meraviglie*, 1865

Operazioni logiche	Operatori C	Operatori Java	Istruzioni RISC-V
Shift a sinistra	<<	<<	sll, slli
Shift a destra	>>	>>>	srl, srli
Shift a destra aritmetico	>>	>>	sra, srai
AND bit a bit	&	&	and, andi
OR bit a bit			or, ori
XOR bit a bit	^	^	xor, xorl
NOT bit a bit	~	~	xori

**Figura 2.8** Operatori logici C e Java e le corrispondenti istruzioni RISC-V. Un modo per implementare l'operatore NOT in RISC-V è utilizzare XOR con uno degli operandi impostato a tutti uni (FFFF FFFF FFFF<sub>bin</sub>).

L'operazione di shift logico a sinistra fornisce un'ulteriore funzionalità; infatti, lo scorrimento di un numero a sinistra di  $i$  cifre produce lo stesso risultato di una moltiplicazione per  $2^i$ , esattamente come lo scorrimento di un numero decimale a sinistra di  $i$  cifre è equivalente a una moltiplicazione per  $10^i$ . Per esempio, la precedente operazione di `slli` di 4 posizioni produce lo stesso risultato che si otterebbe moltiplicando il contenuto del registro per  $2^4$  (cioè 16): la prima sequenza di bit rappresentava il numero 9, e  $9 \times 16 = 144$  è il numero rappresentato nella seconda sequenza di bit.

Il RISC-V fornisce anche un terzo tipo di operazione di scorrimento: *scorrimento a destra aritmetico* (`srai`). Questa variante è simile a `srl1` tranne che invece di riempire con degli zero i bit che si liberano a sinistra, questi bit vengono riempiti copiando il bit del segno. Il RISC-V fornisce anche una variante per ciascuna delle tre operazioni di scorrimento che prendono il numero di posizioni di scorrimento da un registro, invece che dal campo immediato: `ssl`, `srl` e `sra`.

Un'altra utile operazione logica, che permette di isolare i campi di una parola, è **AND** (è scritta in maiuscolo per evitare la confusione con la congiunzione inglese). L'operazione di AND è un'operazione bit a bit che scrive un 1 nel risultato solamente se entrambi gli operandi hanno un 1 in quella stessa posizione. Per esempio, supponiamo che il registro `x11` contenga:

00000000 00000000 00000000 00000000 00000000 00000000 00001101 11000000<sub>due</sub>

e il registro `x10` contenga:

00000000 00000000 00000000 00000000 00000000 00000000 00111100 00000000<sub>due</sub>

dopo l'esecuzione della seguente istruzione RISC-V:

`and x9, X10, X11, // reg x9 = reg x10 & reg x11`

il valore contenuto nel registro `x9` sarebbe:

00000000 00000000 00000000 00000000 00000000 00000000 00001100 00000000<sub>due</sub>

Come si può notare, l'operazione di AND bit a bit può essere utilizzata per forzare a 0 i bit di una parola fornendo in input all'AND una parola che contiene 0 in quelle posizioni, dato che essa "nasconde" alcuni bit (questa parola viene tradizionalmente chiamata *maschera*).

Per inserire un 1 in questo mare di zeri esiste un'operazione bit a bit duale rispetto alla AND, chiamata **OR**. Quest'ultima è un'operazione bit a bit che inserisce un 1 nel risultato laddove *almeno uno dei due operandi* sia uguale a 1. Per aiutarci a capire meglio, supponiamo che i registri `x10` e `x11` contengano gli stessi valori visti in precedenza; il risultato dell'esecuzione dell'istruzione RISC-V:

`or x9, x10, x11 // reg x9 = reg x10 | reg x11`

verrà scritto nel registro `x9` e sarà:

00000000 00000000 00000000 00000000 00000000 00000000 00111101 11000000<sub>due</sub>

L'ultima operazione logica è la negazione: l'operazione di **NOT** prende un operando e pone un 1 nel risultato laddove nell'operando era presente uno 0 e viceversa. Utilizzando la notazione introdotta in precedenza, l'operatore NOT calcola  $\bar{x}$ .

Per mantenere il formato dell'operazione a tre operandi, i progettisti RISC-V hanno deciso di includere l'istruzione **XOR** (OR esclusivo) invece della negazione semplice (NOT). Dato che XOR produce uno 0 quando due bit sono uguali e 1 altrimenti, si può ottenere l'operatore NOT mediante `xor` di un numero con 111...111.

**AND:** un'operazione logica bit a bit su due operandi, che restituisce 1 se entrambi gli operandi sono uguali a 1.

**OR:** un'operazione logica bit a bit su due operandi, che restituisce 1 se almeno uno dei due operandi è uguale a 1.

**NOT:** un'operazione logica bit a bit su un singolo operando che inverte i bit: sostituisce cioè tutti gli 1 con 0, e tutti gli 0 con 1.

**XOR:** un'operazione logica bit a bit su due operandi, che calcola l'OR esclusivo dei due operandi. Restituisce cioè 1 solo se i valori dei due operandi sono diversi.

Supponiamo che il contenuto del registro x10 sia quello degli esempi precedenti e che il registro x12 contenga il valore 0. Il risultato dell'esecuzione dell'istruzione RISC-V:

```
xor x9, x10, x12 // reg x9 = reg x10 ^ reg x12
```

verrà scritto nel registro x9 e sarà:

```
00000000 00000000 00000000 00000000 00000000 00110001 1100000016
```

La Figura 2.8 illustra la relazione tra gli operatori C e Java e le istruzioni RISC-V. Le costanti sono molto utili anche nelle operazioni logiche di AND e OR, come lo erano nelle operazioni aritmetiche. Il RISC-V fornisce quindi anche le istruzioni di *AND immediato* (*andi*), di *OR immediato* (*ori*) e di *OR esclusivo immediato* (*xori*).

**Approfondimento.** Il C consente di definire all'interno di una doppia parola *campi di bit* o *campi*, che consentono di impaccare degli oggetti all'interno della parola doppia e che sono uguali ai campi esposti all'esterno dalle interfacce per esempio dei dispositivi di I/O. Tutti i campi devono essere contenuti all'interno di una singola parola doppia. I campi sono degli interi senza segno che possono essere piccoli fino a occupare 1 solo bit. I compilatori C inseriscono ed estraggono i diversi campi utilizzando le istruzioni logiche che in RISC-V sono: *andi*, *ori*, *slli* e *srls*.

## Autovalutazione

Quali operazioni possono isolare un campo all'interno di una parola doppia?

1. AND.
2. Uno shift a sinistra seguito da uno shift a destra.

## 2.7 Istruzioni per prendere decisioni

Ciò che contraddistingue un calcolatore da una semplice calcolatrice è la sua capacità di prendere delle decisioni: in base ai dati di ingresso e ai valori calcolati durante l'elaborazione, possono essere eseguite istruzioni diverse. Nei linguaggi di programmazione un processo di decisione è comunemente rappresentato con il costrutto *if*, a volte seguito dal costrutto *go to* seguito da un'etichetta. Il linguaggio assembler RISC-V contiene due istruzioni per implementare un processo di decisione, simili a un costrutto *if* seguito da *go to*. La prima istruzione è:

```
beq rs1, rs2, L1
```

e significa: vai all'istruzione etichettata con L1 se il valore contenuto nel registro rs1 è uguale al valore contenuto nel registro rs2; il codice mnemonico *beq* significa *branch if equal* (salta se uguale). La seconda istruzione è:

```
bne rs1, rs2, L1
```

e significa: vai all'istruzione etichettata con L1 se il valore contenuto nel registro rs1 non è uguale al valore contenuto nel registro rs2; il codice mnemonico *bne* indica *branch if not equal* (salta se non uguale). Queste due istruzioni prendono tradizionalmente il nome di **salvi condizionati** (*conditional branches*).

*L'utilità di un calcolatore automatico sta nella possibilità di utilizzare ripetutamente una sequenza di istruzioni, un numero di volte che dipende dal risultato dei calcoli. [...] Questa scelta può dipendere dal segno di un numero (considerando lo zero come un numero positivo per il calcolatore). Di conseguenza, introduciamo una [istruzione] ([l'istruzione] di trasferimento condizionato) che, a seconda del segno di un dato numero, provocherà l'esecuzione di una tra due possibili procedure.*

Burks, Goldstine e von Neumann, 1947.

**Salto condizionato:** istruzione che richiede il confronto fra due valori e che consente il trasferimento del controllo a un altro indirizzo del programma a seconda del risultato del confronto.

## Compilazione di un costrutto di tipo *if-then-else* mediante salti condizionati

### ESEMPIO

Nel seguente frammento di codice C ci sono cinque variabili: f, g, h, i e j. Supponiamo che corrispondano ai cinque registri da x19 a x23. Quale sarà il codice RISC-V prodotto dalla compilazione del seguente codice C?

```
if (i == j) f = g + h; else f = g - h;
```

### SOLUZIONE

La Figura 2.9 mostra il diagramma di flusso dei passaggi logici che il codice RISC-V dovrà rispettare; la prima espressione in C confronta le variabili contenute nei due registri per vedere se sono uguali. L'intuito ci suggerisce di utilizzare l'istruzione `beq` se i e j sono uguali. In generale, però, il codice sarà più efficiente se viene utilizzata la condizione opposta e saltiamo le istruzioni seguenti se le due variabili *non* sono uguali (`bne`). Questo è il codice risultante:

```
bne x22, x23, Else // vai a Else se i ≠ j
```

L'istruzione successiva in C è un assegnamento che svolge una singola operazione e, se tutti gli operandi si trovano all'interno dei registri, può essere tradotta in una singola istruzione RISC-V:

```
add x19, x20, x21 // f = g + h (saltata se i ≠ j)
```

Ora è necessario spostarsi alla fine dell'istruzione `if`. Questo esempio ci consente di introdurre un altro tipo di salto, spesso chiamato *salto incondizionato*: l'istruzione corrispondente indica al calcolatore di eseguire sempre il salto. Un modo per esprimere un salto incondizionato in RISC-V è utilizzare un salto condizionato a una condizione che è sempre vera:

```
beq x0, x0, Esci // se 0 == 0, vai a Esci
```

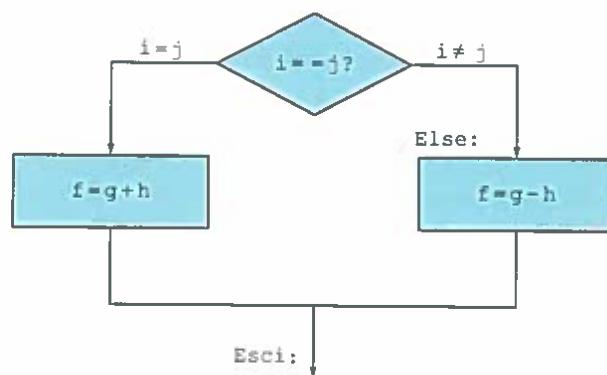
L'assegnamento che si trova nel segmento di codice relativo alla clausola `else` del costrutto `if` può anch'esso essere tradotto in un'unica istruzione RISC-V; occorre solo attaccare l'etichetta `Else` a questa istruzione. Nell'esempio viene mostrata anche l'etichetta `Esci`, posizionata dopo questa istruzione, che segna la fine del costrutto `if-then-else` tradotto in assembler:

```
Else: sub x19, x20, x21 // f = g - h (saltata se i = j)
      Esci:
```

Si noti che l'assembler evita al compilatore o al programmatore il noioso compito di calcolare gli indirizzi dei salti, proprio come nelle istruzioni di trasferimento dati viene evitato il calcolo dell'indirizzo della memoria (vedi par. 2.12).

## Interfaccia hardware/software

I compilatori creano frequentemente salti ed etichette che non appaiono nel linguaggio di programmazione. Uno dei vantaggi dei linguaggi di programmazione ad alto livello è che permettono di evitare di scrivere esplicitamente le etichette e i salti. Questa è una delle ragioni per cui programmare a questo livello è più veloce.



**Figura 2.9** Schema delle operazioni eseguite dal costrutto *if*. Il rettangolo a sinistra corrisponde al *then* del costrutto *if*, quello a destra all'*else*.

## Cicli

Prendere una decisione è importante sia per scegliere tra due alternative, come avviene nei costrutti *if*, sia per iterare un certo calcolo più volte, come avviene nei cicli. Le stesse istruzioni assembler sono utilizzate in entrambi i casi.

### Compilazione di un ciclo while scritto in C

Il seguente codice rappresenta un tipico ciclo in C:

```
while (salva[i] == k)
    i += 1;
```

Supponiamo che *i* e *k* corrispondano ai registri *x22* e *x24* e che l'indirizzo base del vettore *salva* sia contenuto in *x25*. Qual è il codice assembler RISC-V corrispondente a questo frammento di codice C?

Il primo passo consiste nel caricare *salva[i]* in un registro temporaneo. Per poter effettuare questa operazione, però, è necessario prima determinare il suo indirizzo in memoria. E prima di sommare *i* all'indirizzo di base del vettore *salva* per comporre l'indirizzo di *salva[i]*, dobbiamo moltiplicare l'indice *i* per 8 per ottenere l'indirizzamento al byte. Fortunatamente possiamo utilizzare l'operazione di shift logico a sinistra per moltiplicare per  $2^3$ , cioè 8 (vedi par. 2.6). Inoltre, è necessario aggiungere l'etichetta Ciclo alla prima istruzione, in modo che sia possibile ritornare indietro al termine del ciclo:

```
Ciclo: slli x10, x22, 3 // registro temporaneo
      // x10 = i * 8
```

Per ottenere l'indirizzo di *salva[i]* è sufficiente sommare *x10* all'indirizzo base di *salva* contenuto in *x25*:

```
add x10, x10, x25 // x10 = indirizzo di salva[i]
```

A questo punto è possibile usare questo indirizzo per caricare il valore di *salva[i]* in un registro temporaneo:

```
ld x9, 0(x10) // registro temporaneo
                // x9 = salva[i]
```

### ESEMPIO

### SOLUZIONE

(continua)

(continua)

L'istruzione successiva esegue il controllo di uscita dal ciclo: si esce dal ciclo se `salva[i] ≠ k`:

```
bne x9, x24, Esci // vai a Esci se salva[i] ≠ k
```

L'istruzione successiva somma 1 a i:

```
addi x22, x22, 1 // i = i + 1
```

Alla fine del ciclo c'è un salto all'indietro al test del *while* che si trova all'inizio del ciclo. Per completare la scrittura del codice assembler basta aggiungere l'etichetta `Esci` alla fine del codice:

```
beq x0, x0, Ciclo // vai a Ciclo  
Esci:
```

(Vedi gli esercizi sull'ottimizzazione di questa sequenza.)

## Interfaccia hardware/software

Le sequenze di istruzioni che terminano con un salto condizionato (*branch*) sono talmente fondamentali per la compilazione che viene dato loro un nome particolare: **blocchi di base**. Un blocco di base è una sequenza di istruzioni che non contiene né istruzioni di salto (con la possibile eccezione della fine) né etichette di destinazione di salti (con la possibile eccezione dell'inizio). Una delle prime fasi della compilazione consiste nel suddividere il programma nei blocchi di base da cui è composto.

**Blocco di base:** sequenza di istruzioni che non contiene né diramazioni né istruzioni di salto, ad eccezione eventualmente dell'ultima istruzione, e non contiene etichette di destinazione di salti, ad eccezione eventualmente della prima istruzione.

I test per verificare l'uguaglianza o la disuguaglianza fra due variabili sono probabilmente i più diffusi, ma esistono molte altre relazioni tra due numeri. Per esempio, all'interno di un ciclo *for* si può volere controllare se la variabile indice è minore di 0. L'insieme dei confronti possibili prevede: minore di (<), minore di o uguale ( $\leq$ ), maggiore di (>), maggiore di o uguale ( $\geq$ ), uguale (=) e non uguale ( $\neq$ ).

La comparazione di stringhe di bit deve prevedere sia i numeri dotati di segno sia quelli senza: se una stringa di bit con un 1 nel bit più significativo rappresenta un numero negativo, quel numero sarà minore di qualsiasi numero positivo, il quale dovrà avere uno 0 nel bit più significativo. Invece, se la stringa di bit rappresenta un numero senza segno e il bit più significativo è uguale a 1, quel numero sarà *maggiore* di qualsiasi altro numero con uno 0 nel bit più significativo. (Sfrutteremo a breve questo doppio significato del bit più significativo per rendere più semplice il controllo sullo sfioramento dei limiti dei vettori.) Il RISC-V offre istruzioni che gestiscono queste due possibilità. Queste istruzioni hanno la stessa forma di `beq` e `bne`, ma eseguono il confronto in modo diverso. L'istruzione di salta se minore (`blt, branch if less than`) confronta il contenuto nei registri `rs1` e `rs2` e prende il salto quando il valore contenuto in `rs1` è inferiore, trattando il contenuto dei due registri come numeri in complemento a due. L'istruzione di salta se maggiore o uguale (`bge, branch if greater than or equal`) prende il salto nel caso opposto, cioè quando il valore contenuto di `rs1` è maggiore o uguale a quello contenuto in `rs2`. Salta se minore di senza segno (`bltu, branch if less than unsigned`) prende il salto quando il valore contenuto in `rs1` è inferiore di quello in `rs2`, trattando il contenuto dei due registri come numeri senza segno. Infine, l'istruzione salta se maggiore o uguale senza segno (`bgeu, branch if greater than or equal unsigned`) prende il segno nel caso opposto.

Un'alternativa a queste altre istruzioni di salto condizionato è impostare un registro in funzione del risultato del confronto, per poi saltare o meno a

seconda del contenuto di questo registro con le istruzioni di `beq` o `bne`. Questo approccio, utilizzato dall'insieme delle istruzioni dei MIPS, può rendere l'unità di elaborazione di una CPU leggermente più semplice, ma richiede più istruzioni per scrivere un programma.

Un'ulteriore alternativa, utilizzata nell'insieme delle istruzioni ARM, è prevedere dei bit aggiuntivi per registrare quello che succede durante l'esecuzione di un'istruzione. Questi bit aggiuntivi, chiamati *codici di condizione* o *flag* (letteralmente bandiera), indicano, per esempio, se il risultato di un'operazione aritmetica è negativo, o è uguale a zero, oppure se si è verificato un overflow.

Le istruzioni di salto condizionato utilizzano quindi una combinazione di questi codici di condizione per eseguire il confronto desiderato.

Un inconveniente dei codici di condizione è che se le istruzioni che li impostano sono molte, si creano delle dipendenze che rendono difficile l'esecuzione in pipeline del codice (vedi Cap. 4).

## Scorciatoie per il controllo dei confini di vettori e matrici

Considerare un numero come non dotato di segno fornisce un modo semplice per controllare se  $0 \leq x < y$ , dove  $0$  e  $y$  sono i limiti di un dato vettore. La chiave è che i numeri negativi in complemento a 2 equivalgono a numeri molto grandi quando sono considerati nella notazione senza segno. In particolare, il bit più significativo è il bit di segno nel caso in cui il numero sia dotato di segno, mentre costituisce una parte importante del numero nel caso in cui il numero sia considerato senza segno. Perciò il confronto  $x < y$  tra due numeri senza segno valuta non solo se  $x$  sia minore di  $y$ , ma anche se  $x$  sia negativo.

### Scorciatoia per il controllo dello sfornamento dei limiti

Utilizzare la scorciatoia mostrata in precedenza per ridurre il controllo di sfornamento dei limiti di un vettore: si salti a `IndiceFuoriLimiti` se  $x_{20} \geq x_{11}$  o se  $x_{20}$  è negativo.

**ESEMPIO**

Il codice utilizza la sola istruzione di salta se maggiore uguale senza segno, `bgeu` per effettuare entrambi i controlli:

```
bgeu x20, x11, IndiceFuoriLimiti // Se x20 >= x11 o x20 < 0
                                // vai a IndiceFuoriLimiti
```

**SOLUZIONE**

## Costrutto *case/switch*

La maggior parte dei linguaggi di programmazione possiede un costrutto *case* o *switch* che consente al programmatore di scegliere fra più alternative in base al valore che assume un'unica variabile. Il costrutto *switch* può essere realizzato mediante una sequenza di condizioni, trasformandolo in una catena di costrutti *if-then-else*.

Talvolta, però, le diverse alternative di un'istruzione *switch* possono essere specificate in modo più efficiente memorizzandole in una tabella contenente gli indirizzi dell'inizio dei frammenti di codice alternativi; questa tabella viene chiamata **tabella degli indirizzi di salto** (*branch address table*) o **tabella di salto**. Il programma deve solo indicizzare la tabella e saltare all'inizio della sequenza di istruzioni corretta. La tabella è quindi semplicemente un vettore di parole che

**Tabella degli indirizzi di salto:** detta anche **tabella di salto**, contiene gli indirizzi dell'inizio dei frammenti di codice da eseguire in alternativa tra loro.

contiene gli indirizzi corrispondenti alle etichette presenti nel codice. Il programma caricherà la parola corretta dalla tabella di salto in un registro e salterà all'indirizzo caricato nel registro. Per supportare quest'operazione, alcuni calcolatori come il RISC-V hanno inserito un'istruzione di *salto indiretto*, che esegue un salto incondizionato all'indirizzo specificato nel registro. Nel RISC-V viene utilizzata l'istruzione di salta a registro (*jalr, jump and link register*). Vedremo un utilizzo molto più comune di questa versatile istruzione nel prossimo paragrafo.

## Interfaccia hardware/software

Nonostante nei linguaggi ad alto livello come C e Java ci siano molti costrutti per le decisioni o per i cicli, a livello di linguaggio assembler il costrutto base che li implementa tutti è il salto condizionato.

### Autovalutazione

- I. Il linguaggio C ha molti costrutti per le decisioni e per i cicli, mentre il RISC-V ne ha pochi. Quali delle seguenti affermazioni sono una spiegazione di questo sbilanciamento e quali non lo sono? Perché?
  1. Più costrutti decisionali rendono il codice facile da leggere e da capire.
  2. Meno costrutti decisionali semplificano il compito del livello sottostante che è responsabile dell'esecuzione.
  3. Più costrutti decisionali significano meno linee di codice, il che generalmente riduce il tempo di programmazione.
  4. Più costrutti decisionali significano meno linee di codice, il che generalmente comporta l'esecuzione di meno operazioni.
  
- II. Perché il linguaggio C fornisce due tipi di operazioni per l'AND (`&` e `&&`) e due tipi di operazioni per l'OR (`|` e `||`), mentre il RISC-V no?
  1. Le operazioni logiche AND e OR sono implementate con gli operatori `&` e `|` mentre le condizioni multiple verificate nei salti condizionati sono implementate con `&&` e `||`.
  2. L'opposto della precedente affermazione: `&&` e `||` corrispondono alle operazioni logiche mentre `&` e `|` a quelle utilizzate nei salti condizionati.
  3. Sono ridondanti e significano la stessa cosa: `&&` e `||` sono semplicemente state ereditate dal linguaggio di programmazione B, il predecessore del C.

## 2.8 | Supporto hardware alle procedure

**Procedura:** sottoprogramma memorizzato che svolge un compito specifico basandosi sui parametri che gli vengono passati in ingresso.

Le **procedure**, o funzioni, costituiscono uno degli strumenti utilizzati dai programmati per organizzare i programmi in modo da renderli più comprensibili e per permettere il riutilizzo del codice. Le procedure consentono ai programmati di concentrarsi su una parte del problema alla volta; l'interfaccia fra la procedura e il resto del programma e dei dati è costituita dai parametri, i quali permettono di passare dei valori alla procedura e di restituire i risultati alla porzione di programma in cui essa è stata chiamata. L'equivalente Java di una procedura è descritto nel paragrafo 2.15 , ma per questo linguaggio l'analisi non è così diversa, poiché ha bisogno che il calcolatore faccia tutto ciò di cui ha bisogno il C. Le procedure sono un modo per implementare l'astrazione nel software.

Potete pensare alle procedure come a delle spie che lavorano in segreto, acquisiscono risorse, svolgono il loro compito, nascondono le tracce e tornano



ASTRAZIONE

al punto di partenza con i risultati desiderati; nulla viene più modificato dopo che la missione è stata completata. Inoltre, le spie operano soltanto in base a ciò che “devono sapere” e non possono fare ipotesi su chi le abbia assoldate.

In modo simile, per l'esecuzione di una procedura, un programma deve eseguire questi sei passi:

1. mettere i parametri in un luogo accessibile alla procedura;
2. trasferire il controllo alla procedura;
3. acquisire le risorse necessarie per l'esecuzione della procedura;
4. eseguire il compito richiesto;
5. mettere il risultato in un luogo accessibile al programma chiamante;
6. restituire il controllo al punto di origine, dato che la stessa procedura può essere chiamata in diversi punti di un programma.

Come detto, i registri sono il luogo che permette l'accesso ai dati più rapido, quindi è opportuno utilizzarli il più possibile: il software RISC-V, per le chiamate a procedura, utilizza i suoi 32 registri secondo queste convenzioni:

- x10-x17: otto registri argomento per il passaggio dei parametri o la restituzione dei valori calcolati;
- x1: un registro contenente l'indirizzo di ritorno per tornare al punto di origine.

Oltre a utilizzare questi registri, il linguaggio assembler RISC-V comprende un'istruzione apposita per le procedure; questa istruzione salta all'indirizzo della procedura e contemporaneamente salva l'indirizzo dell'istruzione successiva a quella di salto nel registro rd. L'**istruzione jump and link** (jal, salta e crea un collegamento) si scrive come:

```
jal x1, IndirizzoProcedura // salta a IndirizzoProcedura e
                           // scrivi l'indirizzo di ritorno in x1
```

La parola *link* del nome si riferisce al fatto che viene creato un indirizzo, o collegamento (*link*), che punta alla posizione della chiamata, così da permettere alla procedura di tornare all'indirizzo corretto. Questo “collegamento”, memorizzato nel registro x1, è detto **indirizzo di ritorno**. L'indirizzo di ritorno è necessario, perché la stessa procedura può essere chiamata da diversi punti del programma.

Per supportare ciò, i calcolatori come i RISC-V utilizzano un'istruzione di salto indiretto, come l'istruzione *jump-and-link register* (jalr), introdotta in precedenza, per la gestione del costrutto case:

```
jalr x0, 0(x1)
```

L'istruzione jalr salta all'indirizzo memorizzato nel registro x1, che è esattamente quello che desideriamo. Perciò il programma chiamante, detto anche semplicemente **chiamante**, mette i valori dei parametri da passare alla procedura nei registri da x10-x17 e utilizza l'istruzione jal x1, X per saltare alla procedura X (detta anche **chiamata**). Questa esegue le operazioni richieste, memorizza i risultati negli stessi registri dei parametri, x10-x17, e restituisce il controllo al chiamante con l'istruzione jalr x0, 0(x1).

Nell'idea di programma memorizzato è implicita la necessità di mantenere in un registro l'indirizzo dell'istruzione in esecuzione. Per ragioni storiche questo registro viene quasi sempre chiamato **program counter** (contatore di programma), abbreviato in PC nell'architettura RISC-V, anche se il nome più corretto sarebbe *instruction address register* (registro dell'indirizzo dell'istruzione). L'istruzione jal in realtà salva PC + 4 nel registro destinazione designato (di solito x1) per creare il collegamento all'indirizzo dell'istruzione successiva, predisponendo così il ritorno dalla procedura.

**Istruzione jump and link:** istruzione che esegue un salto a un dato indirizzo e che contemporaneamente salva in un registro (x1 nel caso del RISC-V) l'indirizzo dell'istruzione successiva a quella di salto.

**Indirizzo di ritorno:** collegamento che consente a una procedura di restituire il controllo all'indirizzo opportuno del programma chiamante. Nel RISC-V è mantenuto nel registro x1.

**(Programma) chiamante:** programma che invoca una procedura e le fornisce i parametri necessari.

**(Procedura) chiamata:** procedura che esegue una sequenza di istruzioni sulla base dei parametri forniti dal chiamante e che al termine restituisce il controllo al chiamante.

**Program counter (PC):** registro contenente l'indirizzo dell'istruzione in esecuzione.

**Approfondimento.** L'istruzione di "jump-and-link" può essere utilizzata anche per eseguire un salto incondizionato all'interno di una procedura utilizzando `x0` come registro destinazione. Dato che `x0` contiene sempre il valore zero, il risultato è quello di non valutare l'indirizzo di ritorno:

```
jal x0, Etichetta // Salto incondizionato a Etichetta
```

## Utilizzo di più registri

Supponiamo che un compilatore abbia bisogno, all'interno di una procedura, di un numero maggiore di registri rispetto agli otto previsti per il passaggio dei parametri. Dato che quando la missione è compiuta occorre cancellare le tracce, il contenuto dei registri utilizzati dal programma chiamante deve essere ripristinato con il valore *precedente* la chiamata. Questa è una tipica situazione nella quale è necessario copiare il contenuto dei registri in memoria, come descritto nella sezione *Interfaccia hardware/software* del paragrafo 2.3.

La struttura dati ideale per copiare il contenuto dei registri è lo **stack** (pila), una coda di tipo *last-in-first-out* (ultimo-inserito-primo-estratto). Lo stack ha bisogno di un puntatore all'indirizzo dell'ultimo dato introdotto per indicare dove la procedura successiva possa salvare il contenuto dei registri e da dove possa poi recuperare il vecchio valore. Nel RISC-V il puntatore allo stack (**stack pointer**) è il registro `x2`, chiamato anche `sp`. Lo stack pointer viene incrementato o decrementato di una parola doppia ogni volta che si toglie o si inserisce il contenuto di un registro. L'utilizzo dello stack è così comune che sono stati coniati dei termini onomatopeici per indicare il trasferimento dati da e verso lo stack: per **push** si intende l'inserimento di un dato nello stack e per **pop** la sua estrazione.

Per ragioni storiche lo stack "cresce" a partire da indirizzi di memoria alti verso indirizzi di memoria bassi. Questa convenzione implica che quando vengono inseriti dei dati nello stack, il valore dello stack pointer diminuisce; al contrario, quando i dati sono estratti dallo stack, aumenta il valore dello stack pointer e si riduce la dimensione dello stack.

Nell'esempio qui riportato sono stati utilizzati dei registri temporanei ed è stata fatta l'ipotesi che il loro valore originale dovesse essere salvato e poi ripristinato. Per evitare di salvare e ripristinare registri il cui valore non verrà mai

**Stack:** struttura dati contenente una coda di tipo *last-in-first-out* (ultimo-inserito-primo-estratto), utilizzata per salvare il contenuto dei registri.

**Stack pointer:** indirizzo dell'elemento dello stack allocato più di recente; indica il punto in cui i registri potrebbero essere salvati o dove si può trovare il vecchio valore dei registri da ripristinare. Nel RISC-V, lo stack pointer è contenuto nel registro `sp` o `x2`.

**Push:** aggiunta di un elemento nello stack.

**Pop:** estrazione di un elemento dallo stack.

## Compilazione di una procedura C che non chiama altre procedure

### ESEMPIO

Si trasformi il secondo esempio del paragrafo 2.2 in una procedura C:

```
long long int esempio_foglia(long long int g, long long
int h, long long int i, long long int j)
{
    long long int f;
    f = (g + h) - (i + j);
    return f;
}
```

Qual è il codice assembler RISC-V che si otterrebbe dalla compilazione?

### SOLUZIONE

I parametri `g`, `h`, `i` e `j` corrispondono ai registri argomento `x10`, `x11`, `x12` e `x13`, mentre `f` corrisponde a `x20`. Il programma assembler inizia con l'etichetta della procedura:

```
esempio_foglia:
```

(continua)

(continua)

Il passo successivo consiste nel salvare tutti i registri utilizzati dalla procedura. L'assegnazione C nel corpo della procedura è identica a quella del primo esempio del paragrafo 2.3 e utilizza due registri temporanei ( $x5$  e  $x6$ ); è quindi necessario salvare il contenuto dei tre registri  $x5$ ,  $x6$  e  $x20$ . Memorizziamo il contenuto di questi tre registri nello stack (operazione di push) dopo avere allocato loro lo spazio:

La Figura 2.10 mostra lo stack prima, durante e dopo la chiamata alla procedura. Le prossime tre istruzioni corrispondono al corpo della procedura e sono analoghe all'esempio del paragrafo 2.3:

```
add x5, x10, x11 // il registro x5 contiene g + h  
add x6, x12, x13 // il registro x6 contiene i + j  
sub x20, x5, x6 // f = x5 - x6, cioè (g + h) - (i + j)
```

Per restituire il valore di  $f$ , occorre copiarlo in un registro dei parametri:

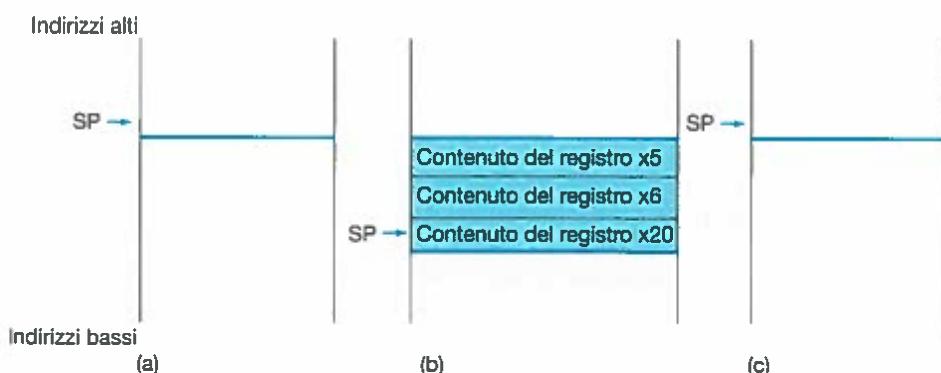
```
addi x10, x20, 0 // restituzione di f (x10 = x20 + 0)
```

Prima del ritorno al programma chiamante, il vecchio valore dei registri deve essere ripristinato, copiandolo (operazione di pop) dallo stack:

```
ld x20, 0(sp)      // ripristino del registro x20 per il chiamante
ld x6, 8(sp)       // ripristino del registro x6 per il chiamante
ld x5, 16(sp)      // ripristino del registro x5 per il chiamante
addi sp, sp, 24     // aggiornamento dello stack pointer con
                    // l'eliminazione di 3 elementi
```

La procedura termina con un'istruzione di salto a registro, utilizzando l'indirizzo di ritorno:

```
jalr x0, 0(x1) // ritorno al programma chiamante
```



**Figura 2.10** Contenuto dello stack pointer e situazione dello stack prima (a), durante (b) e dopo (c) la chiamata a procedura. Lo stack pointer punta sempre alla cima (top) dello stack, cioè all'ultima parola doppia presente nello stack.

utilizzato, cosa che può accadere con i registri temporanei, il software RISC-V suddivide 19 dei registri in due gruppi:

- $x5-x7$  e  $x28-x31$ : registri temporanei, che non sono salvati in caso di chiamata (dalla procedura chiamata) di una procedura;
- $x8-x9$  e  $x18-x27$ : registri da salvare il cui contenuto deve essere preservato in caso di chiamata a procedura (se vengono utilizzati dalla procedura chiamata, il loro contenuto deve essere prima salvato).

Questa semplice convenzione riduce il numero di registri da salvare in memoria. Nell'esempio precedente, dato che il programma chiamante non si aspetta che i registri  $x5$  e  $x6$  siano preservati dalla procedura, possiamo eliminare dal codice della stessa due istruzioni di memorizzazione e due di lettura. Dobbiamo invece salvare e ripristinare il registro  $x20$ , dal momento che la procedura è costretta a supporre che il chiamante abbia bisogno del suo valore originale.

### Procedure annidate

Le procedure che non chiamano altre procedure sono chiamate procedure *foglia*; la vita del programmatore sarebbe più semplice se tutte le procedure fossero di questo tipo, ma non è così. Proprio come una spia può arruolare altre spie per compiere parte di una missione, e queste a loro volta possono arruolarne altre ancora, le procedure possono chiamare altre procedure. Inoltre, le procedure ricorsive invocano perfino "cloni" di se stesse. Se occorre fare attenzione all'utilizzo dei registri nelle procedure, occorre stare ancora più attenti se una procedura non è di tipo foglia.

Si supponga, per esempio, che il programma principale chiami la procedura A con un parametro di valore 3, mettendo quindi il valore 3 nel registro  $x10$  e utilizzando subito dopo l'istruzione `jal x1, A`. Si supponga poi che la procedura A chiami a sua volta la procedura B con una `jal x1, B`, passandole il valore 7, anch'esso posto in  $x10$ . Dato che A non ha ancora terminato il proprio compito, si verifica un conflitto nell'utilizzo del registro  $x10$ ; in modo analogo si verifica un conflitto sull'indirizzo di ritorno di B. Se non si adottassero delle contromisure, questo conflitto farebbe sì che la procedura A non sia più in grado di restituire il controllo al suo chiamante.

Una soluzione consiste nel salvare nello stack tutti gli altri registri, il cui contenuto deve essere conservato, esattamente come abbiamo fatto per i registri da salvare. Il programma chiamante, quindi, salverà nello stack qualunque registro argomento ( $x10-x17$ ) o registro temporaneo ( $x5-x7$  e  $x28-x31$ ) quando gli servirà il contenuto di questi registri dopo la chiamata; il chiamato invece salverà nello stack il registro di ritorno,  $x1$ , e quei registri da preservare ( $x8-x9$  e  $x18-x27$ ) di cui ha bisogno. Lo stack pointer,  $sp$ , sarà aggiornato per tener conto del numero di registri da memorizzare nello stack; al termine della procedura, il contenuto dei registri verrà ripristinato con il contenuto dello stack e lo stack pointer sarà riaggiornato.

### Compilazione di una procedura ricorsiva C come esempio dei collegamenti tra procedure annidate

#### ESEMPIO

Analizziamo una procedura ricorsiva che calcola i numeri fattoriali:

```
Long long int fatt (long long int n)
{
    if (n < 1) return (1);
    else return (n * fatt(n-1));
}
```

(continua)

(continua)

Quale sarà il codice assembler RISC-V?

Il parametro `n` corrisponde al registro argomento `x10`. Il programma compilato inizierà con l'etichetta della procedura seguita dal salvataggio sullo stack di due registri: l'indirizzo di ritorno e `x10`:

```
fatt:
    addi sp, sp, -16 // aggiorna lo stack per fare
                      // posto a 2 elementi
    sd x1, 8(sp)   // salvataggio dell'indirizzo
                      // di ritorno
    sd x10, 0(sp)  // salvataggio del parametro n
```

**SOLUZIONE**

La prima volta che la procedura fattoriale viene chiamata, l'istruzione `sd` salva in stack un indirizzo del programma chiamante. Le due istruzioni successive verificano se  $n < 1$ ; se  $n \geq 1$ , la procedura salta a `L1`.

```
addi x5, x10, 1 // x5 = n - 1
bge x5, x0, L1 // se (n - 1) >= 0, salta a L1
```

Se  $n$  è minore di 1, fattoriale restituisce il valore 1 mettendolo in un registro valore: somma 1 a 0 e memorizza il risultato in `x10`. Elimina quindi dallo stack i due registri e salta all'indirizzo di ritorno:

```
addi x10, x0, 1 // restituisci 1
addi sp, sp, 16 // elimina 2 elementi dallo stack,
jalr x0, 0(x1) // ritorna al programma chiamante
```

Prima di eliminare i due elementi dallo stack, avremmo potuto salvare `x1` e `x10`; tuttavia, dato che `x1` e `x10` non cambiano quando  $n$  è minore di 1, possiamo evitare queste due istruzioni.

Se  $n$  non è minore di 1, il parametro `n` viene decrementato e viene nuovamente chiamata la procedura fattoriale passandole il nuovo valore:

```
L1: addi x10, x10, -1 // n >= 1: l'argomento diventa (n - 1)
    jal x1, fatt      // chiamata a fatt con (n - 1)
```

L'istruzione successiva è quella a cui la procedura fattoriale ritornerà al termine della sua esecuzione; il risultato prodotto si trova in `x10`. Occorre ora ripristinare il vecchio indirizzo di ritorno e il vecchio valore del parametro, e lo stack pointer:

```
addi x6, x10, 0 // ritorno da jal; sposta il risultato
                  // di fatt(n - 1) in x6:
ld x10, 0(sp)   // ripristino dell'argomento, n
ld x1, 8(sp)    // ripristino del registro di ritorno
addi sp, sp, 16 // aggiorna lo stack pointer per
                  // eliminare 2 elementi
```

Successivamente, nel registro argomento `x10` viene memorizzato il prodotto del vecchio argomento con il risultato di `fatt(n - 1)` contenuto in `x6`. Supponiamo di avere a disposizione un'istruzione per la moltiplicazione (anche se questa verrà descritta nel Capitolo 3):

```
mul x10, x10, x6 // restituisci n * fatt(n - 1)
```

Infine, la procedura `fatt` salta nuovamente all'indirizzo di ritorno:

```
jalr x0, 0(x1) // ritorna al chiamante
```

## Interfaccia hardware/software

**Global pointer:** registro riservato all'indirizzamento dei dati statici.

**Record di attivazione:** detto anche **frame della procedura**, è il segmento dello stack che contiene i registri salvati da una procedura e le variabili locali.

**Frame pointer:** valore che individua la posizione dei registri salvati e delle variabili locali di una data procedura.

Una variabile di un programma in C corrisponde in genere a una locazione di memoria e la sua interpretazione dipende sia dal *tipo* sia dalla *modalità di allocazione*. Un tipico esempio sono gli interi e i caratteri (par. 2.9). Il C ha due modalità di memorizzazione: *automatica* e *statica*. Le variabili automatiche sono locali a una procedura e vengono eliminate all'uscita da quest'ultima; le variabili statiche, invece, continuano a esistere anche quando viene chiamata una procedura o la procedura ritorna al chiamante. Le variabili C dichiarate al di fuori di tutte le procedure sono considerate statiche, così come tutte le variabili precedute dalla parola chiave *static*; tutte le altre variabili sono allocate dinamicamente. Per semplificare l'accesso ai dati statici, il software RISC-V utilizza un altro registro, chiamato **global pointer** (x3, puntatore globale).

La **Figura 2.11** riassume ciò che viene preservato nella chiamata a una procedura. Si noti che si utilizzano diversi accorgimenti per salvare il contenuto dello stack, in modo da garantire che il chiamante ottenga dallo stack gli stessi valori che vi aveva memorizzato. La parte di stack che sta sopra lo **sp** viene preservata semplicemente assicurandosi che il chiamato non scriva sopra lo stack pointer; **sp** viene a sua volta preservato dal chiamato, che sommerà sempre al valore corrente di **sp** esattamente lo stesso valore che gli aveva sottratto; gli altri registri vengono preservati salvandoli in stack (nel caso in cui vengano utilizzati) per essere ripristinati in un secondo tempo.

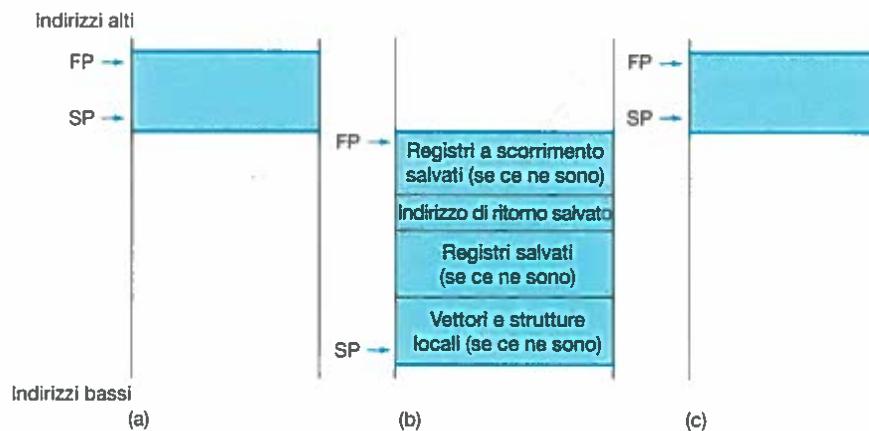
### Allocazione dello spazio nello stack per nuovi dati

Un'ultima complicazione consiste nel fatto che lo stack può essere utilizzato anche per memorizzare le variabili locali della procedura che non trovano spazio nei registri, per esempio i vettori o le strutture locali. Il segmento dello stack che contiene i registri salvati da una procedura e le variabili locali prende il nome di **record di attivazione** (blocco di attivazione), o **frame della procedura** (cornice della procedura). La **Figura 2.12** mostra lo stato dello stack prima, durante e dopo la chiamata di una procedura.

Alcuni programmi RISC-V utilizzano un **frame pointer** (**fp**, puntatore a frame, registro x8) per puntare alla prima parola doppia del frame di una procedura. Lo stack pointer potrebbe cambiare durante l'esecuzione di una procedura; in questo caso il riferimento alle variabili locali in memoria potrebbe assumere offset diversi a seconda della loro posizione nella procedura, rendendo così le procedure più difficili da leggere. In questo caso il frame pointer può essere vantaggiosamente utilizzato come registro di base stabile all'interno della procedura, per l'accesso in memoria alle variabili locali. Si noti che il record di attivazione è presente nello stack indipendentemente dal fatto che il frame pointer sia usato esplicitamente o meno. Nell'esempio precedente non abbiamo utilizzato **fp**, evitando di modificare **sp** all'interno della procedura; la dimensione dello stack veniva modificata solo all'entrata e all'uscita da essa.

Conservato	Non conservato
Registri da salvare: x8-x9, x18-x27	Registri temporanei: x5-x7, x28-x31
Registro stack pointer: x2 (sp)	Registri argomento/risultato: x10-x17
Frame pointer: x8 (fp)	
Registro dell'indirizzo di ritorno: x1 (ra)	
Stack al di sopra dello stack pointer	Stack al di sotto dello stack pointer

**Figura 2.11** Che cosa viene e che cosa non viene conservato in una chiamata a procedura. Se il software utilizza il registro frame pointer o il registro global pointer, che esamineremo in seguito, anche il contenuto di questi registri deve essere preservato.



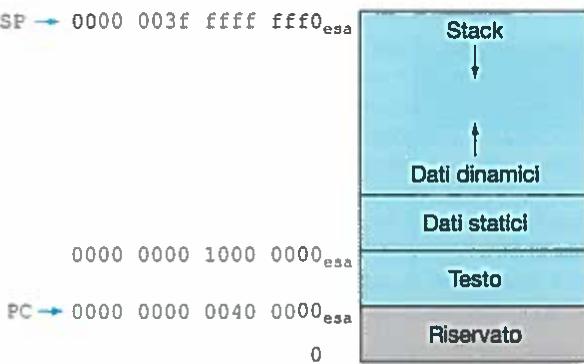
**Figura 2.12 Allocazione dello stack prima (a), durante (b) e dopo (c) la chiamata a procedura.** Il frame pointer (fp o  $\times 8$ ) punta alla prima parola doppia del frame, spesso contenente un registro da preservare, e lo stack pointer (sp) punta alla cima dello stack. Lo stack viene aggiornato per fare spazio a tutti i registri da preservare e alle variabili locali che devono essere salvate in memoria. Dal momento che lo stack pointer può cambiare durante l'esecuzione del programma, è più facile per i programmati indirizzare le variabili in stack tramite il frame pointer, che invece rimane stabile (anche se ciò potrebbe essere fatto con il solo stack pointer e un po' di aritmetica sugli indirizzi). Se lo stack non contiene variabili locali alla procedura, il compilatore risparmia tempo di esecuzione evitando di impostare e ripristinare il frame. Quando viene utilizzato, il frame pointer viene inizializzato con l'indirizzo che ha sp all'alto della chiamata della procedura e sp viene ripristinato al termine della procedura utilizzando il valore di fp. Si possono trovare queste informazioni anche nella quarta colonna della scheda tecnica riassuntiva del RISC-V in fondo al libro.

## Allocazione dello spazio nello heap per nuovi dati

In aggiunta alle variabili ad allocazione dinamica, locali alle procedure, i programmati C hanno bisogno di spazio in memoria per le variabili di tipo statico e per le strutture dati dinamiche esterne alle procedure. La Figura 2.13 mostra la convenzione utilizzata dal RISC-V per l'allocazione della memoria quando viene utilizzato il sistema operativo Linux. Lo stack inizia nella parte alta della memoria dello spazio di indirizzamento utente (*vedi Capitolo 5*) e cresce verso il basso. Dal lato opposto, la prima parte di memoria che si incontra è riservata, seguita dalla zona in cui risiede il codice macchina RISC-V, tradizionalmente chiamata **segmento di testo**. Sopra il codice si trova il *segmento dei dati statici*, che è il luogo in cui risiedono le costanti e le altre variabili statiche. Nonostante i vettori siano di lunghezza fissa e quindi siano adatti a essere allocati nella zona dei dati statici, strutture come le liste concatenate hanno la caratteristica di crescere e di ridursi durante la loro vita. Il segmento di memoria dati utilizzato per queste strutture è chiamato tradizionalmente **heap** ("cumulo") ed è posizionato sopra la porzione di memoria statica. Si noti che con questa modalità di allocazione lo stack e lo heap crescono uno verso l'altro, permettendo un uso efficiente della memoria nonostante la dimensione dei due segmenti sia altalenante.

Il linguaggio C alloca e libera spazio nello heap mediante funzioni esplicite: `malloc()` alloca spazio nello heap e ritorna il puntatore al primo indirizzo dello spazio allocato, e `free()` libera lo spazio nello heap puntato dal puntatore. L'allocazione della memoria in C viene perciò controllata dal programma ed è fonte di molti errori comuni e difficili da localizzare. Dimenticare di liberare lo spazio allocato provoca una "perdita" di memoria che potrebbe diventare tanto rilevante da portare al *crash* (blocco) del sistema operativo; liberare la memoria troppo presto, invece, rende i puntatori "volanti" (*dangling pointers*), ossia questi potrebbero puntare a cose molto diverse da quelle che interessavano al programma. Java utilizza l'allocazione e la liberazione automatica della memoria proprio per evitare questi errori.

**Segmento di testo:** segmento di memoria che contiene il codice in linguaggio macchina di tutte le procedure di un'applicazione, contenute in un file oggetto UNIX.



**Figura 2.13 Allocazione della memoria per programmi e dati nel RISC-V.** Gli indirizzi indicati sono frutto solo di convenzioni software e non costituiscono parte integrante dell'architettura del RISC-V. Lo spazio indirizzamento utente è impostato a  $2^{38}$  rispetto ai teorici  $2^{64}$  valori che si possono raggiungere con un'architettura a 64 bit (vedi Cap. 5). Lo stack pointer viene inizializzato a 0000 003f ffff ffff0<sub>esa</sub> e decresce verso il basso verso il segmento dati. Dall'altra parte, il segmento di testo, contenente il codice del programma, parte dall'indirizzo 0000 0000 0040 0000<sub>esa</sub>. I dati statici nell'esempio partono immediatamente dopo il segmento testo, assumiamo che questo indirizzo sia 0000 0000 1000 0000<sub>esa</sub>. I dati dinamici, che vengono allocati con una malloc in C o con una new in Java, sono memorizzati subito dopo i dati statici: questa area di memoria cresce verso l'alto (verso lo stack), ed è denominata *heap*. Si possono trovare queste informazioni anche nella scheda tecnica riassuntiva del RISC-V in fondo al libro.



La Figura 2.14 riassume la convenzione di utilizzo dei registri nel linguaggio assembler RISC-V. Questa convenzione è un altro esempio di come si possano rendere veloci le situazioni più comuni: la maggiore parte delle procedure richiede non più di otto registri per gli argomenti, dodici registri per variabili permanenti e sette registri temporanei, senza dover ricorrere alla memoria.

**Approfondimento.** Che cosa succederebbe se si volessero passare più di otto parametri a una procedura? Per convenzione i parametri aggiuntivi vengono messi nello stack al di sopra dell'indirizzo puntato dal frame pointer: la procedura si aspetterà quindi i primi otto parametri nei registri x10-x17 e i restanti nell'area di stack, indirizzabili attraverso il frame pointer.

Come citato nella didascalia della Figura 2.12, il frame pointer è comodo perché fa sì che i riferimenti alle variabili in stack contenuti all'interno di una procedura mantengano lo stesso offset. Il frame pointer però non è strettamente necessario: il compilatore per il RISC-V utilizza solo un frame pointer nelle procedure che modificano lo stack pointer all'interno del corpo della procedura.

**Figura 2.14 Convenzioni di utilizzo dei registri RISC-V.** Si possono trovare queste informazioni anche nella seconda colonna della scheda tecnica riassuntiva del RISC-V in fondo al libro.

Nome	Numero del registro	Utilizzo	Da conservare nella chiamata?
x0	0	Costante 0	n.a.
x1 (ra)	1	Registro di ritorno (registro di collegamento)	sì
x2 (sp)	2	Stack pointer	sì
x3 (gp)	3	Global pointer	sì
x4 (tp)	4	Thread pointer	sì
x5-x7	5-7	Variabili temporanee	no
x8-x9	8-9	Variabili da preservare	sì
x10-x17	10-17	Argomenti/risultati	no
x18-x27	18-27	Variabili da preservare	sì
x28-x31	28-31	Variabili temporanee	no

**Approfondimento.** Alcune procedure ricorsive possono essere implementate attraverso iterazioni che non implicano la ricorsione. L'utilizzo delle iterazioni migliora le prestazioni eliminando il sovraccarico di lavoro associato alle chiamate di procedura. Per esempio, si consideri una procedura che effettua delle somme ricorsive:

```
long long int somma (long long int n, long long int acc) {
    if (n > 0)
        return somma(n - 1, acc + n);
    else
        return acc;
}
```

Esaminiamo la procedura somma (3, 0). Si può osservare che essa chiamerà ricorsivamente le procedure somma (2, 3), somma (1, 5) e somma (0, 6); il risultato, 6, verrà quindi restituito alla procedura chiamante quattro volte. Questo tipo di chiamata ricorsiva viene detto anche *tail call* (chiamata ad anello) e in questo esempio può essere implementato in modo molto efficiente come segue (si assuma  $x10 = n$  e  $x11 = acc$  e il risultato in  $x12$ ):

```
somma: ble x10, x0, somma_esci // vai a somma_esci se n <= 0
        add x11, x11, x10      // aggiungi n ad acc
        addi x10, x10, -1      // sottrai 1 a n
        jal x0, somma         // vai a somma
somma_esci:
        addi x12, x11, 0       // restituisci il valore acc
        jalr x0, 0(x1)         // ritorna al chiamante
```

## Autovalutazione

Quali delle seguenti affermazioni su C o su Java sono generalmente vere?

1. I programmati C gestiscono i dati esplicitamente, mentre la gestione dei dati in Java è automatica.
2. Il linguaggio C induce più errori sui puntatori e sull'allocazione dinamica della memoria di quanto faccia Java.

## 2.9 | Comunicare con le persone

I calcolatori furono inventati per effettuare calcoli matematici complessi, ma non appena diventarono commercialmente disponibili iniziarono a essere utilizzati anche per l'elaborazione di testi. Oggi la maggior parte dei calcolatori utilizza 8 bit (1 byte) per rappresentare i caratteri, e la codifica ASCII (*American Standard Code for Information Interchange*) è di gran lunga la più utilizzata. La Figura 2.15 riassume la codifica ASCII.

!(@=> (wow open tab at bar is great - oh! il conto aperto al bar è fantastico)

La quarta riga del poema della tastiera "Hatless Atlas" 1991 (alcuni associano a ogni carattere ASCII una parola: "!" è "wow", "(" è "open", "@" è "at", "]" è "bar" ecc.).

### Confronto tra codifica ASCII e codifica binaria

Potremmo rappresentare i numeri come stringhe di caratteri ASCII invece che come numeri interi. Di quanto aumenterebbe la richiesta di memoria se il numero 1 miliardo venisse rappresentato sotto forma di stringa di caratteri invece che come intero su 32 bit?

### ESEMPIO

(continua)

(continua)

**SOLUZIONE**

Il numero 1 miliardo è rappresentato dalla stringa 1 000 000 000, che quindi richiede 10 caratteri ASCII, ciascuno di 8 bit. La richiesta di memoria sarebbe pari a  $(10 \times 8)/32$ , ovvero 2,5 parole. Oltre a richiedere più memoria, l'hardware diventerebbe molto più complesso e assorbirebbe una maggiore quantità di energia per compiere le operazioni di somma, sottrazione, moltiplicazione e divisione in base 10. Questi problemi spiegano come mai per un calcolatore la codifica binaria è naturale, mentre quella in base 10 risulta non funzionale.

Valore ASCII	Carattere										
32	Spazio	48	0	64	@	80	P	096	'	112	p
33	!	49	1	65	A	81	Q	097	a	113	q
34	"	50	2	66	B	82	R	098	b	114	r
35	#	51	3	67	C	83	S	099	c	115	s
36	\$	52	4	68	D	84	T	100	d	116	t
37	%	53	5	69	E	85	U	101	e	117	u
38	&	54	6	70	F	86	V	102	f	118	v
39	'	55	7	71	G	87	W	103	g	119	w
40	(	56	8	72	H	88	X	104	h	120	x
41	)	57	9	73	I	89	Y	105	i	121	y
42	*	58	:	74	J	90	Z	106	j	122	z
43	+	59	;	75	K	91	[	107	k	123	{
44	,	60	<	76	L	92	\	108	l	124	
45	-	61	=	77	M	93	]	109	m	125	}
46	.	62	>	78	N	94	^	110	n	126	~
47	/	63	?	79	O	95	_	111	o	127	DEL

**Figura 2.15 Codifica ASCII dei caratteri.** Si noti che le lettere maiuscole e minuscole differiscono esattamente di 32; ciò consente di utilizzare una scorciatoia per controllare e convertire le minuscole in maiuscole e viceversa. I valori non mostrati sono associati ai caratteri di formattazione: per esempio, il carattere 8 rappresenta lo spostamento di uno spazio a sinistra, il 9 la tabulazione, mentre il 13 rappresenta il ritorno a capo. Un altro carattere particolarmente utile è lo 0, che rappresenta il carattere "null", che viene utilizzato nel linguaggio di programmazione C per segnalare la fine di una stringa.

Una serie di istruzioni consente di estrarre un byte da una parola doppia, quindi le istruzioni di trasferimento tra memoria e registri (caricamento di un registro e salvataggio in memoria) sono sufficienti anche per trasferire byte. L'elaborazione di testi è però talmente diffusa che il RISC-V offre istruzioni apposite per trasferire i singoli byte: *load byte unsigned* (1bu) prende un byte dalla memoria mettendolo negli 8 bit di un registro, collocati più a destra, mentre *store byte* (sb) prende il byte corrispondente agli 8 bit di un registro, collocati più a destra, e lo salva in memoria. Perciò, possiamo copiare un byte con la seguente coppia di istruzioni:

```
lbu x12, 0(x10) // Leggi un byte dall'indirizzo sorgente
sb  x12, 0(x11) // Scrivi il byte all'indirizzo di
                  // destinazione
```

Normalmente i caratteri vengono raggruppati in stringhe, composte da un numero variabile di caratteri. Ci sono tre modi per rappresentare una stringa: (1) il primo elemento della stringa viene utilizzato per definirne la lunghezza, (2) una variabile di appoggio contiene la lunghezza della stringa (come in una struttura) oppure (3) l'ultimo elemento della stringa è individuato da un carattere particolare di fine stringa. Il linguaggio C utilizza questo terzo modo e termina le stringhe con un byte che contiene il valore 0 (carattere "null" in ASCII). In C, quindi, la stringa "Cal" viene rappresentata da 4 byte, i cui valori decimali sono rispettivamente 67, 97, 108 e 0 (come vedremo in seguito, Java invece utilizza il primo modo).

### Compilazione di una procedura per copiare una stringa come esempio di utilizzo delle stringhe nel C

La procedura `copiaStringa` copia la stringa `y` nella stringa `x`, utilizzando il byte null come carattere di fine stringa:

```
void copiaStringa(char x[], char y[])
{
    size_t i;
    i = 0;
    while ((x[i] = y[i]) != '\0') /* copia & controllo del byte
        i += 1;
}
```

Qual è il codice assembler RISC-V corrispondente?

Il frammento di codice assembler RISC-V viene mostrato di seguito. Si assuma che gli indirizzi di base dei vettori `x` e `y` si trovino rispettivamente in `x10` e `x11`, mentre `i` sia contenuto in `x19`. `copiaStringa` aggiorna lo stack pointer e poi salva il registro da preservare `x19` nello stack:

```
copiaStringa:
    addi sp, sp, -8 // aggiorna lo stack
                  // per inserire un elemento
    sd x19, 0(sp) // salva x19
```

Per inizializzare `i` a 0, l'istruzione successiva imposta `x19` a 0 sommando 0 con 0 e mettendo il risultato in `x19`:

```
add x19, x0, x0 // i = 0 + 0
```

Questo è l'inizio del ciclo. Come prima cosa, l'indirizzo di `y[i]` viene creato sommando `i` a `y[]` e ponendo il risultato in `x5`:

```
I1: add x5, x19, x11 // indirizzo di y[i] in x5
```

Si noti che non è necessario moltiplicare `i` per 8, dato che `y` è un vettore di byte e non di parole doppie come nei casi precedenti.

Per caricare il carattere presente in `y[i]` utilizziamo l'istruzione "*load byte unsigned*" e mettiamo il carattere letto in `x6`:

```
lbu x6, 0(x5) // x6 = y[i]
```

Analogamente calcoliamo l'indirizzo di `x[i]` e lo scriviamo in `x7`; poniamo quindi il carattere che si trova in `x6` nella locazione di memoria

### ESEMPIO

### SOLUZIONE

(continua)

(continua)

a quell'indirizzo:

```
add x7, x19, x10 // indirizzo di x[i] in x7
sb x6, 0(x7) // x[i] = y[i]
```

Successivamente possiamo uscire dal ciclo se il carattere letto è 0, cioè se si tratta del carattere di fine stringa:

```
beq x6, x10, L2
```

Altrimenti si incrementa *i* di 1 e si torna all'inizio del ciclo:

```
addi x19, x19, 1 // i = i + 1
jal x0, L1 // salta a L1
```

Se non si ritorna all'inizio del ciclo vuol dire che la stringa è terminata; in questo caso viene ripristinato il valore di *x19* e dello stack pointer e si esce dalla procedura:

```
L2: ld x19, 0(sp) // ripristina il vecchio x19
    addi sp, sp, 8 // aggiorna lo stack eliminando
                    // un elemento
    jral x0, 0(x1) // ritorna al chiamante
```

In C la copia di stringhe di solito avviene tramite puntatori e non tramite vettori per evitare di effettuare operazioni sull'indice *i*. Vedi il paragrafo 2.14 per un confronto tra l'utilizzo di vettori e puntatori.

Dal momento che la procedura *copiaStringa* appena descritta è una procedura di tipo foglia, il compilatore può allocare *i* in un registro temporaneo evitando di salvare e ripristinare *x19*. Pertanto, invece di considerare questi registri soltanto come registri temporanei, è più corretto considerarli come registri che le procedure possono utilizzare a loro piacimento. Quando il compilatore incontra una procedura foglia esaurisce tutti i registri temporanei prima di utilizzare quelli di cui deve salvare il contenuto.

### Caratteri e stringhe in Java

L'*Unicode* è una codifica universale degli alfabeti di gran parte delle lingue utilizzate sul nostro pianeta. La Figura 2.16 riporta un elenco di alfabeti inclusi nell'*Unicode*. Come si può notare, il numero di alfabeti è quasi uguale a quello dei simboli nel codice ASCII. Java utilizza Unicode per rappresentare i caratteri e quindi normalmente utilizza 16 bit per rappresentare un singolo carattere.

L'insieme RISC-V contiene istruzioni esplicite per caricare e memorizzare questi gruppi di 16 bit, che sono detti *halfword* (mezze parole). L'istruzione *load half unsigned* (*lh*, carica mezza parola senza segno) legge 16 bit dalla memoria e li posiziona nei bit più a destra di un registro e riempie i rimanenti 48 bit di zero. Come *load byte*, *load half* tratta le mezze parole come numeri dotati di segno e perciò estende il segno per riempire i 48 bit più significativi del registro. *Store half* (*sh*) prende i 16 bit meno significativi di un registro e li scrive in memoria. Possiamo copiare mezza parola con la seguente coppia di istruzioni:

```
lhu x19, 0(x10) // Leggi mezza parola (16 bit)
                  // dall'indirizzo sorgente
sh x19, 0(x11) // Scrivi mezza parola (16 bit)
                  // all'indirizzo di destinazione
```

Latino	Malayalam	Tagbanwa	Punteggiatura generale
Greco	Sinhala	Khmer	Lettere che modificano la spaziatura
Cirillico	Thailandese	Mongolo	Simboli delle monete
Armeno	Laotiano	Limbu	Combinazione di segni diacritici
Ebraico	Tibetano	Tai Le	Combinazione di segni per ottenere simboli
Arabo	Myanmar	Kangxi Radicals	Apici e pedici
Siriaco	Georgiano	Hiragana	Forme numeriche
Thaana	Hangul Jamo	Katakana	Operatori matematici
Devanagari	Etiope	Bopomofo	Simboli matematici alfanumerici
Bengalese	Cherokee	Kanbun	Sequenze Braille
Gurmukhi	Sillabico Unificato degli Aborigeni Canadesi	Shavian	Caratteri per il riconoscimento ottico
Gujarati	Ogham	Osmanya	Simboli musicali bizantini
Oriya	Runico	Sillabico Cipriota	Simboli musicali
Tamil	Tagalog	Thailandese Simbolico Xuan Jing	Frecce
Telugu	Hanunoo	Simboli Yijing Hexagram	Disegno di riquadri
Kannada	Buhid	Numeri Aegean	Forme geometriche

**Figura 2.16 Esempi di alfabeti rappresentati con l'Unicode.** La versione 4.0 di Unicode contiene più di 160 "blocchi", che è il nome dato a un insieme di simboli. Ogni blocco è un multiplo di 16. Per esempio, l'alfabeto greco parte dalla posizione 0370<sub>ess</sub>, il cirillico da 0400<sub>ess</sub>. Le prime tre colonne mostrano il contenuto di 48 blocchi che corrispondono alle lingue rappresentate dall'Unicode, seguendo approssimativamente l'ordine numerico dell'Unicode. L'ultima colonna contiene 16 blocchi, considerati multilingua e che non sono in ordine. La codifica a 16 bit, chiamata UTF-16, è quella di default, mentre una codifica a lunghezza variabile, denominata UTF-8, contiene come sottosistema la codifica ASCII su 8 bit e utilizza 16-32 bit per gli altri caratteri. Infine, UTF-32 utilizza 32 bit per ogni carattere. Per saperne di più *vedi* il sito web [www.unicode.org](http://www.unicode.org).

Le stringhe costituiscono una classe standard Java, con speciali supporti pre-costituiti e metodi predefiniti per la loro concatenazione, comparazione e conversione. A differenza del C, Java prevede una parola che contiene la lunghezza della stringa, in modo simile a quanto fa per i vettori.

**Approfondimento.** Il software RISC-V cerca di mantenere lo stack allineato agli indirizzi di quattro parole ("quad word" o parola quadrupla, 16 byte) per ottenere prestazioni migliori. Questa convenzione comporta che una variabile di tipo char nello stack possa occupare fino a 16 byte, anche se ne servirebbero di meno. In ogni caso, le variabili di tipo stringa o i vettori di byte in C vengono compattati in modo tale da inserire 16 byte per ogni parola quadrupla, mentre le stringhe Java e i vettori di valori interi codificati su 16 bit (short) vengono compattati in 8 elementi di mezza parola ciascuno per ogni parola quadrupla.

**Approfondimento.** La maggior parte delle pagine web oggigiorno utilizza il codice Unicode invece del codice ASCII, riflettendo così la natura internazionale del web. Quindi Unicode può essere considerato oggigiorno ancora più noto del codice ASCII.

**Approfondimento.** Il RISC-V contiene anche istruzioni per trasferire dati su 32 bit da e verso la memoria. *Load word unsigned* (**lwu**, caricamento di una parola senza segno) carica una parola di 32 bit dalla memoria nei 32 bit più a destra di un registro, riempiendo i 32 bit più a sinistra di zero. *Load word* (**lw**, caricamento di una parola) invece riempie i 32 bit più a sinistra con una copia del bit 31 (bit di segno). *Store word* (**sw**, salvataggio in memoria) estrae una parola dai 32 bit più a destra di un registro e la salva in memoria.

## Autovalutazione

- I. Quali delle seguenti affermazioni su caratteri e stringhe in C e Java sono vere?
  1. Una stringa in C occupa metà dello spazio di memoria rispetto alla stessa stringa in Java.
  2. Le stringhe sono solo un nome informale per indicare vettori di caratteri sia in C sia in Java.
  3. Le stringhe in C e in Java utilizzano il carattere null (0) per segnalare il termine della stringa.
  4. Le operazioni sulle stringhe, come il calcolo della loro lunghezza, sono più veloci in C che in Java.
  
- II. Quali tipi di variabili in grado di memorizzare il numero 1 000 000 000<sub>dec</sub> richiedono più spazio in memoria?
  1. long long int in C
  2. string in C
  3. string in Java

## 2.10 Indirizzamento RISC-V di un campo immediato e di un indirizzo ampio

Anche se mantenere tutte le istruzioni RISC-V a 32 bit semplifica l'hardware, ci sono casi in cui sarebbe utile avere una costante o un indirizzo a 32 bit. Questo paragrafo inizia illustrando la soluzione generale al problema delle costanti molto grandi e prosegue mostrando le ottimizzazioni adottate per comporre gli indirizzi nelle istruzioni di salto.

### Operandi immediati ampi

Sebbene le costanti siano molto spesso piccole e possano quindi trovare spazio all'interno del campo di 12 bit a loro assegnato, qualche volta sono più grandi.

Per questo motivo l'insieme di istruzioni RISC-V include l'istruzione *load upper immediate* (*lui*) che permette di caricare i 20 bit più significativi di una costante nei bit da 12 a 31 di un registro, imposta i 32 bit più a sinistra al valore del bit 31 (bit di segna) e riempie i 12 bit più a destra del registro con zero. Questa istruzione consente, per esempio, di creare una costante su 32 bit con due istruzioni. L'istruzione *lui* utilizza un nuovo formato di istruzione, tipo U, dato che gli altri formati non possono ospitare una costante così ampia.

### Caricamento di una costante su 32 bit

#### ESEMPIO

Qual è il codice assembler RISC-V utilizzato per caricare la seguente costante a 64 bit nel registro x19?

```
00000000 00000000 00000000 00000000 00000000 00111101 00000101 00000000
```

#### SOLUZIONE

Per prima cosa attraverso l'istruzione *lui* si caricano i bit dalla posizione 12 alla 31, che costituiscono il numero 976 in decimale:

```
lui x19, 976 // 976 in decimale = 0000 0000 0011 1101 0000
```

(continua)

A questo punto il registro x19 contiene:

```
00000000 00000000 00000000 00000000 00111101 00000000 00000000
```

Il passo successivo consiste nell'inserire nei 12 bit meno significativi della costante, il cui valore è 1280 in notazione decimale:

```
addi x19, x19, 1280 // 1280dec = 00000101 00000000
```

Il valore finale contenuto nel registro x19 è il valore desiderato:

```
00000000 00000000 00000000 00000000 00000000 00111101 00000101 00000000
```

(continua)

**Approfondimento.** Nell'esempio precedente, il bit 11 della costante era uguale a 0. Se il bit 11 fosse stato uguale a 1 ci sarebbe stata una complicazione aggiuntiva: viene esteso automaticamente il bit di segno del campo immediato su 12 bit, per cui l'addendo della `addi` sarebbe risultato negativo. Questo significa che oltre a sommare gli 11 bit più a destra della costante, sottrarremmo al numero finale  $2^{12}$ . Per correggere questo errore, è sufficiente aggiungere 1 alla costante caricata con la `lui`, dato che la costante caricata è scalata per  $2^{12}$ .

I compilatori o gli assemblatori devono spezzare le costanti molto grandi e successivamente assemblarle di nuovo in un registro. Come ci si può aspettare, la dimensione limitata del campo riservato alle costanti nelle istruzioni di tipo immediato può rappresentare un problema, sia per gli indirizzi di memoria nelle operazioni di trasferimento dati sia per le costanti.

In conclusione, la rappresentazione simbolica delle istruzioni in linguaggio macchina RISC-V non è più limitata dall'hardware, ma da ciò che il progettista dell'assemblatore ha scelto di supportare (par. 2.12). Noi abbiamo scelto di rimanere legati all'hardware per spiegare l'architettura del calcolatore, rendendo esplicito l'eventuale utilizzo di istruzioni avanzate dell'assemblatore che non fanno però parte dell'architettura del processore.

## Interfaccia hardware/software

### Indirizzamento nei salti

Le istruzioni di salto condizionato RISC-V utilizzano il formato chiamato di *tipo-SB*. Questo formato può rappresentare indirizzi di salto da -4096 a 4094 in multipli di 2. Per i motivi che spiegheremo fra poco, è possibile saltare solo a indirizzi pari. Il formato di tipo-SB consiste in 7 bit di codice operativo, 3 bit di codice funzione, due registri operandi su 5 bit (rs1 e rs2) e un campo immediato di indirizzo. L'indirizzo è implementato con una codifica insolita, che semplifica l'elaborazione da parte della CPU ma complica l'assembler. L'istruzione

```
bne x10, x11, 2000 // se x10 != x11, vai  
// all'indirizzo 2000dec = 0111 1101 0000
```

può essere tradotta secondo il seguente formato (in realtà le cose sono un po' più complicate come vedremo a breve):

0	111110	01011	01010	001	1000	0	1100111
imm[12]	imm[10:5]	rs2	rs1	funz3	imm[4:1]	imm[11]	codop

dove il codice operativo dei salti condizionati è 1100111<sub>due</sub> e il codice funz3 della bne è 001<sub>due</sub>.

L'istruzione di salto incondizionato *jump-and-link* (*jal*) è l'unica istruzione che utilizza il formato di tipo *UJ*. Questa istruzione consiste in un codice operativo su 7 bit, un registro operando di destinazione su 5 bit e un indirizzo immediato su 20 bit. L'indirizzo di collegamento, che è l'indirizzo dell'istruzione successiva alla *jal*, viene scritto nel campo *rd*.

Come per il formato di tipo *SB*, l'operando che contiene l'indirizzo nel formato di tipo *UJ*, utilizza una codifica insolita, e non può codificare gli indirizzi dispari. Quindi,

*jal x0, 2000 // vai all'indirizzo 2000<sub>dec</sub> = 0111 1101 0000*

viene assemblata in questo formato:

0	1111101000	0	00000000	00000	1101111
imm[20]	imm[10:1]	imm[11]	imm[19:12]	rd	codop

Se gli indirizzi del programma trovassero posto in questo campo a 20 bit, risulterebbe che nessun programma potrebbe avere una dimensione superiore a  $2^{20}$ , troppo piccola per essere utilizzata nelle applicazioni reali. Una valida alternativa consiste nello specificare un registro il cui contenuto deve essere sommato all'indirizzo del salto; l'istruzione di salto dovrebbe quindi effettuare il seguente calcolo:

$$\text{Program counter} = \text{Registro} + \text{Spiazzamento del salto}$$

Questa somma consentirebbe a un programma di indirizzare  $2^{64}$  posizioni pur continuando a utilizzare i salti condizionati, risolvendo il problema della dimensione dell'indirizzo di salto. La domanda che dobbiamo porci è: quale registro si deve utilizzare?

La risposta proviene dall'analisi di come i salti condizionati vengono utilizzati: essi si trovano tipicamente nei cicli e nei costrutti *if*, e quindi di solito eseguono salti a istruzioni vicine. Per fare un esempio, quasi la metà delle istruzioni di salto condizionato nei programmi appartenenti ai benchmark SPEC saltano a indirizzi che distano meno di 16 istruzioni.

Dal momento che il program counter (PC) contiene l'indirizzo dell'istruzione corrente, si può saltare con un salto condizionato fino a una distanza di  $\pm 2^{10}$  parole rispetto all'istruzione corrente, e con un salto incondizionato fino a  $\pm 2^{18}$  parole rispetto all'istruzione corrente, utilizzando il PC come registro da sommare all'indirizzo di salto. Dato che quasi tutti i cicli e i costrutti *if* hanno dimensione inferiore alle  $2^{10}$  parole, il PC è la scelta ideale. Questo metodo di indirizzamento per i salti si chiama **indirizzamento relativo al program counter** (*PC-relative addressing*).

Come molti dei calcolatori più recenti, il RISC-V utilizza l'indirizzamento relativo al PC per tutti i tipi di salto condizionato e incondizionato, dato che l'indirizzo di destinazione ha un'alta probabilità di essere vicino a quello dell'istruzione di salto. D'altra parte le chiamate a procedura facilmente richiedono di saltare più di  $2^{18}$  parole, dato che non è garantito che il chiamante si trovi vicino in memoria al chiamato. Quindi il RISC-V consente di effettuare salti molto lunghi a uno qualsiasi tra  $2^{32}$  indirizzi utilizzando una sequenza di due istruzioni: lui scrive i bit da 12 a 31 in un registro temporaneo e *jalr* somma i 12 bit meno significativi all'indirizzo contenuto nel registro temporaneo e salta all'indirizzo ottenuto con la somma.

Dato che le istruzioni RISC-V sono ampie 4 byte, le istruzioni di salto RISC-V sono state progettate per ampliare il loro spazio di indirizzamento definendo l'indirizzo relativo al PC in termini di numero di parole tra l'istruzione corrente

**Indirizzamento relativo al program counter:** una modalità di indirizzamento secondo la quale l'indirizzo è ottenuto mediante la somma del program counter (PC) e il campo costante contenuto nell'istruzione.

di salto e l'istruzione di destinazione del salto, invece che in termini di numero di byte. Tuttavia gli architetti RISC-V hanno voluto supportare la possibilità che le istruzioni siano ampie solo 2 byte, per cui lo spiazzamento viene definito nelle istruzioni di salto in termini di *mezze parole* che intercorrono tra l'indirizzo dell'istruzione corrente di salto e quello di destinazione del salto. Quindi il campo di indirizzi di 20 bit nell'istruzione `jal` può codificare una distanza di  $\pm 2^{19}$  mezze parole, o  $\pm 1$  MiB a partire dal valore attuale del PC. Analogamente, il campo di 12 bit delle istruzioni di salto condizionato è espresso anch'esso in termini di mezze parole; questo vuol dire che rappresenta un indirizzo di 13 bit in termini di byte.

### Calcolo dello spiazzamento dei salti condizionati in linguaggio macchina

Il ciclo `while` riportato nel paragrafo 2.7 è stato compilato nelle seguenti istruzioni assembler:

```
Ciclo:slli x10, x22, 3 // registro temporaneo x10 = i * 8
    add x10, x10, x25 // x10 = indirizzo di salva[i]
    ld x9, 0(x10) // registro temporaneo x9 = salva[i]
    bne x9, x24, Esci // vai a Esci se salva[i] ≠ k
    addi x22, x22, 1 // i = i + 1
    beq x0, x0, Ciclo // vai a Ciclo
Esci:
```

Se si suppone che il ciclo sia contenuto in memoria a partire dalla locazione 80000, quale sarà il codice macchina RISC-V corrispondente?

Le istruzioni in linguaggio macchina e gli indirizzi corrispondenti vengono riportati di seguito.

Indirizzo	Istruzione					
80000	0000000	00011	10110	001	01010	0010011
80004	0000000	11001	01010	000	01010	0110011
80008	0000000	00000	01010	011	01001	0000011
80012	0000000	11000	01001	001	01100	1100011
80016	0000000	00001	10110	000	10110	0010011
80020	1111111	00000	00000	000	01101	1100011

Si ricordi che le istruzioni RISC-V indirizzano il byte, quindi gli indirizzi di parole consecutive differiscono di quattro unità, cioè il numero di byte in una parola. L'istruzione `bne` della quarta riga somma 12 byte (e quindi 3 parole) all'indirizzo dell'istruzione; specifica quindi l'indirizzo di destinazione del salto come distanza dall'istruzione di salto ( $12 + 80012$ ) e non come indirizzo assoluto dell'istruzione di arrivo (80024). L'istruzione di branch dell'ultima riga utilizza un calcolo simile per saltare all'indietro ( $-20 + 80020$ ), che corrisponde all'etichetta `Ciclo`.

### ESEMPIO

### SOLUZIONE

## Interfaccia hardware/software

La quasi totalità delle istruzioni di salto condizionato ha come destinazione una locazione di memoria vicina, ma occasionalmente può accadere di dover eseguire un salto condizionato a una distanza tale che non può essere rappresentata nei 12 bit del campo indirizzo di questo tipo di istruzioni. L'assemblatore risolve il problema nello stesso modo in cui risolveva il problema degli indirizzi o delle costanti molto grandi: inserisce un salto incondizionato all'istruzione di arrivo del salto, e inverte il test sui registri in modo tale che l'istruzione di salto condizionato decida se saltare o meno l'istruzione di salto incondizionato.

### ESEMPIO

#### Salto a indirizzi lontani

Si consideri il seguente salto condizionato che viene preso quando il contenuto di  $x10$  è uguale a zero:

```
beq x10, x0, L1
```

Si sostituisca l'istruzione `beq` con una coppia di istruzioni che permetta di saltare a una distanza molto maggiore.

### SOLUZIONE

Queste istruzioni sostituiscono il salto condizionato a un indirizzo vicino:

```
bne x10, x0, L2
jal x0, L1
```

L2:

**Modalità di indirizzamento:** uno dei possibili modi di indirizzamento definito dal diverso utilizzo degli operandi e/o degli indirizzi.

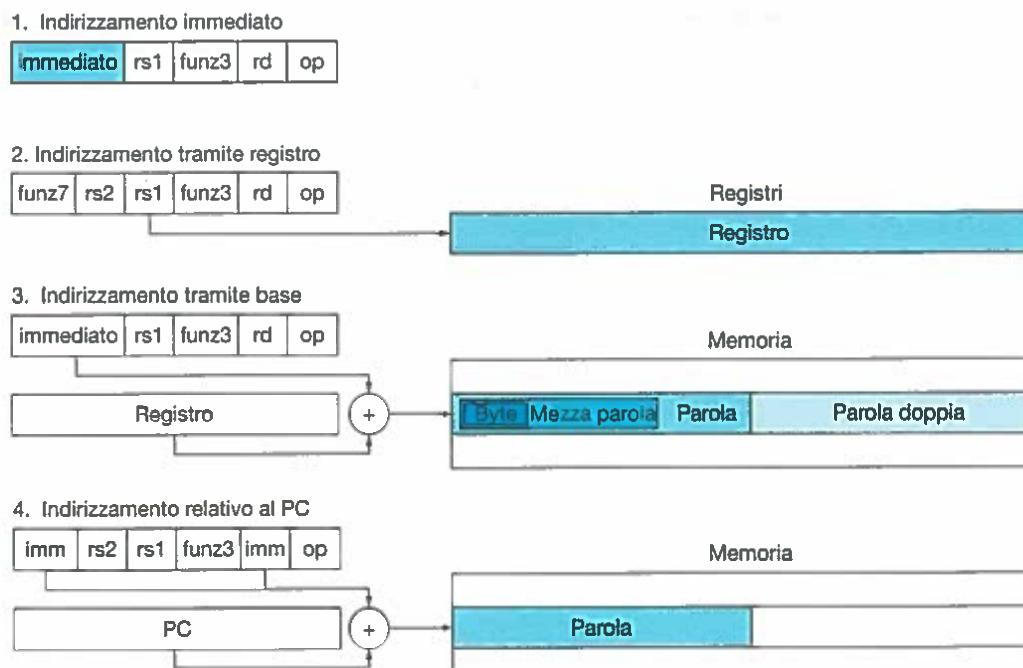
#### Riassunto delle modalità di indirizzamento del RISC-V

I diversi formati di indirizzamento sono in generale chiamati **modalità di indirizzamento**. La Figura 2.17 mostra come sono utilizzati gli operandi in ciascuna modalità di indirizzamento, che nel RISC-V sono le seguenti:

1. *indirizzamento immediato (immediate addressing)*, in cui l'operando è una costante contenuta nell'istruzione stessa;
2. *indirizzamento tramite registro (register addressing)*, in cui l'operando è un registro;
3. *indirizzamento tramite base e spiazzamento (base and displacement addressing)*, in cui l'operando è in una locazione di memoria individuata dalla somma del contenuto di un registro e di una costante contenuta nell'istruzione stessa;
4. *indirizzamento relativo al program counter (PC-relative addressing)*, in cui l'indirizzo di salto è la somma del contenuto del program counter e di una costante contenuta nell'istruzione stessa.

#### Come decodificare il linguaggio macchina

In alcuni casi è utile ricostruire il codice assembler a partire dal codice macchina, per esempio per analizzare i *core dump*, ossia i file contenenti una copia del contenuto dei registri e della memoria generati da un errore in un programma. La Figura 2.18 mostra la codifica dei diversi campi delle istruzioni in linguaggio macchina RISC-V e può essere utilizzata per tradurre manualmente da codice macchina a codice assembler e viceversa.



**Figura 2.17** Illustrazione delle quattro modalità di indirizzamento del RISC-V. Gli operandi sono evidenziati in blu. L'operando della modalità 3 si trova in memoria, mentre quello della modalità 2 si trova in un registro. Si noti che le varianti delle istruzioni di load e store possono accedere al byte, alla mezza parola, alla parola o alla parola doppia. Nella modalità 1, l'operando è contenuto nell'istruzione stessa. La modalità 4 viene utilizzata per indirizzare le istruzioni in memoria, aggiungendo un indirizzo ampio al PC. Si noti che un'operazione può utilizzare diverse modalità di indirizzamento: la somma, per esempio, può avere sia un operando immediato (addi) sia tutti gli operandi nei registri (add).

## Come decodificare le istruzioni in linguaggio macchina

Qual è l'istruzione assembler corrispondente alla seguente istruzione in linguaggio macchina?

00578833<sub>esa</sub>

Il primo passo consiste nel convertire il valore esadecimale in binario:

0000 0000 0101 0111 1000 1000 0011 0011

Per capire come interpretare i diversi bit, occorre determinare il formato dell'istruzione, e per fare ciò esaminiamo per prima cosa il codice operativo. Questo è definito dai sette bit più a destra dell'istruzione, cioè 0110011. Cercando questo numero in Figura 2.18, vediamo che il codice operativo corrisponde alle istruzioni aritmetiche di tipo R. Possiamo quindi suddividere la stringa dell'istruzione nei campi riportati in Figura 2.19.

funz7	rs2	rs1	funz3	rd	codop
0000000	00101	01111	000	10000	0110011

Decodifichiamo il resto dell'istruzione analizzando i valori contenuti nei diversi campi. I campi funz7 e funz3 contengono entrambi zero, indicando che l'istruzione è una add. Il valore decimale dei registri degli operandi è 5 per il campo rs2, 15 per rs1 e 16 per rd. Questi numeri rappresentano i registri x5, x15 e x16. Siamo ora in grado di scrivere l'istruzione assembler:

add x16, x15, x5

ESEMPIO

SOLUZIONE

Formato	Istruzione	Codice Operativo	funz3	funz6/7
Tipo R	add	0110011	000	0000000
	sub	0110011	000	0100000
	sll	0110011	001	0000000
	xor	0110011	100	0000000
	srl	0110011	101	0000000
	sra	0110011	101	0000000
	or	0110011	110	0000000
	and	0110011	111	0000000
	lr.d	0110011	011	0001000
	sc.d	0110011	011	0001100
Tipo I	lb	0000011	000	n.a.
	lh	0000011	001	n.a.
	lw	0000011	010	n.a.
	ld	0000011	011	n.a.
	lbu	0000011	100	n.a.
	lhu	0000011	101	n.a.
	lwu	0000011	110	n.a.
	addi	0010011	000	n.a.
	slli	0010011	001	000000
	xori	0010011	100	n.a.
	srl	0010011	101	000000
	srai	0010011	101	010000
	ori	0010011	110	n.a.
	andi	0010011	111	n.a.
	jalr	1100111	000	n.a.
Tipo S	sb	0100011	000	n.a.
	sh	0100011	001	n.a.
	sw	0100011	010	n.a.
	sd	0100011	111	n.a.
Tipo SB	beq	1100111	000	n.a.
	bne	1100111	001	n.a.
	blt	1100111	100	n.a.
	bge	1100111	101	n.a.
	bltu	1100111	110	n.a.
	bgeu	1100111	111	n.a.
Tipo U	lui	0110111	n.a.	n.a.
Tipo UJ	jal	1101111	n.a.	n.a.

**Figura 2.18 Codifica delle istruzioni RISC-V.** Tutte le istruzioni hanno un codice operativo, e tutti i formati tranne il tipo U e il tipo UJ utilizzano il campo funz3. Le istruzioni di tipo R utilizzano il campo funz7, e le istruzioni di scorrimento mediante costante (slli, srl, srai) utilizzano il campo funz6.

La Figura 2.19 mostra tutti i formati delle istruzioni RISC-V. In Figura 2.1 avevamo mostrato le istruzioni assembler RISC-V che sono descritte in questo capitolo. Il prossimo capitolo tratta le istruzioni RISC-V per moltiplicare e dividere, e l'aritmetica per i numeri decimali.

Nome (dimensione del campo)	Campi						Commenti
	7 bit	5 bit	5 bit	3 bit	5 bit	7 bit	
Tipo R	funz7	rs2	rs1	funz3	rd	codop	Istruzioni aritmetiche
Tipo I	Immediato[11:0]		rs1	funz3	rd	codop	Istruzioni di caricamento dalla memoria e aritmetica con costanti
Tipo S	immed[11:5]	rs2	rs1	funz3	immed[4:0]	codop	Istruzioni di trasferimento alla memoria (store)
Tipo SB	immed[12, 10:5]	rs2	rs1	funz3	immed[4:1,11]	codop	Istruzioni di salto condizionato
Tipo UJ	immediato[20, 10:1, 11, 19:12]				rd	codop	Istruzioni di salto incondizionato
Tipo U	immediato[31:12]				rd	codop	Formato caricamento stringhe di bit più significativi

Figura 2.19 Formatì delle istruzioni RISC-V.

## Autovalutazione

- I. Qual è l'intervallo degli indirizzi raggiungibili da un salto condizionato nel RISC-V ( $K = 1024$ )?
  1. Indirizzi tra 0 e  $4K - 1$
  2. Indirizzi tra 0 e  $8K - 1$
  3. Indirizzi fino a  $2K$  prima e fino a  $2K$  dopo l'istruzione di salto
  4. Indirizzi fino a  $4K$  prima e fino a  $4K$  dopo l'istruzione di salto
- II. Qual è l'intervallo degli indirizzi raggiungibili da un'istruzione di jump e di jump-and-link RISC-V ( $M = 1024K$ )?
  1. Indirizzi tra 0 e  $512K - 1$
  2. Indirizzi tra 0 e  $1M - 1$
  3. Indirizzi fino a  $512K$  prima e fino a  $512K$  dopo l'istruzione di salto
  4. Indirizzi fino a  $1M$  prima e fino a  $1M$  dopo l'istruzione di salto

## 2.11 Parallelismo e istruzioni: la sincronizzazione

L'esecuzione parallela è più semplice quando i task sono indipendenti, ma spesso essi devono cooperare e la cooperazione di solito significa che uno dei task deve scrivere dei valori che devono essere letti da altri task. Per sapere quando un task ha terminato la scrittura, permettendo a un altro task di leggere i valori scritti, è necessaria una sincronizzazione. Se i task non fossero sincronizzati si potrebbe creare una **competizione sui dati** (*data race*), a causa della quale i risultati di un programma possono cambiare a seconda dell'ordine degli eventi.

Riprendiamo l'analogia con gli otto giornalisti che lavorano contemporaneamente a un articolo, descritta nel paragrafo 1.6. Si supponga che uno dei giornalisti debba leggere il testo scritto dagli altri prima di potere scrivere la conclusione: dovrà sapere quando gli altri giornalisti hanno terminato di scrivere, in modo tale che queste parti non rischino di essere modificate dopo che lui le ha lette. Cioè, è opportuno che i giornalisti sincronizzino scrittura e lettura delle varie parti dell'articolo in modo tale che la conclusione sia consistente con il contenuto del resto dell'articolo.

Nelle architetture i meccanismi di sincronizzazione sono tipicamente integrati in procedure software di livello utente basate su istruzioni hardware di



**Competizione sui dati:** due accessi a memoria entrano in competizione per un dato se vengono effettuati da due thread diversi, riguardano lo stesso indirizzo di memoria, uno dei due accessi è in scrittura e si verificano uno subito dopo l'altro.

sincronizzazione. In questo paragrafo ci focalizzeremo sull'implementazione delle operazioni di sincronizzazione mediante *lock* (blocco) e *unlock* (sblocco). Lock e unlock possono essere utilizzate in modo semplice per creare delle regioni di dati, chiamate di *mutua esclusione*, sulle quali un solo processore alla volta può operare, ma anche per implementare meccanismi di sincronizzazione più complessi.

La funzionalità critica richiesta per implementare meccanismi di sincronizzazione nei multiprocessori è fornita da un insieme di primitive hardware in grado di leggere e modificare una locazione della memoria *atomicamente*, cioè in modo tale che non si possano compiere operazioni tra la lettura e la scrittura di quella locazione di memoria. Senza queste primitive hardware il costo di costruzione di primitive di sincronizzazione di base sarebbe elevato e crescerebbe in modo insostenibile con il numero dei processori.

Esistono diverse formulazioni possibili di queste primitive hardware di base, ma tutte consentono di leggere e modificare atomicamente una locazione di memoria, e a tutte è associato un modo per riconoscere se lettura e scrittura siano state eseguite in modo veramente atomico. In generale, i progettisti non si aspettano che siano gli utenti a implementare le primitive hardware di base, ma si aspettano che le primitive vengano utilizzate dai programmati di sistema per costruire una libreria di sincronizzazione, processo che si rivela spesso complesso e insidioso.

Iniziamo a esaminare una di queste primitive hardware e vediamo in che modo utilizzarla per costruire una primitiva di sincronizzazione di base. Una tipica operazione con cui costruire processi di sincronizzazione è lo *scambio atomico* (*atomic swap*), che scambia il contenuto di un registro con il contenuto di una parola di memoria.

Per vedere come utilizzare questa operazione, supponiamo di volere costruire un lock semplice in cui il valore 0 viene utilizzato per indicare che il lock è aperto e 1 per indicare che la memoria associata al lock non è disponibile (lock chiuso). Un processore cerca di bloccare una locazione di memoria scambiando un 1, contenuto in un registro, con l'indirizzo di memoria che si vuole bloccare. Il valore restituito dalla funzione che esegue lo scambio sarà 1 se un altro processore avrà già richiesto e ottenuto l'accesso a quella locazione di memoria, altrimenti sarà 0. In questo secondo caso il lock verrà modificato e impostato a 1, in modo da evitare ulteriori competizioni con altri processori, che vedrebbero come valore del lock ancora 0.

Per esempio, si considerino due processori che cercano di eseguire uno scambio simultaneamente: questa competizione sui dati viene interrotta, poiché solamente uno dei due processori potrà eseguire lo scambio per primo e la sua funzione di scambio restituirà uno 0, mentre la funzione del secondo processore restituirà 1 perché non potrà completare l'operazione. La chiave per utilizzare le primitive di scambio per implementare scambi sincronizzati è che l'operazione sia atomica: lo scambio deve essere indivisibile e due scambi richiesti simultaneamente devono essere messi in sequenza dall'hardware. È impossibile che due processori possano impostare la variabile di sincronizzazione a 1.

L'implementazione di una singola operazione atomica sulla memoria introduce alcune sfide nella progettazione dei processori, dato che richiede che la lettura e la scrittura della memoria avvengano in un'unica istruzione non interrompibile.

Una possibile alternativa è avere una coppia di istruzioni nella quale la seconda restituisca un valore che indica che le due istruzioni sono state eseguite come una singola istruzione atomica. La coppia di istruzioni è effettivamente eseguita in modo atomico se si ha la certezza che tutte le altre operazioni effettuate dai processori sono state eseguite prima di essa. In tal caso nessun altro processore può modificare il valore della memoria tra le due istruzioni.

Nel RISC-V questa coppia di istruzioni è realizzata da un'istruzione di load speciale, chiamata *load-reserved doubleword* (lr.d, load riservata di una parola doppia) e da un'istruzione di store speciale chiamata *store-conditioned doubleword* (sc.d, store condizionata di una parola doppia). Queste due istruzioni sono utilizzate in sequenza: se il contenuto di una locazione di memoria associata a una load riservata (linked) venisse alterato prima che la store condizionata associata abbia salvato in quella cella di memoria il dato, l'istruzione di store condizionata fallirebbe e non scriverebbe il dato in memoria. L'istruzione di store condizionata ha quindi due funzioni: salva il contenuto di un registro (presumibilmente diverso) in memoria e imposta il contenuto di quel registro a 1 se la scrittura ha avuto successo, e a 0 se invece è fallita. Quindi l'istruzione sc.d specifica tre registri: uno che contiene l'indirizzo, uno per indicare se l'operazione atomica fallisce o ha successo, e uno per il dato da scrivere in memoria se l'operazione ha avuto successo. Dato che l'istruzione load riservata restituisce il valore letto e l'istruzione store condizionata restituisce 1 solamente se la scrittura ha avuto successo, la sequenza di istruzioni seguente implementa uno scambio atomico relativo alla locazione di memoria specificata dal contenuto di x20:

```

tentativo:lr.d x10, (x20)          // carica il dato e
                                    // blocca la cella di memoria
    sc.d x11, x23, (x20)          // store condizionata
    bne x11, x0, tentativo // ritorna a tentativo se la
                            // store fallisce
    addi x23, x10, 0           // copia il valore letto in x23

```

Se il processore intervenisse per modificare il contenuto della memoria in x20 tra l'esecuzione dell'istruzione lr.d e della sc.d, l'istruzione sc.d scriverebbe un valore non 0 nel registro x11, forzando la ripetizione della sequenza di istruzioni. Al termine di questa sequenza di istruzioni il contenuto di x23 e della locazione di memoria indicata da x20 sono stati scambiati.

**Approfondimento.** Abbiamo presentato lo scambio atomico per la sincronizzazione dei sistemi multiprocessore, ma questo è utilizzato anche dal sistema operativo che gestisce più processi in esecuzione su un singolo processore. Per assicurarsi che niente interferisca con l'esecuzione sul processore, l'istruzione di store condizionata fallisce ogni volta che il processore esegue un cambio di contesto tra le due istruzioni (vedi Cap. 5).

**Approfondimento.** Un vantaggio del meccanismo della coppia load linked/store condizionata è che può essere utilizzato per costruire altre primitive di sincronizzazione, quali *atomic compare and swap* e *atomic fetch and increment*, utilizzate in alcuni modelli di programmazione parallela. Queste primitive richiedono un certo numero di istruzioni tra la lr.d e la sc.d, ma non tante.

Dato che l'istruzione di store condizionata potrà fallire sia dopo la richiesta di un altro tentativo di scrittura nell'indirizzo da cui una load linked ha letto il dato, sia dopo un'eccezione, molta attenzione deve essere posta nella scelta delle istruzioni da inserire tra questa coppia di istruzioni. In particolare, si possono inserire in modo sicuro solamente le istruzioni di aritmetica con gli interi, di salto in avanti e di salto all'indietro al di fuori del blocco compreso tra l'istruzione di load linked e store condizionata.

Viceversa, è possibile creare delle situazioni di stallo (*deadlock*) nelle quali il processore non può terminare l'esecuzione della sc.d perché continua a ricevere un'eccezione di page fault. Inoltre, il numero di istruzioni inserite tra la load linked e la store condizionata deve essere ridotto al minimo, in modo tale da minimizzare la probabilità che un evento scorrelato o un altro processore causino il fallimento frequente della store condizionata.

**Approfondimento.** Mentre il codice precedente eseguiva uno scambio atomico, il codice seguente blocca in modo più efficace la locazione di memoria contenuta nel registro x20, dove il valore 0 indica che la locazione è stata sbloccata mentre 1 che il blocco è stato eseguito:

```
addi x12, x0, 1      // copia del valore associato al blocco
ancora: lr.d x10, (x20) // carica con il blocco per vedere se
                         // la locazione è bloccata
bne x10, x0, ancora // controlla se è ancora 0
sc.d x11, x12,(x20) // prova a memorizzare il nuovo valore
bne x11, x0, ancora // salta se la memorizzazione ha fallito
```

Per rilasciare il blocco è sufficiente utilizzare una normale istruzione di store per scrivere 0 nella locazione di memoria:

```
sd x0, 0(x20) // elimina il blocco scrivendo 0
```

### Autovalutazione

Quando si utilizzano primitive come load linked riservata e store condizionata?

1. Quando processi cooperanti di un programma parallelo devono essere sincronizzati per ottenere un comportamento corretto nella lettura e nella scrittura di dati condivisi.
2. Quando processi cooperanti su un'architettura a singolo processore devono essere sincronizzati per leggere e scrivere dati condivisi.

## 2.12 Tradurre e avviare un programma

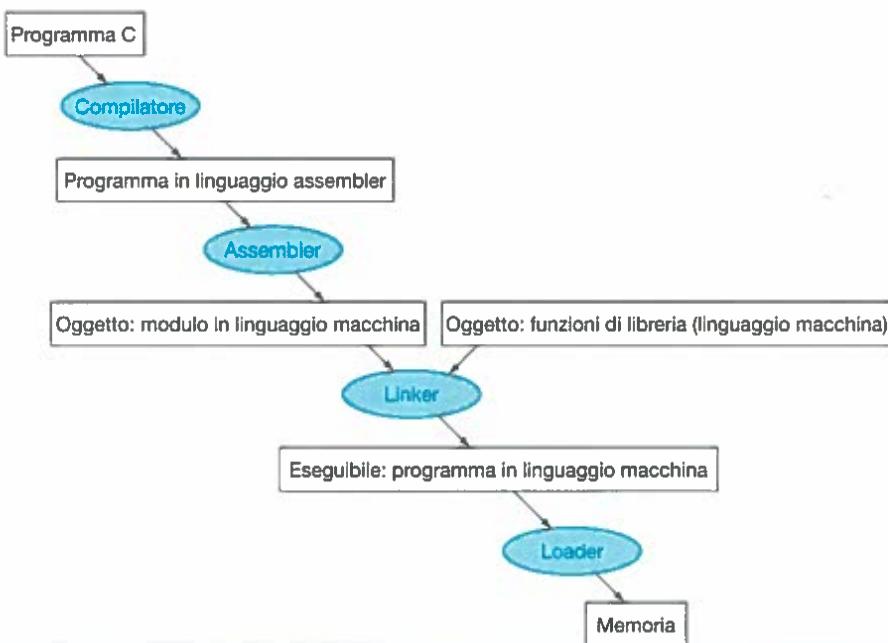
Questo paragrafo descrive i quattro passi necessari a trasformare un programma C, contenuto in un file in memoria (disco o memoria flash), in un programma pronto per essere eseguito su un calcolatore. La Figura 2.20 mostra la sequenza delle trasformazioni; alcuni calcolatori possono raggruppare questi passi per ridurre il tempo di traduzione, ma concettualmente queste sono le quattro fasi attraverso le quali tutti i programmi devono passare. Noi seguiremo questa sequenza di trasformazioni.

### Compilatore

**Linguaggio assembler:** un linguaggio simbolico che può essere tradotto nel linguaggio macchina binario.

Il compilatore trasforma il programma C in un *programma in linguaggio assembler*, cioè in una forma simbolica di ciò che il calcolatore è in grado di capire. I programmi scritti in linguaggio ad alto livello sono costituiti da un numero di righe di codice minore rispetto al linguaggio assembler, e questo fa sì che la produttività di un programmatore sia più elevata.

Nel 1975 molti sistemi operativi e molti assemblatori erano scritti in **linguaggio assembler**, dal momento che la capacità della memoria era ridotta e i compilatori erano inefficienti. Oggi la capacità delle memorie per ogni singolo chip di DRAM è cresciuta di un milione di volte, riducendo il problema delle dimensioni dei programmi; inoltre i compilatori ottimizzanti attuali sono in grado di produrre codice assembler quasi con la medesima qualità di quello scritto da programmatore assembler esperti e qualche volta addirittura migliore, nel caso di programmi di grandi dimensioni.



**Figura 2.20 Sequenza di passi di traduzione per il C.** Un programma scritto in linguaggio ad alto livello viene per prima cosa compilato in linguaggio assembleur e poi assemblato per ottenere un modulo oggetto in linguaggio macchina. Il linker unisce uno o più moduli tra loro e con le procedure contenute nelle librerie, e risolve tutti i riferimenti incrociati. Il loader, infine, carica il codice macchina nell'opportuna area di memoria, in modo che possa essere eseguito dal processore. Per accelerare il processo di traduzione, alcuni passi possono essere saltati o eseguiti insieme ad altri. Alcuni compilatori producono direttamente il codice macchina mentre altri utilizzano dei linker-loader, i quali eseguono gli ultimi due passi contemporaneamente. Per identificare il tipo di file, il sistema operativo UNIX segue una precisa convenzione sull'estensione dei file: i file sorgente C vengono chiamati `x.c`, i file assembler `x.s`, i file oggetto `x.o`, le librerie collegate staticamente `x.a`, le librerie collegate dinamicamente `x.so` e i file eseguibili sono solitamente denominati `a.out`. Il sistema operativo MS-DOS utilizza come estensioni rispettivamente `.C`, `.ASM`, `.OBJ`, `.LIB`, `.DLL` e `.EXE` con lo stesso significato.

## Assemblatore

Dato che il linguaggio assembler rappresenta l'interfaccia verso il software di livello più alto, l'assemblatore può anche trattare varianti delle istruzioni in linguaggio macchina come se fossero istruzioni vere e proprie. Non è richiesto che l'hardware implementi queste istruzioni, ma esse consentono di semplificare la traduzione e la programmazione, aumentando la leggibilità del codice. Queste istruzioni sono chiamate **pseudoistruzioni**. Come detto in precedenza, l'hardware RISC-V si assicura che il registro `x0` contenga sempre il valore 0; quindi il registro `x0` fornisce il valore 0 ogni volta che viene letto e se il programmatore cerca di cambiarne il valore di `x0`, il nuovo valore viene semplicemente scartato. Il registro `x0` viene utilizzato per creare l'istruzione assembler che copia il contenuto di un registro in un altro. L'assemblatore RISC-V accetta l'istruzione seguente anche se non fa parte del linguaggio macchina RISC-V:

```
li x9, 123 // carica il valore 123 nel registro x9
```

L'assemblatore converte questa istruzione assembler nell'istruzione equivalente in linguaggio macchina, che è:

```
addi x9, x0, 123 // il registro x9 assume il valore x0 + 123
```

L'assemblatore RISC-V converte anche la pseudoistruzione `mv` (*move*, sposta) in un'istruzione `addi`. Quindi:

```
mv x10, x11 // il registro x10 assume il valore del registro x11
```

diventa

```
addi x10, x11, 0 // Il registro x10 riceve il contenuto
                  // del registro x11 + 0
```

**Pseudoistruzione:** una variante di un'istruzione in linguaggio assembler frequentemente utilizzata; viene spesso trattata come una vera e propria istruzione.

L'assemblatore accetta anche la pseudoistruzione `j Etichetta` come istruzione di salto incondizionato a Etichetta, che corrisponde all'istruzione `jal x0, Etichetta`. Converte anche i salti condizionati a locazioni lontane in un'istruzione di salto condizionato seguita da una di salto incondizionato. Come accennato in precedenza, l'assemblatore RISC-V consente di caricare delle costanti ampie nei registri nonostante la dimensione limitata dei campi immediati. Quindi, la pseudoistruzione *carica immediata* (`li, load immediate`) introdotta in precedenza può creare delle costanti più ampie di quelle che possono essere contenute nel campo immediato della addi; la macro `load address` (`la`) lavora in modo simile per gli indirizzi simbolici. Infine può semplificare l'insieme delle istruzioni determinando quale variante di un'istruzione desideri il programmatore. Per esempio, l'assemblatore RISC-V non richiede che il programma specifichi la versione immediata dell'istruzione utilizzando una costante per le istruzioni aritmetiche e logiche; genera semplicemente il codice operativo adeguato. Quindi:

```
and x9, x10, 15 // registro x9 assume il valore
                  // di x10 AND 15
```

diventa

```
andi x9, x10, 15 // registro x9 assume il valore
                  // di x10 AND 15
```

Abbiamo aggiunto la “i” nell'istruzione per ricordare al lettore che andi produce un codice operativo diverso con un formato di istruzione differente rispetto all'istruzione and che non ha operandi immediati.

Riassumendo, le pseudoistruzioni consentono all'assembler RISC-V di avere un insieme di istruzioni più ricco di quello implementato in hardware. Quando si scrivono programmi in assembler, l'utilizzo delle pseudoistruzioni rende il compito più semplice; tuttavia per capire l'architettura RISC-V e per ottenere le migliori prestazioni è bene conoscere le vere istruzioni RISC-V, riportate nelle Figure 2.1 e 2.18.

Gli assemblatori, inoltre, accettano numeri espressi in basi diverse. Oltre alla base due e a quella decimale, solitamente accettano anche una base più concisa di quella binaria, che può essere facilmente convertita in una sequenza di bit: gli assemblatori RISC-V utilizzano la base esadecimale e ottale.

Tutte le caratteristiche dell'assemblatore appena illustrate sono molto utili, ma l'assemblatore ha come compito principale la generazione del linguaggio macchina: converte un programma assembler in un *file oggetto*, cioè una sequenza di istruzioni in linguaggio macchina, di dati e di informazioni necessarie a collocare le istruzioni in memoria nella posizione opportuna.

Per produrre la versione binaria di ogni istruzione di un programma assembler, l'assemblatore deve determinare gli indirizzi corrispondenti a tutte le etichette: esso tiene traccia di tutte le etichette utilizzate nei salti e nei trasferimenti di dati scrivendole in una tabella, detta **tabella dei simboli** (*symbol table*). Come si può immaginare, questa tabella conterrà coppie di tipo simbolo-indirizzo.

I file oggetto nei sistemi UNIX tipicamente contengono sei segmenti distinti: vediamoli.

- L'*intestazione del file oggetto* (*object file header*) descrive la dimensione e la posizione degli altri segmenti del file oggetto stesso.
- Il *segmento di testo* (*text segment*) contiene il codice in linguaggio macchina.
- Il *segmento di dati statici* (*static data segment*) contiene tutti i dati allocati per tutta la durata del programma. UNIX consente l'utilizzo sia di *dati statici*, che rimangono allocati per tutta la durata dell'esecuzione del programma, sia di *dati dinamici*, che possono crescere o diminuire di dimensione a seconda delle esigenze del programma (Figura 2.13).

**Tabella dei simboli:** tabella che fa corrispondere i nomi delle etichette agli indirizzi delle parole di memoria contenenti le istruzioni.

- Le *informazioni di rilocazione* (*relocation information*) identificano le istruzioni e i dati che, quando il programma è posto in memoria, dipendono da indirizzi assoluti.
- La *tabella dei simboli* (*symbol table*) contiene le rimanenti etichette di cui non è stata trovata una definizione, per esempio quelle che fanno riferimento a moduli esterni.
- Le *informazioni per il debugger* (*debugging information*) contengono una descrizione concisa di come sono stati compilati i moduli, in modo che il debugger (programma per il supporto della ricerca degli errori) possa associare le istruzioni in linguaggio macchina al codice sorgente C e rendere leggibili le strutture dati.

Il paragrafo seguente, relativo al linker, mostrerà come vengono collegate procedure già assemblate, per esempio quelle contenute nelle librerie.

## Linker

Tutto ciò che abbiamo visto finora suggerisce che la modifica anche di una sola linea di codice richieda di ricompilare e riassemblare l'intero programma. Ma ritradurre tutto dall'inizio è un terribile spreco di risorse computazionali, soprattutto per quanto riguarda le librerie standard, che verrebbero ricompilate e riassemblate anche se, per definizione, esse non cambiano praticamente mai. Una valida alternativa consiste nel compilare e assemblare ciascuna procedura indipendentemente dalle altre, in modo che la modifica di una linea di codice renda necessario ricompilare e riassemblare solo la procedura a cui la linea di codice appartiene. Tale alternativa necessita di un nuovo programma di sistema, chiamato **link editor** o **linker**; esso prende tutti i programmi (le procedure) in codice macchina che sono stati assemblati indipendentemente e li "cuce" insieme. Il motivo per cui il linker è particolarmente utile è che risulta molto più veloce "correggere" il codice piuttosto che ricompilarlo e assemblarlo di nuovo.

Il linker esegue i seguenti tre passi:

1. inserisce in memoria in modo simbolico il codice e i moduli dati;
2. determina gli indirizzi dei dati e delle etichette che compaiono nelle istruzioni;
3. corregge i riferimenti interni ed esterni.

Il linker utilizza le informazioni di rilocazione e la tabella dei simboli di ciascun modulo oggetto per risolvere tutte le etichette non definite. Queste si trovano nelle istruzioni di salto e negli indirizzi dei dati; il compito del linker è quindi più o meno quello di un editor: trovare gli indirizzi vecchi e sostituirli con quelli nuovi. Questa attività di "editing" è all'origine del nome "link editor", abbreviato in linker.

Dopo che tutti i riferimenti esterni sono stati risolti, il linker determina le locazioni di memoria che ciascun modulo dovrà occupare. In Figura 2.13 avevamo mostrato la convenzione RISC-V per l'allocazione dei programmi e dei dati in memoria. Poiché tutti i file vengono assemblati in modo indipendente, l'assemblatore non può conoscere la posizione relativa delle istruzioni e dei dati di un modulo rispetto a un altro; perciò, quando il linker inserisce un modulo in memoria, tutti i riferimenti *assoluti*, cioè gli indirizzi di memoria definiti non in relazione a un registro, devono essere *rilocati* in modo da poter riflettere la loro reale posizione.

Il linker produce un **file eseguibile** che può essere eseguito su un calcolatore. Di norma questo file ha lo stesso formato di un file oggetto, ma non contiene più riferimenti non risolti. È possibile eseguire questo passo anche in maniera parziale, come è il caso delle librerie, che contengono ancora indirizzi non risolti e quindi in realtà sono assimilabili a file oggetto.

**Link editor o linker:** programma di sistema che unisce varie procedure in linguaggio macchina assemblate indipendentemente e risolve tutte le etichette non definite, producendo un programma eseguibile.

**File eseguibile:** programma in grado di essere eseguito, scritto nel formato del file oggetto ma che non contiene riferimenti non risolti. Può contenere tabelle di simboli e informazioni per il debugger. Un "eseguibile spogliato" (*stripped executable*) non contiene queste ultime informazioni. Le informazioni di rilocazione possono venire incluse per il loader.

## Eseguire il link di due file oggetto

### ESEMPIO

Si esegua il link dei due file oggetto riportati qui sotto, mostrando gli indirizzi aggiornati delle prime istruzioni del file eseguibile completo. Le istruzioni sono mostrate in linguaggio assembler solo per rendere l'esempio più leggibile, ma nella realtà le istruzioni sono rappresentate da numeri.

Si noti che nel file oggetto sono evidenziati gli indirizzi e i simboli che devono essere aggiornati nel processo di link, cioè le istruzioni che fanno riferimento agli indirizzi delle procedure A e B e quelle che fanno riferimento agli indirizzi delle parole doppie dei dati X e Y.

Intestazione del file oggetto			
	Nome	Procedura A	
	Dimensione del testo	100 <sub>esa</sub>	
	Dimensione dei dati	20 <sub>esa</sub>	
Segmento testo	Indirizzo	Istruzione	
	0	ld x10, 0(x3)	
	4	jal x1, 0	
	...	...	
Segmento dati	0	(X)	
	...	...	
Informazioni di rilocazione	Indirizzo	Tipo di istruzione	Dipendenza
	0	ld	X
	4	jal	B
Tabella dei simboli	Etichetta	Indirizzo	
	X	-	
	B	-	
Procedura B			
	Nome	Procedura B	
	Dimensione del testo	200 <sub>esa</sub>	
	Dimensione dei dati	30 <sub>esa</sub>	
Segmento testo	Indirizzo	Istruzione	
	0	sd x11, 0(x3)	
	4	jal x1, 0	
	...	...	
Segmento dati	0	(Y)	
	...	...	
Informazioni di rilocazione	Indirizzo	Tipo di istruzione	Dipendenza
	0	sd	Y
	4	jal	A
Tabella dei simboli	Etichetta	Indirizzo	
	Y	-	
	A	-	

(continua)

La procedura A deve trovare l'indirizzo della variabile etichettata X per inserirlo nell'istruzione di load e l'indirizzo della procedura B per inserirlo nell'istruzione jal. La procedura B ha bisogno dell'indirizzo della variabile etichettata Y per inserirlo nell'istruzione di store e dell'indirizzo della procedura A per la sua istruzione jal.

Dalla Figura 2.14 si ricava che il segmento di testo inizia all'indirizzo 0000 0000 0040 0000<sub>esa</sub> e che il segmento dati inizia in 0000 0000 1000 0000<sub>esa</sub>. Il testo della procedura A è inserito a partire dal primo dei due indirizzi e i suoi dati a partire dal secondo. L'intestazione del file oggetto della procedura A indica che il suo segmento testo è lungo 100<sub>esa</sub> byte e che i suoi dati occupano 20<sub>esa</sub> byte; quindi l'indirizzo di partenza del testo della procedura B sarà 40 0100<sub>esa</sub> e i suoi dati partiranno da 1000 0020<sub>esa</sub>.

(continua)

**SOLUZIONE**

Intestazione del file eseguibile	Dimensione del testo	300 <sub>esa</sub>
Segmento testo	Indirizzo	Istruzione
	0000 0000 0040 0000 <sub>esa</sub>	ld x10, 0(x3)
	0000 0000 0040 0004 <sub>esa</sub>	jal x1, 252 <sub>dec</sub>
	...	...
	0000 0000 0040 0100 <sub>esa</sub>	sd x11, 32(x3)
	0000 0000 0040 0104 <sub>esa</sub>	jal x1, -260 <sub>dec</sub>
	...	...
Segmento dati	Indirizzo	
	0000 0000 1000 0000 <sub>esa</sub>	(X)
	...	...
	0000 0000 1000 0020 <sub>esa</sub>	(Y)
	...	...

A questo punto il linker aggiorna il campo indirizzo delle istruzioni, utilizzando il campo del codice operativo per riconoscere il formato dell'indirizzo. In questo caso sono presenti due tipi di istruzioni:

- Le istruzioni jal utilizzano l'indirizzamento relativo al PC. Il campo indirizzo dell'istruzione jal, che si trova all'indirizzo 40 0004<sub>esa</sub> e salta all'indirizzo 40 0100<sub>esa</sub> (l'indirizzo della procedura B), conterrà quindi il valore (40 0100<sub>esa</sub> - 40 0004<sub>esa</sub>) o 252<sub>dec</sub> nel suo campo indirizzo. Analogamente, dato che 40 0000<sub>esa</sub> è l'indirizzo della procedura A, la jal a 40 0104<sub>esa</sub> conterrà il numero negativo -260<sub>dec</sub> (40 0000<sub>esa</sub> - 40 0104<sub>esa</sub>) nel proprio campo indirizzo.
- Gli indirizzi delle istruzioni di trasferimento dati sono un po' più complessi da ricavare, dal momento che sono relativi a un registro di base. In questo esempio viene utilizzato il registro x3 come registro base, che supponiamo inizializzato a 0000 0000 1000 0000<sub>esa</sub>. Per ottenere l'indirizzo della parola doppia X (0000 0000 1000 0000<sub>esa</sub>), inseriamo 0<sub>dec</sub> nel campo indirizzo dell'istruzione ld all'indirizzo 40 0000<sub>esa</sub>. Allo stesso modo inseriamo 20<sub>esa</sub> nel campo indirizzo dell'istruzione

(continua)

(continua)

- sd all'indirizzo 40 0100<sub>esa</sub> per ottenere l'indirizzo 0000 0000 1000 0020<sub>esa</sub> (l'indirizzo della parola doppia Y).
3. Gli indirizzi associati alle operazioni di store vengono gestiti esattamente come nelle operazioni di load, ad eccezione del fatto che il formato delle istruzioni di tipo S rappresenta le costanti in modo differente dal formato I delle load. Inseriamo 32<sub>dec</sub> nel campo indirizzo della sd all'indirizzo 40 0100<sub>esa</sub> per ottenere l'indirizzo 0000 0000 1000 0020<sub>esa</sub> (l'indirizzo della parola doppia Y).

## Loader

**Loader:** programma di sistema che carica un programma oggetto nella memoria principale in modo che sia pronto per essere eseguito.

Una volta che il file eseguibile è stato memorizzato su disco, il sistema operativo può leggerlo e trasferirlo in memoria per avviare l'esecuzione. Nei sistemi UNIX il **loader** (programma di caricamento) esegue i seguenti passi.

1. Lettura dell'intestazione del file eseguibile per determinare la lunghezza del segmento testo e del segmento dati.
2. Creazione di uno spazio di indirizzamento sufficiente a contenere testo e dati.
3. Copia delle istruzioni e dei dati dal file eseguibile in memoria.
4. Copia nello stack degli eventuali parametri passati al programma principale.
5. Inizializzazione dei registri del calcolatore e impostazione dello stack pointer affinché punti alla prima locazione libera.
6. Salto a una procedura di *start-up* (avviamento), la quale copia i parametri nei registri argomento e chiama la procedura principale del programma. Quando la procedura principale restituisce il controllo, la procedura di start-up termina il programma con una chiamata alla funzione di sistema `exit`.

## Librerie a caricamento dinamico

Praticamente tutti i problemi dell'informatica possono essere risolti con un altro livello di indirizzamento.

David Wheeler

Nella prima parte di questo paragrafo abbiamo descritto l'approccio tradizionale, detto *statico*, al collegamento delle librerie, che vengono mappate prima che il programma sia eseguito. Nonostante l'approccio statico sia il modo più veloce per chiamare le funzioni di libreria, ci sono alcuni svantaggi:

- Le funzioni di libreria diventano parte del codice eseguibile. Se viene rilasciata una nuova versione della libreria, che per esempio può contenere la correzione di alcuni errori o supportare nuovi dispositivi hardware, un programma che carica staticamente le librerie continua a utilizzare la vecchia versione.
- Vengono caricate tutte le funzioni di libreria utilizzate dal programma, qualsiasi sia il punto del programma nel quale le funzioni vengono chiamate, e anche quelle funzioni che non vengono utilizzate. La libreria può essere molto più grande del programma; per esempio la libreria standard C di un sistema RISC-V su cui è montato il sistema operativo Linux è di 1,5 MiB.

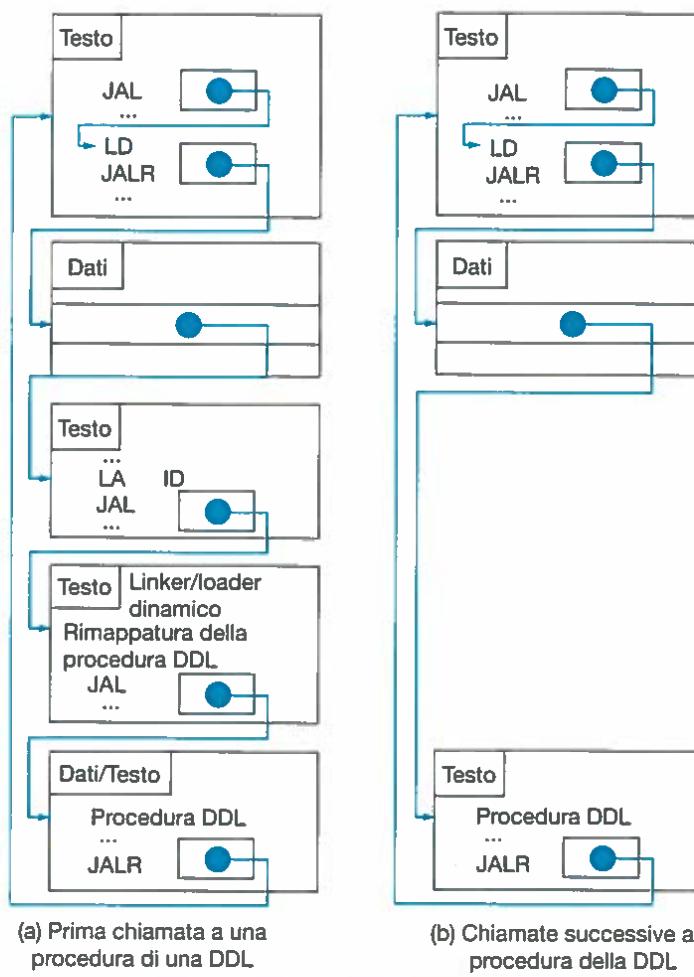
**Librerie a caricamento dinamico (DLL):** procedure di libreria che vengono collegate a un programma durante la sua esecuzione.

Questi svantaggi hanno portato allo sviluppo delle **librerie a caricamento dinamico (DLL, Dynamically Linked Libraries)**, per cui le funzioni di libreria non vengono collegate e caricate finché non inizia l'esecuzione del programma. Sia il programma sia le librerie contengono ulteriori informazioni riguardanti la posizione delle procedure non locali e il loro nome. Nella versione originale delle DLL, il loader lanciava un linker dinamico, utilizzando queste ulteriori informazioni presenti nel file per cercare le librerie appropriate e per aggiornare tutti i riferimenti esterni. L'inconveniente era che anche in questo caso venivano

caricate tutte le funzioni di libreria disponibili, invece di caricare solamente quelle utilizzate effettivamente durante l'esecuzione del programma. Questa osservazione ha portato alla versione delle DLL con collegamento *lazy* (pigro), in cui ogni procedura viene caricata solo *dopo* la sua chiamata.

Come per molti altri casi nel nostro campo, questo trucco è basato sulla rimappatura, come mostrato in Figura 2.21. Ciascuna procedura non locale inizia chiamando una procedura fasulla (*dummy*), inserita subito dopo la fine del codice del programma; questa procedura contiene un'istruzione di salto indiretto: si inserisce quindi una procedura fasulla diversa per ognuna delle procedure non locali.

La prima volta che la funzione di libreria viene chiamata, il programma salta all'istruzione fasulla ad essa associata ed esegue il salto indiretto indicato. L'istruzione di salto porta a una porzione di codice che inserisce un numero in un registro, il quale identifica la funzione di libreria desiderata, e quindi salta a sua volta al linker-loader dinamico. Il linker-loader dinamico trova la funzione di libreria che si sta cercando, la rimappa in memoria e cambia l'indirizzo contenuto nell'istruzione di salto indiretto, in modo che la procedura fasulla salti ora alla prima istruzione della funzione richiesta. A questo punto la funzione è disponibile e si ritornerà al punto originale di chiamata una volta completata.



**Figura 2.21** Procedura per il caricamento dinamico delle librerie a collegamento *lazy*. (a) Passi seguiti per la prima chiamata a una procedura di una DLL. (b) I passi associati alla ricerca della procedura, alla sua rimappatura e al suo collegamento vengono saltati nelle chiamate successive. Come vedremo nel Capitolo 5, il sistema operativo può evitare di copiare la funzione desiderata, rimappandola mediante l'utilizzo della memoria virtuale.

la sua esecuzione. In seguito la chiamata alla funzione di libreria prevederà un salto indiretto unico alla sua prima istruzione, senza salti addizionali.

Riassumendo, le DLL richiedono spazio aggiuntivo per memorizzare le informazioni necessarie per il collegamento dinamico, ma non richiedono di caricare o collegare l'intera libreria. Queste librerie pagano un costo computazionale significativo la prima volta che una funzione viene chiamata, ma solamente il costo di un salto indiretto in tutte le chiamate successive. Si noti che il ritorno da queste funzioni di libreria non richiede alcun costo aggiuntivo. Il sistema operativo Windows di Microsoft è basato pesantemente sulle librerie a caricamento dinamico, che costituiscono anche lo standard negli attuali sistemi UNIX.

### Come avviare un programma Java

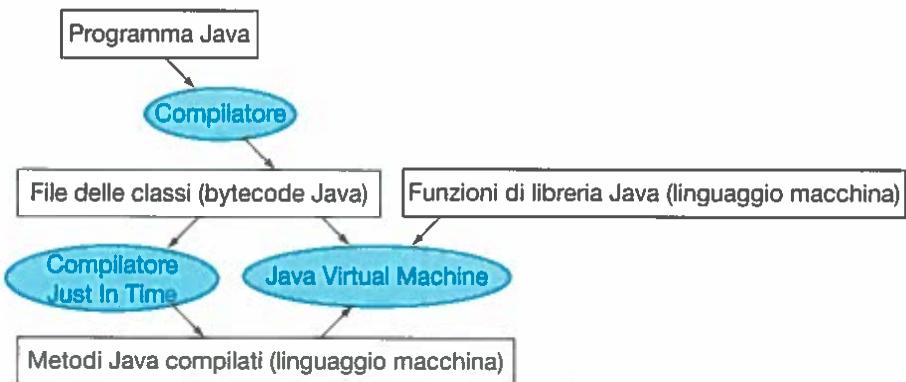
La discussione precedente ha riguardato il modello tradizionale di esecuzione di un programma, dove si dà importanza alla velocità di esecuzione per una specifica architettura dell'insieme di istruzioni o per una sua specifica implementazione. Anche se è possibile eseguire un programma Java come se fosse un programma C, Java è stato creato con obiettivi differenti: uno dei più importanti è quello di poter eseguire correttamente una applicazione su qualsiasi calcolatore, eventualmente con velocità di esecuzione diverse.

La Figura 2.22 mostra i passaggi tipici della traduzione ed esecuzione di un programma Java. Diversamente dal caso in cui le applicazioni vengono compilate traducendole nel linguaggio assembler del calcolatore, i programmi Java vengono tradotti, in prima istanza, in un insieme di istruzioni universali e semplici da interpretare, il **bytecode Java** (par. 2.15). Questo insieme di istruzioni è stato progettato per essere simile al linguaggio Java, in modo da rendere molto semplice la fase di compilazione; in pratica non vengono effettuate ottimizzazioni sul codice. Come per i compilatori C, anche il compilatore Java controlla i tipi di dato e genera le operazioni necessarie per ogni tipo. I programmi Java vengono distribuiti in forma binaria, come bytecode.

Un interprete software, chiamato **Java Virtual Machine (JVM)**, può eseguire il bytecode Java. Un interprete è un programma che simula una certa architettura di un insieme di istruzioni; per esempio, il simulatore RISC-V a cui si fa riferimento in questo libro è un interprete. Non c'è bisogno di una fase separata per la traduzione del codice in linguaggio macchina (assemblaggio), dato che

**Bytecode Java:** istruzione appartenente all'insieme di istruzioni progettato per interpretare i programmi Java.

**Java Virtual Machine (JVM):** programma che interpreta il bytecode Java.



**Figura 2.22 Gerarchia della traduzione del codice Java.** Il primo passo per un programma Java è la compilazione in una versione binaria, detta bytecode Java, nella quale tutti gli indirizzi vengono definiti dal compilatore. A questo punto il programma Java è pronto per essere eseguito da un interprete, chiamato **Java Virtual Machine (JVM)** che collega i metodi desiderati della libreria Java durante l'esecuzione stessa del programma. Per ottenere prestazioni migliori, la Java Virtual Machine può invocare il compilatore **Just In Time (JIT)**, che selettivamente compila i metodi del linguaggio Java nel linguaggio macchina del calcolatore sul quale è in esecuzione.

questa viene effettuata direttamente dal compilatore o dalla stessa JVM durante l'esecuzione del codice.

Il principale vantaggio nell'uso di un interprete è la portabilità. La disponibilità della Java Virtual Machine per diverse piattaforme ha comportato che gran parte dei programmatore abbia potuto scrivere ed eseguire i programmi Java subito dopo la creazione del linguaggio. Oggi la Java Virtual Machine si può trovare su centinaia di milioni di dispositivi, dai cellulari ai browser Internet.

Lo svantaggio nell'uso di un interprete è dato dalle basse prestazioni. L'incredibile miglioramento dell'hardware in termini di prestazioni tra gli anni '80 e gli anni '90 è riuscito a rendere sensato l'utilizzo di un interprete per molte applicazioni importanti; tuttavia, quando si confrontano i programmi C compilati tradizionalmente con quelli Java interpretati, si trova spesso un fattore 10 di differenza nelle prestazioni, il che rende i programmi interpretati poco interessanti per alcune tipologie di applicazioni.

Per mantenere la portabilità e migliorare la velocità di esecuzione, la successiva fase dello sviluppo di Java ha visto la nascita di compilatori che traducono il codice interpretato in codice macchina del calcolatore durante l'esecuzione del programma stesso. Questi programmi, chiamati **compilatori Just In Time (JIT)**, tracciano un profilo dell'applicazione durante la sua esecuzione in modo da sapere quali siano i metodi "caldi" (quelli che durante l'esecuzione consumano più tempo CPU), e quindi li compilano, traducendoli in istruzioni nel linguaggio macchina proprio del calcolatore. La porzione compilata viene salvata su disco per la successiva esecuzione del programma, ottenendo così che lo stesso programma possa essere eseguito più velocemente in tutte le esecuzioni successive. Il rapporto tra codice compilato e interpretato evolve quindi nel tempo, e i programmi Java eseguiti molto di frequente soffrono solo di un piccolo costo addizionale di interpretazione.

Dato che l'incremento di prestazioni dei calcolatori permette di creare compilatori sempre più potenti e che i ricercatori sviluppano tecniche sempre più efficienti per compilare il codice Java durante l'esecuzione, la differenza di prestazioni tra Java e C/C++ sta assottigliandosi sempre più. Il paragrafo 2.15 approfondisce maggiormente la descrizione di Java, del bytecode Java, della JVM e dei compilatori JIT.

#### **Compilatore Just In Time (JIT):**

il nome comunemente dato a un compilatore che opera nel corso dell'esecuzione di un programma, traducendo il codice interpretato nel codice macchina del calcolatore su cui il programma è in esecuzione.

### Autovalutazione

Quale dei seguenti vantaggi di un interprete rispetto a un compilatore è più importante per un programmatore Java?

1. La facilità di scrivere un interprete.
2. I messaggi di errore più chiari.
3. Il codice oggetto più piccolo.
4. L'indipendenza dal calcolatore.

## 2.13 | Un esempio riassuntivo in linguaggio C

Il rischio che si corre nel mostrare solo piccoli frammenti di codice assembler è di non dare un'idea di quale aspetto abbia un programma completo. Per ovviare a ciò, in questo paragrafo presentiamo il codice RISC-V corrispondente a due procedure scritte in linguaggio C: la prima inverte la posizione di due elementi consecutivi di un vettore e la seconda li ordina.

```
void scambia(long long int v[], int k)
{
    long long int temp;
    temp = v[k];
    v[k] = v[k + 1];
    v[k+1] = temp;
}
```

**Figura 2.23** Procedura C che scambia il contenuto di due locazioni di memoria. Questa procedura sarà utilizzata all'interno della procedura di ordinamento descritta in seguito.

### Procedura scambia

Iniziamo esaminando il codice C della procedura **scambia**, descritta in Figura 2.23, che scambia semplicemente il contenuto di due locazioni di memoria. Per tradurre da C ad assembler occorre seguire i seguenti passaggi:

1. vengono allocati i registri per le variabili del programma;
2. si scrive il codice relativo al corpo della procedura;
3. viene salvato il contenuto dei registri dalle modifiche che possono avvenire nel corso della procedura.

Nel seguito le tre parti della procedura verranno prima descritte separatamente e poi ricomposte.

#### Allocazione dei registri per la procedura scambia

Come descritto in precedenza, la convenzione RISC-V per quanto riguarda il passaggio dei parametri prevede l'utilizzo dei registri da x10 a x17. Dal momento che la procedura **scambia** ha solo due parametri, **v** e **k**, questi si troveranno nei registri x10 e x11. L'unica variabile della procedura è **temp**, che associamo al registro x5. Questo è possibile perché **scambia** è una procedura foglia. Lo schema di allocazione dei registri corrisponde alla dichiarazione delle variabili contenuta nella prima riga della procedura **scambia** di Figura 2.23.

#### Codice del corpo della procedura scambia

Le restanti righe di codice C della procedura **scambia** sono:

```
temp = v[k];
v[k] = v[k + 1];
v[k+1] = temp;
```

Ricordando che gli indirizzi di memoria delle istruzioni RISC-V fanno riferimento ai byte e quindi che due parole doppie consecutive differiscono tra loro di 8 byte, occorre moltiplicare l'indice **k** per 8 prima di sommarlo all'indirizzo base del vettore. *Dimenticare che gli indirizzi di due parole doppie consecutive differiscono di 8 byte, invece che di 1, è un errore molto comune nella programmazione assembler.* Perciò, il primo passo di esecuzione consiste nell'ottenere l'indirizzo di **v[k]** moltiplicando **k** per 8 attraverso un'operazione di shift a sinistra di 3 posizioni:

```
slli x6, x11, 3 // registro x6 = k * 8
add x6, x10, x6 // registro x6 = v + (k * 8)
```

Ora si può caricare **v[k]** utilizzando **x6** e, successivamente, **v[k+1]** aggiungendo 8 a **x6**:

```
ld x5, 0(x6) // registro x5 (temp) = v[k]
ld x7, 8(x6) // registro x7 = v[k + 1]
// si riferisce all'elemento di v successivo
```

Il passo successivo consiste nel memorizzare  $x_9$  e  $x_{11}$  scambiando la loro posizione in memoria:

```
sd x7, 0(x6)    // v[k] = registro x7
sd x5, 8(x6)    // v[k + 1] = registro x5 (temp)
```

A questo punto abbiamo allocato i registri e scritto il corpo di scambia; non resta che aggiungere il codice per salvare il contenuto dei registri eventualmente utilizzati dalla procedura scambia. Dato che non vengono utilizzati registri che devono essere preservati ed essendo scambia una procedura foglia, non dobbiamo salvare il contenuto di alcun registro.

### Procedura scambia completa

Siamo ora pronti per la procedura scambia completa, che comprende l'etichetta della procedura e l'istruzione di ritorno.

```
scambia:
    slli  x6, x11, 3    // reg x6 = k * 8
    add   x6, x10, x6  // reg x6 = v + (k * 8)
    ld    x5, 0(x6)    // reg x5 (temp) = v[k]
    ld    x7, 8(x6)    // reg x7 = v[k + 1]
    sd    x7, 0(x6)    // v[k] = reg x7
    sd    x5, 8(x6)    // v[k+1] = reg x5 (temp)
    jalr  x0, 0(x1)    // ritorno alla procedura chiamante
```

### Procedura ordina

Per capire meglio il rigore richiesto dalla programmazione in linguaggio assembler, descriveremo un secondo esempio, un po' più lungo del primo. Scriveremo una procedura che richiama la procedura scambia, il cui scopo è quello di ordinare gli elementi di un vettore di numeri interi. Utilizzeremo l'algoritmo *Bubble Sort*, uno degli algoritmi di ordinamento più semplici e veloci. La Figura 2.24 mostra la versione C del programma. Anche in questo caso la procedura verrà costruita attraverso passi successivi fino ad arrivare al codice assembler completo.

```
void scambia (long long int v[], int n)
{
    int i, j;
    for (i = 0; i < n; i += 1) {
        for (j = i - 1; j >= 0 && v[j] > v[j + 1]; j -= 1) {
            scambia(v, j);
        }
    }
}
```

**Figura 2.24** Procedura C che opera un ordinamento degli elementi di un vettore v.

### Allocazione dei registri per la procedura ordina

I due parametri della procedura ordina, v e n, saranno contenuti nei registri dei parametri  $x_{10}$  e  $x_{11}$ ; assegniamo i registri  $x_{19}$  e  $x_{20}$  rispettivamente alle variabili i e j.

### Codice del corpo della procedura ordina

Il corpo della procedura consiste in due cicli for annidati e nella chiamata alla procedura scambia, a cui vengono passati i parametri opportuni. Sviluppere-

mo il codice procedendo dall'esterno all'interno. Il primo passo consiste nella traduzione del ciclo *for* più esterno:

```
for (i = 0; i < n; i += 1) {
```

Si ricordi che il ciclo *for* in C è composto da tre parti: inizializzazione, test di fine ciclo e incremento associato a ogni iterazione. È sufficiente una sola istruzione per inizializzare *i* a 0, contenuta nella prima parte dell'istruzione *for*:

```
li x19, 0
```

(Si ricordi che *li* è una pseudoistruzione fornita dall'assemblatore per comodità di programmazione; par. 2.12.) Viene richiesta una seconda istruzione per incrementare il contatore *i*, traducendo così l'ultima parte del ciclo *for*:

```
addi x19, x19, 1 // i += 1
```

Il ciclo termina se la condizione *i* < *n* non è vera o, in altre parole, se *i* ≥ *n*. Questo confronto richiede solamente una istruzione:

```
for1tst: bge x19, x11, escil // vai a escil se
           // x19 ≥ x1 (i ≥ n)
```

L'ultima istruzione del ciclo sarà un salto incondizionato all'inizio del ciclo stesso:

```
j for1tst // Torna al test contenuto
           // nel ciclo più esterno
escil:
```

Lo scheletro del codice per il ciclo *for* più esterno sarà quindi:

```
li x19, 0          // i = 0
for1tst:
    bge x19, x11, escil // vai a escil se
                           // x19 ≥ x1 (i ≥ n)
    ...
    (corpo del ciclo for più esterno)
    ...
    addi x19, x19, 1 // i += 1
    j for1tst        // torna al test del
                       // ciclo più esterno
escil:
```

Ecco fatto. (Negli esercizi potrete sperimentare come scrivere una procedura che esegua più velocemente cicli di questo tipo.)

Il ciclo *for* più interno ha la seguente forma in linguaggio C:

```
for (j = i - 1; j >= 0 && v[j] > v[j + 1]; j -= 1) {
```

La parte di inizializzazione è anche in questo caso composta da una istruzione:

```
addi x20, x19, -1 // j = i - 1
```

Anche il decremento di *j* alla fine del ciclo può essere tradotto in una singola istruzione:

```
addi x20, x20, -1 // j -= 1
```

Il test di fine ciclo è composto da due parti. Si esce dal ciclo se almeno una delle condizioni non è verificata; il primo test deve quindi provocare l'uscita dal ciclo se *j* < 0:

```
for2tst:
    blt x20, x0, esci2 // vai a esci2 se
                           // x20 < 0 (j < 0)
```

L'ultima istruzione fa sì che si salti il secondo test se non risulta  $j \geq 0$ .

Il secondo test impone che il ciclo termini se  $v[j] > v[j + 1]$  non è vera o, in altre parole, se  $v[j] \leq v[j + 1]$ . Per prima cosa occorre ottenere l'indirizzo dell'elemento del vettore, moltiplicando  $j$  per 8 (dato che gli indirizzi sono sempre espressi in byte) e sommarlo all'indirizzo base di  $v$ :

```
slli x5, x20, 3 // registro x5 = j * 8
add x5, x10, x5 // registro x5 = v + (j * 8)
```

Ora è possibile caricare  $v[j]$ :

```
ld x6, 0(x5) // registro x6 = v[j]
```

Dal momento che l'elemento successivo del vettore è semplicemente la parola successiva in memoria, basta sommare 8 all'indirizzo contenuto in  $x5$  per ottenere l'indirizzo di  $v[j + 1]$ :

```
ld x7, 8(x5) // registro x7 = v[j + 1]
```

Verifichiamo se  $v[j] \leq v[j + 1]$  per uscire dal ciclo:

```
ble x6, x7, esci2 // vai a esci2 se x6 \leq x7
```

L'ultima istruzione del ciclo sarà un salto incondizionato all'inizio del ciclo interno stesso:

```
j for2tst // torna al test di ingresso del ciclo più interno
```

Unendo i vari frammenti di codice, lo scheletro del secondo ciclo *for* diventa:

```
addi x20, x19, -1 // j = i - 1
for2tst: blt x20, x0, esci2 // vai a esci2 se
           // x20 < 0 (j < 0)
           slli x5, x20, 3 // registro x5 = j * 8
           add x5, x10, x5 // registro x5 = v + (j * 8)
           ld x6, 0(x5) // registro x6 = v[j]
           ld x7, 8(x5) // registro x7 = v[j + 1]
           ble x6, x7, esci2 // vai a esci2 se x6 \leq x7
           ...
           (corpo del ciclo for più interno)
           ...
           addi x20, x20, -1 // j -= 1
           j for2tst // salta al test del ciclo
           // più interno
esci2:
```

### Chiamata alla procedura contenuta in ordina

Il passo successivo consiste nel tradurre il corpo del ciclo *for* più interno:

```
scambia(v, j);
```

La chiamata alla procedura *scambia* è molto semplice:

```
jal x1, scambia
```

### Passaggio dei parametri in ordina

Il problema nasce dal fatto che la procedura *ordina* ha bisogno di conservare i propri parametri, contenuti nei registri  $x10$  e  $x11$ , mentre la procedura *scambia* ha bisogno di inserire i propri parametri in questi stessi due registri. Una possibile soluzione consiste nel copiare i parametri della procedura *ordina*

in altri registri nella fase iniziale della procedura stessa, in modo da rendere disponibili per la procedura scambia i registri x10 e x11 (la copia in un registro è più veloce del salvataggio e della lettura dei dati dallo stack). Quindi all'inizio della procedura si copia il contenuto di x10 e x11 rispettivamente in x21 e x22:

```
mv x21, x22      // copia del parametro x10 in x21
mv x22, x11      // copia del parametro x11 in x22
```

A questo punto si possono passare i parametri alla procedura scambia con le due istruzioni seguenti:

```
mv x10, x21      // il primo parametro
                  // di scambia è v
mv x11, x20      // il secondo parametro
                  // di scambia è j
```

### Salvataggio dei registri in ordina

Rimane ancora da scrivere la parte di codice relativa al salvataggio e al ripristino dei registri. Chiaramente è necessario salvare l'indirizzo di ritorno, contenuto nel registro x1, dal momento che ordina è una procedura e quindi verrà chiamata a sua volta da un'altra procedura. Inoltre ordina utilizza i registri x19, x20, x21 e x22 che devono quindi essere salvati. La parte iniziale della procedura sarà quindi:

```
addi sp, sp, -40    // crea spazio nello stack per 5 registri
sd x1, 32(sp)      // salva x1 in stack
sd x22, 24(sp)     // salva x22 in stack
sd x21, 16(sp)     // salva x21 in stack
sd x20, 8(sp)      // salva x20 in stack
sd x19, 0(sp)      // salva x19 in stack
```

La parte finale della procedura semplicemente inverte tutte queste istruzioni e le conclude con una jalr all'indirizzo di ritorno.

### Procedura ordina completa

Possiamo ora mettere assieme tutti i diversi frammenti della procedura e ottenere il codice mostrato in Figura 2.25; occorre avere cura di sostituire i riferimenti ai registri x10 e x11 all'interno dei cicli *for* con i riferimenti ai registri x21 e x22. Per rendere il codice più leggibile, a ciascun blocco di codice viene associato il suo ruolo. In questo esempio le 9 linee di codice C della procedura ordina sono diventate 34 righe in codice assembler RISC-V.

#### Salvataggio dei registri

ordina:	addi sp, sp, -40	// crea spazio per 5 registri in stack
	sd x1, 32(sp)	// salva l'indirizzo di ritorno in stack
	sd x22, 24(sp)	// salva x22 in stack
	sd x21, 16(sp)	// salva x21 in stack
	sd x20, 8(sp)	// salva x20 in stack
	sd x19, 0(sp)	// salva x19 in stack

#### Corpo della procedura

Spostamento dei parametri	mv x21, x10	// copia il parametro x10 in x21
	mv x22, x11	// copia il parametro x11 in x22
Ciclo esterno	li x19, 0	// i = 0
	forltst:bge x19, x22, escil	// vai a escil se i >= n

Corpo della procedura		
Ciclo interno	addi x20, x19, -1	// j = i - 1
	for2tst: blt x20, x0, esci2	// vai a esci2 se j < 0
	slli x5, x20, 3	// x5 = j * 8
	add x5, x21, x5	// x5 = v + (j * 8)
	ld x6, 0(x5)	// x6 = v[j]
	ld x7, 8(x5)	// x7 = v[j + 1]
	ble x6, x7, esci2	// vai a esci2 se x6 < x7
Passaggio dei parametri e chiamata	mv x10, x21	// il primo parametro di ordina è v
	mv x11, x20	// il secondo parametro di ordina è j
	jal x1, scambia	// chiama scambia
Ciclo interno	addi x20, x20, -1	j for2tst
		// vai a for2tst
Ciclo esterno	esci2: addi x19, x19, 1	// i += 1
	j for1tst	// vai a for1tst
Ripristino dei registri		
esci1:	ld x19, 0(sp)	// ripristina x19 dallo stack
	ld x20, 8(sp)	// ripristina x20 dallo stack
	ld x21, 16(sp)	// ripristina x21 dallo stack
	ld x22, 24(sp)	// ripristina x22 dallo stack
	ld x1, 32(sp)	// ripristina l'indirizzo di ritorno dallo stack
	addi sp, sp, 40	// ripristina lo stack pointer
Ritorno alla procedura		
	jalr x0, 0(x1)	// ritorno alla procedura chiamante

**Figura 2.25** Versione in linguaggio assembler RISC-V della procedura `ordina` di Figura 2.24.

**Approfondimento.** Un tipo di ottimizzazione che in questo caso può dare dei buoni risultati è l'*inserimento in linea (in-lining) della procedura*. Invece di passare gli argomenti sotto forma di parametri e di richiamare il codice della procedura `scambia` con un'istruzione `jal`, il compilatore copia il codice del corpo della procedura `scambia` nel punto in cui compare la sua chiamata. In questo esempio l'inserimento in linea del codice eliminerebbe quattro istruzioni. Lo svantaggio di questa ottimizzazione è che il codice compilato sarà di dimensioni maggiori quando la procedura inserita in linea viene chiamata da più posizioni. Questa espansione del codice potrebbe peggiorare le prestazioni se aumentasse la frequenza di miss della cache (vedi Cap. 5).

La Figura 2.26 mostra l'impatto delle ottimizzazioni effettuate dal compilatore sulle prestazioni di un programma di ordinamento in termini di tempo di compilazione, cicli di clock, numero di istruzioni e CPI. Si può notare come il codice non ottimizzato abbia il CPI migliore, l'ottimizzazione di tipo O1 produca il minor numero di istruzioni, ma l'ottimizzazione O3 sia la più veloce. Si ricordi che, in ultima analisi, il tempo è l'unica misura accurata delle prestazioni di un programma.

La Figura 2.27 confronta l'impatto dei linguaggi di programmazione, dell'utilizzo di codice compilato o interpretato, e dell'algoritmo sulle prestazioni del programma di ordinamento. Le quattro colonne mostrano che il programma non ottimizzato C è 8,3 volte più veloce di quello scritto in codice Java interpretato per quanto riguarda l'algoritmo Bubble Sort. Utilizzando il compilatore JIT Java, il codice Java diventa 2,1 volte più veloce del codice C non ottimizzato e risulta più lento solamente di un fattore 1,13 rispetto al codice C con il massimo livello di ottimizzazione. (Il paragrafo 2.15  fornisce maggiori dettagli sul confronto tra

## Capire le prestazioni dei programmi

(continua)

(continua)

compilazione e interpretazione di Java, e riporta il codice Java della procedura di ordinamento tramite Bubble Sort.) Le prestazioni non sono altrettanto vicine quando si considera l'algoritmo di ordinamento Quicksort (colonna 5), probabilmente perché non si riesce ad ammortizzare il costo della compilazione JIT con un tempo di esecuzione così breve. L'ultima colonna mostra l'impatto sulle prestazioni di un algoritmo migliore: si può ottenere un aumento della velocità di esecuzione addirittura di tre ordini di grandezza per l'ordinamento di 100 000 elementi. Anche confrontando le prestazioni del codice Java interpretato per Quicksort, riportate nella colonna 5, con quelle del codice C compilato per Bubble Sort riportate nella colonna 4, l'algoritmo Quicksort batte Bubble Sort di un fattore 50 ( $0,05 \times 2468$ , ovvero il codice è 123 volte più veloce del codice C non ottimizzato, contro un aumento di velocità di 2,41 volte).

Ottimizzazione gcc	Prestazioni relative	Cicli di clock (milioni)	Numero di istruzioni (milioni)	CPI
Nessuna	1,00	158 615	114 938	1,38
01 (media)	2,37	66 990	37 470	1,79
02 (completa)	2,38	66 521	39 993	1,66
03 (integrazione delle procedure)	2,41	65 747	44 993	1,46

**Figura 2.26** Confronto delle prestazioni in termini di numero di istruzioni e di CPI, utilizzando diverse ottimizzazioni del compilatore per l'implementazione dell'algoritmo Bubble Sort. Il programma ordina 100 000 valori di una parola, contenuti in un vettore inizializzato con valori casuali. I programmi sono stati eseguiti su un Pentium 4 con una frequenza di clock di 3,06 GHz, con bus di sistema a 533 MHz e con 2 GB di memoria SDRAM (PC2100 DDR). È stato utilizzato il kernel di Linux 2.4.20.

Linguaggio	Metodo esecuzione	Ottimizzazione	Prestazioni relative per Bubble Sort	Prestazioni relative per Quicksort	Speedup di Quicksort rispetto a Bubble Sort
C	Compilato	Nessuna	1,00	1,00	2468
	Compilato	01	2,37	1,50	1562
	Compilato	02	2,38	1,50	1555
	Compilato	03	2,41	1,91	1955
Java	Interpretato	-	0,12	0,05	1050
	Compilato Just In Time	-	2,13	0,29	338

**Figura 2.27** Confronto delle prestazioni di due programmi di ordinamento ottenute eseguendo un codice C ottenuto da un compilatore con diversi livelli di ottimizzazione e interpretando il corrispondente codice Java. Le prestazioni sono riferite a quelle ottenute con codice C non ottimizzato. Nell'ultima colonna è riportato l'incremento di prestazioni di Quicksort rispetto a Bubble Sort per ciascuna delle versioni di codice. I programmi sono stati eseguiti sullo stesso sistema con cui sono stati ottenuti i dati riportati in Figura 2.26. La JVM è la versione 1.3.1 della Sun, e il compilatore JIT è Hotspot, versione 1.3.1, della Sun.

## 2.14 | Confronto tra vettori e puntatori

La comprensione dei puntatori è una delle maggiori sfide per chi inizia a programmare in C. Il confronto tra un programma assembler che utilizza vettori e indici e un programma che utilizza i puntatori aiuta a capire meglio il funzionamento di questi ultimi. Questo paragrafo mostra la versione C e assembler di due procedure che azzerano il contenuto di una sequenza di parole in memoria:

```

azzeral(long long int vettore[], int dim)
{
    int i;
    for (i = 0; i < dim; i += 1)
        vettore[i] = 0;
}
azzerar2(long long int *vettore, int dim)
{
    long long int *p;
    for (p = &vettore[0]; p < &vettore[dim]; p = p + 1)
        *p = 0;
}

```

**Figura 2.28 Due procedure C che azzerano tutti gli elementi di un vettore.** Azzeral utilizza un indice, mentre azzerar2 utilizza un puntatore. La seconda procedura richiede qualche spiegazione per chi non ha familiarità con il linguaggio C. L'indirizzo di una variabile viene indicato con &, mentre l'oggetto puntato da un puntatore viene indicato con \*. Dalla dichiarazione delle variabili si può vedere che sia vett sia p sono puntatori a interi. La prima parte dell'istruzione for della procedura azzerar2 assegna l'indirizzo del primo elemento di vett al puntatore p, mentre la seconda parte effettua un controllo per vedere se il puntatore ha superato l'ultimo elemento di vett. L'incremento di un'unità del puntatore nell'ultima parte dell'istruzione for equivale a spostare il puntatore all'elemento successivo, incrementando l'indirizzo di un numero di byte pari a quello dichiarato: dato che p è un puntatore a interi, il compilatore genera delle istruzioni RISC-V che incrementano p di 8, che corrisponde al numero di byte che costituiscono un intero nei RISC-V. L'istruzione di assegnamento scrive uno 0 nell'oggetto puntato da p.

la prima utilizza l'indice di un vettore, la seconda i puntatori. La versione C delle due procedure è mostrata in Figura 2.28.

Lo scopo del paragrafo è mostrare come i puntatori si possano tradurre in istruzioni RISC-V e non di raccomandare uno stile di programmazione oramai datato. Al termine vedremo l'impatto delle ottimizzazioni dei moderni compilatori su queste due procedure.

### Versone della procedura azzerar che utilizza vettore e indice

Iniziamo dalla prima versione della procedura, azzeral, che utilizza un vettore, concentrando sul corpo del ciclo e ignorando il codice di collegamento. Supponiamo che i due parametri vett e dim si trovino nei registri x10 e x11, e che i sia allocato nel registro x5.

L'inizializzazione di i, cioè la prima parte dell'istruzione for, è elementare:

```
li      x5, 0 // i = 0 (registro x5 = 0)
```

Per impostare vett[i] a 0 occorre per prima cosa determinare il suo indirizzo; per fare ciò si moltiplica i per 8 e si ottiene l'indirizzo in byte:

```
ciclo1: slli x6, x5, 3 // x6 = i * 8
```

Dal momento che l'indirizzo di partenza del vettore è in un registro, lo si deve sommare al valore dell'indice per ottenere vett[i]:

```
add x7, x10, x6 // x7 = indirizzo di vett[i]
```

Infine, possiamo scrivere 0 in questo indirizzo:

```
sd x0, 0(x7) // vett[i] = 0
```

Questa istruzione si trova alla fine del corpo del ciclo for, quindi il passo successivo consiste nell'incrementare i:

```
addi x5, x5, 1 // i = i + 1
```

Il test di fine ciclo verifica se i è minore della dimensione del vettore:

```
blt x5, x11, ciclo1 // se (i < dim) vai a ciclo1
```

Abbiamo così visto tutte le parti della procedura. Il codice RISC-V completo che azzerà tutti gli elementi del vettore (utilizzando un indice) sarà quindi il seguente:

```
lv x5, 0           // i = 0
ciclol: slli x6, x5, 3    // x6 = i * 8
        add x7, x10, x6    // x7 = indirizzo di vett[i]
        sd x0, 0(x7)      // vett[i] = 0
        addi x5, x5, 1     // i = i + 1
        blt x5, x11, ciclol // se (i < dim) vai a ciclol
```

(Si noti che questo frammento di codice funziona solo se `dim` è maggiore di 0; il C ANSI richiede che venga effettuato un test su `dim` prima del ciclo, ma in questo esempio non considereremo questo vincolo.)

### Versione della procedura azzerà che utilizza i puntatori

La seconda procedura, che utilizza i puntatori, alloca i due parametri `vett` e `dim` nei registri `x10` e `x11` e il puntatore `p` in `x5`. Il codice di questa seconda procedura inizia con l'assegnazione dell'indirizzo del primo elemento del vettore al puntatore `p`:

```
mv x5, x10 // p = indirizzo di vett[0]
```

La parte successiva è relativa al codice del corpo del ciclo `for`, che semplicemente scrive uno 0 nella posizione puntata da `p`:

```
ciclo2: sd x0, 0(x5) // Memoria[p] = 0
```

Questa istruzione costituisce il corpo del ciclo, mentre l'istruzione seguente aggiorna `p` per fare sì che punti alla parola successiva:

```
addi x5, x5, 8 // p = p + 8
```

In C incrementare un puntatore di 1 significa farlo puntare all'oggetto successivo della sequenza. Dato che `p` è un puntatore a interi, ciascuno dei quali occupa 4 byte, il compilatore incrementa `p` di 8 unità.

La fase seguente consiste nel test di fine ciclo. Per prima cosa occorre calcolare l'indirizzo dell'ultimo elemento del vettore, moltiplicando `dim` per 8 al fine di ottenere l'indirizzo in byte:

```
slli x6, x11, 3 // x6 = dim * 8
```

Il prodotto ottenuto viene quindi sommato all'indirizzo di partenza del vettore, per ottenere l'indirizzo della prima parola che segue il vettore:

```
add x7, x10, x6 // x7 = indirizzo di vett[dim]
```

Il test di fine ciclo controlla semplicemente che `p` sia minore dell'ultimo elemento di `vett`:

```
bltu x5, x7 ciclo2 // se (p < &vett[dim]) vai a ciclo2
```

Abbiamo così completato tutte le parti del programma. Possiamo ora vedere la versione completa della procedura che azzerà un vettore, basata sui puntatori:

```
mv x5, x10 // p = indirizzo di vett[0]
ciclo2: sd x0, 0(x5) // memoria[p] = 0
        addi x5, x5, 8 // p = p + 8
        slli x6, x11, 3 // x6 = dim * 8
        add x7, x10, x6 // x7 = indirizzo di vett[dim]
        bltu x5, x7 ciclo2 // se (p < &vett[dim]) vai a ciclo2
```

Come nel caso precedente, si supponga che `dim` sia maggiore di 0.

Si noti che questa versione del programma calcola l'indirizzo corrispondente alla fine del vettore durante ogni iterazione del ciclo, anche se esso non cambia.

Spostando questo calcolo fuori dal ciclo si ottiene una versione più veloce:

```
mv x5, x10      // p = indirizzo di vett[0]
slli x6, x11, 3  // x6 = dim * 8
add x7, x10, x6  // x7 = indirizzo di vett[dim]
ciclo2: sd x0, 0(x5) // memoria[p] = 0
        addi x5, x5, 8    // p = p + 8
        bltu x5, x7 ciclo2 // se (p < &vett[dim]) vai a ciclo2
```

## Confronto tra le due versioni di azzera

Il confronto delle due versioni di codice, riportate una a fianco dell'altra, consente di evidenziare le differenze fra indici e puntatori (i cambiamenti introdotti dall'uso dei puntatori sono discussi nel testo):

<pre>li x5, 0          // i = 0 ciclo1: slli x6, x5, 3   // x6 = i * 8         add x7, x10, x6  // x7 = indirizzo di vett[i]         sd x0, 0(x7)     // vett[i] = 0         addi x5, x5, 1    // i = i + 1         blt x5, x11, ciclo1 // se (i &lt; dim) vai a ciclo1</pre>	<pre>mv x5, x10      // p = indirizzo di vett[0] slli x6, x11, 3  // x6 = dim * 8 add x7, x10, x6  // x7 = indirizzo di vett[dim] ciclo2: sd x0, 0(x5) // memoria[p] = 0         addi x5, x5, 8    // p = p + 8         bltu x5, x7, ciclo2 // se (p &lt; &amp;vett[dim]) vai a ciclo2</pre>
---	--

La versione di sinistra deve necessariamente contenere le operazioni di moltiplicazione e somma all'interno del ciclo, dato che `i` viene incrementato e quindi l'indirizzo di ciascun elemento del vettore deve essere calcolato a partire dal nuovo valore dell'indice, mentre la versione di destra incrementa direttamente il puntatore alla memoria `p`. La versione con i puntatori sposta le istruzioni che determinano la posizione in memoria del primo e dell'ultimo elemento del vettore fuori dal ciclo, riducendo così da cinque a tre il numero di istruzioni eseguite a ogni ciclo.

Questa ottimizzazione effettuata manualmente corrisponde a due ottimizzazioni effettuate automaticamente dai compilatori: riduzione del peso delle istruzioni (operazione di shift invece che di moltiplicazione) ed eliminazione di variabili per induzione (eliminazione del calcolo dell'indirizzo degli elementi del vettore all'interno del ciclo). Il paragrafo 2.15  descrive queste e altre ottimizzazioni.

**Approfondimento.** Come abbiamo detto in precedenza, un compilatore C aggiungerebbe un test per assicurarsi che `dim` sia maggiore di 0; un modo per farlo sarebbe saltare all'istruzione dopo il ciclo mediante `blt x0, 11, dopo Ciclo`.

Tradizionalmente veniva raccomandato agli aspiranti programmatori di utilizzare i puntatori in C per ottenere un codice più efficiente di quello ottenibile utilizzando indici e vettori. I compilatori ottimizzanti moderni sono in grado di produrre, partendo da un codice scritto con i vettori, un codice altrettanto efficiente. Molti programmati attualmente preferiscono che sia il compilatore a effettuare le ristrutturazioni pesanti del codice.

**Capire le prestazioni  
dei programmi**

## 2.15 Approfondimento: compilazione del C e interpretazione di Java

Questo paragrafo contiene una breve introduzione sul funzionamento di un compilatore C e su come viene eseguito un programma Java. Dato che il compilatore ha un grande impatto sulle prestazioni, comprendere la tecnologia dei compilatori moderni è un elemento critico per valutarle. La descrizione di come si costruisce un compilatore è oggetto di corsi della durata di uno o due semestri interi, per cui in questa introduzione possiamo solamente descrivere i concetti fondamentali.

La seconda parte del paragrafo è rivolta ai lettori interessati a comprendere come un **linguaggio orientato agli oggetti**, quale Java, venga eseguito su un'architettura RISC-V. Vengono mostrati i bytecode Java utilizzati per l'interpretazione del codice e il RISC-V che si ottiene dalla versione Java di alcuni frammenti di procedura riportati nei paragrafi precedenti, tra i quali quello relativo all'algoritmo di Bubble Sort. Viene descritto, inoltre, il funzionamento sia della Java Virtual Machine sia dei compilatori Just In Time. Troverete il testo di questo paragrafo online .

**Linguaggio orientato agli oggetti:** un linguaggio di programmazione che è orientato agli oggetti piuttosto che alle azioni, in altre parole che è orientato ai dati piuttosto che alla logica di funzionamento.

## 2.16 Un caso reale: le istruzioni dell'architettura MIPS

L'insieme di istruzioni più vicino al RISC-V è quello del MIPS. Anch'esso è nato nell'Università ma è ora proprietà della Imagination Technologies. I MIPS e i RISC-V condividono la stessa filosofia di progetto, anche se il MIPS ha 25 anni più del RISC-V. La buona notizia è che se conoscete il RISC-V, è veramente facile conoscere anche il MIPS. Per illustrare il loro grado di somiglianza, la Figura 2.29 compara i formati delle istruzioni dei RISC-V e dei MIPS.

L'ISA dei MIPS ha una versione a 32 bit e una a 64 bit, saggiamente denominate MIPS-32 e MIPS-64. Questi insiemi di istruzioni sono virtualmente

Registro-registro									
RISC-V	31	25 24	20 19	15 14	12 11	7	6	0	
		funz7(7)	rs2(5)	rs1(5)	funz3(3)	rd(5)			codop(7)
MIPS	31	26 25	21 20	16 15	11 10	6	5	0	
		Op(6)	Rs1(5)	Rs2(5)	Rd(5)	Cost(5)			Opx(6)
Trasferimento dalla memoria									
RISC-V	31	immed(12)	20 19	15 14	12 11	7	6	0	
				rs1(5)	funz3(3)	rd(5)			codop(7)
MIPS	31	26 25	21 20	16 15				0	
		Op(6)	Rs1(5)	Rs2(5)		Cost(16)			
Trasferimento alla memoria									
RISC-V	31	immed(7)	25 24	20 19	15 14	12 11	7	6	0
			rs2(5)	rs1(5)	funz3(3)	immed(5)			codop(7)
MIPS	31	26 25	21 20	16 15				0	
		Op(6)	Rs1(5)	Rs2(5)		Cost(16)			
Salto condizionato									
RISC-V	31	immed(7)	25 24	20 19	15 14	12 11	7	6	0
			rs2(5)	rs1(5)	funz3(3)	immed(5)			codop(7)
MIPS	31	26 25	21 20	16 15				0	
		Op(6)	Rs1(5)	Opx/Rs2(5)		Cost(16)			

**Figura 2.29** Formati delle istruzioni RISC-V e MIPS. Le similità provengono in parte dal fatto che entrambi gli insiemi di istruzioni si appoggiano su 32 registri.

identici ad eccezione della dimensione maggiore dello spazio di indirizzamento che richiede registri a 64 bit invece che a 32 bit. Riportiamo ora le caratteristiche comuni a RISC-V e MIPS:

- Tutte le istruzioni sono ampie 32 bit in entrambe le architetture.
- Entrambe le architetture hanno 32 registri a uso generale, con un registro impostato fisso a 0.
- L'unico modo per accedere alla memoria è attraverso le istruzioni di load e store in entrambe le architetture.
- A differenza di alcune architetture, non ci sono istruzioni nei MIPS e nei RISC-V che possano trasferire alla o dalla memoria più di un registro.
- Entrambe le architetture hanno istruzioni che possono saltare se il contenuto di un registro è uguale a zero o è diverso da zero.
- Le modalità di indirizzamento in entrambe le architetture possono trasferire dati di tutte le dimensioni.

Una delle differenze principali tra i RISC-V e i MIPS è sui salti condizionati basati su condizioni diverse dall'uguaglianza o disuguaglianza. Mentre il RISC-V fornisce semplicemente istruzioni di salto condizionato che confrontano il contenuto di due registri, i MIPS si devono basare su un'istruzione di comparazione che imposta il contenuto di un registro a 0 o a 1 a seconda che il risultato del confronto sia vero o falso. I programmati devono fare seguire questa istruzione di comparazione con un'istruzione di saldo condizionato basato su uguaglianza o disuguaglianza, a seconda del risultato desiderato della comparazione. Obbedendo alla sua filosofia minimalista, i MIPS eseguono soltanto il confronto di tipo "minore di", lasciando al programmatore l'onere di scambiare l'ordine degli operandi o la condizione che deve verificare il salto condizionato per ottenere l'uscita desiderata. I MIPS hanno sia la versione con segno che senza segno dell'istruzione di imposta se minore: `slt` e `sltu`.

Guardando oltre il nucleo delle istruzioni che vengono utilizzate più comunemente, l'altra differenza principale è che il MIPS completo ha un insieme di istruzioni molto più ampio del RISC-V, come vedremo nel paragrafo 2.18.

## 2.17 | Un caso reale: le istruzioni dell'architettura x86

I progettisti di insiemi di istruzioni talvolta forniscono istruzioni più potenti di quelle che si trovano nel RISC-V e nel MIPS. Lo scopo è di ridurre il numero di istruzioni eseguite da un programma, ma il rischio è che tale riduzione vada a scapito della semplicità aumentando il tempo di esecuzione dei programmi, perché le istruzioni diventano più complesse. Ciò può derivare da un periodo del ciclo di clock più lungo o dalla necessità di un maggior numero di cicli rispetto alle istruzioni più semplici.

La strada che porta alla creazione di istruzioni più complesse è dunque irta di pericoli. Alcuni trabocchetti legati alla complessità delle istruzioni sono descritti nel paragrafo 2.19.

### Evoluzione dell'Intel x86

RISC-V e MIPS nacquero dalla visione di singoli gruppi di ricercatori che lavorarono nello stesso periodo; i diversi blocchi di queste architetture si adattano bene l'uno all'altro e l'intera architettura può essere descritta in modo conciso. Questo non è il caso delle architetture Intel x86, che sono il risultato del lavoro indipendente di molti gruppi di progettisti che hanno continuato a sviluppare l'architettura nel corso di quasi 40 anni, raggiungendo progressivamente nuove

*La bellezza è tutta negli occhi  
di chi guarda.*

Margaret Wolfe Hungerford,  
*Molly Bawn*, 1877

**Registro di utilizzo generale (GPR):** un registro che può essere utilizzato praticamente da tutte le istruzioni, per contenere indirizzi o dati.

caratteristiche all'insieme di istruzioni originario. Le pietre miliari dell'evoluzione degli x86 sono le seguenti:

- **1978:** l'architettura Intel 8086 è introdotta come estensione, compatibile a livello di linguaggio assembler, dell'allora popolarissimo Intel 8080, un microprocessore a 8 bit. L'8086 è un'architettura a 16 bit con tutti i registri interni a 16 bit. A differenza del MIPS, i registri erano dedicati a usi specifici, quindi l'8086 non può essere considerata un'architettura con **registri di utilizzo generale (GPR, General-Purpose Register)**.
- **1980:** viene annunciato il coprocessore matematico Intel 8087 che estende l'architettura 8086 con 60 nuove istruzioni in virgola mobile. Invece di utilizzare registri, si serve di uno stack per le proprie operazioni (*vedi* parr. 2.21 e 3.7).
- **1982:** l'80286 estende l'architettura 8086 allargando lo spazio di indirizzamento a 24 bit mediante la creazione di un complesso sistema di traduzione e protezione degli indirizzi di memoria (*vedi* Cap. 5). Vengono anche aggiunte alcune istruzioni per completare l'insieme di istruzioni e per gestire il modello di protezione.
- **1985:** l'80386 estende a 32 bit l'architettura 80286. Oltre a fornire un'architettura a 32 bit, con registri da 32 bit e uno spazio di indirizzamento di 32 bit, l'80386 contiene nuove modalità di indirizzamento e nuove operazioni. Le istruzioni aggiunte rendono l'80386 praticamente un calcolatore a registri di utilizzo generale. L'80386 fornisce un supporto alla paginazione, oltre che alla segmentazione della memoria (*vedi* Cap. 5). Come l'80286, anche l'80386 possiede una modalità che permette di eseguire i programmi scritti per l'8086 senza apportare modifiche.
- **1989-1995:** i processori successivi, l'80486 del 1989, il Pentium del 1992 e il Pentium Pro del 1995, mirano a ottenere prestazioni sempre più elevate; vengono aggiunte solamente quattro nuove istruzioni all'insieme delle istruzioni visibile all'utente: tre per la gestione del *multiprocessing* (*vedi* Cap. 6) e una per il trasferimento dati condizionato.
- **1997:** dopo l'introduzione del Pentium e del Pentium Pro, l'Intel annuncia l'espansione di queste architetture con l'insieme di istruzioni MMX (*Multi Media Extensions*). Questo nuovo insieme di 57 istruzioni utilizza lo stack dei dati in virgola mobile per accelerare le applicazioni legate alla multimedialità e alla comunicazione. Le istruzioni MMX tipicamente operano su più elementi alla volta, ciascuno di piccola dimensione, secondo la tradizione delle architetture SIMD (*Single Instruction Multiple Data*; *vedi* Cap. 6). Il Pentium II non aggiunge nuove istruzioni.
- **1999:** Intel aggiunge altre 70 istruzioni, chiamate SSE (*Streaming SIMD Extensions*), come parte del Pentium III. Le modifiche principali consistono nell'introdurre otto nuovi registri separati, nel raddoppiarne la dimensione a 128 bit e nell'introdurre come tipo di dato i numeri in virgola mobile in singola precisione. Ciò permette di effettuare fino a quattro operazioni in virgola mobile su 32 bit in parallelo. Per migliorare le prestazioni della memoria, l'insieme SSE comprende istruzioni per il precaricamento della cache, oltre a istruzioni di memorizzazione dei dati in streaming che evitano il passaggio dalla cache e scrivono direttamente in memoria.
- **2001:** Intel estende l'insieme di istruzioni con altre 144 istruzioni, chiamate questa volta SSE2. Vengono introdotti i numeri in virgola mobile in doppia precisione, il che permette di effettuare due operazioni in virgola mobile su 64 bit in parallelo. Quasi tutte le 144 istruzioni sono estensioni di istruzioni presenti negli insiemi MMX e SSE, che operano in parallelo su dati a 64 bit. Questo cambiamento non solo permette più operazioni multimediali, ma dà anche la possibilità al compilatore di scegliere un registro destinazione dif-

ferente per le operazioni in virgola mobile rispetto all'architettura a singolo stack: infatti, il compilatore può scegliere di utilizzare gli otto registri SSE come registri in virgola mobile, come negli altri calcolatori. Queste modifiche hanno aumentato molto le prestazioni sulle operazioni in virgola mobile del Pentium 4, il primo microprocessore a includere le istruzioni SSE2.

- 2003: questa volta è un'azienda diversa da Intel a migliorare l'architettura x86. AMD annuncia una serie di estensioni per aumentare lo spazio di indirizzamento da 32 a 64 bit. In maniera del tutto simile alla transizione dello spazio di indirizzamento da 16 a 32 bit avvenuta nel 1985 con l'80386, l'AMD64 allarga tutti i registri a 64 bit. Inoltre sale a 16 unità sia il numero di registri di tipo generale sia il numero dei registri a 128 bit di tipo SSE. Il principale cambiamento nell'ISA è l'aggiunta di una nuova modalità, chiamata *long mode* (modalità lunga), che ridefinisce l'esecuzione di tutte le istruzioni x86 con indirizzi e dati a 64 bit. Inoltre, per poter indirizzare un maggior numero di registri, viene aggiunto alle istruzioni un nuovo prefisso. Questa modalità lunga aggiunge da 4 a 10 nuove istruzioni, a seconda di come viene effettuato il conto, ed elimina 27 delle istruzioni introdotte in precedenza. Un'ulteriore estensione riguarda l'indirizzamento relativo al PC. L'AMD64 mantiene ancora una modalità di funzionamento identica a quella dell'x86 (*legacy mode*, modalità ereditaria), ma la affianca a una modalità che vincola i programmi utente a utilizzare le caratteristiche dell'x86, e permette anche al sistema operativo di utilizzare tutte le caratteristiche dell'AMD64 (*compatibility mode*, modalità compatibile). Queste modalità permettono una transizione all'indirizzamento a 64 bit più morbida delle architetture HP/Intel IA-64.
- 2004: Intel capitolà e adotta la soluzione AMD64, rinominandola *tecnologia di estensione della memoria a 64 bit* (EM64T, Extended Memory 64 Technology). La differenza principale tra le due architetture è che Intel ha aggiunto un'operazione atomica di confronto e scambio a 128 bit che probabilmente avrebbe dovuto essere inclusa anche nell'AMD64. Allo stesso tempo Intel annuncia un'altra generazione di estensioni multimediali: SSE3, che aggiunge 13 istruzioni per supportare operazioni aritmetiche con numeri complessi, operazioni grafiche su vettori di strutture, codifica video, conversione in virgola mobile e sincronizzazione di thread (par. 2.11). AMD ha aggiunto l'estensione SSE3 e l'operazione atomica di swap che mancava nell'AMD64 per poter mantenere la compatibilità binaria con Intel.
- 2006: Intel annuncia 54 nuove istruzioni come parte dell'estensione dell'insieme di istruzioni SSE4. Questa estensione consente di eseguire funzioni particolari come la somma di differenze, prese in valore assoluto, i prodotti vettore su vettori di strutture, l'aggiunta di zeri o l'estensione del segno di dati codificati su un numero ridotto di bit, il conteggio degli elementi di una popolazione e così via. L'SSE4 fornisce anche il supporto alle macchine virtuali (vedi Cap. 5).
- 2007: AMD annuncia 170 istruzioni come parte dell'SSE5, comprendenti 46 istruzioni con tre operandi che completano l'insieme di base, come nel RISC-V.
- 2011: Intel rilascia l'estensione *Advanced Vector Extension*, che espande la larghezza dei registri da 128 a 256 bit, ridefinendo di conseguenza circa 250 istruzioni e aggiungendone 128 nuove.

Come si vede da questa cronologia, la compatibilità ha influito come una "gabbia dorata" sull'architettura x86: il software già in circolazione era troppo importante perché si potesse metterlo a rischio con modifiche architettoniche significative.

Qualunque siano i difetti architettonici dell'x86, si ricordi che lo sviluppo dei PC è largamente basato su questo insieme di istruzioni, che domina ancora

nei sistemi cloud dell'era post-PC. La produzione annuale di 350 milioni di chip x86 può sembrare piccola rispetto ai 14 miliardi di chip ARM, ma molte società sarebbero felici di avere il controllo di questo mercato. Tuttavia questo mosaico di modelli diversi ha prodotto un'architettura difficile da spiegare e impossibile da apprezzare.

Un consiglio al lettore su quanto segue: non cerchi di leggere il prossimo paragrafo con l'attenzione necessaria a scrivere programmi assembler per l'x86: lo scopo qui è mostrare la forza e la debolezza dell'architettura per calcolatori desktop più diffusa al mondo.

Invece di mostrare l'intero insieme di istruzioni a 16, 32 e 64 bit, in questo paragrafo ci concentreremo sul sottoinsieme delle istruzioni a 32 bit nato con l'80386. Inizieremo la trattazione con i registri e le modalità di indirizzamento, per passare poi alle operazioni sui numeri interi e concludere con l'analisi della codifica delle istruzioni.

## Registri e modalità di indirizzamento dell'x86

L'evoluzione dell'insieme di istruzioni si nota chiaramente osservando il formato dei registri dell'80386 (Figura 2.30): questo processore estende tutti i registri, precedentemente a 16 bit, a 32 bit (ad eccezione dei registri di segmento). La nuova versione dei registri viene indicata anteponendo al loro nome il prefisso *E*. I registri, nel loro complesso, saranno qui considerati a utilizzo generale (GPR). L'80386 ha solo 8 registri GPR; questo significa che un programma RISC-V e MIPS ne può utilizzare il quadruplo.

Le istruzioni aritmetiche, logiche e di trasferimento dati hanno due operandi e le combinazioni consentite sono mostrate in Figura 2.31. Ci sono due importanti differenze rispetto al RISC-V e al MIPS. La prima è che nell'x86 le istruzioni aritmetico-logiche hanno sempre un operando che funge sia da sorgente sia da destinazione, mentre il RISC-V e il MIPS consentono di definire registri distinti per operandi sorgente e operando destinazione. Questa restrizione rende più complesso il problema del numero limitato di registri, dato che il contenuto di uno dei registri sorgente viene necessariamente modificato. La seconda importante differenza consiste nel fatto che uno degli operandi può essere contenuto in memoria; perciò praticamente tutte le istruzioni possono prendere uno dei due operandi dalla memoria, a differenza del RISC-V e del MIPS.

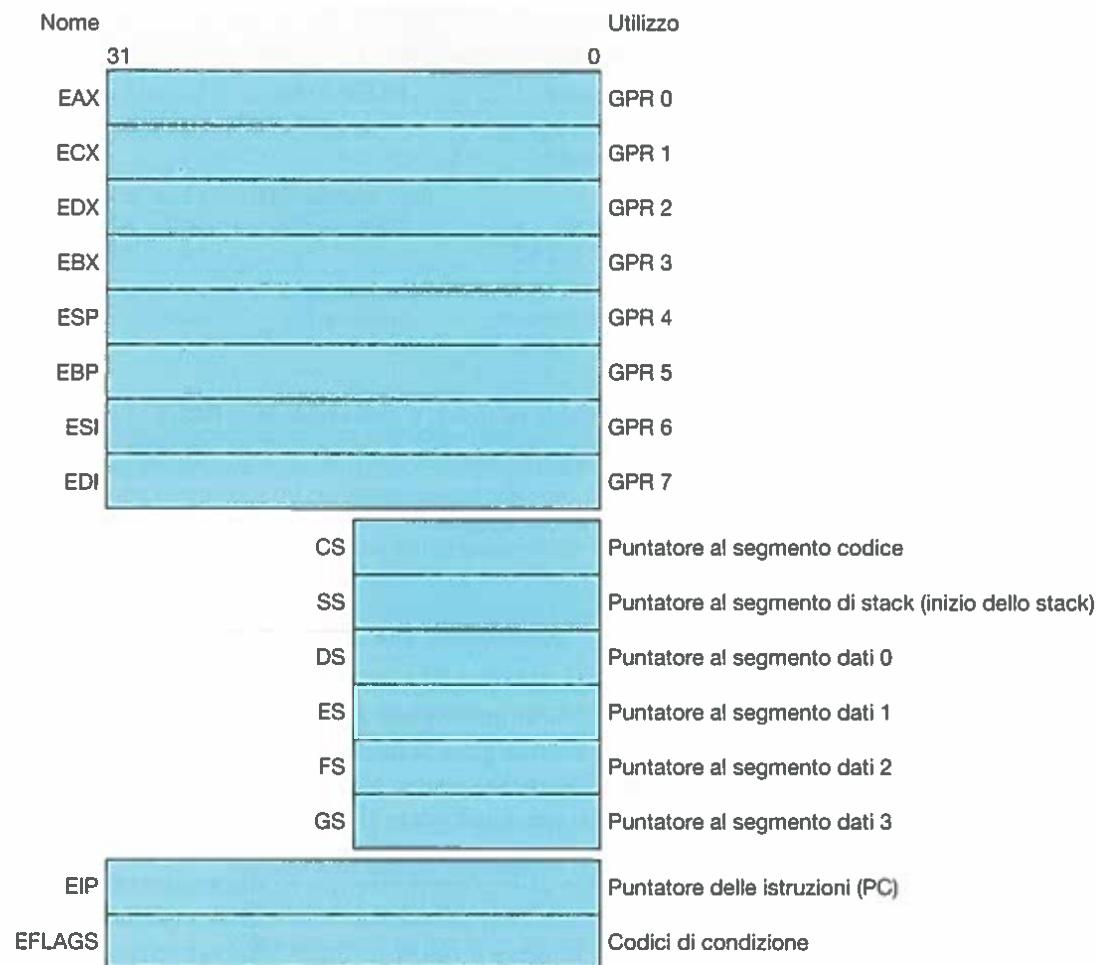
I modi di indirizzamento dei dati in memoria, che descriveremo in dettaglio fra poco, consentono di inserire nella stessa istruzione indirizzi di due dimensioni diverse: il cosiddetto *spiazzamento* può essere di 8 o 32 bit.

Nonostante per un operando in memoria si possa usare qualunque modo di indirizzamento, ci sono alcune restrizioni sui *registri* utilizzabili nei vari modi. La Figura 2.32 mostra le modalità di indirizzamento dell'x86 e i registri GPR che non possono essere utilizzati nelle diverse modalità. Sono riportate anche le istruzioni RISC-V che consentono di ottenere lo stesso risultato.

## Operazioni su numeri interi dell'x86

L'8086 consente di utilizzare sia dati a 8 bit (*byte*) che a 16 bit (*parola*); come abbiamo visto l'80386 aggiunge indirizzi e dati a 32 bit (*doppia parola*) all'architettura x86. L'AMD64 prevede anche la possibilità di utilizzare indirizzi e dati a 64 bit, chiamati *quad word* (*parola quadrupla*), ma in questo paragrafo ci limiteremo all'analisi dell'80386. La differenza tra i diversi tipi di dati si applica sia alle operazioni sui registri sia agli accessi alla memoria.

Quasi tutte le istruzioni possono operare sia su dati a 8 bit che su dati di dimensione più ampia. Questa viene definita nel campo modalità (*mode*) e può essere di 16 o 32 bit.



**Figura 2.30** Insieme di registri dell'80386. A partire dall'80386, gli 8 registri in alto sono stati estesi a 32 bit e possono essere impiegati come registri a utilizzo generale.

Tipo di operando sorgente/destinazione	Secondo operando sorgente
Registro	Registro
Registro	Immediato
Registro	Memoria
Memoria	Registro
Memoria	Immediato

**Figura 2.31** Combinazioni di operandi per le istruzioni aritmetiche, logiche e di trasferimento dati. L'x86 permette le combinazioni di operandi mostrate in figura: l'unica restrizione riguarda l'assenza di una modalità memoria-memoria. Gli operandi immediati possono avere una lunghezza di 8, 16 o 32 bit; il termine registro indica uno dei 14 registri principali di Figura 2.30 (tutti ad eccezione di EIP e EFLAGS).

Chiaramente alcuni programmi devono operare su dati di tutte e tre le dimensioni, quindi l'80386 fornisce un modo pratico per specificare la diversa dimensione dei dati senza estendere le dimensioni del codice in modo significativo. I progettisti pensarono che la maggior parte dei programmi avrebbe usato prevalentemente dati a 16 o 32 bit e che quindi aveva senso impostare come default una di queste due dimensioni. La dimensione di default viene specificata in uno dei bit del registro di segmento del codice; per sovrascriverla, occorre

Modalità	Descrizione	Restrizione sui registri	Codice RISC-V equivalente
Indiretta tramite registro	L'indirizzo è in un registro	No ESP o EBP	ld x10, 0(x11)
Base più spiazzamento a 8 o 32 bit	L'indirizzo è il contenuto del registro base più lo spiazzamento	No ESP	ld x10, 40(x11)
Base più indice scalato	L'indirizzo è Base + $(2^{\text{Scala}} \times \text{Indice})$ , dove Scala può valere 0, 1, 2 o 3	Base: qualsiasi GPR Indice: no ESP	slli x12, x12, 3 add x11, x11, x12 ld x10, 0(x11)
Base più indice scalato con spiazzamento a 8 o 32 bit	L'indirizzo è Base + $(2^{\text{Scala}} \times \text{Indice}) + \text{spiazzamento}$ , dove Scala può valere 0, 1, 2 o 3	Base: qualsiasi GPR Indice: no ESP	slli x12, x12, 3 add x11, x11, x12 ld x10, 40(x11)

**Figura 2.32 Modalità di indirizzamento a 32 bit dell'x86, con restrizioni sui registri e codice equivalente RISC-V.** La modalità di indirizzamento base più indice scalato, che non si trova né nel RISC-V né nel MIPS, è stata inserita per evitare la moltiplicazione per otto (costante di scorrimento pari a 3) necessaria a trasformare l'indice contenuto in un registro in un indirizzo in byte (Figure 2.26 e 2.28). Una costante di scorrimento pari a 1 viene utilizzata per i dati a 16 bit mentre una costante di scorrimento pari a 2 per quelli a 32 bit. Una costante di scorrimento pari a 0 significa che l'indirizzo non viene modificato. Se nella seconda o nella quarta modalità di indirizzamento lo spiazzamento fosse maggiore di 12 bit, allora il codice RISC-V equivalente richiederebbe più istruzioni: una lui per caricare i bit dello spiazzamento da 12 a 31 in un registro seguita da una add per sommare questi bit con l'indirizzo base. L'Intel dà due nomi differenti alla modalità di indirizzamento tramite registro base, cioè indirizzamento tramite base (Based) o indicizzato (Indexed) ma, dato che sono sostanzialmente identici, in questa figura sono riportati assieme.

precedere l'istruzione con un prefisso di 8 bit, che specifica al calcolatore di utilizzare l'altra dimensione grande dei dati (16 o 32 bit) per quella istruzione.

La soluzione di adottare un prefisso fu presa in prestito dall'8086, che prevedeva prefissi diversi per modificare il comportamento delle istruzioni. I tre prefissi originari servivano per modificare i registri di segmento predefiniti, per impedire l'accesso al bus permettendo la sincronizzazione dei processi (par. 2.11) o per ripetere l'istruzione successiva fino a quando il registro contatore, ECX, non raggiungeva il valore 0. Quest'ultimo prefisso era stato pensato per essere utilizzato con un'istruzione di move e aveva lo scopo di trasferire un numero variabile di byte. L'80386 prevede anche un altro prefisso da utilizzare per modificare la dimensione di default degli indirizzi.

Le operazioni sui numeri interi dell'x86 si possono suddividere in quattro classi principali:

1. le istruzioni per lo spostamento dei dati, che comprendono move, push e pop;
2. le istruzioni aritmetico-logiche, comprese quelle di confronto, e le operazioni aritmetiche su numeri interi e decimali;
3. le istruzioni che controllano il flusso del programma e che comprendono i salti condizionati e incondizionati, le chiamate a procedura e il ritorno da procedura;
4. le istruzioni per operazioni su stringhe, che comprendono lo spostamento e il confronto tra due stringhe.

Le prime due classi non meritano particolare attenzione, ad eccezione del fatto che le istruzioni aritmetico-logiche prevedono che la destinazione sia un registro o una locazione di memoria. La Figura 2.33 mostra alcune istruzioni tipiche dell'x86 e il loro funzionamento.

Nell'x86 i salti condizionati si basano sui bit del risultato dei test, detti *condition codes* o *flags*. Questi bit vengono posti a 1 come conseguenza di un'operazione: la maggior parte di essi viene utilizzata per confrontare un certo risultato con 0. I salti condizionati, quindi, effettuano un controllo su un determinato flag. L'indirizzo di salto relativo ai salti condizionati, come spiazzamento rispetto al PC, viene specificato in numero di byte, dato che le istruzioni di un 80386 non sono tutte lunghe 4 byte come nel RISC-V e nel MIPS.

Istruzione	Funzione
je nome	se uguale(codice di condizione) {EIP=nome}; EIP-128 ≤ nome < EIP+128
jmp nome	EIP=nome
call nome	SP=SP-4; M[SP]=EIP+5; EIP=nome;
movw EBX, [EDI+45]	EBX=M[EDI+45]
push ESI	SP=SP-4; M[SP]=ESI
pop EDI	EDI=M[SP]; SP=SP+4
add EAX, #6765	EAX= EAX+6765
test EDX, #42	Imposta i flag con il risultato dell'operazione con EDX e 42
movsl	M[EDI]=M[ESI]; EDI=EDI+4; ESI=ESI+4

**Figura 2.33** Alcune tipiche istruzioni dell'x86 con relativa funzione. Un elenco delle operazioni utilizzate più di frequente compare in Figura 2.34. L'istruzione CALL salva nello stack il puntatore (EIP) all'istruzione successiva (EIP è il PC dell'Intel).

Istruzione	Significato
<b>Controllo</b>	<b>Salti condizionati e incondizionati</b>
jnj, jz	Salto condizionato a EIP + offset su 8 bit; nomi alternativi sono JNE (per JNZ) e JE (per JZ)
jmp	Salto incondizionato; offset su 8 o 16 bit
call	Chiamata a procedura, con offset di 16 bit; l'indirizzo di ritorno viene inserito nello stack
ret	Prende dallo stack l'indirizzo di ritorno ed esegue il salto a quell'indirizzo
loop	Salto all'inizio di un ciclo: decrementa ECX e salta a EIP + offset su 8 bit, se ECX è diverso da 0
<b>Trasferimento dati</b>	<b>Spostamento dati fra registri o fra registri e memoria</b>
move	Sposta un dato da un registro a un altro o da registro a memoria
push, pop	Mette nello stack un operando, recupera dalla cima dello stack un operando e lo mette in un registro
les	Carica nel registro ES o in uno dei GPR il contenuto della memoria
<b>Operazioni aritmetico-logiche</b>	<b>Operazioni aritmetico-logiche su dati nei registri o in memoria</b>
add, sub	Somma il sorgente alla destinazione; sottrae il sorgente dalla destinazione; formato registro-memoria
cmp	Confronta il sorgente con la destinazione; formato registro-memoria
shl, shr, rcr	Scorrimento a sinistra, a destra e rotazione verso destra con impostazione del flag di riporto
cbw	Converte il byte che si trova negli 8 bit meno significativi di EAX in una parola che occupa i 16 bit meno significativi di EAX
test	Mette i valori opportuni nel registro dei flag facendo l'AND logico fra sorgente e destinazione
inc, dec	Incrementa, decrementa il registro destinazione
or, xor	OR logico, OR esclusivo; formato registro-memoria
<b>Stringhe</b>	<b>Trasferimento del contenuto di stringhe; lunghezza fornita nel prefisso di ripetizione</b>
movs	Copia una stringa sorgente in una stringa destinazione incrementando ESI ed EDI; può essere ripetuta
lod	Carica nel registro EAX un byte, una parola singola o doppia contenenti una stringa

**Figura 2.34** Alcune tipiche operazioni dell'x86. Molte operazioni utilizzano il formato registro-memoria, dove uno dei due operandi (sorgente o destinazione) si trova in memoria e l'altro può essere contenuto in un registro o essere un operando immediato.

Le istruzioni che operano su stringhe risalgono all'antenato 8080 e nella maggior parte dei programmi non vengono utilizzate; in generale, infatti, sono più lente delle procedure software equivalenti (vedi la sezione "Errori" riportata nel par. 2.19).

Alcune istruzioni su numeri interi dell'x86 sono riportate in Figura 2.34; molte di esse sono disponibili sia per dati in formato word sia in formato byte.

## Codifica delle istruzioni x86

La codifica delle istruzioni dell'80386 è assai complessa, essendoci per le istruzioni molti formati diversi che possono variare da 1 byte quando non ci sono operandi fino a 15 byte.

La Figura 2.35 mostra il formato di molte istruzioni campione presentate in Figura 2.33. Il byte del codice operativo solitamente contiene un bit che dice se l'operando è a 8 o 32 bit. Il codice operativo di alcune istruzioni può comprendere la modalità di indirizzamento e il registro; questo si verifica in molte istruzioni che hanno la forma "registro = registro operazione-operando immediato". Altre istruzioni utilizzano un "post-byte" o byte aggiuntivo per il codice operativo, chiamato "mod, reg, r/m", che contiene le informazioni relative alla modalità di indirizzamento; questo byte aggiuntivo è presente in molte istruzioni che accedono alla memoria. La modalità di indirizzamento "base più indice scalato" utilizza un secondo byte aggiuntivo, che può assumere i valori "sc, indice, base".

a. JE EIP + spiazzamento

	4	4	8
JE	Condizione	Spiazzamento	

b. CALL

	8	32
CALL		Spiazzamento

c. MOV EBX, [EDI + 45]

	6	1	1	8	8
MOV	d	w		r/m Post-byte	Spiazzamento

d. PUSH ESI

	5	3
PUSH		Reg

e. ADD EAX, #6765

	4	3	1	32
ADD	Reg	w		Operando immediato

f. TEST EDX, #42

	7	1	8	32
TEST	w	Post-byte		Operando immediato

**Figura 2.35** Format tipici delle istruzioni x86. La codifica del post-byte è mostrata in Figura 2.36. Molte istruzioni contengono il campo *w*, su un bit, che specifica se l'operazione agisce su un byte o su una parola doppia. Il campo *d* all'interno della MOV viene utilizzato nelle operazioni di trasferimento da e verso la memoria per specificare la direzione del trasferimento. L'istruzione ADD prevede 32 bit per il campo immediato, dato che nella modalità a 32 bit gli operandi immediati possono essere a 8 o 32 bit. Il campo immediato dell'istruzione TEST è lungo 32 bit, perché nella modalità a 32 bit non si possono utilizzare operandi immediati a 8 bit. Complessivamente la lunghezza di un'istruzione può variare da 1 a 15 byte. Le istruzioni più lunghe devono la loro dimensione alla presenza dei byte di prefisso, associati a un operando immediato su 4 byte e a un indirizzo con spiazzamento su 4 byte, di un codice operativo su 2 byte e di un byte che specifica quale base più indice scalato costituisca la modalità di indirizzamento utilizzata.

reg	w = 0	w = 1		r/m	mod = 0		mod = 1		mod = 2		mod = 3
		16b	32b		16b	32b	16b	32b	16b	32b	
0 AL	AX	EAX	0	ind=BX+SI	=EAX		stesso ind di mod=0 + disp8	stesso ind di mod=0 + disp8	stesso ind di mod=0 + disp16	stesso ind di mod=0 + disp32	come il campo reg
1 CL	CX	ECX	1	ind=BX+DI	=ECX						
2 DL	DX	EDX	2	ind=BP+SI	=EDX						
3 BL	BX	EBX	3	ind=BP+SI	=EBX						
4 AH	SP	ESP	4	ind=SI	= <i>(sib)</i>	SI+disp8	<i>(sib)</i> +disp8	SI+disp8	<i>(sib)</i> +disp32	"	
5 CH	BP	EBP	5	ind=DI	=disp32	DI+disp8	EBP+disp8	DI+disp16	EBP+disp32	"	
6 DH	SI	ESI	6	ind=disp16	=ESI	BP+disp8	ESI+disp8	BP+disp16	ESI+disp32	"	
7 BH	DI	EDI	7	ind=BX	=EDI	BX+disp8	EDI+disp8	BX+disp16	EDI+disp32	"	

**Figura 2.36** La codifica del primo campo utilizzato per la specifica degli indirizzi negli x86: "mod, reg, r/m". Le prime quattro colonne mostrano la codifica su 3 bit del campo *reg* (registro), il nome del registro dipende dal bit *w* del codice operativo e dal fatto che il calcolatore si trovi nella modalità a 16 bit (8086) o in quella a 32 bit (80386). Le restanti colonne illustrano i campi *mod* e *r/m*. Il significato del campo di 3 bit *r/m* dipende dal valore contenuto nel campo da 2 bit, *mod*, e dalle dimensioni dell'indirizzo. In pratica i registri utilizzati nel calcolo dell'indirizzo sono elencati nella sesta e settima colonna sotto *mod = 0*. Viene aggiunto uno spiazzamento (*disp*) su 8 bit se *mod = 1*, e uno spiazzamento su 16 bit o 32 bit, secondo la modalità di indirizzamento corrente, se *mod = 2*. Le seguenti tre combinazioni costituiscono una eccezione: 1) *r/m = 6* e *mod = 1* oppure *mod = 2*; in questa configurazione, se la CPU si trova nella modalità a 16 bit, viene selezionato il registro BP a cui viene aggiunto lo spiazzamento; 2) *r/m = 5* e *mod = 1* oppure *mod = 2*; in questa configurazione, se la CPU si trova nella modalità a 32 bit, viene selezionato il registro EBP a cui viene aggiunto lo spiazzamento; 3) *r/m = 4* e *mod = 3*; in questa configurazione, se la CPU si trova nella modalità a 32 bit, viene utilizzata la modalità di indirizzamento a indice scalato (*sib*) mostrata in Figura 2.35. Quando *mod = 3*, il campo *r/m* indica un registro, seguendo lo stesso tipo di codifica ottenuta combinando il campo *reg* con il bit *w*.

La Figura 2.36 mostra la codifica dei due post-byte per la specifica degli indirizzi, sia per la modalità a 16 bit sia per quella a 32 bit. Sfortunatamente, per capire a fondo quali registri e quali modalità di indirizzamento siano disponibili, occorre considerare la codifica di tutte le modalità di indirizzamento e a volte persino quella delle singole istruzioni.

### Conclusioni sull'x86

L'Intel ha costruito un microprocessore a 16 bit due anni prima di quelli, più eleganti, costruiti dai suoi concorrenti, come il Motorola 68000, e questo ha fatto sì che l'8086 fosse scelto come CPU per il PC IBM. In generale gli ingegneri dell'Intel riconoscono che è più difficile costruire un x86 piuttosto che un MIPS o un RISC-V, ma l'ampiezza del mercato delle loro CPU nell'era dei PC ha consentito a Intel e ad AMD di avere più risorse da investire per superare i problemi posti dalla crescente complessità. Ciò che all'x86 manca in stile è compensato dalla sua diffusione sul mercato, rendendolo apprezzabile se visto nella giusta prospettiva.

Inoltre, i componenti architettonici dell'x86 utilizzati più di frequente non sono troppo difficili da implementare, come Intel e AMD hanno dimostrato migliorando rapidamente le prestazioni dei programmi che lavorano sui numeri interi, sin dal 1978. Per ottenere queste prestazioni i compilatori devono evitare le parti dell'architettura difficili da implementare in modo veloce.

Ma nell'era post-PC, nonostante la grande esperienza nell'architettura e nella sua produzione, l'x86 non è ancora diventato competitivo nel campo dei dispositivi personali mobili (PMD).

## 2.18 | Un caso reale: le altre istruzioni dell'architettura RISC-V

Con l'obiettivo di realizzare un'architettura di un insieme di istruzioni adatto a una grande varietà di calcolatori, i progettisti del RISC-V hanno suddiviso l'in-

Istruzione	Nome	Formato	Descrizione
Somma i bit più significativi al PC	auipc	U	Somma i 20 bit più significativi al PC; scrivi la somma in registro
Imposta a 1 se minore	slt	R	Confronta 2 registri; scrivi il risultato come booleano in registro
Imposta a 1 se minore, senza segno	sltu	R	Confronta 2 registri; scrivi il risultato come booleano in registro
Imposta a 1 se minore, immediato	slti	I	Confronta 1 registro e 1 immediato; scrivi il risultato come booleano in registro
Imposta a 1 se minore, immediato, senza segno	sltiu	I	Confronta 1 registro e 1 immediato; scrivi il risultato come booleano in registro
Somma parole	addw	R	Somma 2 numeri su 32 bit
Sottrai parole	subw	R	Sottrai 2 numeri su 32 bit
Somma parola e immediato	addiw	I	Somma 1 costante e 1 numero su 32 bit
Scorrimento logico di una parola a sinistra	sllw	R	Scorrimento di un numero a 32 bit a sinistra secondo un registro
Scorrimento logico di una parola a destra	srlw	R	Scorrimento di un numero a 32 bit a destra secondo un registro
Scorrimento aritmetico di una parola a destra	sraw	R	Scorrimento aritmetico di un numero a 32 bit a destra
Scorrimento logico di una parola a sinistra, immediato	slliw	I	Scorrimento di un numero a 32 bit a sinistra secondo un immediato
Scorrimento logico di una parola a destra, immediato	srliw	I	Scorrimento di un numero a 32 bit a destra secondo un immediato
Scorrimento aritmetico di una parola a destra immediato	sraiw	I	Scorrimento aritmetico di un numero a 32 bit a destra, secondo un immediato

Figura 2.37 Le restanti 14 istruzioni dell'architettura dell'insieme delle istruzioni di base del RISC-V.

sieme delle istruzioni in un'architettura di base e diverse estensioni. Ciascuna di esse è denominata con le lettere dell'alfabeto e l'architettura di base è chiamata I che sta per interi. L'architettura di base contiene poche istruzioni rispetto ad altri insiemi di istruzioni popolari oggigiorno, e questo capitolo le ha già introdotte quasi tutte. Questo paragrafo riassume le istruzioni dell'architettura base e descrive le cinque estensioni standard.

La Figura 2.37 riporta le istruzioni non ancora trattate dell'architettura di base RISC-V. La prima istruzione, auipc, viene utilizzata per l'indirizzamento della memoria relativo al PC. Come l'istruzione lui, contiene una costante su 20 bit nei bit da 12 a 31 di un intero. L'effetto di auipc è quello di aggiungere questo numero al PC e scrivere la somma in un registro. Combinata con un'altra istruzione come addi, è possibile indirizzare un qualsiasi byte della memoria entro 4 GiB dall'indirizzo corrente nel PC. Questa caratteristica è utile per il codice indipendente dalla posizione, che viene così eseguito correttamente indipendentemente dalla posizione della memoria in cui viene caricato. Viene utilizzata soprattutto nelle librerie a collegamento dinamico.

Le quattro istruzioni successive confrontano due interi e scrivono il risultato come variabile booleana in un registro. slt e sltu confrontano il contenuto di due registri come numeri con e senza segno, rispettivamente, e scrivono 1 in un registro se il primo valore è minore del secondo, o 0 nel caso opposto. slti e sltiu eseguono lo stesso confronto, ma con un valore immediato come secondo operando.

Le istruzioni restanti dovrebbero suonare familiari, dato che il loro nome è simile a quello di altre istruzioni introdotte in questo capitolo, ma con la lettera w che sta per parola (word), aggiunta in coda. Queste istruzioni eseguono le stesse

Nome simbolico	Descrizione	Numero di istruzioni
I	Architettura di base	51
M	Moltiplicazione/divisione intera	13
A	Operazioni atomiche	22
F	Virgola mobile in precisione singola	30
D	Virgola mobile in doppia precisione	32
C	Istruzioni compresse	36

**Figura 2.38** L'architettura dell'insieme delle istruzioni RISC-V è suddivisa in ISA di base, denominata I, e in cinque estensioni M, A, F, D e C.

operazioni di quelle con lo stesso nome senza la lettera *w* delle quali abbiamo già parlato in questo capitolo tranne che lavorano sui 32 bit meno significativi degli operandi e ignorano i bit dal 32 al 63. Inoltre, producono un risultato con il segno esteso: viene cioè ricoppiato il bit 31 nei bit da 32 a 63. I progettisti del RISC-V hanno incluso queste istruzioni con suffisso *w* perché le operazioni su numeri a 32 bit rimangono molto comuni nei calcolatori con indirizzamento a 64 bit. La ragione principale è che il diffuso tipo di dati `int` rimane a 32 bit in Java e nella maggior parte delle implementazioni del linguaggio C.

E questo è tutto per l'architettura di base! La Figura 2.38 riporta le cinque estensioni standard. La prima, M, aggiunge le istruzioni per moltiplicare e dividere i numeri interi. Il Capitolo 3 introdurrà diverse istruzioni dell'estensione M.

La seconda estensione, A, supporta le operazioni atomiche sulla memoria per la sincronizzazione dei multiprocessori. Le istruzioni di load riservata (*load reserved*, `lr.r.d`) e store condizionata (*store conditioned*, `sc.d`), introdotte nel paragrafo 2.11, fanno parte dell'estensione A. Sono anche comprese le versioni che operano su dati a 32 bit (`lr.w` e `sc.w`). Le restanti 18 istruzioni sono ottimizzazioni di schemi tipici di sincronizzazione, quali scambio atomico e somma atomica, ma non aggiungono funzionalità alla load riservata e alla store condizionata.

La terza e quarta estensione, F e D, forniscono operazioni su numeri in virgola mobile, che sono descritte nel Capitolo 3.

L'ultima estensione, C, non fornisce nuove funzionalità, ma prende le istruzioni più comuni del RISC-V, come la addi, e fornisce le istruzioni equivalenti che sono lunghe solo 16 bit, invece di 32. Consente così di esprimere i programmi con meno byte, cosa che può ridurre i costi e, come vedremo nel Capitolo 5, può migliorare le prestazioni. Per potere essere contenute in 16 bit, queste nuove istruzioni hanno alcune restrizioni sugli operandi: per esempio, alcune istruzioni possono accedere solamente ad alcuni dei 32 registri, e i campi immediati sono più stretti.

Considerando tutte le istruzioni, il RISC-V di base e le sue estensioni hanno 184 istruzioni, più 13 istruzioni di sistema che verranno introdotte alla fine del Capitolo 5.

## 2.19 Errori e trabocchetti

*Errore: più le istruzioni sono potenti, migliori sono le prestazioni.*

Parte della potenza dell'architettura Intel x86 consiste nell'utilizzo dei prefissi che possono modificare l'esecuzione dell'istruzione. Per esempio, esiste un prefisso che causa la ripetizione dell'istruzione fino a che un contatore non

raggiunge lo 0; quindi per eseguire trasferimenti sequenziali di dati da 32 bit da una porzione di memoria a un'altra, può sembrare che la soluzione migliore sia un'istruzione di move, preceduta dal prefisso di ripetizione.

Un metodo alternativo, che utilizza le istruzioni standard presenti in tutti i calcolatori, consiste nel caricare i dati nei registri (load) per poi memorizzare il contenuto dei registri in memoria (store). Questa seconda versione del programma, con il codice ripetuto in modo da evitare il costo addizionale dei cicli, è 1,5 più veloce. Una terza versione, che utilizza i registri in virgola mobile invece dei registri interi dell'x86, trasferisce i dati 2 volte più velocemente dell'istruzione complessa.

*Errore: programmare in linguaggio assembler permette di ottenere le prestazioni migliori.*

Una volta i compilatori producevano sequenze di istruzioni poco efficienti, ma il loro crescente miglioramento rende la distanza fra il codice generato dal compilatore e quello scritto a mano sempre minore. In effetti, per essere in grado di competere con gli attuali compilatori, un programmatore assembler dovrebbe avere perfettamente chiari i concetti descritti nei Capitoli 4 e 5 (pipeline del processore e gerarchia delle memorie).

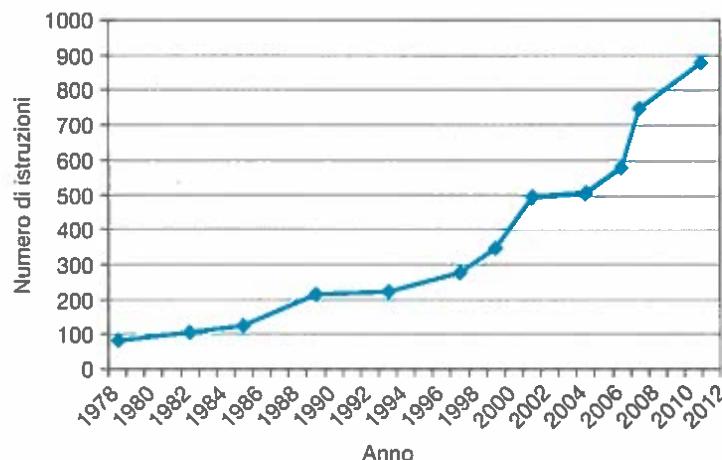
La battaglia fra compilatori e programmatore assembler è in una fase in cui l'uomo perde terreno ogni giorno. Per esempio, il C offre al programmatore la possibilità di suggerire al compilatore quali variabili mantenere nei registri e quali copiare in memoria. Quando i compilatori non erano molto efficienti nell'allocazione dei registri, questi suggerimenti avevano un'importanza vitale, e alcuni vecchi libri di testo sul C dedicavano un buon numero di pagine a esempi sul modo di dare suggerimenti efficaci al compilatore. Attualmente i compilatori C di solito ignorano questi suggerimenti, dato che sono in grado di allocare i registri meglio di quanto sappia fare il programmatore.

Anche se scrivere direttamente in assembler consentisse di ottenere un codice più efficiente, gli svantaggi sarebbero la quantità maggiore di tempo impiegata nella codifica e nella ricerca degli errori, la mancanza di portabilità e la difficoltà di manutenzione del codice. Uno dei pochi assiomi dell'ingegneria del software ampiamente accettati è che il tempo di codifica è proporzionale al numero di linee di codice e, ovviamente, un programma assembler contiene molte più linee di codice di un programma C o Java. Una volta terminato, inoltre, un programma può diventare assai diffuso: la vita delle applicazioni è sempre più lunga del previsto e ciò significa che qualcun altro dovrà necessariamente aggiornare il codice a distanza di anni, facendolo funzionare su nuove versioni del sistema operativo o su nuovi modelli di calcolatore. Scrivere in linguaggio ad alto livello non solo permette ai compilatori futuri di ottimizzare il codice per i calcolatori che man mano entrano in commercio, ma rende anche più semplice la manutenzione del software e garantisce una maggiore portabilità su architetture diverse.

*Errore: l'importanza commerciale della compatibilità binaria implica che gli insiemi di istruzioni che hanno successo non cambino.*

Anche se la compatibilità dei nuovi sistemi rispetto ai precedenti è sacrosanta, la Figura 2.39 mostra come l'architettura dell'x86 sia cresciuta enormemente, con un tasso di crescita medio di un'istruzione al mese per 35 anni!

*Trabocchetto: dimenticare che gli indirizzi di parole consecutive nei calcolatori con indirizzamento al byte non differiscono di una unità.*



**Figura 2.39** Crescita negli anni dell'insieme di istruzioni dell'x86. Sebbene alcune estensioni dell'insieme di istruzioni abbiano un evidente contenuto tecnico, il rapido incremento dell'insieme di istruzioni rende difficile per altre società cercare di costruire processori compatibili.

Molti programmatore assembler hanno faticato a lungo per trovare errori commessi supponendo che l'indirizzo della parola successiva fosse ottenibile incrementando di una unità l'indirizzo contenuto in un registro, invece che della dimensione della parola o della parola doppia in byte. Uomo avvisato, mezzo salvato!

**Trabocchetto:** utilizzare un puntatore a una variabile locale al di fuori della procedura in cui la variabile è definita.

Un errore comune, quando si ha a che fare con i puntatori, consiste nel passare come risultato di una procedura un puntatore a un vettore definito localmente nella procedura stessa. Secondo le regole di utilizzo dello stack descritte in Figura 2.12, la parte di memoria che contiene il vettore locale verrà riutilizzata non appena la procedura sarà terminata e l'esecuzione tornerà alla procedura chiamante. I puntatori a variabili locali possono generare il caos.

## 2.20 | Note conclusive

I calcolatori a *programma memorizzato* sono basati su due assunti fondamentali: le istruzioni sono indistinguibili dai numeri e la memoria in cui sono scritti i programmi può essere modificata. Questi principi fanno sì che lo stesso calcolatore possa essere utile a scienziati che si occupano di ambiente, a operatori finanziari e a romanzi. La scelta dell'insieme di istruzioni che il calcolatore è in grado di eseguire è il risultato di un delicato compromesso fra il numero di istruzioni necessarie per eseguire un programma, il numero di cicli di clock per istruzione e il periodo di clock. Come è stato illustrato in questo capitolo, sono tre i criteri che guidano i progettisti degli insiemi di istruzioni nel delicato compito di trovare il bilanciamento giusto: vediamoli.

*Less is more* (Meno è meglio).  
Robert Browning, *Andrea del Sarto*, 1855

1. *La semplicità favorisce la regolarità.* La ricerca della regolarità è alla base di molte caratteristiche dell'insieme di istruzioni RISC-V: stessa dimensione per tutte le istruzioni, richiesta di operandi di tipo registro in tutte le istruzioni aritmetiche, mantenimento dei campi registro nella stessa posizione all'interno di tutti i formati delle istruzioni.

2. *Minori sono le dimensioni, maggiore è la velocità.* La velocità di esecuzione è il motivo per cui il RISC-V ha 32 registri anziché molti di più.
3. *Un buon progetto richiede buoni compromessi.* Un esempio relativo al RISC-V è il compromesso raggiunto tra fornire indirizzi e costanti di dimensioni maggiori e mantenere la stessa lunghezza per tutte le istruzioni.



Abbiamo anche visto all'opera la grande idea introdotta nel primo capitolo di rendere veloci le situazioni più comuni applicata agli insiemi di istruzioni e all'architettura del calcolatore. Esempi di questo tipo che riguardano il RISC-V sono la modalità di indirizzamento relativo al PC per i salti condizionati e l'indirizzamento immediato per gli operandi costanti di dimensioni più ampie.

Al di sopra del linguaggio macchina c'è il linguaggio assembler, leggibile anche dagli uomini, che viene tradotto in numeri binari comprensibili al calcolatore dall'assemblatore. Questo può addirittura "estendere" l'insieme di istruzioni creando istruzioni simboliche non implementate a livello hardware. Per esempio, le costanti o gli indirizzi troppo grandi vengono suddivisi in frammenti di dimensione adeguata; le varianti più comuni di determinate istruzioni prendono un nome che li contraddistingue e così via. Le istruzioni RISC-V viste finora, sia quelle vere e proprie sia le pseudoistruzioni, sono riportate in Figura 2.40. Nascondere i dettagli delle istruzioni ai livelli superiori è un altro esempio della grande idea dell'astrazione.

I diversi tipi di istruzioni RISC-V sono associati a costrutti contenuti nei linguaggi di programmazione ad alto livello:

- le istruzioni aritmetiche corrispondono alle operazioni che si possono trovare negli assegnamenti;
- le istruzioni di trasferimento dati sono più frequenti quando si utilizzano strutture dati, come i vettori o le strutture;
- le istruzioni di salto condizionato vengono utilizzate nei costrutti *if* e nei cicli;
- le istruzioni di salto incondizionato vengono utilizzate nelle chiamate a procedura, nel ritorno da procedura e nei costrutti *case/switch*.

Non tutte le istruzioni sono uguali: alcune sono utilizzate con più frequenza rispetto ad altre. Per esempio, la Figura 2.41 mostra la frequenza di utilizzo di ciascun tipo di istruzione nei programmi dello SPEC2000. Questa differenza gioca un ruolo importante nella progettazione, come verrà discusso nei capitoli che riguardano unità di elaborazione, unità di controllo e pipeline.

Dopo avere illustrato l'aritmetica dei calcolatori, a cui è dedicato il Capitolo 3, presenteremo le altre istruzioni dell'architettura dell'insieme di istruzioni RISC-V.

## 2.21 Inquadramento storico e approfondimenti

Questo paragrafo contiene una rassegna dell'evoluzione delle *architetture degli insiemi di istruzioni* (ISA), assieme a una breve storia dei linguaggi di programmazione e dei compilatori.

Le architetture di insiemi di istruzioni comprendono architetture ad accumulatore, architetture a registri a utilizzo generale, architetture a stack e una breve storia dell'architettura x86 e ARM a 32 bit. Verrà anche esaminata la tematica controversa del confronto tra architetture di calcolatori basate su linguaggi ad alto livello e architetture di calcolatori basate su un insieme ridotto di istruzioni. La rassegna dei linguaggi di programmazione tratterà Fortran, Lisp, Algol, C, Cobol, Pascal, Simula, Smalltalk, C++ e Java, men-

tre quella dei compilatori comprenderà le pietre miliari del loro sviluppo e i pionieri che hanno permesso di raggiungerle. Il testo del paragrafo 2.21 è disponibile online .

Istruzione RISC-V	Nome	Formato	Pseudoistruzione RISC-V	Nome	Istruzione reale
Somma	add	R	Sposta	mv	addi
Sottrazione	sub	R	Carica immediato	li	addi
Add immediato	addi	I	Salta	j	Jal
Lettura parola doppia	ld	I	Carica indirizzo	la	lui + addi
Memorizzazione parola doppia	sd	S			
Lettura parola	lw	I			
Lettura parola senza segno	lwu	I			
Memorizzazione parola	sw	S			
Lettura mezza parola	lh	I			
Lettura mezza parola senza segno	lhu	I			
Memorizzazione mezza parola	sh	S			
Lettura byte	lb	I			
Lettura byte senza segno	lbu	I			
Memorizzazione byte	sb	S			
Lettura riservata	lr.d	R			
Memorizzazione condizionata	sc.d	R			
Lettura dei bit superiori	lui	U			
And	and	R			
Or inclusivo	or	R			
Or esclusivo	xor	R			
And immediato	andi	I			
Or inclusivo immediato	ori	I			
Or esclusivo immediato	xori	I			
Scorrimento logico sinistra	sll	R			
Scorrimento logico destra	srl	R			
Scorrimento aritmetico destra	sra	R			
Scorrimento logico sinistra immediato	slli	I			
Scorrimento logico destra immediato	srli	I			
Scorrimento aritmetico destra immediato	srai	I			
Salta se uguale	beq	SB			
Salta se non uguale	bne	SB			
Salta se minore	blt	SB			
Salta se maggiore o uguale	bge	SB			
Salta se minore senza segno	bltu	SB			
Salta se maggiore o uguale senza segno	bgeu	SB			
Jump and link	jal	UJ			
Jump and link a registro	jalr	I			

**Figura 2.40** Insieme delle istruzioni RISC-V esaminate finora; le istruzioni vere e proprie sono riportate a sinistra mentre a destra sono riportate le pseudoistruzioni. La Figura 2.1 mostra ulteriori dettagli sulla parte di architettura RISC-V presentata in questo capitolo. Si possono trovare queste informazioni anche nella prima e seconda colonna della scheda tecnica riassuntiva del RISC-V in fondo al libro.

Tipi di istruzioni	Esempi RISC-V	Corrispondenza con linguaggi ad alto livello	Frequenza	
			Interi	Virgola mobile
Aritmetiche	add, sub, addi	Operazioni in istruzioni di assegnamento	16%	48%
Trasferimento dati	ld, sd, lw, sw, lh, sh, lb, sb, lui	Riferimento a strutture dati, per esempio a vettori	35%	36%
Logiche	and, or, xor, sll, srl, sra	Operazioni in istruzioni di assegnamento	12%	4%
Salti condizionati	beq, bne, blt, bge, bltu, bgeu	Costrutti if e cicli	34%	8%
Salti incondizionati	jal, jalr	Chiamate a procedura, ritorno da procedura e costrutti switch	2%	0%

**Figura 2.41** Tipi di istruzioni RISC-V, esempi, corrispondenza con i costrutti dei linguaggi di programmazione ad alto livello e percentuale di istruzioni RISC-V di ognuno dei diversi tipi, eseguite in media nei programmi di benchmark dello SPEC CPU2006. In Figura 3.24 del Capitolo 3 verranno mostrate le percentuali relative alle singole istruzioni RISC-V eseguite.

## 2.22 | Esercizi

**2.1** [5] < 2.2> Qual è il codice assembler RISC-V corrispondente alla seguente istruzione C? Si supponga che le variabili f, g e h siano già state memorizzate nei registri x5, x6 e x7 rispettivamente. Si utilizzi il numero minimo di istruzioni assembler.

$$f = g + (h - 5);$$

**2.2** [5] < 2.2> Quali sono le istruzioni C che corrispondono alle seguenti due istruzioni in assembler RISC-V?

```
add f, g, h
add f, i, f
```

**2.3** [5] < 2.2, 2.3> Qual è il codice assembler RISC-V corrispondente alla seguente istruzione C? Si supponga che le variabili f, g, h, i e j siano assegnate rispettivamente ai registri x5, x6, x7, x28 e x29. Si assuma che l'indirizzo base dei vettori A e B sia contenuto nei registri x10 e x11 rispettivamente nei registri x10 e x11.

$$B[8] = A[i - j];$$

**2.4** [10] < 2.2, 2.3> Qual è il codice assembler RISC-V corrispondente al seguente spezzone di codice C? Si supponga che le variabili f, g, h, i e j siano assegnate rispettivamente ai registri x5, x6, x7, x28 e x29. Si assuma che l'indirizzo base dei vettori A e B sia contenuto nei registri x10 e x11 rispettivamente nei registri x10 e x11.

```
slli x30, x5, 3      // x30 = f * 8
add x30, x10, x30   // x30 = &A[f]
slli x31, x6, 3      // x31 = g * 8
add x31, x11, x31   // x31 = &B[g]
ld  x5, 0(x30)      // f = A[f]
```

```
addi x12, x30, 8
ld  x30, 0(x12)
add x30, x30, x5
sd  x30, 0(x31)
```

**2.5** [5] < 2.3> Mostrare come il numero esadecimale 0abcdef12 viene disposto in memoria in una macchina con codifica little-endian e in una con codifica big-endian. Si supponga che i dati vengano memorizzati a partire dall'indirizzo 0 e che la dimensione della parola sia di 4 byte.

**2.6** [5] < 2.4> Tradurre 0abcdef12 in decimale.

**2.7** [5] < 2.2, 2.3> Tradurre la seguente istruzione C in assembler RISC-V. Si supponga che le variabili f, g, h, i e j siano assegnate rispettivamente ai registri x5, x6, x7, x28 e x29. Si assuma che l'indirizzo base dei vettori A e B sia contenuto nei registri x10 e x11 rispettivamente. Si assuma che gli elementi dei vettori A e B siano parole su 8 byte:

$$B[8] = A[i] + A[j];$$

**2.8** [10] < 2.2, 2.3> Tradurre il seguente spezzone di codice assembler RISC-V in C. Si supponga che le variabili f, g, h, i e j siano assegnate rispettivamente ai registri x5, x6, x7, x28 e x29. Si assuma che l'indirizzo base dei vettori A e B sia contenuto nei registri x10 e x11 rispettivamente.

```
addi x30, x10, 8
addi x31, x10, 0
sd  x31, 0(x30)
ld  x30, 0(x30)
add x5, x30, x31
```

**2.9** [20] < 2.2, 2.5> Per ciascuna istruzione RISC-V dell'Esercizio 2.8, definire il valore del codice operativo (OP – opcode) e assegnare un valore ai campi registro sorgente (rs1) e registro destinazione (rd). Per le istruzioni di tipo I, assegnare un valore al contenuto del campo immediato e, per le istruzioni di tipo R, al contenuto del campo del registro del secondo operando (rs2). Per le istruzioni non di tipo U o UJ, mostrare il contenuto del campo funz3, e per le istruzioni di tipo R e S, mostrare il contenuto anche del campo funz7.

**2.10** Supponiamo che i registri x5 e x6 contengano i numeri 0x8000000000000000 e 0xD000000000000000 rispettivamente.

**2.10.1** [5] < 2.4> Determinare quale sarà il contenuto di x30 dopo l'esecuzione di questa istruzione assembler:

```
add x30, x5, x6
```

**2.10.2** [5] < 2.4> Il contenuto di x30 è corretto, o si è verificato un overflow?

**2.10.3** [5] < 2.4> Si supponga che i registri x5 e x6 contengano i valori riportati sopra. Determinare quale sarà il contenuto di x30 dopo l'esecuzione di questa istruzione assembler:

```
sub x30, x5, x6
```

**2.10.4** [5] < 2.4> Il contenuto di x30 è corretto o si è verificato un overflow?

**2.10.5** [5] < 2.4> Si supponga che i registri x5 e x6 contengano i valori riportati sopra. Determinare quale sarà il contenuto di x30 dopo l'esecuzione di queste due istruzioni assembler:

```
add x30, x5, x6
add x30, x30, x5
```

**2.10.6** [5] < 2.4> Il contenuto di x30 è corretto o si è verificato un overflow?

**2.11** Si supponga che x5 contenga il valore 128<sub>dec</sub>.

**2.11.1** [5] < 2.4> Determinare qual è l'intervallo dei valori di x6 che produrrebbero un overflow quando viene eseguita l'istruzione add x30, x5, x6.

**2.11.2** [5] < 2.4> Determinare qual è l'intervallo dei valori di x6 che produrrebbero un overflow quando viene eseguita l'istruzione sub x30, x5, x6.

**2.11.3** [5] < 2.4> Determinare qual è l'intervallo dei valori di x6 che produrrebbero un overflow quando viene eseguita l'istruzione sub x30, x6, x5.

**2.12** [5] < 2.2, 2.5> Determinare quale istruzione assembler corrisponde alla seguente stringa binaria:

0000 0000 0001 0000 1000 0000 10111 0011<sub>due</sub>

Suggerimento: la Figura 2.20 può essere utile.

**2.13** [5] < 2.2, 2.5> Definire il tipo di istruzione e la rappresentazione esadecimale della seguente istruzione assembler:

```
sd x5, 32(x30)
```

**2.14** [5] < 2.5> Definire il tipo di istruzione, l'istruzione in linguaggio assembler e la rappresentazione binaria dell'istruzione descritta dai seguenti campi RISC-V:

```
codop = 0x33, funz3 = 0x0, funz7 = 0x20, rs2 = 5, rs1 = 7, rd = 6
```

**2.15** [5] < 2.5> Definire il tipo di istruzione, l'istruzione in linguaggio assembler e la rappresentazione binaria dell'istruzione descritta dai seguenti campi RISC-V:

```
codop = 0x3, funz3 = 0x3, rs1 = 27, rd = 3, imm = 0x4
```

**2.16** Si supponga di voler espandere il register file del RISC-V a 128 registri ed espandere l'insieme delle istruzioni in modo che contenga quattro volte il numero di istruzioni.

**2.16.1** [5] < 2.5> Quali saranno le dimensioni risultanti, in numero di bit, dei diversi campi di un'istruzione di tipo R?

**2.16.2** [5] < 2.5> Quali saranno le dimensioni risultanti, in numero di bit, dei diversi campi di un'istruzione di tipo I?

**2.16.3** [5] < 2.5, 2.8, 2.10> In che modo le due modifiche proposte sopra potrebbero consentire di diminuire le dimensioni dei programmi assembler RISC-V? In che modo le modifiche suggerite sopra potrebbero fare aumentare le dimensioni dei programmi assembler RISC-V?

**2.17** Si supponga che i registri seguenti contengano i valori:

```
x5 = 0x00000000AAAAAAA, x6 = 0x1234567812345678
```

**2.17.1** [5] < 2.6> Supponendo che i registri x5 e x6 contengano i valori riportati sopra, determinare il contenuto di x7 dopo l'esecuzione delle seguenti istruzioni:

```
slli x7, x5, 4
or x7, x7, x6
```

**2.17.2** [5] < 2.6> Supponendo che i registri x5 e x6 contengano i valori riportati sopra, determinare il contenuto di x7 dopo l'esecuzione delle seguenti istruzioni:

```
slli x7, x6, 4
```

**2.17.3 [5] < 2.6>** Supponendo che i registri x5 e x6 contengano i valori riportati sopra, determinare il contenuto di x7 dopo l'esecuzione delle seguenti istruzioni:

```
srl x7, x5, 3
andi x7, x7, 0xFF
```

**2.18 [10] < 2.6>** Determinare la sequenza più corta di istruzioni RISC-V che consente di estrarre i bit da 11 a 16 dal registro x5 e li sostituisce ai bit da 26 a 31 del registro x6 senza modificare gli altri bit del registro x5 e x6. (Assicuratevi di verificare il vostro codice utilizzando  $x5 = 0$  e  $x6 = 0xffffffffffffffffff$ . Questo potrebbe rivelare una svista comune.)

**2.19 [5] < 2.6>** Determinare la sequenza più corta di istruzioni RISC-V che si possono utilizzare per implementare la pseudoistruzione seguente:

```
not x5, x6 // inversione bit a bit
```

**2.20 [5] < 2.6>** Determinare la sequenza più corta di istruzioni RISC-V che si possono utilizzare per tradurre la seguente istruzione C. Si supponga  $x6 = A$  e che  $x17$  contenga l'indirizzo di base di C.

```
A = C[0] << 4;
```

**2.21 [5] < 2.7>** Si supponga che x5 contenga il valore  $0x000000001010000$ . Determinare il contenuto di x6 dopo l'esecuzione delle seguenti istruzioni:

```
bge x5, x0, ELSE
jal x0, FINE
ELSE: ori x6, x0, 2
FINE:
```

**2.22 Si supponga che il program counter (PC) sia impostato a 0x20000000.**

**2.22.1 [5] < 2.10>** Qual è l'intervallo degli indirizzi a cui si può saltare utilizzando l'istruzione RISC-V di *jump-and-link* (jal)? (In altre parole, qual è l'insieme dei possibili valori che può assumere il PC dopo l'esecuzione dell'istruzione jal?).

**2.22.2 [5] < 2.10>** Qual è l'intervallo degli indirizzi a cui si può saltare utilizzando l'istruzione RISC-V di *branch se uguale* (beq)? (In altre parole, qual è l'insieme dei possibili valori che può assumere il PC dopo l'esecuzione dell'istruzione beq?).

**2.23 Si consideri la proposta di una nuova istruzione chiamata rpt.** Questa istruzione combina in una sola istruzione il controllo di una condizione di fine ciclo

e il decremento dell'indice di ciclo. Per esempio rpt x29, ciclo avrebbe questo effetto:

```
if (x29 > 0) {
    x29 = x29 - 1;
    goto ciclo
}
```

**2.23.1 [5] < 2.7, 2.10>** Quale sarebbe il formato di istruzione più appropriato per implementare l'istruzione riportata sopra nell'insieme di istruzioni RISC-V?

**2.23.2 [5] < 2.7>** Determinare la sequenza minima di istruzioni RISC-V che consenta di implementare l'istruzione riportata sopra.

**2.24 Si consideri il seguente ciclo in assembler RISC-V:**

```
CICLO: beq x6, x0, FINE
        addi x6, x6, -1
        addi x5, x5, 2
        jal x0, CICLO
FINE:
```

**2.24.1 [5] < 2.7>** Si supponga che il registro x6 venga inizializzato al valore 10. Quale sarà il contenuto finale di x5 supponendo che x5 venga inizializzato a 0?

**2.24.2 [5] < 2.7>** Scrivere la procedura C equivalente al ciclo assembler riportato sopra. Si supponga che i registri x5 e x6 contengano rispettivamente le variabili intere acc e i rispettivamente.

**2.24.3 [5] < 2.7>** Si supponga che il registro x6 sia inizializzato con il valore numerico N. Quante istruzioni RISC-V verranno eseguite dal ciclo scritto in linguaggio assembler RISC-V e riportato sopra?

**2.24.4 [5] < 2.7>** Sostituire l'istruzione "beq x6, x0, FINE" con l'istruzione "blt x6, x0, FINE" nel ciclo scritto in linguaggio assembler RISC-V e riportato sopra. Scrivere l'equivalente codice C.

**2.25 [10] < 2.7>** Tradurre il seguente frammento di codice C in codice assembler RISC-V. Si utilizzi il minor numero possibile di istruzioni. Si supponga che le variabili a, b, i e j siano contenute rispettivamente nei registri x5, x6, x7 e x29. Si supponga anche che il registro x10 contenga l'indirizzo base del vettore D.

```
for (i=0; i<a; i++)
    for (j=0; j<b; j++)
        D[4*j] = i+j;
```

**2.26 [5] < 2.7>** Quante istruzioni RISC-V sono necessarie per implementare il frammento di codice C

dell'Esercizio 2.25? Supponendo che le variabili *a* e *b* vengano inizializzate a 10 e 1, e che tutti gli elementi di *D* contengano 0 all'inizio del ciclo, quante istruzioni RISC-V verranno eseguite per completare il ciclo?

**2.27** [5] < 2.7> Tradurre il seguente ciclo in C. Si supponga che il registro *x5* contenga la variabile intera *i*, *x6* contenga la variabile intera *Ris* e *x10* l'indirizzo base del vettore MemVett.

```
addi x6, x6, 0
addi x29, x0, 100
CICLO: ld x7, 0(x10)
       add x5, x5, x7
       addi x10, x10, 8
       addi x6, x6, 1
       blt x6, x29, CICLO
```

**2.28** [10] < 2.7> Riscrivere il ciclo in assembler RISC-V dell'Esercizio 2.27 in modo da ridurre il numero totale di istruzioni RISC-V eseguite. Suggerimento: si noti che la variabile *i* viene utilizzata solo per il controllo del ciclo.

**2.29** [30] < 2.8> Tradurre il seguente frammento di codice C in codice assembler RISC-V. Suggerimento: ricordate che lo stack pointer deve rimanere allineato a valori multipli di 16.

```
int fib(int n) {
    if (n==0)
        return 0;
    else if (n==1)
        return 1;
    else
        return(fib(n-1) + fib(n-2));
}
```

**2.30** [20] < 2.8> Mostrare il contenuto dello stack dopo ciascuna chiamata di funzione dell'Esercizio 2.29. Si supponga che lo stack pointer contenga inizialmente l'indirizzo 0x7fffffff. Seguire le convenzioni sull'utilizzo dei registri specificate in Figura 2.11.

**2.31** [20] < 2.8> Tradurre la funzione *f* in codice assembler RISC-V. Si supponga che la dichiarazione della funzione *funz* sia `int funz(int a, int b);`. Il codice della funzione *f* è il seguente:

```
int f(int a, int b, int c, int d) {
    return funz(funz(a,b),c+d);
}
```

**2.32** [5] < 2.8> È possibile utilizzare l'ottimizzazione relativa alle chiamate ad anello (*tail call*) in questa funzione? Se la risposta fosse negativa, spiegare perché; se fosse affermativa, calcolare la differenza nel

numero di istruzioni eseguite per calcolare *f* con e senza l'ottimizzazione.

**2.33** [5] < 2.8> Un attimo prima che la funzione *f* dell'Esercizio 2.31 ritorni alla funzione chiamante, che cosa si sa sul contenuto dei registri *x10-x14*, *x8*, *x1* e *sp*? Si ricordi che conosciamo l'intera funzione *f*, ma della funzione *funz* abbiamo solamente la dichiarazione.

**2.34** [30] < 2.9> Scrivere un programma in linguaggio assembler RISC-V che converta un numero positivo o negativo rappresentato come sequenza di cifre codificate come caratteri ASCII in un numero intero. Il programma dovrà contenere nel registro *x10* l'indirizzo della stringa, terminata con il carattere *null*. La stringa conterrà una sequenza di cifre comprese tra 0 e 9 eventualmente preceduta dal segno: “-” o “+”. Il programma dovrà determinare il numero associato a questa sequenza di cifre e scrivere questo numero nel registro *x10*. Se all'interno della stringa comparisse un carattere non numerico, in una qualunque posizione, il programma dovrà terminare e restituire il valore -1 nel registro *x10*. Per esempio, se il registro *x10* puntasse alla sequenza di tre byte costituita dai numeri  $50_{dec}$ ,  $52_{dec}$ ,  $0_{dec}$  (corrispondenti alla stringa “24” seguita dal carattere terminatore null), il registro *x10* dovrà contenere il numero  $24_{dec}$  al termine del programma. L'istruzione RISC-V *mul* prende due registri in ingresso. Non esiste l'istruzione “*mult*”. Quindi, memorizzare semplicemente la costante 10 in un registro.

**2.35** Si consideri il seguente frammento di codice:

```
lb x6, 0(x7)
sd x6, 8(x7)
```

Si supponga che il registro *x7* contenga l'indirizzo 0x10000000 e che l'indirizzo dei dati sia 0x112334455667788.

**2.35.1** [5] < 2.3., 2.9> Quale valore è contenuto all'indirizzo 0x10000008 in una macchina con codifica “big-endian”?

**2.35.2** [5] < 2.3., 2.9> Quale valore è contenuto all'indirizzo 0x10000008 in una macchina con codifica “little-endian”?

**2.36** [5] < 2.10> Scrivere il frammento di codice assembler RISC-V che crea la costante su 64 bit: 0x1122334455667788<sub>due</sub> e la memorizza nel registro *x10*.

**2.37** [10] < 2.11> Scrivere il codice assembler RISC-V che implementa il seguente frammento di codice C come operazione atomica di “imposta il massimo” utilizzando la coppia di istruzioni *lr.d/sc.d*. In questo caso l'argomento *shvar* contiene l'indirizzo di una

variabile condivisa che dovrebbe essere sostituita da  $x$  se  $x$  è maggiore del valore a cui punta:

```
void setmax(int *shvar, int x= {
    // Inizia la sezione critica
    if (x > *shvar)
        *shvar = x;
    // Fine della sezione critica
}
```

**2.38** [5] < 2.11> Utilizzando il codice sviluppato per risolvere l'Esercizio 2.37 come esempio, spiegare cosa succede quando due processori iniziano a eseguire questa sezione critica contemporaneamente, ipotizzando che ciascun processore esegua esattamente una sola istruzione per ciclo di clock.

**2.39** Si supponga che per un certo processore, il CPI delle istruzioni aritmetiche sia 1, il CPI delle istruzioni di load/store sia 10 e il CPI delle istruzioni di salto condizionato (branch) sia 3. Si supponga che il programma sia costituito dal seguente numero di istruzioni: 500 milioni di istruzioni aritmetiche, 300 milioni di istruzioni di load/store e 100 milioni di istruzioni di salto condizionato.

**2.39.1** [5] < 1.6, 2.13> Si supponga che vengano aggiunte nuove istruzioni aritmetiche, più potenti, all'insieme delle istruzioni. Utilizzando queste nuove istruzioni, in media si può ridurre il numero di istruzioni aritmetiche richieste per eseguire un programma del 25%, con un aumento della durata del ciclo di clock del 10%. È una buona scelta di progetto? Perché?

**2.39.2** [5] < 1.6, 2.13> Si supponga che si trovi un modo per raddoppiare le prestazioni delle istruzioni aritmetiche. Quale sarebbe lo speedup complessivo del calcolatore? Quale sarebbe lo speedup se si migliorassero le prestazioni delle istruzioni aritmetiche di 10 volte?

**2.40** Si supponga che per un certo programma il 70% delle istruzioni eseguite siano istruzioni aritmetiche, il 10% sia load/store e il 20% sia salti condizionati.

**2.40.1** [5] < 1.6, 2.13> Data la distribuzione delle istruzioni riportata in precedenza e assumendo che un'istruzione aritmetica richieda 2 cicli di clock, un'istruzione di lettura/scrittura richieda 6 cicli e un'istruzione di salto condizionato richieda 3 cicli, determinare il CPI medio.

**2.40.2** [5] < 1.6, 2.13> Quanti cicli di clock deve durare un'istruzione aritmetica, in media, per ottenere un miglioramento delle prestazioni del 25%, se le istruzioni di lettura/scrittura e di salto condizionato non vengono rese più veloci?

**2.40.3** [5] < 1.6, 2.13> Quanti cicli di clock deve durare un'istruzione aritmetica, in media, per ottenere un miglioramento delle prestazioni del 50%, se le istruzioni di lettura/scrittura e di salto condizionato non vengono velocizzate?

**2.41** [10] < 2.19> Si supponga che l'ISA del RISC-V comprenda una modalità di indirizzamento mediante offset scalato simile a quella dell'x86 descritta nel paragrafo 2.17 (Figura 2.35). Descrivere come si potrebbe utilizzare l'offset scalato nelle istruzioni di load per ridurre ulteriormente il numero di istruzioni assembler necessarie per eseguire la funzione dell'Esercizio 2.4.

**2.42** [10] < 2.19> Si supponga che l'ISA del RISC-V comprenda una modalità di indirizzamento mediante offset scalato simile a quella dell'x86 descritta nel paragrafo 2.17 (Figura 2.35). Descrivere come si potrebbe utilizzare l'offset scalato nelle istruzioni di load per ridurre ulteriormente il numero di istruzioni assembler necessarie per implementare il codice C dell'Esercizio 2.7.

## Risposte alle domande di autovalutazione

Paragrafo 2.2, pagina 58 – RISC-V, C, Java.

Paragrafo 2.3, pagina 64 – 2. Molto lento.

Paragrafo 2.4, pagina 71 – 2.  $-8_{dec}$ .

Paragrafo 2.5, pagina 78 – 3. sub  $x11, x10, x9$

Paragrafo 2.6, pagina 81 – Entrambe. L'operazione di AND può essere vista come una maschera di 1, scrive 0 dappertutto tranne che nei campi desiderati. Lo shift a sinistra di un numero di posizioni adeguato rimuove i bit alla sinistra del campo. Lo shift a destra di un numero di posizioni adeguato posiziona il campo nei bit meno significativi della parola doppia e inserisce 0 nel resto della parola doppia. Si noti che AND lascia il campo dov'era originariamente, mentre la coppia di operazioni di shift sposta il campo nei bit meno significativi della parola doppia.

Paragrafo 2.7, pagina 86 – I. Sono tutte vere. II. 1.

Paragrafo 2.8, pagina 95 – Sono entrambe vere.

Paragrafo 2.9, pagina 100 – I. 1. e II. 2. 3.

Paragrafo 2.10, pagina 107 – I. 4)  $\pm 4K$ . II. 4)  $\pm 1M$ .

Paragrafo 2.11, pagina 110 – Sono entrambe vere.

Paragrafo 2.12, pagina 119 – Indipendenza dal calcolatore.



# 3

## L'aritmetica dei calcolatori

*La precisione numerica è la vera anima della scienza.*

Sir D'Arcy Wentworth Thompson,  
*On Growth and Form*, 1917

### 3.1 | Introduzione

Le parole dei calcolatori sono composte da bit e pertanto possono essere rappresentate come numeri binari. Nel Capitolo 2 abbiamo visto che i numeri interi possono essere rappresentati in forma decimale o in forma binaria; più complessa è la rappresentazione di altri tipi di numeri di uso comune. Per esempio:

- Come vengono trattati le frazioni e gli altri numeri reali?
- Che cosa succede se un'operazione genera un numero più grande di quello rappresentabile?
- Dietro a queste domande si nasconde un mistero: come fa l'hardware a eseguire le operazioni di moltiplicazione e divisione sui numeri?

Scopo di questo capitolo è svelare il mistero della rappresentazione dei numeri reali, degli algoritmi che implementano le operazioni, dell'hardware che implementa questi algoritmi e delle implicazioni che tutto questo ha sull'insieme delle istruzioni. I temi trattati possono chiarire anche alcuni dubbi sul funzionamento dei calcolatori. Vi mostreremo anche come sfruttare queste conoscenze per rendere molto più veloci i programmi ad alta intensità aritmetica.

### 3.2 | Somme e sottrazioni

L'operazione di somma in un calcolatore è svolta in maniera del tutto simile a come la si svolge a mano. Le cifre vengono sommate bit a bit da destra verso

*La sottrazione: l'amica complicata della somma.*

N. 10, Top Ten Courses for Athletes at a Football Factory, David Letterman et al., *Book of Top Ten Lists*, 1990

sinistra, con il riporto passato alla cifra a sinistra. La sottrazione è basata sulla somma: l'operando opportuno (sottraendo) viene semplicemente negato prima di essere sommato.

### Somma e sottrazione binaria

#### ESEMPIO

Proviamo a sommare  $6_{dec}$  a  $7_{dec}$  e poi a sottrarre  $6_{dec}$  da  $7_{dec}$  utilizzando la codifica binaria.

$$\begin{array}{r}
 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000111_{due} = 7_{dec} \\
 + 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000110_{due} = 6_{dec} \\
 \hline
 - 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00001101_{due} = 13_{dec}
 \end{array}$$

Tutto si svolge sui 4 bit a destra; la Figura 3.1 mostra le operazioni di somma e di riporto. I riporti sono scritti in parentesi e le frecce mostrano la loro propagazione.

#### SOLUZIONE

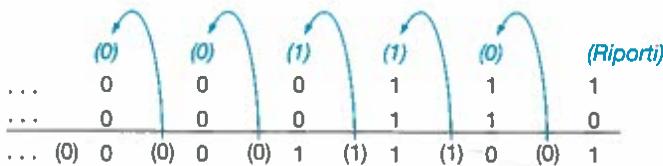
La sottrazione di  $6_{dec}$  da  $7_{dec}$  può essere eseguita direttamente:

$$\begin{array}{r}
 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000111_{due} = 7_{dec} \\
 - 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000110_{due} = 6_{dec} \\
 \hline
 - 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000001_{due} = 1_{dec}
 \end{array}$$

oppure attraverso la somma utilizzando la rappresentazione in complemento a 2 di  $-6$ :

$$\begin{array}{r}
 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000111_{due} = 7_{dec} \\
 + 11111111 11111111 11111111 11111111 11111111 11111111 11111111 11111010_{due} = -6_{dec} \\
 \hline
 - 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000001_{due} = 1_{dec}
 \end{array}$$

Come già detto in precedenza, l'overflow si verifica quando il risultato di un'operazione non può essere rappresentato con l'hardware a disposizione, in questo caso una parola di 64 bit. In quali casi può verificarsi l'overflow in una somma? Se si sommano operandi con segno opposto non può verificarsi overflow: la ragione è che la somma non può essere più grande, in valore assoluto, di uno degli operandi; per esempio,  $-10 + 4 = -6$ . Poiché gli operandi sono rappresentabili con 64 bit e la somma non può essere più grande di un operando, anche la somma deve essere contenuta in 64 bit. Quindi non può verificarsi overflow quando si sommano un operando positivo e uno negativo.



**Figura 3.1** Somma binaria: sono visibili i riporti da destra verso sinistra. Nella posizione del bit più a destra si somma 1 a 0 e il risultato della somma di questo bit è 1, mentre il riporto generato vale 0. Quindi l'operazione da effettuare sul secondo bit a destra sarà 0 + 1 + 1. Questa genera uno 0 come somma e un riporto pari a 1. Sul terzo bit si esegue la somma 1 + 1 + 1, producendo una somma pari a 1 e un riporto pari a 1. Sul quarto bit si esegue 1 + 0 + 0, generando 1 come bit di somma e nessun riporto.

Operazione	Operando A	Operando B	Risultato che indica un overflow
$A + B$	$\geq 0$	$\geq 0$	$< 0$
$A + B$	$< 0$	$< 0$	$\geq 0$
$A - B$	$\geq 0$	$< 0$	$< 0$
$A - B$	$< 0$	$\geq 0$	$\geq 0$

Figura 3.2 Condizioni di overflow per somma e sottrazione.

Vi sono regole analoghe sull'overflow nel caso di una sottrazione, ma il principio è opposto al precedente: quando il segno degli operandi è lo stesso, l'overflow non può verificarsi. Per comprendere ciò, ricordiamo che  $c - a = c + (-a)$ , poiché si può eseguire la sottrazione negando il secondo operando e poi sommandolo. Quindi, quando si sottraggono operandi con lo stesso segno si finisce per sommare operandi con segno opposto. Sulla base di quanto detto nel paragrafo precedente, è chiaro che in questo caso l'overflow non può verificarsi.

Abbiamo visto quando l'overflow non può verificarsi in caso di somma e sottrazione, ma non abbiamo ancora dato una risposta a come rilevarlo quando si verifica. Chiaramente, sommando o sottraendo due numeri su 64 bit si può ottenere un risultato che richiede 65 bit per essere completamente rappresentato; la mancanza del 65-esimo bit implica che quando si verifica l'overflow il bit di segno assume il valore del risultato, anziché il segno del risultato. Dal momento che sarebbe necessario un bit in più, soltanto il bit di segno può essere sbagliato. Perciò, l'overflow si può verificare quando si sommano due numeri positivi e la loro somma appare negativa e viceversa: in questi casi viene generato un riporto sul bit di segno.

Nelle operazioni di sottrazione, l'overflow si verifica quando si sottrae un numero negativo da un numero positivo e si ottiene un risultato negativo, o quando si sottrae un numero positivo da un numero negativo e si ottiene un numero positivo: questo significa che si è verificato un riporto sul bit di segno. La Figura 3.2 mostra le combinazioni di operazioni, operandi e risultati che indicano il verificarsi di un overflow.

Abbiamo visto come rilevare l'overflow nella rappresentazione dei numeri in complemento a 2. Che cosa possiamo dire dell'overflow quando abbiamo a che fare con numeri interi senza segno? Gli interi senza segno sono comune-mente utilizzati per l'indirizzamento della memoria, operazione per la quale l'overflow viene ignorato.

Fortunatamente, il compilatore può facilmente controllare l'overflow associato a numeri senza segno (*unsigned*) utilizzando un'istruzione di salto condizionato. Un'addizione genera overflow se la somma è inferiore a uno dei due addendi, mentre una sottrazione genera overflow quando la differenza è maggiore del minuendo.

L'Appendice A  descrive l'hardware che esegue le addizioni e le sottrazioni, chiamato **unità aritmetico-logica** o **ALU** (*Arithmetic Logic Unit*).

**Unità aritmetico-logica (ALU):**  
hardware che esegue le addizioni, le sottrazioni e di solito anche le operazioni logiche quali AND e OR.

Il progettista di un calcolatore deve decidere come gestire l'overflow aritmetico. Anche se alcuni linguaggi come C e Java ignorano l'overflow dei numeri interi, ci sono linguaggi di programmazione come Ada e Fortran che richiedono che la condizione di overflow sia segnalata al programma. Sarà compito del programmatore, o dell'ambiente di programmazione, decidere poi che cosa fare in caso di overflow.

**Interfaccia hardware/software**

## Riepilogo

Il concetto principale espresso in questo paragrafo è che, indipendentemente dalla rappresentazione, la dimensione finita delle parole dei calcolatori implica che le operazioni aritmetiche possano generare risultati troppo grandi per poter stare nella dimensione fissata delle parole. È facile rilevare l'overflow quando si lavora con i numeri senza segno, anche se in questo caso l'overflow viene quasi sempre ignorato, poiché i programmi non richiedono di rilevarlo quando si fanno operazioni sugli indirizzi, utilizzo più comune dei numeri interi senza segno. I numeri in complemento a 2 rappresentano un problema maggiore e alcuni sistemi software richiedono la rilevazione dell'overflow; perciò oggi tutti gli elaboratori possiedono un meccanismo progettato per questo scopo.

## Autovalutazione

Alcuni linguaggi di programmazione consentono di eseguire operazioni aritmetiche su numeri interi codificati in complemento a 2 in variabili dichiarate come byte o mezza parola, mentre i RISC-V offrono solamente operazioni aritmetiche su interi larghi l'intera parola. Come abbiamo visto nel Capitolo 2, i RISC-V hanno delle istruzioni apposite per trasferire dalla memoria singoli byte o singole mezze parole. Quale sequenza di istruzioni RISC-V è necessaria per eseguire delle operazioni aritmetiche su byte o mezze parole?

1. Caricare i dati dalla memoria mediante l'istruzione `lb` o `lh`; eseguire l'operazione aritmetica desiderata mediante `add`, `sub`, `mul` o `div`, utilizzando l'operatore AND per applicare una maschera di 8 o 16 bit al risultato di ciascuna operazione; al termine salvare il risultato in memoria utilizzando l'istruzione `sb` o `sh`.
2. Caricare i dati dalla memoria mediante l'istruzione `lb` o `lh`; eseguire l'operazione aritmetica desiderata mediante `add`, `sub`, `mul` o `div`; al termine salvare il risultato in memoria utilizzando l'istruzione `sb` o `sh`.

**Approfondimento.** Una caratteristica che in genere non si trova nei microprocessori a utilizzo generale è la saturazione. *Saturazione* significa che quando si verifica un overflow, il risultato assume il valore del numero positivo o negativo più grande rappresentabile, invece del valore del risultato calcolato nell'aritmetica in complemento a 2. La saturazione è ampiamente utilizzata nelle operazioni multimediali. Per esempio, sarebbe fastidioso se girando la manopola del volume della radio il volume aumentasse man mano per diventare poi improvvisamente bassissimo. Con un meccanismo di saturazione, il dispositivo, una volta raggiunto il volume massimo, continuerebbe a riprodurre il suono a quel volume anche se si continua a girare la manopola. Le estensioni multimediali degli insiemi standard di istruzioni prevedono spesso l'aritmetica con saturazione.

**Approfondimento.** La velocità di un'addizione viene aumentata determinando il riporto del bit più significativo in anticipo rispetto alla sua propagazione. Ci sono diversi schemi per farlo, in modo tale che nel peggior dei casi il riporto sia una funzione  $\log_2$  del numero di bit di ingresso dell'addizionatore. I bit di riporto vengono calcolati da questi schemi più velocemente perché i segnali attraversano un numero minore di porte logiche che lavorano in parallelo, anche se i circuiti corrispondenti richiedono un numero maggiore di porte. Lo schema più diffuso è chiamato *ad anticipazione di riporto (carry lookahead)* e viene descritto nel paragrafo A.6 dell'Appendice A .

### 3.3 | Moltiplicazione

Dopo aver illustrato somma e sottrazione, siamo ora pronti per affrontare la moltiplicazione, operazione più complessa.

Per prima cosa, ripassiamo la moltiplicazione dei numeri in rappresentazione decimale, con l'obiettivo di rivedere i passi elementari dell'algoritmo e il nome degli operandi.

Per ragioni che diventeranno chiare tra poco, nell'esempio seguente ci limiteremo a utilizzare solamente le cifre decimali 0 e 1. Moltiplichiamo  $1000_{\text{dec}}$  per  $1001_{\text{dec}}$ :

$$\begin{array}{r}
 \text{Moltiplicando} & \times & 1000_{\text{dec}} \\
 \text{Moltiplicatore} & & 1001_{\text{dec}} \\
 & & \hline
 & & 1000 \\
 & & 0000 \\
 & & 0000 \\
 & & 1000 \\
 & & \hline
 \text{Prodotto} & & 1001000_{\text{dec}}
 \end{array}$$

Il primo operando è chiamato *moltiplicando*, il secondo operando *moltiplicatore* e il risultato finale *prodotto*. Come ricorderete, l'algoritmo imparato alle scuole elementari consiste nel considerare le cifre del moltiplicatore una alla volta da destra a sinistra, nel moltiplicare il moltiplicando per la cifra del moltiplicatore considerata, e nello scalare ogni prodotto intermedio di una cifra verso sinistra rispetto al prodotto intermedio precedente.

La prima osservazione è che il numero di cifre del prodotto è considerevolmente maggiore rispetto al numero di cifre sia del moltiplicando sia del moltiplicatore. Infatti, ignorando il carattere riservato al segno, la lunghezza del prodotto fra un moltiplicando con  $n$  cifre decimali e un moltiplicatore con  $m$  cifre sarà in generale un numero di  $n + m$  cifre. Sono quindi necessarie  $n + m$  cifre per rappresentare tutti i possibili prodotti.

Di conseguenza, come per la somma, anche nella moltiplicazione si deve tenere conto della possibilità che si verifichi un overflow, dal momento che spesso si vuole ottenere un prodotto su 64 bit come risultato della moltiplicazione di due numeri su 64 bit.

Nell'esempio precedente, le cifre decimali erano limitate a 0 e a 1. Se queste sono le uniche possibilità, ciascun passo della moltiplicazione può essere semplificato:

1. si mette una copia del moltiplicando ( $1 \times$  moltiplicando) nella posizione opportuna se la cifra del moltiplicatore è 1;
2. oppure si mette 0 ( $0 \times$  moltiplicando) nella posizione opportuna se la cifra del moltiplicatore è 0.

Mentre nell'esempio precedente l'utilizzo delle sole due cifre decimali 0 e 1 era stato imposto, la moltiplicazione di numeri binari deve utilizzare per forza solo 0 e 1, e quindi possono verificarsi solamente le due possibilità riportate sopra.

Dopo aver ripassato le basi della moltiplicazione, il passo successivo consiste di solito nel presentare i circuiti logici altamente ottimizzati che implementano questa operazione. Qui non seguiremo l'approccio tradizionale, perché riteniamo che possiate capire meglio i moltiplicatori analizzando l'evoluzione dell'hardware e dell'algoritmo per la moltiplicazione attraverso le generazioni che si sono susseguite. Per il momento, assumiamo di moltiplicare solamente numeri positivi.

*La moltiplicazione è una vessazione. La divisione è altrettanto terribile. La regola del tre mi confonde, e la pratica mi rende pazzo.*

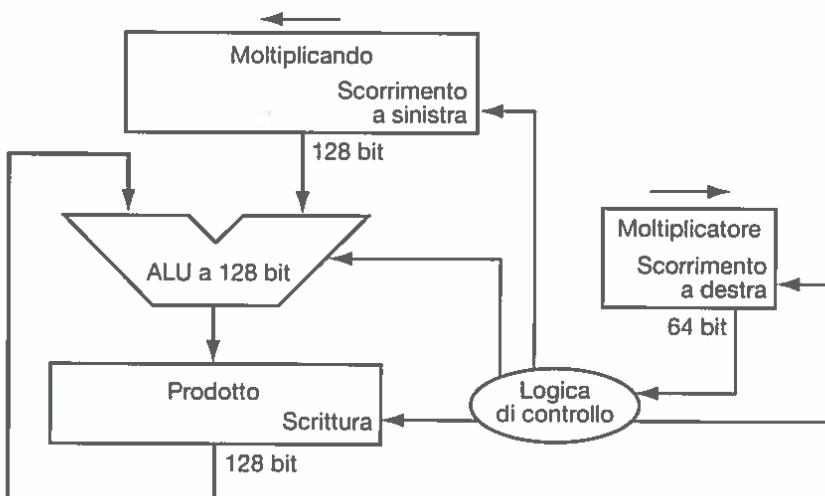
Anonimo, Manoscritto Elisabettiano, 1570

## Versione sequenziale dell'algoritmo della moltiplicazione e sua implementazione hardware

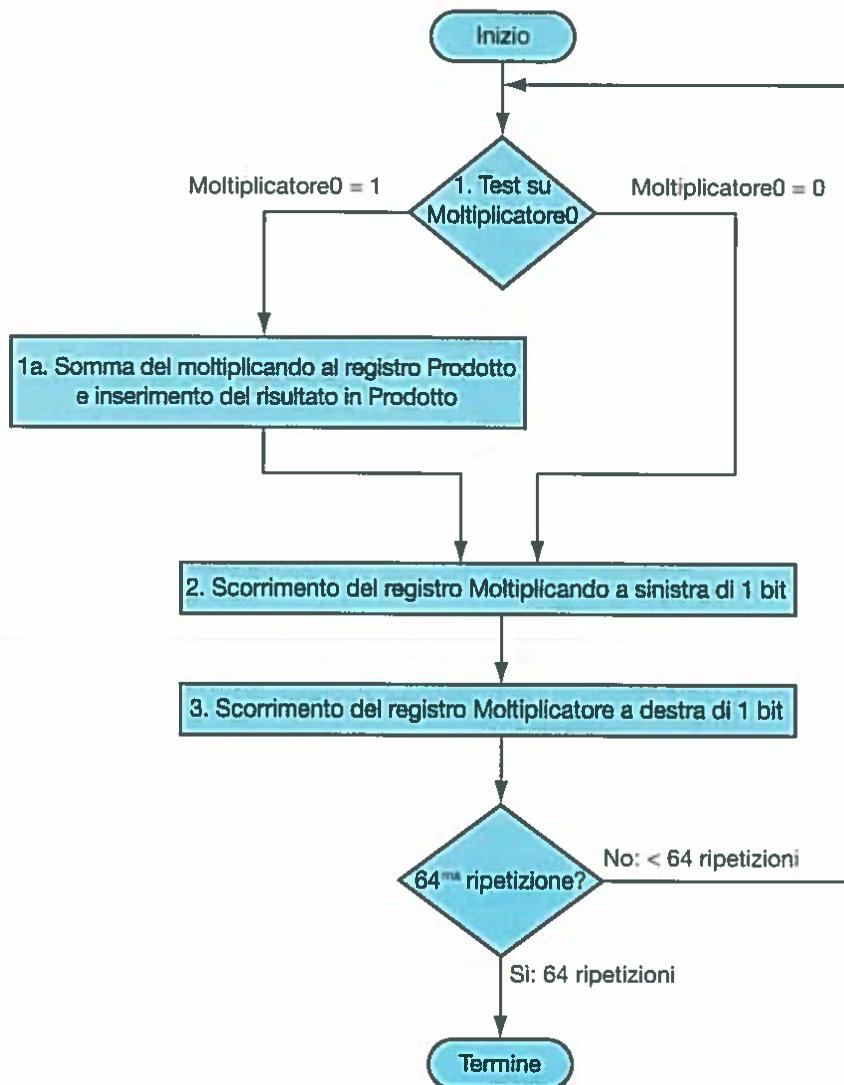
Il circuito sequenziale della moltiplicazione riprende l'algoritmo imparato a scuola; l'hardware corrispondente è mostrato in Figura 3.3. Il disegno è stato fatto in maniera tale che i dati scorrono dall'alto verso il basso, così da evidenziare meglio la similitudine con il metodo "carta e penna".

Si supponga che il moltiplicatore si trovi nel registro Moltiplicatore di 64 bit e che il registro Prodotto di 128 bit sia inizializzato a 0. Sulla base dell'esempio fatto con il metodo "carta e penna", è chiaro che sarà necessario spostare il moltiplicando a sinistra di una cifra a ogni passo, in maniera che possa essere sommato in modo corretto agli altri prodotti intermedi. Dopo 64 passi un moltiplicando su 64 bit si sarà spostato di 64 bit verso sinistra. Ci sarà quindi bisogno di un registro Moltiplicando a 128 bit; esso verrà inizializzato inserendo il moltiplicando nei 64 bit della metà di destra e una sequenza di 0 nella metà di sinistra. Il registro verrà poi scalato a sinistra di una posizione a ogni passo per allineare il moltiplicando con la somma parziale che viene accumulata nel registro Prodotto di 128 bit.

La Figura 3.4 mostra i tre passi fondamentali necessari per il calcolo di ogni bit. Il bit meno significativo del moltiplicatore (Moltiplicatore0) determina se Moltiplicando debba essere sommato al registro Prodotto. Lo scorrimento a sinistra nel passo 2 ha l'effetto di spostare il moltiplicando a sinistra, come si fa quando si esegue la moltiplicazione a mano. Lo scorrimento a destra nel passo 3 fornisce il prossimo bit del moltiplicatore, quello che verrà esaminato nell'iterazione successiva. Questi tre passi vengono ripetuti 64 volte per ottenere il prodotto. Se ogni passo corrispondesse a un ciclo di clock, l'algoritmo richiederebbe quasi 200 colpi di clock per la moltiplicazione di due numeri a 64 bit. L'importanza relativa delle operazioni aritmetiche come la moltiplicazione varia da programma a programma, ma, in generale, la somma e la sottrazione sono tra 5 e 100 volte più frequenti della moltiplicazione. Conseguentemente, in molte applicazioni è accettabile che la moltiplicazione richieda un certo



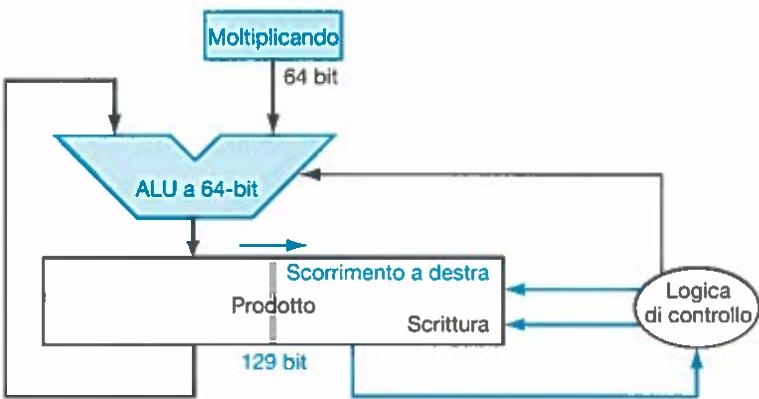
**Figura 3.3** Prima versione dell'hardware che realizza la moltiplicazione. Il registro Moltiplicando, la ALU (Arithmetic Logic Unit) e il registro Prodotto sono tutti lunghi 128 bit; solo il registro Moltiplicatore ne contiene 64 (l'Appendice A descrive le ALU). Il moltiplicando su 64 bit risiede inizialmente nella metà destra del registro Moltiplicando e scorre a sinistra di 1 bit a ogni passo. Il moltiplicatore scorre nella direzione opposta di 1 bit a ogni passo. All'inizio dell'algoritmo, il registro Prodotto viene inizializzato a 0. La logica di controllo decide quando fare scorrere il contenuto dei registri Moltiplicando e Moltiplicatore e quando scrivere il nuovo valore nel registro Prodotto.



**Figura 3.4** Prima versione dell'algoritmo della moltiplicazione, basata sull'hardware di Figura 3.3. Se il bit meno significativo del moltiplicatore è 1, si somma il moltiplicando al contenuto del registro Prodotto. Altrimenti, si procede al passo successivo. Nei due passi successivi si fanno scorrere, rispettivamente, il moltiplicando verso sinistra e il moltiplicatore verso destra. Questi tre passi vengono ripetuti 64 volte.

numero di colpi di clock in più, dato che questo non influenza significativamente le prestazioni. Tuttavia, la Legge di Amdahl (vedi par. 1.10) ricorda come un'operazione lenta, seppure utilizzata con una frequenza bassa, possa limitare le prestazioni.

L'algoritmo e l'hardware che lo implementa possono essere facilmente raffinati per eseguire la moltiplicazione in un ciclo di clock per passo. Questo miglioramento è dovuto all'esecuzione delle operazioni in parallelo: il moltiplicatore e il moltiplicando vengono fatti scorrere e contemporaneamente il moltiplicando viene sommato alla somma parziale, somma che avviene solo se il bit corrente del moltiplicatore è pari a 1. La logica di controllo deve assicurarsi che il test avvenga sul bit corretto del moltiplicatore e che prenda il moltiplicando prima che venga fatto scorrere. Tenendo conto che parte dei registri non viene utilizzata, l'hardware può essere ulteriormente ottimizzato per dimezzare la larghezza del sommatore e dei registri. La Figura 3.5 mostra la versione ottimizzata dell'hardware.



**Figura 3.5** Versione raffinata dell'hardware che realizza la moltiplicazione. Si confronti questo circuito con la prima versione, riportata in Figura 3.3. Il registro Moltiplicando e la ALU sono su 64 bit. Ora il contenuto del registro Prodotto viene fatto scorrere verso destra. Ora il contenuto del registro Prodotto viene fatto scorrere verso destra. La larghezza del registro Prodotto è diventata di 129 bit per poter contenere anche il riporto dell'addizionatore. Si noti che il registro Moltiplicatore è scomparso; il moltiplicatore, infatti, viene inserito nella metà destra del registro Prodotto. I cambiamenti sono evidenziati in colore.

## Interfaccia hardware/software

La sostituzione delle operazioni aritmetiche con quelle di scorrimento può essere utilizzata anche nel caso delle moltiplicazioni per valori costanti. Alcuni compilatori sostituiscono le moltiplicazioni per piccole costanti con una sequenza di operazioni di scorrimento, di somme e di sottrazioni. Dal momento che uno scorrimento a sinistra di un bit produce un numero due volte più grande in base 2, eseguire uno scorrimento a sinistra ha lo stesso effetto di una moltiplicazione per una potenza di 2. Perciò, come mostrato nel Capitolo 2, quasi tutti i compilatori sostituiscono la moltiplicazione per una potenza di 2 con uno scorrimento a sinistra, realizzando un'ottimizzazione basata sulla riduzione del peso delle istruzioni.

### Un algoritmo di moltiplicazione

#### ESEMPIO

Utilizzando numeri a 4 bit, per risparmiare spazio, eseguire la moltiplicazione  $2_{\text{dec}} \times 3_{\text{dec}}$ , ovvero  $0010_{\text{due}} \times 0011_{\text{due}}$ .

#### SOLUZIONE

La Figura 3.6 mostra il contenuto di ogni registro per ciascuno dei passi evidenziati in Figura 3.4 e il valore calcolato alla fine, che è pari a  $0000\ 0110_{\text{due}}$ , ovverosia  $6_{\text{dec}}$ . Il colore blu viene utilizzato per indicare le cifre contenute nei registri che vengono modificate in quel particolare passo, e il bit cerchiato è il bit che viene esaminato per determinare l'operazione del passo successivo.

## Moltiplicazione di numeri dotati di segno

Sinora abbiamo considerato solo numeri positivi. Il modo più semplice per gestire i numeri dotati di segno è quello di convertire prima moltiplicatore e moltiplicando in numeri positivi, effettuare la moltiplicazione, e poi considerare i segni originali. L'algoritmo dovrebbe quindi lavorare per 31 iterazioni, lasciando i bit di segno fuori dal calcolo. Come avete imparato a scuola, si deve negare il prodotto solo se i segni originali sono discordi.

Ne consegue che l'ultimo algoritmo funziona anche per i numeri dotati di segno, purché si ricordi che i numeri con cui si ha a che fare hanno un numero infinito di cifre di cui si rappresentano solo le prime 64. Di conseguenza, per

Iterazione	Passo	Moltiplicatore	Moltiplicando	Prodotto
0	Valori iniziali	0011	0000 0010	0000 0000
1	1a: $1 \Rightarrow$ Prod = Prod + Mcando 2: Scorr. sinistra Mcando 3: Scorr. destra Mcatore	0011	0000 0010 <b>0000 0100</b>	0000 0010
2	1a: $1 \Rightarrow$ Prod = Prod + Mcando 2: Scorr. sinistra Mcando 3: Scorr. destra Mcatore	0001	0000 0100 <b>0000 1000</b>	0000 0110
3	1: $0 \Rightarrow$ Nessuna operazione 2: Scorr. sinistra Mcando 3: Scorr. destra Mcatore	0000	0000 1000 0001 0000	0000 0110
4	1: $0 \Rightarrow$ Nessuna operazione 2: Scorr. sinistra Mcando 3: Scorr. destra Mcatore	0000	0001 0000 <b>0010 0000</b>	0000 0110
			0010 0000	0000 0110

**Figura 3.6** Esempio di moltiplicazione con l'impiego dell'algoritmo di Figura 3.4. Il bit esaminato per determinare il passo successivo è cerchiato.

i numeri dotati di segno, nelle operazioni di scorrimento occorre estendere il segno del numero contenuto in Prodotto. Al termine dell'esecuzione dell'algoritmo la parola meno significativa conterrà il prodotto su 64 bit.

## Moltiplicazione veloce

La Legge di Moore ha consentito di avere sempre più risorse a disposizione, tanto che i progettisti hardware possono ora costruire moltiplicatori molto più veloci rispetto al passato. Guardando ognuno dei 64 bit del moltiplicatore è possibile sapere, già all'inizio della moltiplicazione, se il moltiplicando debba essere sommato o meno. Moltiplicatori veloci si possono realizzare essenzialmente fornendo un sommatore a 64 bit per ogni bit del moltiplicatore: per ogni stadio del sommatore, uno dei due ingressi sarà il moltiplicando posto in AND con il bit corrispondente del moltiplicatore, mentre l'altro ingresso sarà l'uscita del sommatore precedente.

L'implementazione più immediata del circuito dei sommatori prevede di connettere l'uscita di ciascun addizionatore all'input dell'addizionatore alla sua sinistra, creando così una pila di addizionatori di 64 elementi. Un modo alternativo di organizzare queste 64 addizioni è in un albero parallelo, come mostrato in Figura 3.7. Invece di dover aspettare il tempo di 64 addizioni su 64 bit, dobbiamo aspettare solamente  $\log_2(64)$ , ovvero il tempo di 6 addizioni su 64 bit.

In realtà, la moltiplicazione può essere ulteriormente velocizzata se si utilizzano *sommatori a salvataggio di riporto* (vedi par. A.6 dell'Appendice A). Inoltre questo schema può essere facilmente implementato in modalità pipeline, così da poter supportare più moltiplicazioni simultaneamente (vedi Cap. 4).

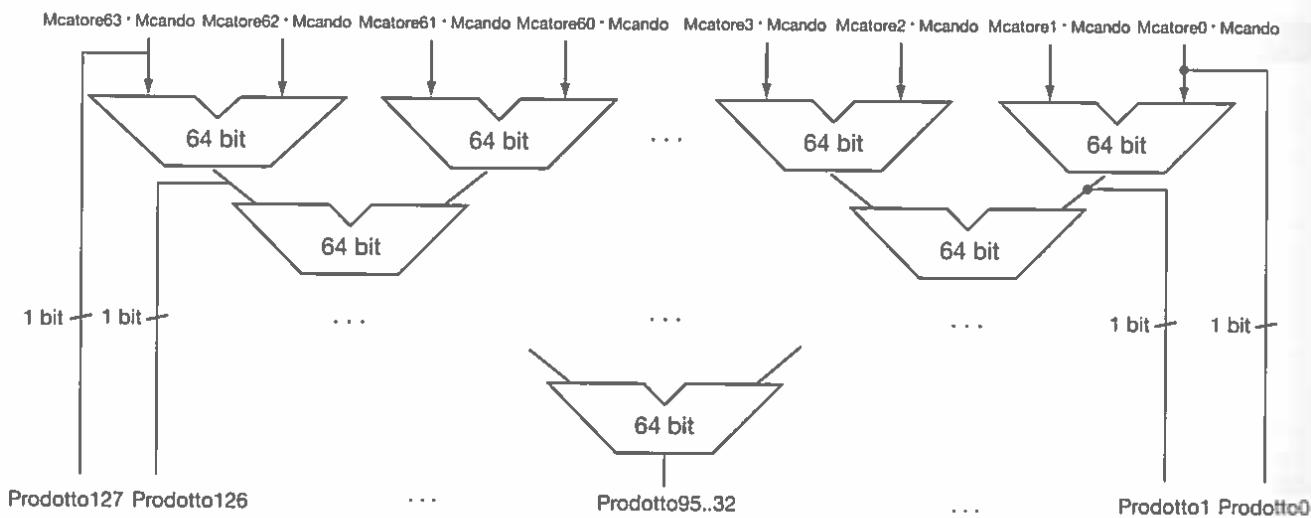


LEGGE DI MOORE



## Moltiplicazione nel RISC-V

Per produrre un risultato corretto, con numeri dotati o non dotati di segno, il RISC-V ha quattro istruzioni di moltiplicazione: *multiply* (`mul`), *multiply high* (`mulh`), *multiply high unsigned* (`mulhu`) e *multiply high signed-unsigned* (`mulhsu`).



**Figura 3.7** Hardware per la moltiplicazione veloce. Invece di utilizzare un sommatore a 64 bit per 64 volte, il ciclo viene scomposto in modo da utilizzare 63 sommatori a 64 bit connessi tra loro, minimizzando così i ritardi.

Per calcolare il prodotto di due numeri interi su 64 bit, il programmatore utilizza l'istruzione `mul`. Per recuperare i 64 bit più significativi, il programmatore utilizza `mulh` se entrambi gli operandi sono dotati di segno, `mulhu` se entrambi gli operandi sono non dotati di segno, e `mulhsu` se uno dei due operandi è dotato di segno e l'altro no.

### Riepilogo



L'hardware della moltiplicazione è basato semplicemente sulle operazioni di somma e scorrimento, le stesse utilizzate nella moltiplicazione effettuata a mano imparata a scuola. Anche i compilatori utilizzano le istruzioni di scorrimento per implementare le operazioni di moltiplicazione per potenze di 2. Con l'aggiunta di più hardware possiamo eseguire le somme in parallelo, velocizzando così la moltiplicazione.

## Interfaccia hardware/software

Il software può utilizzare l'istruzione `multiply high` per verificare se si è verificato un overflow in una moltiplicazione di numeri su 64 bit. Non si verifica overflow nella moltiplicazione di numeri su 64 bit non dotati di segno se il contenuto di `mulhu` è zero dopo avere calcolato il prodotto. Non si verifica overflow nella moltiplicazione di numeri su 64 bit dotati di segno se tutti i bit di `mulh` sono la replica del bit di segno di `mul` dopo avere calcolato il prodotto.

## 3.4 Divisione

*Divide et impera.*  
“Dividi e comanda”, antica massima latina

L'operazione inversa della moltiplicazione è la divisione, ancora meno frequente e più complessa. Essa offre anche la possibilità di eseguire un'operazione matematicamente non valida, la divisione per 0.

Iniziamo mostrando un esempio di divisione di due numeri, scritti in codifica decimale, per ripassare i nomi degli operandi e l'algoritmo della divisione imparati a scuola. Per ragioni simili a quelle del precedente paragrafo, ci limiteremo a considerare solamente le cifre decimali 0 e 1. L'esempio prevede la divisione

di  $1001010_{dec}$  per  $1000_{dec}$ :

	$1001_{dec}$							
Divisore	$1000_{dec}$							
	Quoziente							
	Dividendo							
	-1000							
	10							
	101							
	1010							
	-1000							
	10 <sub>dec</sub>							
	Resto							

I due operandi (**dividendo** e **divisore**) di una divisione, oltre a produrre il risultato (**quoziente**), producono anche un secondo risultato chiamato **resto**. Un altro modo per esprimere la relazione tra i componenti della divisione è il seguente:

$$\text{Dividendo} = \text{Quoziente} \times \text{Divisore} + \text{Resto}$$

dove il resto è più piccolo del divisore. In alcuni casi, poco frequenti, i programmi utilizzano l'istruzione di divisione soltanto per ottenere il resto, ignorando il valore del quoziente.

L'algoritmo di base per la divisione, appreso alle scuole elementari, verifica quanto è grande il numero che può essere sottratto, costruendo il quoziente una cifra alla volta, a ogni tentativo. L'esempio proposto (accuratamente scelto) utilizza soltanto i numeri 0 e 1, in modo che sia più facile verificare quante volte il divisore sta nella porzione di dividendo considerata: sarà sempre 0 volte oppure 1 volta. I numeri binari contengono soltanto 0 e 1, quindi anche la divisione binaria è limitata a queste due possibilità, risultando così relativamente semplice.

Si supponga che sia il dividendo sia il divisore siano numeri positivi e quindi che il quoziente e il resto siano non negativi; si supponga anche che gli operandi della divisione e ambedue i risultati siano numeri su 64 bit. Ignoriamo, per ora, il segno.

**Dividendo:** il numero che deve essere diviso.

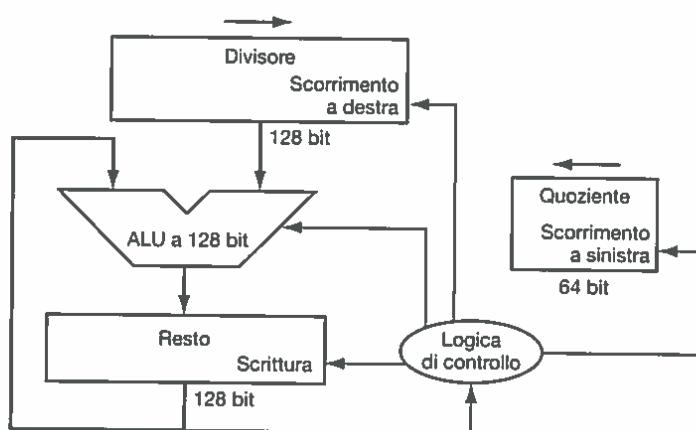
**Divisore:** il numero per cui il dividendo deve essere diviso.

**Quoziente:** il risultato principale di una divisione; il numero che quando viene moltiplicato per il divisore e poi sommato al resto produce il dividendo.

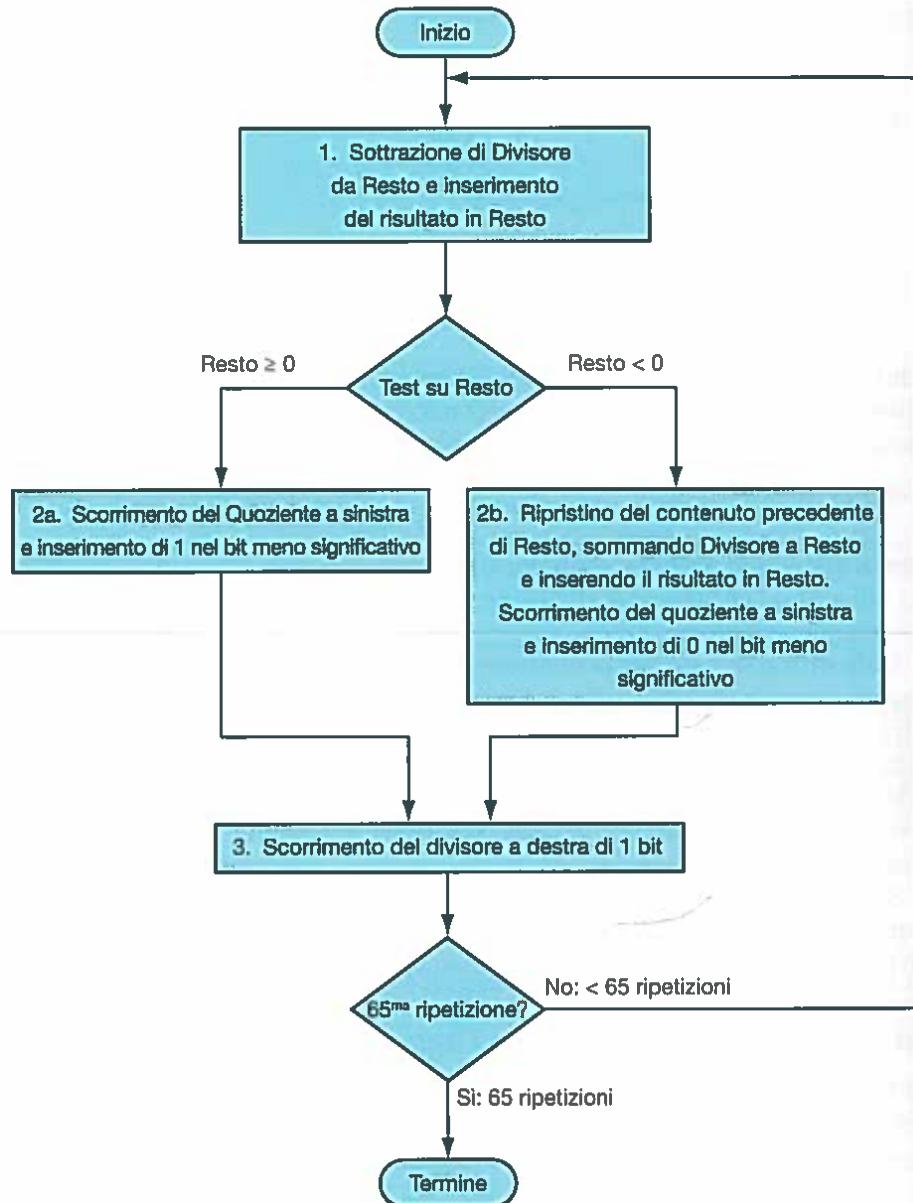
**Resto:** il risultato secondario di una divisione; il numero che, sommato al prodotto tra il quoziente e il divisore, produce il dividendo.

## Un algoritmo della divisione e l'hardware che lo implementa

La Figura 3.8 mostra l'hardware che implementa qualcosa di simile all'algoritmo imparato a scuola. Si comincia ponendo a 0 il registro Quoziente di 64 bit. A ogni passo l'algoritmo fa scorrere il divisore a destra di una posizione e quindi,



**Figura 3.8** Prima versione del circuito hardware che implementa la divisione. Il registro Divisore, la ALU e il registro Resto sono tutti di 128 bit, mentre il registro Quoziente è l'unico a essere di 64 bit. Il divisore su 64 bit si trova inizialmente nella metà sinistra del registro Divisore e viene fatto scorrere verso destra di 1 bit a ogni passo. Il resto viene inizializzato con il valore del dividendo. La logica di controllo decide quando eseguire lo scorrimento dei registri Divisore e Quoziente e quando scrivere il nuovo valore nel registro Resto.



**Figura 3.9** Un algoritmo di divisione basato sull'hardware di Figura 3.8. Se il resto è positivo, il divisore è contenuto nel dividendo e il passo 2a genera un 1 nel quoziente. Un resto negativo al passo 1 significa che il divisore non è contenuto nel dividendo; in questo caso, il passo 2b genera uno 0 nel quoziente e somma il divisore al resto, annullando così la sottrazione eseguita nel passo 1. L'operazione finale di scorrimento, nel passo 3, allinea opportunamente il divisore rispetto al dividendo per la successiva iterazione. I passi elencati vengono ripetuti 65 volte.

all'inizio, il divisore deve essere inserito nella metà sinistra del registro Divisore di 128 bit. A ogni passo viene eseguito uno scorrimento del divisore di 1 bit a destra per allineararlo con il dividendo. Il registro Resto viene inizializzato con il valore del dividendo.

La Figura 3.9 mostra i tre passi di questo primo algoritmo di divisione. A differenza degli esseri umani, il calcolatore non è abbastanza furbo da capire in anticipo se il divisore è più piccolo del dividendo. Deve prima sottrarre il divisore (passo 1); si ricordi che questo è il modo in cui si effettua il confronto nell'istruzione *set less than*. Se il risultato è positivo, il divisore è più piccolo o uguale al dividendo, e quindi viene generato un 1 nel quoziente (passo 2a). Se il risultato è negativo, il passo successivo consiste nel ripristinare il valore

precedente del resto, sommando il divisore al resto, e nell'inserire uno 0 nel quoziente (passo 2b). Il divisore viene quindi fatto scorrere a destra di una posizione e si ripete l'iterazione. Al termine delle iterazioni, il resto e il quoziente si troveranno nei registri corrispondenti.

### Un primo algoritmo di divisione

Utilizzando numeri a 4 bit, per risparmiare spazio, utilizzare l'algoritmo illustrato in precedenza per dividere  $7_{dec}$  per  $2_{dec}$ , ovvero  $0000\ 0111_{due}$  per  $0010_{due}$ .

La Figura 3.10 mostra il contenuto di ogni registro per ogni passo dell'algoritmo; al termine il quoziente vale  $3_{dec}$  e il resto  $1_{dec}$ . Si noti che il controllo eseguito al passo 2 per verificare se il resto sia positivo o negativo valuta semplicemente se il bit di segno del registro Resto sia uguale a 0 o a 1. La caratteristica sorprendente di questo algoritmo è che richiede  $n + 1$  passi per produrre il valore corretto di quoziente e resto.

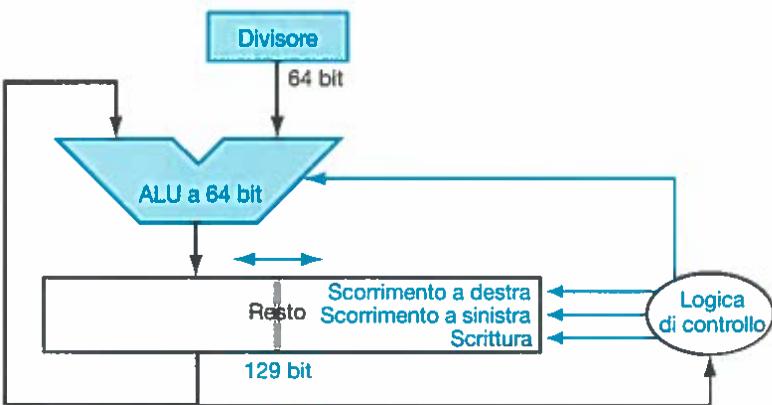
### ESEMPIO

### SOLUZIONE

Iterazione	Passo	Quoziente	Divisore	Resto
0	Valori iniziali	0000	0010 0000	0000 0111
1	1. Resto = Resto – Div	0000	0010 0000	0110 0111
	2b. Resto < 0 $\Rightarrow$ +Div, SLL Q, Q0 = 0	0000	0010 0000	0000 0111
	3: Scorr destra Div	0000	0001 0000	0000 0111
2	1. Resto = Resto – Div	0000	0001 0000	0111 0111
	2b. Resto < 0 $\Rightarrow$ +Div, SLL Q, Q0 = 0	0000	0001 0000	0000 0111
	3: Scorr destra Div	0000	0000 1000	0000 0111
3	1. Resto = Resto – Div	0000	0000 1000	0111 1111
	2b. Resto < 0 $\Rightarrow$ +Div, SLL Q, Q0 = 0	0000	0000 1000	0000 0111
	3: Scorr destra Div	0000	0000 0100	0000 0111
4	1. Resto = Resto – Div	0000	0000 0100	0000 0011
	2a. Resto $\geq$ 0 $\Rightarrow$ SLL Q, Q0 = 1	0001	0000 0100	0000 0011
	3: Scorr destra Div	0001	0000 0010	0000 0011
5	1. Resto = Resto – Div	0001	0000 0010	0000 0001
	2a. Resto $\geq$ 0 $\Rightarrow$ SLL Q, Q0 = 1	0011	0000 0010	0000 0001
	3: Scorr destra Div	0011	0000 0001	0000 0001

Figura 3.10 Esempio di divisione utilizzando l'algoritmo di Figura 3.9. Il bit che si esamina per determinare il passo successivo è cerchiato in colore.

L'algoritmo e l'hardware possono essere migliorati e diventare più veloci e più economici. Identificando le porzioni dei registri non utilizzate si può incrementare la velocità effettuando lo scorrimento degli operandi e del quoziente contemporaneamente alla sottrazione. Quindi questa strategia dimezza la larghezza del sommatore e dei registri. La versione così rivisitata dell'hardware è mostrata in Figura 3.11.



**Figura 3.11 Versione migliorata del circuito della divisione.** Registro Divisore, ALU e registro Quoziente sono tutti a 64 bit. Confrontata con la versione di Figura 3.8, l'ALU e il registro Divisore sono dimezzati, mentre il resto, all'inizio, viene inserito nei 64 bit più significativi del registro Resto. In questa versione i 64 bit meno significativi del registro Resto contengono il quoziente. (Come in Figura 3.5, il registro Resto è lungo 129 bit, per evitare che vada perso il riporto in uscita del sommatore del bit più significativo.)

### Divisione di numeri dotati di segno

Finora i numeri dotati di segno sono stati ignorati. La soluzione più semplice per gestirli consiste nel memorizzare il segno del divisore e del dividendo, e nel negare il quoziente se i segni sono discordi al termine dell'esecuzione.

**Approfondimento.** La complicazione nella divisione con segno è che si deve anche decidere il segno del resto. Si ricordi che deve sempre valere la seguente equazione:

$$\text{Dividendo} = \text{Quoziente} \times \text{Divisore} + \text{Resto}$$

Per capire come decidere il segno del resto, si consideri il seguente esempio, in cui si dividono tutte le combinazioni possibili di  $\pm 7_{\text{dec}}$  per  $\pm 2_{\text{dec}}$ . Il primo caso è facile:

$$+7 \div +2: \text{Quoziente} = +3, +\text{Resto} = +1$$

Verificando l'equazione della divisione si ha:

$$+7 = 3 \times 2 + (+1) = 6 + 1$$

Se si cambia il segno del dividendo, anche quello del quoziente deve cambiare:

$$-7 \div +2: \text{Quoziente} = -3$$

Riscrivendo la formula base per calcolare il resto, si ottiene:

$$\text{Resto} = (\text{Dividendo} - \text{Quoziente} \times \text{Divisore}) = -7 - (-3 \times +2) = -7 - (-6) = -1$$

Quindi:

$$-7 \div +2: \text{Quoziente} = -3, \text{Resto} = -1$$

Verificando l'equazione della divisione si ha:

$$-7 = -3 \times 2 + (-1) = -6 - 1$$

La ragione per cui il risultato non è un quoziente pari a  $-4$  e un resto pari a  $+1$ , che soddisfarebbe comunque la formula, è che il valore assoluto del quoziente

cambierebbe, a quel punto, a seconda del segno del dividendo e del divisore! Chiaramente, se fosse valida la relazione:

$$-(x \div y) \neq (-x) \div y$$

la programmazione sarebbe un compito ancora più impegnativo. Questo comportamento anomalo viene evitato seguendo la regola per cui il dividendo e il resto devono avere lo stesso segno, indipendentemente dal segno del divisore e del quoziente. Seguendo la stessa regola, si possono calcolare le divisioni delle altre combinazioni:

$$\begin{aligned} +7 \div -2: & \text{ Quoziente} = -3, \text{ Resto} = +1 \\ -7 \div -2: & \text{ Quoziente} = +3, \text{ Resto} = -1 \end{aligned}$$

In conclusione, un corretto algoritmo di divisione di numeri dotati di segno produce un quoziente negativo se i segni degli operandi sono opposti, e fa sì che il segno del resto, quando questo non è nullo, corrisponda a quello del dividendo.

## Una divisione più veloce

La Legge di Moore si applica anche all'hardware delle divisioni e delle moltiplicazioni, per cui vorremmo potere velocizzare queste operazioni aggiungendo altro hardware. In precedenza per velocizzare la moltiplicazione abbiamo adottato il trucco di utilizzare molti sommatori, ma non possiamo applicare la stessa tecnica alla divisione. Il motivo è che abbiamo bisogno di conoscere il segno della differenza prima di poter eseguire il passo successivo dell'algoritmo, mentre con la moltiplicazione potevamo calcolare immediatamente i 64 prodotti parziali.

Ci sono tecniche che producono più di un bit del quoziente per volta. La tecnica di *divisione SRT* tenta di predire il valore di più bit del quoziente a ogni passo, utilizzando una tabella predefinita che calcola i valori di alcuni bit del quoziente a partire da alcuni bit del dividendo e del resto calcolati prima. La tecnica fa affidamento sui passi successivi per correggere eventuali errori di predizione. Il numero tipico di bit previsti attualmente è 4 e la chiave è la previsione del valore da sottrarre. Nella divisione binaria c'è un'unica scelta; l'algoritmo utilizza 6 bit del resto e 4 bit del divisore per indicizzare una tabella che determina la predizione a ogni passo.

L'accuratezza di questo metodo veloce dipende dall'opportunità di avere i valori corretti nella tabella predefinita. L'errore tipico che si può verificare nel caso in cui la tabella di predizione non sia corretta viene illustrato al termine del capitolo, nel paragrafo 3.8.



LEGGE DI MOORE



PREDIZIONE

## Divisione nel RISC-V

Avrete probabilmente notato, confrontando le Figure 3.5 e 3.11, che lo stesso circuito logico può essere utilizzato sia per la moltiplicazione sia per la divisione. L'unico vincolo è che vi sia un registro di 128 bit in grado di eseguire lo scorrimento dei bit a sinistra o a destra, oltre a una ALU a 64 bit che possa eseguire somma e sottrazione. Per questo motivo, il RISC-V utilizza due istruzioni per la divisione e due per calcolare il resto: *divisione* (div), *divisione senza segno* (divu), *resto* (rem), *resto senza segno* (remu).

## Riepilogo

Lo stesso hardware implementa sia la moltiplicazione che la divisione e permette al RISC-V di fornire una coppia di registri a 64 bit che vengono utilizzati sia per la moltiplicazione sia per la divisione. La divisione viene accelerata predicendo

più bit del quoziente e correggendo gli errori di predizione nei passi successivi. La Figura 3.12 riassume quanto è stato aggiunto all'architettura RISC-V negli ultimi due paragrafi.

## Interfaccia hardware/software

Le istruzioni RISC-V per la divisione ignorano l'overflow, e quindi è il software che deve determinare se il quoziente è troppo grande. Oltre all'overflow, la divisione può anche ridursi a un calcolo improprio: la divisione per 0. Alcuni calcolatori distinguono questi due eventi anomali. Il software del RISC-V deve invece controllare il divisore per rilevare il verificarsi della divisione per 0, così come avviene per l'overflow.

### Linguaggio assembler MIPS

Categoria	Istruzioni	Esempi	Significato	Commenti
Aritmetica	Somma	add x5, x6, x7	$x5 = x6 + x7$	Operandi in tre registri
	Sottrazione	sub x5, x6, x7	$x5 = x6 - x7$	Operandi in tre registri
	Somma immediata	addi x5, x6, 20	$x5 = x6 + 20$	Somma di costanti
	Poni uguale a 1 se minore	slt x5, x6, x7	Se $(x5 < x6)$ $x5 = 1$ ; else $x5 = 0$	Operandi in tre registri
	Poni uguale a 1 senza segno	sltu x5, x6, x7	Se $(x5 < x6)$ $x5 = 1$ ; else $x5 = 0$	Operandi in tre registri
	Poni uguale a 1 se minore, immediato	slti x5, x6, x7	Se $(x5 < x6)$ $x5 = 1$ ; else $x5 = 0$	Comparazione con una costante
	Poni uguale a 1 se minore, immediato e senza segno	sltiu x5, x6, x7	Se $(x5 < x6)$ $x5 = 1$ ; else $x5 = 0$	Comparazione con una costante
	Moltiplicazione	mul x5, x6, x7	$x5 = x6 \times x7$	64 bit meno significativi del prodotto su 128 bit
	Moltiplicazione, parte alta	mulh x5, x6, x7	$x5 = (x6 \times x7) \gg 64$	64 bit più significativi del prodotto con segno su 128 bit
	Moltiplicazione senza segno, parte alta	mulhu x5, x6, x7	$x5 = (x6 \times x7) \gg 64$	64 bit più significativi del prodotto senza segno su 128 bit
	Moltiplicazione con o senza segno, parte alta	mulhsu x5, x6, x7	$x5 = (x6 \times x7) \gg 64$	64 bit più significativi del prodotto con o senza segno su 128 bit
	Divisione	div x5, x6, x7	$x5 = x6 / x7$	Divisione di numeri con segno su 64 bit
	Divisione senza segno	divu x5, x6, x7	$x5 = x6 / x7$	Divisione di numeri senza segno su 64 bit
	Resto	rem x5, x6, x7	$x5 = x6 \% x7$	Resto di una divisione di numeri con segno su 64 bit
	Resto senza segno	remu x5, x6, x7	$x5 = x6 \% x7$	Resto di una divisione di numeri senza segno su 64 bit
Trasferimento dati	Lettura parola doppia	ld x5, 40(x6)	$x5 = \text{Memoria}[x6 + 40]$	Trasferimento di una parola doppia da memoria a registro
	Memorizzazione parola doppia	sd x5, 40(x6)	$\text{Memoria}[x6 + 40] = x5$	Trasferimento di una parola doppia da registro a memoria
	Lettura parola	lw x5, 40(x6)	$x5 = \text{Memoria}[x6 + 40]$	Trasferimento di una parola da memoria a registro
	Lettura parola senza segno	lwu x5, 40(x6)	$x5 = \text{Memoria}[x6 + 40]$	Trasferimento di una parola senza segno da memoria a registro
	Memorizzazione parola	sw x5, 40(x6)	$\text{Memoria}[x6 + 40] = x5$	Trasferimento di una parola da registro a memoria
	Lettura mezza parola	lh x5, 40(x6)	$x5 = \text{Memoria}[x6 + 40]$	Trasferimento di una mezza parola da memoria a registro
	Lettura mezza parola senza segno	lhu x5, 40(x6)	$x5 = \text{Memoria}[x6 + 40]$	Trasferimento di una mezza parola senza segno da memoria a registro

## Linguaggio assembler MIPS

Categoria	Istruzioni	Esempi	Significato	Commenti
Trasferimento dati (continua)	Memorizzazione mezza parola	sh x5, 40(x6)	Memoria[x6 + 40] = x5	Trasferimento di una mezza parola da registro a memoria
	Lettura byte	lb x5, 40(x6)	x5 = Memoria[x6 + 40]	Trasferimento di un byte da memoria a registro
	Lettura byte senza segno	lbu x5, 40(x6)	x5 = Memoria[x6 + 40]	Trasferimento di un byte da memoria a registro
	Memorizzazione byte	sb x5, 40(x6)	Memoria[x6 + 40] = x5	Trasferimento di un byte da registro a memoria
	Lettura riservata	lr.d x5, (x6)	x5 = Memoria[x6]	Lettura: primo passo di uno scambio atomico
	Memorizzazione condizionata	sc.d x7, x5, (x6)	Memoria[x6] = x5; x7 = 0/1	Trasferimento: secondo passo di uno scambio atomico
	Caricamento dei bit superiori	lui x5, 0x12345	x5 = 0x12345000	Caricamento di una costante su 20 bit nei 20 bit superiori di una parola
Istruzioni logiche	Somma dei bit superiori di una costante al PC	auipc x5, 0x12345	x5 = PC + 0x12345000	Utilizzata per calcolare un indirizzo relativo al PC
	And	and x5, x6, x7	x5 = x6 & x7	Operandi in tre registri; AND bit a bit
	Or inclusivo	or x5, x6, x8	x5 = x6   x8	Operandi in tre registri; OR bit a bit
	Or esclusivo	xor x5, x6, x9	x5 = x6 ^ x9	Operandi in tre registri; XOR bit a bit
	And immediato	andi x5, x6, 20	x5 = x6 & 20	AND bit a bit tra un operando in registro e una costante
	Or inclusivo immediato	ori x5, x6, 20	x5 = x6   20	OR bit a bit tra un operando in registro e una costante
	Or esclusivo immediato	xori x5, x6, 20	x5 = x6 ^ 20	XOR bit a bit tra un operando in registro e una costante
Scorrimento	Scorrimento logico a sinistra	sll x5, x6, x7	x5 = x6 << x7	Scorrimento a sinistra del numero di bit specificato da un registro
	Scorrimento logico a destra	srl x5, x6, x7	x5 = x6 >> x7	Scorrimento a destra del numero di bit specificato da un registro
	Scorrimento aritmetico a destra	sra x5, x6, x7	x5 = x6 >> x7	Spostamento a destra del numero di bit specificato da un registro
	Scorrimento logico a sinistra immediato	slli x5, x6, 3	x5 = x6 << 3	Scorrimento a sinistra del numero di bit specificato dalla costante
	Scorrimento logico a destra immediato	srli x5, x6, 3	x5 = x6 >> 3	Scorrimento a destra del numero di bit specificato dalla costante
	Scorrimento aritmetico a sinistra immediato	srai x5, x6, 3	x5 = x6 >> 3	Scorrimento aritmetico a sinistra specificato dalla costante
	Salto condizionato	beq x5, x6, 100	Se (x5 == x6) vai a PC+100	Test di uguaglianza; salto relativo al PC
Salto condizionato	Salta se non è uguale	bne x5, x6, 100	Se (x5 != x6) vai a PC+100	Test di disuguaglianza; salto relativo al PC
	Salta se minore	bltx5, x6, 100	Se (x5 < x6) vai a PC+100	Test di minoranza; salto relativo al PC
	Salta se maggiore o uguale	bge x5, x6, 100	Se (x5 >= x6) vai a PC+100	Test di maggioranza; salto relativo al PC
	Salta se minore senza segno	bltu x5, x6, 100	Se (x5 < x6) vai a PC+100	Test di minoranza; salto relativo al PC
	Salta se maggiore o uguale senza segno	bgeu x5, x6, 100	Se (x5 >= x6) vai a PC+100	Test di maggioranza; salto relativo al PC
	Salto incondizionato	jal x1, 100	x1 = PC+4; vai all'indirizzo PC+100	Chiamata a procedura relativa al PC
	Salta e collega a registro	jalr x1, 100(x5)	x1 = PC+4; vai all'indirizzo x5+100	Ritorno da procedura; salto indiretto

Figura 3.12 Architettura del RISC-V. Si possono trovare le istruzioni in linguaggio macchina RISC-V nella scheda tecnica riassuntiva del RISC-V in fondo al libro.

**Approfondimento.** Vi è un algoritmo ancora più veloce che non esegue subito la somma del dividendo con il resto se il resto risulta negativo. L'algoritmo si limita a sommare il dividendo al resto scalato nel passo successivo, dal momento che:

$$(r + d) \times 2 - d = r - 2 + d \times 2 - d = r \times 2 + d$$

Questo algoritmo di divisione, denominato *senza ripristino (non restoring)*, richiede un solo ciclo di clock per ogni passo, e verrà analizzato più in dettaglio negli esercizi. Esiste anche un terzo algoritmo che non salva il risultato della sottrazione se questo è negativo ed è chiamato algoritmo di divisione *senza esecuzione (non performing)*. In media richiede un terzo in meno di operazioni aritmetiche.

### 3.5 Numeri in virgola mobile

*La velocità non ti porta da nessuna parte se stai andando nella direzione sbagliata.*

Proverbo americano

**Notazione scientifica:** notazione che prevede che il numero venga scritto con una sola cifra alla sinistra della virgola.

**Normalizzato:** un numero, in notazione a virgola mobile, che non ha zeri iniziali.

**Virgola mobile:** aritmetica dei calcolatori e rappresentazione dei numeri nei quali la virgola non è fissata in una certa posizione.

Oltre agli interi con o senza segno, i linguaggi di programmazione supportano anche i numeri con una parte frazionaria (ossia con alcune cifre dopo la virgola), cioè quelli che in linguaggio matematico sono chiamati *numeri reali*. Ecco alcuni esempi di numeri reali:

$3,14159265\dots_{\text{dec}}$  (pi)

$2,71828\dots_{\text{dec}}$  (e)

$0,000000001_{\text{dec}}$  o  $1,0_{\text{dec}} \times 10^{-9}$  (numero di secondi in un nanosecondo)

$3\,155\,760\,000_{\text{dec}}$  o  $3,15576_{\text{dec}} \times 10^9$  (numero di secondi in un secolo medio)

Si noti che l'ultimo numero non rappresenta una piccola frazione, ma è maggiore del massimo numero rappresentabile con un intero dotato di segno su 32 bit. Per gli ultimi due numeri si ricorre a una notazione alternativa, chiamata **notazione scientifica**, che prevede un'unica cifra a sinistra della virgola. Un numero in notazione scientifica senza zeri davanti alla virgola viene detto **normalizzato**, ed è questo il modo usuale per rappresentare i numeri in tale notazione. Per esempio,  $1,0_{\text{dec}} \times 10^{-9}$  è in notazione scientifica, ma  $0,1_{\text{dec}} \times 10^{-8}$  e  $10,0_{\text{dec}} \times 10^{-10}$  non lo sono.

Come avviene per i numeri in base dieci, anche i numeri binari si possono rappresentare in notazione scientifica:

$$1,0_{\text{due}} \times 2^{-1}$$

Per poter scrivere un numero binario in forma normalizzata, si deve disporre di una base da poter incrementare o decrementare dello stesso numero di bit di cui il numero deve essere scalato, arrivando ad avere una sola cifra non nulla alla sinistra della virgola. Solo una base pari a 2 soddisfa questo vincolo. Poiché la base non è 10, è anche necessario un nuovo nome per la virgola: *virgola binaria* è adatto allo scopo.

L'aritmetica dei calcolatori che supporta tali numeri è detta aritmetica in **virgola mobile** perché rappresenta numeri in cui la virgola binaria non è fissa, come per gli interi. Il linguaggio di programmazione C usa il termine *float* ("mobile") per tali numeri. Come nella notazione scientifica, i numeri vengono rappresentati con una singola cifra non nulla a sinistra della virgola binaria, quindi un numero normalizzato binario avrà la forma:

$$1,xxxxxxxxx_{\text{due}} \times 2^{yyy}$$

(Benché i calcolatori rappresentino anche l'esponente come numero in base 2, come il resto del numero, al fine di semplificare la notazione scriveremo l'esponente in notazione decimale.)

Una notazione scientifica standard in forma normalizzata per i numeri reali offre tre vantaggi: semplifica lo scambio di dati che includono numeri in virgola mobile; semplifica gli algoritmi aritmetici per la virgola mobile, in quanto questi sapranno che i numeri vengono sempre rappresentati in una certa forma; accresce l'accuratezza dei numeri memorizzabili in una parola, in quanto gli 0 alla sinistra del numero possono essere sostituiti da cifre aggiuntive della parte frazionaria, che consentono quindi di avvicinare di più il numero rappresentato al numero vero.

## Rappresentazione in virgola mobile

Il progettista che debba definire una rappresentazione in virgola mobile ha bisogno di trovare un compromesso tra la dimensione della **mantissa** e quella dell'**esponente**, dal momento che, essendo fissa la dimensione della parola, se si aggiunge un bit all'una lo si deve togliere all'altro. Questo compromesso si riflette sull'accuratezza della rappresentazione e della dimensione dell'intervallo dei numeri rappresentabili: se si aumenta la dimensione della mantissa si migliora l'accuratezza del numero, mentre aumentando la dimensione dell'esponente si aumenta l'intervallo dei numeri rappresentabili. Come prescritto dalle linee guida di progetto introdotte nel Capitolo 2, un buon progetto richiede buoni compromessi.

I numeri in virgola mobile sono codificati solitamente in un numero di bit multiplo rispetto alla dimensione della parola. Mostriamo ora la rappresentazione di un numero in virgola mobile nel RISC-V:  $s$  è il segno del numero in virgola mobile (1 se il segno è negativo) **esponente** è il valore del campo esponente su 8 bit (incluso il segno dell'esponente), e **mantissa** è un numero su 23 bit che rappresenta la parte frazionaria del numero. Come abbiamo già visto nel Capitolo 2, questo tipo di rappresentazione è chiamata in *modulo e segno*, dal momento che il segno costituisce un bit separato dal resto del numero.

**Mantissa:** il valore, generalmente compreso tra 0 e 1, che viene posto nel campo dopo la virgola.

**Esponente:** il valore assunto dall'esponente nella rappresentazione in virgola mobile dei numeri.

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
<b>s</b>	esponente							Mantissa																							
1 bit	8 bit							23 bit																							

In generale i numeri in virgola mobile vengono espressi nella forma:

$$(-1)^s \times F \times 2^E$$

dove  $F$  è il contenuto del campo mantissa ed  $E$  del campo esponente; l'esatta relazione tra questi campi verrà chiarita tra poco. (Vedremo tra breve che il RISC-V fa in realtà qualcosa di un po' più sofisticato.)

La dimensione scelta per l'esponente e la mantissa danno all'aritmetica del RISC-V un intervallo di rappresentabilità straordinario: si possono rappresentare numeri frazionari piccoli quanto  $2,0_{dec} \times 10^{-38}$  e numeri grandi quanto  $2,0_{dec} \times 10^{38}$ . Tuttavia straordinario non vuol dire infinito, per cui esiste la possibilità che un numero sia comunque troppo grande per essere rappresentato. Quindi, interrupt generati dall'overflow si possono verificare anche nell'aritmetica in virgola mobile e non solo in quella intera. Si noti che qui **overflow** significa che l'esponente, positivo, è troppo grande per poter essere rappresentato nel campo esponente.

La virgola mobile offre anche un nuovo tipo di situazione particolare. I programmati vogliono sapere quando hanno generato un numero troppo grande per essere rappresentato e, allo stesso modo, desiderano sapere se la parte frazionaria calcolata non sia così piccola da non poter essere rappresentata. In

**Overflow:** situazione in cui un esponente positivo diventa troppo grande per essere contenuto nel campo riservato all'esponente.

**Underflow:** situazione in cui un esponente negativo diventa troppo grande per essere contenuto nel campo riservato all'esponente.

**Doppia precisione:** un numero in virgola mobile rappresentato mediante due parole a 64 bit.

**Singola precisione:** un numero in virgola mobile rappresentato mediante una parola a 32 bit.

ambedue i casi, un programma potrebbe fornire dei risultati scorretti. Questa seconda situazione viene chiamata **underflow** per distinguerla dall'overflow; essa si verifica quando un esponente negativo è troppo grande per poter essere rappresentato nel campo esponente.

Un modo per ridurre la possibilità che si verifichino underflow e overflow è utilizzare una notazione che preveda un esponente più grande. In C questa codifica corrisponde al tipo di variabile *double* e le operazioni sui double sono chiamate operazioni aritmetiche in virgola mobile in **doppia precisione**, mentre la codifica corrispondente al formato introdotto precedentemente si chiama virgola mobile in **singola precisione**.

La rappresentazione di un numero in virgola mobile in doppia precisione richiede due parole RISC-V, come mostrato sotto; s è sempre il segno del numero, *esponente* è il valore del campo esponente su 11 bit e *mantissa* è un numero su 52 bit che rappresenta la parte frazionaria del numero.

63	62	61	60	59	58	57	56	55	54	53	52	51	50	49	48	47	46	45	44	43	42	41	40	39	38	37	36	35	34	33	32				
s	esponente											Mantissa																							
1 bit	11 bit											20 bit																							
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0				
mantissa																																			
32 bit																																			

La doppia precisione del RISC-V permette di rappresentare numeri piccoli quanto  $2,0_{\text{dec}} \times 10^{-308}$  e grandi quanto  $2,0_{\text{dec}} \times 10^{308}$ . Benché la doppia precisione incrementi l'intervallo dell'esponente, il suo principale vantaggio sta nella maggiore accuratezza della mantissa, dovuta alla sua maggiore ampiezza.

## Eccezioni e Interrupt

**Eccezione:** detta anche **interrupt**, è un evento non previsto che interrompe l'esecuzione di un programma. Può essere utilizzata per segnalare le condizioni di overflow.

**Interrupt:** eccezione che ha origine all'esterno del processore. Alcune architetture utilizzano il termine **interrupt** per indicare tutti i tipi di eccezioni.

Cosa deve succedere se si verifica un overflow o un underflow, perché l'utente possa venire a sapere che si è verificato un problema? Alcuni calcolatori segnalano questi eventi sollevando un'eccezione, talvolta chiamata interrupt. Un'eccezione o un interrupt sono essenzialmente una chiamata a procedura non prevista. L'indirizzo dell'istruzione che ha causato l'overflow viene salvato in un registro, e il programma salta a un indirizzo predefinito dove si trova il programma per la risposta appropriata a questa eccezione. L'indirizzo dell'istruzione interrotta viene salvato perché in alcune situazioni il programma interrotto può riprendere l'esecuzione dopo che è stato eseguito del codice di correzione. (Il paragrafo 4.9 descrive le eccezioni con maggiori dettagli: il Capitolo 5 descrive altre situazioni in cui si verificano eccezioni e interrupt.) I calcolatori RISC-V *non* sollevano un'eccezione quando si verifica un overflow o un underflow, ma il software può leggere il *registro di controllo e stato in virgola mobile* (fcsr) per controllare se si è verificata un'eccezione di overflow o underflow.

## Standard IEEE 754 per la virgola mobile

Questi formati hanno una validità più generale del RISC-V. Essi fanno parte dello standard **IEEE 754 per la virgola mobile**, utilizzato praticamente in ogni calcolatore a partire dal 1980. Tale standard ha significativamente aumentato sia la portabilità dei programmi in virgola mobile sia la qualità dell'aritmetica dei calcolatori. Per far stare ancora più bit all'interno della mantissa, l'IEEE 754 rende implicito il primo bit dei numeri binari normalizzati, che sappiamo essere sempre uguale a 1. Quindi, la mantissa è in realtà lunga 24 bit in singola precisione (l'1 sottinteso e la parte frazionaria vera e propria su 23 bit) e 53 bit

in doppia precisione (1 + 52). Per essere rigorosi, utilizzeremo il termine *significando* per rappresentare i 24 o 53 bit della mantissa costituiti dall'1 prima della virgola più la parte frazionaria, mentre con *mantissa* indicheremo i 23 o 52 bit della parte frazionaria vera e propria. Dato che il numero 0 non ha 1 iniziali, gli è stato riservato il valore uguale a 0 dell'esponente, così che l'hardware non inserisce un 1 iniziale.

Quindi 00 ... 00<sub>due</sub> rappresenta il numero 0; la rappresentazione degli altri numeri utilizza la formula riportata in precedenza, a cui va aggiunto l'1 reso implicito:

$$(-1)^S \times (1 + \text{Mantissa}) \times 2^E$$

dove i bit della mantissa rappresentano un numero compreso tra 0 e 1, ed E specifica il valore nel campo esponente, che verrà descritto in dettaglio tra poco. Se si considerano i bit della mantissa da *sinistra a destra* ( $s_1, s_2, s_3, \dots$ ), il numero contenuto al suo interno sarà dato da:

$$(-1)^S \times (1 + (s_1 \times 2^{-1}) + (s_2 \times 2^{-2}) + (s_3 \times 2^{-3}) + (s_4 \times 2^{-4}) + \dots) \times 2^E$$

La Figura 3.13 mostra la codifica IEEE 754 dei numeri in virgola mobile. Altre caratteristiche dello standard IEEE 754 sono le configurazioni speciali utilizzate per rappresentare situazioni particolari. Per esempio, invece di generare un interrupt quando viene richiesta una divisione per 0, il software può impostare la combinazione di bit associata a  $+\infty$  o  $-\infty$ ; l'esponente più grande (255 o 2047) viene riservato per rappresentare i simboli di infinito. In questi casi, un programma invece di stampare il risultato dell'operazione stamperà uno di questi due simboli. (Per gli studenti che hanno conoscenze avanzate di matematica, l'infinito rappresenta la chiusura topologica dell'insieme dei numeri reali.)

L'IEEE 754 prevede anche la situazione in cui il risultato provenga dall'esecuzione di un'operazione non valida, come 0/0 oppure la sottrazione di infinito da infinito; il simbolo associato a questa situazione è *NaN* (*Not a Number*). Lo scopo dei NaN è quello di permettere al programmatore di spostare alcuni test e decisioni in un punto del programma dove risultati convenienti.

I progettisti dell'IEEE 754 hanno introdotto, inoltre, una rappresentazione efficiente per implementare il confronto tra interi, avendo in mente in particolare l'ordinamento. Questa implementazione si riflette sulla posizione del segno, che si trova nel bit più significativo, permettendo di eseguire velocemente un test di maggiore, minore o uguale a 0. (In realtà le cose sono un po' più complicate di un semplice ordinamento di interi, poiché questa notazione è essenzialmente del tipo modulo e segno, anziché in complemento a 2.)

Singola precisione		Doppia precisione		Situazioni rappresentate
Esponente	Mantissa	Esponente	Mantissa	
0	0	0	0	0
0	Diverso da 0	0	Diverso da 0	$\pm$ numero denormalizzato
1-254	qualsiasi numero	1-2046	qualsiasi numero	$\pm$ numero in virgola mobile
255	0	2047	0	$\pm$ infinito
255	Diverso da 0	2047	Diverso da 0	NaN ( <i>Not a Number</i> )

**Figura 3.13** La codifica IEEE 754 dei numeri in virgola mobile. Un bit di segno, separato, determina il segno del numero. I numeri denormalizzati verranno descritti più avanti nell'approfondimento di pagina 191. Si possono trovare queste informazioni anche nella quarta colonna della scheda tecnica riassuntiva del RISC-V in fondo al libro.

Il posizionamento dell'esponente prima della mantissa semplifica inoltre l'ordinamento dei numeri in virgola mobile attraverso l'utilizzo delle istruzioni di confronto tra interi. Infatti, in questo modo i numeri con esponente maggiore appaiono più grandi dei numeri con esponente minore, purché ambedue gli esponenti abbiano lo stesso segno. Gli esponenti negativi rappresentano un problema per l'ordinamento. Se si utilizzasse il complemento a 2 o un'altra notazione in cui gli esponenti negativi hanno un 1 nel bit più significativo del campo esponente, un esponente negativo produrrebbe un numero grande. Per esempio  $1,0_{\text{due}} \times 2^{-1}$  verrebbe rappresentato in singola precisione come:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0			
0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(Si ricordi che l'1 iniziale è隐式 nel significando.) Il valore  $1,0_{\text{due}} \times 2^{-1}$  sarebbe invece un numero binario più piccolo:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0			
0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

La notazione desiderata deve quindi rappresentare l'esponente più negativo come 00 ... 00<sub>due</sub> e quello più positivo come 11 ... 11<sub>due</sub>. Questa notazione è chiamata *notazione polarizzata* e la polarizzazione è il numero sottratto alla rappresentazione normale senza segno per determinare il valore reale.

Lo standard IEEE 754 prevede una polarizzazione pari a 127 per la singola precisione, così che -1 viene rappresentato dalla combinazione di bit corrispondente a  $-1 + 127_{\text{dec}}$ , ovvero  $126_{\text{dec}} = 0111\ 1110_{\text{due}}$ , e +1 viene rappresentato da  $1 + 127$ , ovvero  $128_{\text{dec}} = 1000\ 0000_{\text{due}}$ . La polarizzazione dell'esponente per la doppia precisione è pari a 1023. La presenza di un esponente polarizzato implica che il valore rappresentato da un numero in virgola mobile nella realtà sia dato da:

$$(-1)^S \times (1 + \text{Mantissa}) \times 2^{(\text{Esponente} - \text{Polarizzazione})}$$

L'intervallo dei numeri in virgola mobile rappresentabili in singola precisione va dal numero più piccolo, che è uguale a:

$$\pm 1,0000000000000000000000000000000_{\text{due}} \times 2^{-126}$$

al numero più grande, pari a:

$$\pm 1,111111111111111111111111111111_{\text{due}} \times 2^{+127}$$

Vediamo ora alcuni esempi di questa rappresentazione.

### ESEMPIO

Si mostri la rappresentazione binaria IEEE 754 del numero  $-0,75_{\text{dec}}$  in singola e doppia precisione.

### SOLUZIONE

Il numero  $-0,75_{\text{dec}}$  è anche esprimibile come:

$$-3/4_{\text{dec}} \text{ o } -3/2^2_{\text{dec}}$$

(continua)

Può essere rappresentato dalla seguente frazione binaria:

$$-11_{\text{due}} / 2^2_{\text{dec}} \text{ o } -0,11_{\text{due}}$$

In notazione scientifica la sua rappresentazione è:

$$-0,11_{\text{due}} \times 2^0$$

e in notazione scientifica normalizzata:

$$-1,1_{\text{due}} \times 2^{-1}$$

In generale la rappresentazione per un numero in singola precisione è la seguente:

$$(-1)^s \times (1 + \text{Mantissa}) \times 2^{(\text{Esponente} - 127)}$$

e quindi, sottraendo la polarizzazione di 127 all'esponente di  $-1,1_{\text{due}} \times 2^{-1}$ , si ottiene:

$$(-1)^1 \times (1 + 0,1000\ 0000\ 0000\ 0000\ 0000_{\text{due}}) \times 2^{(126 - 127)}$$

La rappresentazione binaria in singola precisione di  $-0,75_{\text{dec}}$  è quindi:

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
1	0	1	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

1 bit                    8 bit                    23 bit

La rappresentazione in doppia precisione è:

$$(-1)^1 \times (1 + 0,1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_{\text{due}}) \times 2^{(1022 - 1023)}$$

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
1	0	1	1	1	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1 bit	11 bit	20 bit	32 bit																													

Proviamo ora a procedere nella direzione opposta.

### Conversione dalla rappresentazione binaria nella rappresentazione decimale in virgola mobile

Quale numero decimale è rappresentato dal seguente numero codificato in virgola mobile?

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
1	1	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

ESEMPIO

(continua)

(continua)

**SOLUZIONE**

Il bit di segno è uguale a 1, il campo esponente contiene 129 e il campo mantissa contiene  $1 \times 2^{-2} = 1/4$ , cioè 0,25. Utilizzando l'equazione di base:

$$\begin{aligned} (-1)^s \times (1 + \text{Mantissa}) \times 2^{(\text{Esponente} - \text{Polarizzazione})} &= (-1)^1 \times (1 + 0,25) \times 2^{(129 - 127)} \\ &= -1 \times 1,25 \times 2^2 \\ &= -1,25 \times 4 \\ &= -5,0 \end{aligned}$$

Nei paragrafi seguenti verranno descritti gli algoritmi per la somma e la moltiplicazione in virgola mobile. Nella loro parte centrale essi utilizzano le corrispondenti operazioni intere sui significandi, ma è necessaria una parte di calcolo aggiuntiva per gestire gli esponenti e normalizzare il risultato. Verrà dapprima mostrata una derivazione intuitiva dell'algoritmo, utilizzando numeri decimali, per poi arrivare a una versione binaria più dettagliata, mostrata nelle figure.

**Approfondimento.** Seguendo le linee guida dell'IEEE, il comitato dell'IEEE 754 si riunì vent'anni dopo la pubblicazione dello standard per analizzare quali modifiche si fossero rese necessarie. Lo standard aggiornato, IEEE 754-2008, comprende quasi tutte le caratteristiche dell'IEEE 754-1985 e aggiunge i formati a 16 bit (mezza precisione) e 128 bit (quadrupla precisione). Non sono ancora stati costruiti sistemi che supportino la quadrupla precisione, ma sicuramente lo saranno in futuro. Lo standard aggiornato prevede anche la somma in virgola mobile decimale.

**Approfondimento.** Nel tentativo di aumentare l'intervallo di rappresentabilità senza togliere bit al significando, alcuni calcolatori prima dello standard IEEE 754 utilizzavano una base diversa da 2: per esempio, i mainframe IBM 360 e 370 utilizzavano la base 16. Poiché la variazione di 1 nell'esponente IBM implicava lo scorrimento di 4 bit della mantissa, i numeri "normalizzati" con base 16 potevano avere fino a 3 bit iniziali a 0! Quindi l'utilizzo di cifre esadecimali implicava che fino a 3 bit dovevano essere scartati dal significando, il che creava notevoli problemi nell'accuratezza delle operazioni in virgola mobile. Alcuni mainframe IBM ora supportano sia lo standard IEEE 754 sia il vecchio formato esadecimale.

### Addizione in virgola mobile

Per illustrare i problemi che si incontrano con la somma in virgola mobile, si consideri la somma a mano di due numeri in notazione scientifica:

$$9,999_{\text{dec}} \times 10^1 + 1,610_{\text{dec}} \times 10^{-1}$$

Si assume di poter memorizzare soltanto quattro cifre decimali per la mantissa e due cifre decimali per l'esponente.

**Passo 1.** Per sommare correttamente i due numeri, si deve spostare la virgola del numero che ha l'esponente più piccolo in modo da allinearla al numero più grande. Occorre quindi trovare una forma per il numero minore, ossia  $1,610_{\text{dec}} \times 10^{-1}$ , che abbia lo stesso esponente del numero più grande. Osserviamo che esistono diverse rappresentazioni di un numero in virgola mobile non normalizzato in notazione scientifica:

$$1,610_{\text{dec}} \times 10^{-1} = 0,1610_{\text{dec}} \times 10^0 = 0,01610_{\text{dec}} \times 10^1$$

Il numero più a destra è la versione desiderata, poiché il suo esponente coincide con quello del numero più grande,  $9,999_{\text{dec}} \times 10^1$ . Quindi il primo passo da compiere è lo spostamento a destra della mantissa del numero più piccolo, fino a che l'esponente del numero più piccolo non coincide con quello del numero più grande. Si possono però rappresentare solo quattro cifre decimali, quindi dopo lo scorrimento il numero effettivamente rappresentato sarà:

$$0,1610 \times 10^1$$

**Passo 2.** A questo punto si esegue la somma dei significandi:

$$\begin{array}{r} 9,999_{\text{dec}} \\ + 0,016_{\text{dec}} \\ \hline 10,015_{\text{dec}} \end{array}$$

La somma è quindi  $10,015_{\text{dec}} \times 10^1$ .

**Passo 3.** La somma ottenuta non è in notazione scientifica normalizzata, per cui è necessario correggerla:

$$10,015_{\text{dec}} \times 10^1 = 1,0015_{\text{dec}} \times 10^2$$

Quindi, a valle dell'addizione, può essere necessario eseguire uno scorrimento per portare il risultato nella forma normalizzata, con un opportuno aggiustamento dell'esponente. L'esempio mostra uno scorrimento a destra, ma se un numero fosse stato positivo e l'altro negativo, allora la somma avrebbe potuto avere diversi 0 iniziali, e avrebbe quindi richiesto uno scorrimento a sinistra. Ogni volta che l'esponente viene incrementato o decrementato, si deve verificare se ci sia stato un overflow o un underflow, cioè bisogna accertarsi che l'esponente continui a essere rappresentabile nel campo ad esso riservato.

**Passo 4.** Poiché si è assunto che il significando possa essere lungo soltanto quattro cifre (escludendo il segno), occorre arrotondare il numero. Nell'algoritmo imparato a scuola si utilizza la regola per cui il numero viene troncato se la cifra alla destra del punto desiderato è compresa tra 0 e 4, mentre si somma 1 alla cifra così ottenuta se il numero alla destra è compreso tra 5 e 9. Il significando del numero:

$$1,0015_{\text{dec}} \times 10^2$$

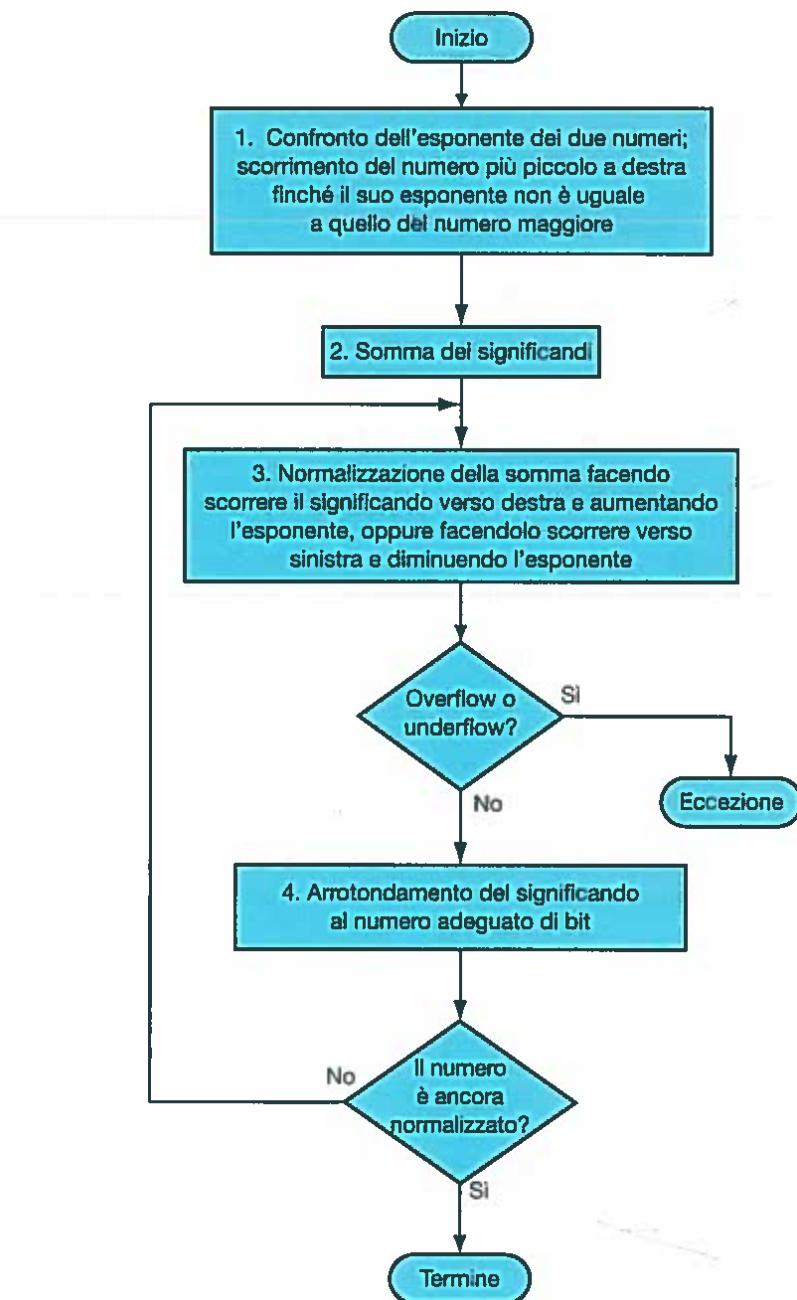
diventa quindi:

$$1,002_{\text{dec}} \times 10^2$$

essendo la quarta cifra alla destra della virgola compresa tra 5 e 9. Si noti che esiste un caso sfortunato per l'arrotondamento: quando dobbiamo sommare 1 a una stringa di 9, la somma può non essere più normalizzata; in questo caso occorre eseguire di nuovo il passo 3.

La Figura 3.14 mostra l'algoritmo per la somma binaria in virgola mobile che ricalca l'esempio decimale. I passi 1 e 2 sono simili all'esempio appena discusso.

so: si modifica il significando del numero con esponente minore e si sommano i due significandi. Il passo 3 normalizza il risultato eseguendo anche un controllo su overflow e underflow; il test del passo 3 dipende dalla precisione degli operandi. Si ricordi che una combinazione di bit tutti a 0 nell'esponente è utilizzata per la rappresentazione in virgola mobile dello 0. Inoltre, la combinazione di tutti 1 nell'esponente è riservata a valori e situazioni particolari, al di fuori dell'insieme dei normali numeri in virgola mobile (*vedi* la sezione *Approfondimento* al termine di questo paragrafo). Per l'esempio seguente, ricordate che per la singola precisione l'esponente più grande è 127 e quello più piccolo -126.



**Figura 3.14 Somma in virgola mobile.** La sequenza normale prevede l'esecuzione dei passi 3 e 4 una sola volta, ma se l'arrotondamento fa sì che la somma risultante sia denormalizzata, si deve ripetere il passo 3.

## Addizione binaria in virgola mobile

Provate a sommare i numeri  $0,5_{\text{dec}}$  e  $-0,4375_{\text{dec}}$  in forma binaria utilizzando l'algoritmo mostrato in Figura 3.14.

Anzitutto si consideri la versione binaria dei due numeri in notazione scientifica normalizzata, assumendo di adottare una precisione di 4 bit:

$$\begin{array}{lll} 0,5_{\text{dec}} & = 1/2_{\text{dec}} & = 1/2^1_{\text{dec}} \\ & = 0,1_{\text{due}} & = 0,1_{\text{due}} \times 2^0 & = 1,000_{\text{due}} \times 2^{-1} \\ -0,4375_{\text{dec}} & = -7/16_{\text{dec}} & = -7/2^4_{\text{dec}} \\ & = -0,0111_{\text{due}} & = -0,0111_{\text{due}} \times 2^0 & = 1,110_{\text{due}} \times 2^{-2} \end{array}$$

Ora seguiamo i passi indicati.

**Passo 1.** Il significando del numero con l'esponente minore ( $-1,11_{\text{due}} \times 2^{-2}$ ) viene fatto scorrere verso destra fino a che il suo esponente non coincide con quello del numero maggiore:

$$-1,110_{\text{due}} \times 2^{-2} = -0,111_{\text{due}} \times 2^{-1}$$

**Passo 2.** Si sommano i significandi:

$$1,000_{\text{due}} \times 2^{-1} + (-0,111_{\text{due}} \times 2^{-1}) = 0,001_{\text{due}} \times 2^{-1}$$

**Passo 3.** Si normalizza la somma, controllando se si è verificato un overflow o un underflow:

$$\begin{aligned} 0,001_{\text{due}} \times 2^{-1} &= 0,010_{\text{due}} \times 2^{-2} = 0,100_{\text{due}} \times 2^{-3} \\ &= 1,000_{\text{due}} \times 2^{-4} \end{aligned}$$

Dato che  $127 \geq -4 \geq -126$ , non c'è né underflow né overflow. L'esponente polarizzato sarebbe  $-4 + 127$ , ossia 123, che è compreso tra 1 e 254, che sono l'esponente più piccolo e quello più grande tra quelli utilizzabili.

**Passo 4.** Si arrotonda la somma:

$$1,0000_{\text{due}} \times 2^{-4}$$

La somma è contenuta esattamente in 4 bit, e quindi non ci sono modifiche dovute all'arrotondamento.

La somma è quindi:

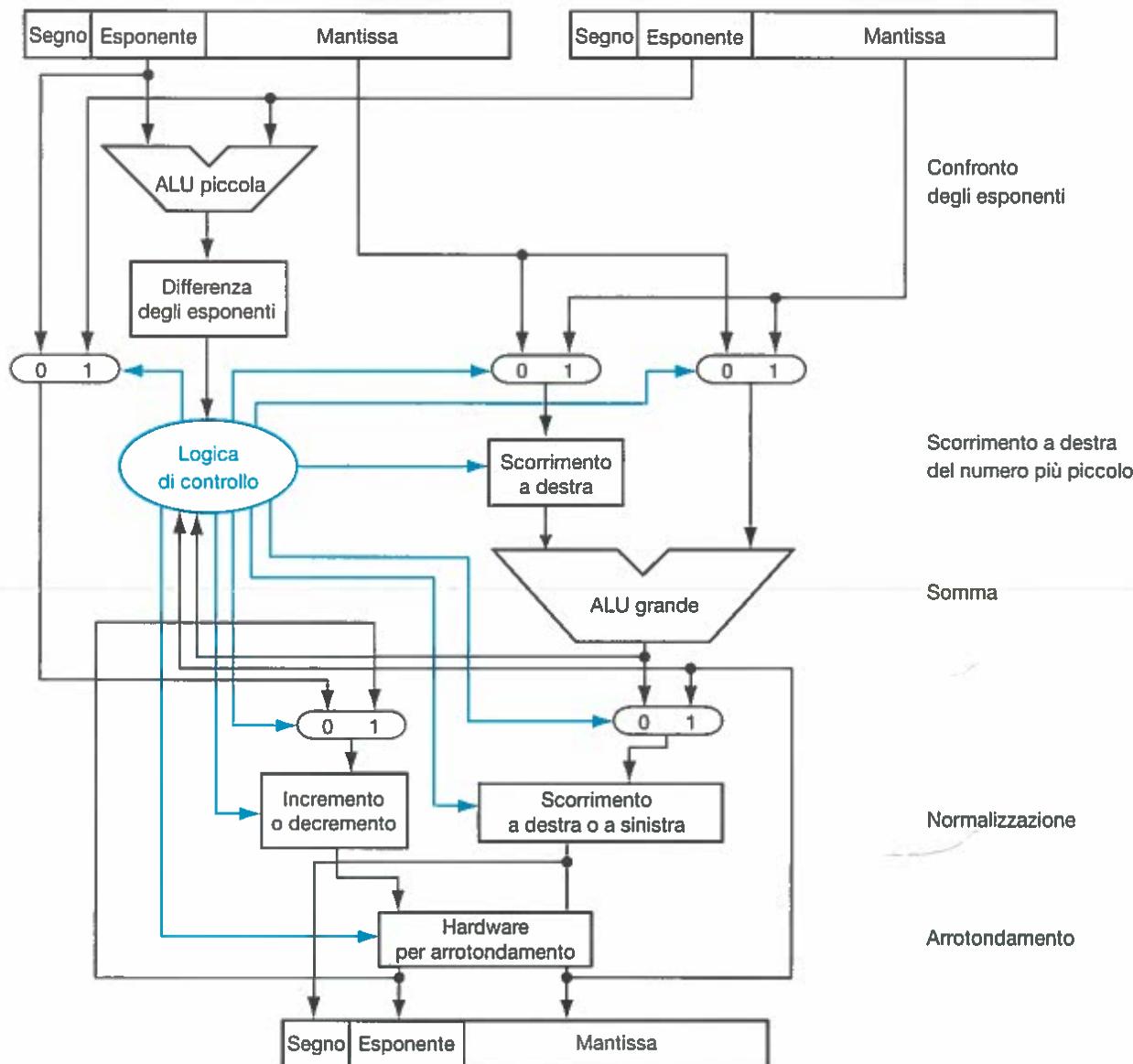
$$\begin{aligned} 1,000_{\text{due}} \times 2^{-4} &= 0,0001000_{\text{due}} &= 0,0001_{\text{due}} \\ &= 1/2^4_{\text{dec}} &= 1/16_{\text{dec}} &= 0,0625_{\text{dec}} \end{aligned}$$

Tale valore coincide con il valore che ci saremmo aspettati sommando

$0,5_{\text{dec}}$  e  $-0,4375_{\text{dec}}$

### ESEMPIO

### SOLUZIONE



**Figura 3.15** Diagramma a blocchi di un'unità aritmetica dedicata alla somma in virgola mobile. I passi indicati in Figura 3.14 corrispondono ai vari blocchi, dall'alto verso il basso. Anzitutto l'esponente di un operando viene sottratto da quello dell'altro, utilizzando la ALU piccola per determinare quale sia l'esponente maggiore e di quanto. Il risultato della differenza controlla i tre multiplexer; da sinistra a destra questi selezionano: l'esponente maggiore, il significando del numero minore e quello del numero maggiore. Il significando minore viene fatto scorrere verso destra e i significandi vengono poi sommati utilizzando la ALU grande. Nel passo di normalizzazione la somma viene fatta scorrere verso sinistra o verso destra, incrementando o decrementando così l'esponente. L'arrotondamento porta al risultato finale, che può richiedere una nuova normalizzazione.

Molti calcolatori possiedono componenti hardware preposti a eseguire le operazioni in virgola mobile nel modo più rapido possibile. La Figura 3.15 delinea l'organizzazione di base di un circuito per la somma in virgola mobile.

### Moltiplicazione in virgola mobile

Dopo aver illustrato la somma, possiamo affrontare la moltiplicazione in virgola mobile. Si supponga di moltiplicare a mano i seguenti numeri decimali, espressi in notazione scientifica:  $(1,110_{dec} \times 10^{10}) \times (9,200_{dec} \times 10^{-5})$ . Si supponga anche di poter memorizzare soltanto quattro cifre della mantissa e due cifre dell'esponente.

**Passo 1.** Diversamente da quanto si fa per la somma, si calcola l'esponente del prodotto semplicemente sommando gli esponenti degli operandi:

$$\text{Nuovo esponente} = 10 + (-5) = 5$$

Ripetiamo ora l'operazione con gli esponenti polarizzati, per verificare che si ottiene lo stesso risultato:  $10 + 127 = 137$  e  $-5 + 127 = 122$ , e quindi:

$$\text{Nuovo esponente} = 137 + 122 = 259$$

Il risultato è troppo grande per il campo esponente di 8 bit, quindi abbiamo dimenticato qualcosa! Il problema nasce dalla polarizzazione; infatti, assieme agli esponenti stiamo sommando anche le polarizzazioni:

$$\text{Nuovo esponente} = (10 + 127) + (-5 + 127) = (5 + 2 \times 127) = 259$$

*Quando si sommano due numeri polarizzati, per ottenere la somma corretta si deve sottrarre la polarizzazione dalla somma:*

$$\text{Nuovo esponente} = 137 + 122 - 127 = 259 - 127 = 132 = (5 + 127)$$

e 5 è effettivamente l'esponente che era stato calcolato inizialmente.

**Passo 2.** Nel passo successivo si esegue la moltiplicazione dei significandi:

$$\begin{array}{r} & 1,000_{\text{dec}} \\ \times & 9,200_{\text{dec}} \\ \hline & 0000 \\ & 0000 \\ & 2220 \\ & 9990 \\ \hline & 1110000_{\text{dec}} \end{array}$$

Ci sono tre cifre a destra della virgola per ciascun operando, per cui la virgola va inserita dopo la sesta cifra del prodotto:

$$10,212000_{\text{dec}}$$

Se si suppone di poter conservare soltanto tre cifre a destra della virgola, il prodotto sarà  $10,212 \times 10^5$ .

**Passo 3.** Il prodotto non è normalizzato, per cui diventa:

$$10,212_{\text{dec}} \times 10^5 = 1,0212_{\text{dec}} \times 10^6$$

Quindi, dopo la moltiplicazione, può succedere che il prodotto debba essere fatto scorrere verso destra di una posizione per metterlo in forma normalizzata e, parallelamente, occorre sommare 1 all'esponente. A questo punto si può controllare se si sia verificato un overflow o un underflow. L'underflow può verificarsi se ambedue gli operandi sono piccoli, e quindi se ambedue hanno degli esponenti negativi con valore assoluto grande.

- Passo 4. Era stata fatta l'ipotesi che la mantissa fosse lunga soltanto 4 bit (escludendo il segno), per cui occorre effettuare l'arrotondamento. Il numero:

$$1,0212_{\text{dec}} \times 10^6$$

viene arrotondato a quattro cifre:

$$1,021_{\text{dec}} \times 10^6$$

- Passo 5. Il segno del prodotto dipende dal segno degli operandi di partenza. Se questi sono concordi il segno del prodotto è positivo, altrimenti è negativo. Il risultato finale, quindi, è:

$$+1,021_{\text{dec}} \times 10^6$$

Il segno della somma nell'algoritmo per l'addizione veniva determinato dalla somma delle mantisse, mentre nella moltiplicazione il segno del prodotto è determinato dal segno degli operandi.

Di nuovo, come mostrato in Figura 3.16, i passi dell'algoritmo per la moltiplicazione di numeri binari in virgola mobile sono simili a quelli appena descritti. Si inizia con il calcolo del nuovo esponente del prodotto sommando gli esponenti polarizzati, facendo attenzione a sottrarre una volta la polarizzazione in maniera da ottenere il risultato corretto. Quindi, si esegue la moltiplicazione dei significandi, seguita dall'eventuale normalizzazione.

La dimensione dell'esponente viene controllata per rilevare eventuali situazioni di overflow o underflow e il prodotto viene poi arrotondato. Se l'arrotondamento porta a un'ulteriore normalizzazione, occorre di nuovo verificare la dimensione dell'esponente. Infine, si forza il bit di segno a 1 se i segni degli operandi erano diversi (prodotto negativo), o a 0 se erano uguali (prodotto positivo).

### Moltiplicazione decimale in virgola mobile

#### ESEMPIO

Si moltiplichino i numeri  $0,5_{\text{dec}}$  e  $-0,4375_{\text{dec}}$  utilizzando i passi riportati in Figura 3.16.

#### SOLUZIONE

In binario il compito corrisponde a moltiplicare  $1,000_{\text{due}} \times 2^{-1}$  per  $-1,110_{\text{due}} \times 2^{-2}$ .

- Passo 1. Si sommano gli esponenti senza polarizzazione:

$$-1 + (-2) = -3$$

o, utilizzando la rappresentazione polarizzata:

$$\begin{aligned} (-1 + 127) + (-2 + 127) - 127 &= (-1 - 2) + (127 + 127 - 127) \\ &= -3 + 127 = 124 \end{aligned}$$

(continua)

(continua)

Passo 2. Si moltiplicano i significandi:

$$\begin{array}{r} 1,000_{\text{due}} \\ \times 1,110_{\text{due}} \\ \hline 0000 \\ 1000 \\ 1000 \\ 1000 \\ \hline 1110000_{\text{due}} \end{array}$$

Il prodotto è  $1,110000_{\text{due}} \times 2^{-3}$ , ma bisogna ridurlo a 4 bit, per cui si ottiene  $1,110_{\text{due}} \times 2^{-3}$ .

Passo 3. Ora si verifica il prodotto per essere sicuri che sia normalizzato, e poi si controlla l'esponente per rilevare eventuali condizioni di overflow o underflow. Il prodotto è già normalizzato e, dato che  $127 \geq -3 \geq -126$ , non vi è né overflow né underflow. (Utilizzando la rappresentazione normalizzata si avrebbe  $254 \geq 124 \geq 1$ , per cui l'esponente può essere contenuto nel campo ad esso riservato.)

Passo 4. L'arrotondamento del prodotto non produce modifiche:

$$1,110_{\text{due}} \times 2^{-3}$$

Passo 5. Dato che il segno degli operandi è discorda, il segno del prodotto è negativo. Quindi il risultato finale della moltiplicazione è:

$$-1,110_{\text{due}} \times 2^{-3}$$

Convertiamolo in decimale per verificare il risultato:

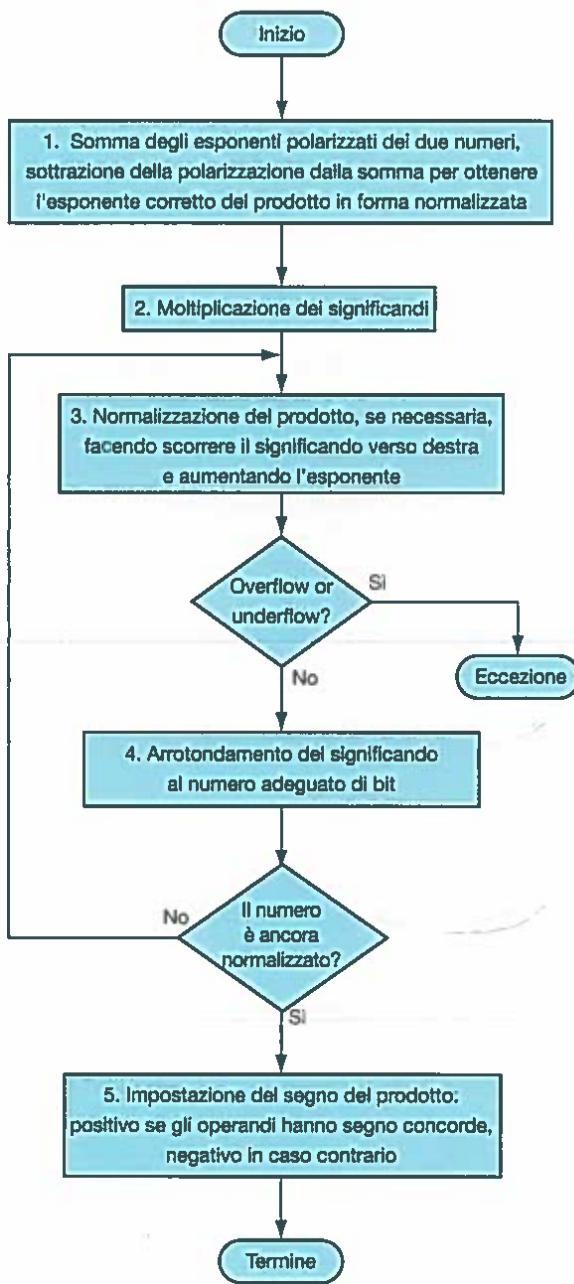
$$\begin{aligned} -1,110_{\text{due}} \times 2^{-3} &= -0,001110_{\text{due}} = -0,00111_{\text{due}} \\ &= -7/2^5_{\text{dec}} = -7/32_{\text{dec}} = -0,21875_{\text{dec}} \end{aligned}$$

Il prodotto di  $0,5_{\text{dec}}$  per  $-0,4375_{\text{dec}}$  è effettivamente  $-0,21875_{\text{dec}}$ .

## Istruzioni in virgola mobile nel RISC-V

Il RISC-V supporta i formati IEEE 754 per la singola e la doppia precisione con le seguenti istruzioni:

- Addizione in virgola mobile: *addition single* (fadd.s) e *addition double* (fadd.d);
- Sottrazione in virgola mobile: *subtraction single* (fsub.s) e *subtraction double* (fsub.d);
- Moltiplicazione in virgola mobile: *multiplication single* (fmul.s) e *multiplication double* (fmul.d);
- Divisione in virgola mobile: *division single* (fdiv.s) e *division double* (fdiv.d);



**Figura 3.16 Moltiplicazione in virgola mobile.** La sequenza normale di esecuzione prevede l'esecuzione dei passi 3 e 4 una sola volta, ma se l'arrotondamento fa sì che la somma risultante sia denormalizzata si deve ripetere il passo 3.

- Radice quadrata in virgola mobile, singola precisione (fsqrt.s) e doppia (fsqrt.d);
- Uguaglianza in virgola mobile, in singola precisione (feq.s) e doppia (feq.d);
- Test di minoranza in virgola mobile, singola precisione (flt.s) e doppia (flt.d)
- Test di minore uguale, singola precisione (fle.s) e doppia (fle.d)

Le istruzioni di comparazione feq, flt e fle impostano il contenuto di un registro intero a 0 se il risultato è falso e a 1 se il risultato è vero. Il software può quindi prendere o meno il salto mediante un'istruzione beq o bne, a seconda del risultato della comparazione.

I progettisti del RISC-V decisero di aggiungere dei registri in virgola mobile separati, chiamati  $f_0, f_1, \dots, f_{31}$ . Di conseguenza, hanno incluso operazioni separate di load e store per i registri in virgola mobile:  $f1d$  e  $f1s$  per la doppia precisione e  $f1w$  e  $f1s$  per la precisione singola. I registri base per il trasferimento di dati in virgola mobile che vengono utilizzati per generare l'indirizzo della memoria rimangono i registri interi. Il codice RISC-V per caricare due numeri in singola precisione dalla memoria, sommarli e poi memorizzare il risultato della somma potrebbe essere qualcosa del tipo:

```
flw    f0, 0(x10) // carica in f0 il numero in VM su 32 bit
flw    f1, 4(x10) // carica in f1 il numero in VM su 32 bit
fadd.s f2, f0,f1   // f2 = f0 + f1 in singola precisione
fsw    f2, 8(x10) // memorizza il numero in VM su 32 bit
                  // contenuto in f2
```

Un registro in singola precisione in realtà corrisponde alla metà inferiore di un registro per la doppia precisione. Si noti che il registro  $f_0$  *non* è bloccato sul valore 0 come il registro intero  $x_0$ .

La Figura 3.17 riassume la parte di architettura RISC-V relativa alla virgola mobile introdotta in questo capitolo; le istruzioni aggiunte per il supporto della virgola mobile sono evidenziate in blu. Le istruzioni in virgola mobile utilizzano lo stesso formato delle corrispondenti istruzioni intere. Le istruzioni di lettura adottano il formato di tipo I, le istruzioni di memorizzazione il formato di tipo S e le istruzioni aritmetiche il formato di tipo R.

### Operandi in virgola mobile del RISC-V

Nome	Esempio	Commenti
32 registri virgola mobile	$f_0, f_1, f_2, \dots, f_{31}$	I registri in virgola mobile del RISC-V possono contenere un numero in virgola mobile in precisione singola o in doppia precisione
$2^{61}$ parole doppie di memoria	Memoria[0], Memoria[8], ..., Memoria[18 446 744 073 709 551 608]	Si accede solo tramite le istruzioni di trasferimento dati. Il RISC-V utilizza l'indirizzamento al byte, per cui gli indirizzi sequenziali delle parole doppie differiscono di 8. La memoria mantiene le strutture dati, come i vettori, e il contenuto scaricato dai registri

### Istruzioni assembler RISC-V in virgola mobile

Categoria	Istruzione	Esempio	Significato	Commenti
Aritmetica	Somma VM singola prec.	fadd.s f0,f1,f2	$f_0 = f_1 + f_2$	somma VM (singola precisione)
	Sottrazione VM singola prec.	fsub.s f0,f1,f2	$f_0 = f_1 - f_2$	sottraz VM (singola precisione)
	Moltiplicaz VM singola prec.	fmul.s f0,f1,f2	$f_0 = f_1 \times f_2$	moltiplicaz VM (singola precisione)
	Divisione VM singola prec.	fdiv.s f0,f1,f2	$f_0 = f_1 / f_2$	divisione VM (singola precisione)
	Radice quadrata VM singola prec.	fsqrt.s f0,f1	$f_0 = \sqrt{f_1}$	Radice quadrata VM (singola precisione)
	Somma VM doppia prec.	fadd.d f0,f1,f2	$f_0 = f_1 + f_2$	somma VM (doppia precisione)
	Sottrazione VM doppia prec.	fsub.d f0,f1,f2	$f_0 = f_1 - f_2$	sottraz VM (doppia precisione)
	Moltiplicaz VM doppia prec.	fmul.d f0,f1,f2	$f_0 = f_1 \times f_2$	moltiplicaz VM (doppia precisione)
	Divisione VM doppia prec.	fdiv.d f0,f1,f2	$f_0 = f_1 / f_2$	divisione VM (doppia precisione)
	Radice quadrata VM doppia prec.	fsqrt.d f0,f1	$f_0 = \sqrt{f_1}$	radice quadrata VM (doppia precisione)

### Istruzioni assembler RISC-V in virgola mobile

Categoria	Istruzione	Esempio	Significato	Commenti
Confronto	Uguaglianza VM singola prec.	feq.s x5,f0,f1	x5 = 1 if (f0 == f1), else x5 = 0	Confronto in VM (precisione singola)
	Test di minoranza VM singola prec.	flt.s x5,f0,f1	x5 = 1 if (f0 < f1), else x5 = 0	Confronto in VM (precisione singola)
	Test di minore uguale VM singola prec.	fle.s x5,f0,f1	x5 = 1 if (f0 <= f1), else x5 = 0	Confronto in VM (precisione singola)
	Uguaglianza VM doppia prec.	feq.d x5,f0,f1	x5 = 1 if (f0 == f1), else x5 = 0	Confronto in VM (doppia precisione)
	Test di minoranza VM doppia prec.	flt.d x5,f0,f1	x5 = 1 if (f0 < f1), else x5 = 0	Confronto in VM (doppia precisione)
	Test di minore uguale VM doppia prec.	fle.d x5,f0,f1	x5 = 1 if (f0 <= f1), else x5 = 0	Confronto in VM (doppia precisione)
Trasferimento dati	Lettura di una parola in VM	flw f0,4(x5)	f0 = Memoria[x5+4]	Leggi da memoria dato in singola precisione
	Lettura di una parola doppia in VM	fld f0,8(x5)	f0 = Memoria[x5+8]	Leggi da memoria dato in doppia precisione
	Scrittura di una parola in VM	fsw f0,4(x5)	Memoria[x5+4] = f0	Trasferisci in memoria dato in singola precisione
	Scrittura di una parola doppia in VM	fsd f0,8(x5)	Memoria[x5+8] = f0	Trasferisci in memoria dato in doppia precisione

**Figura 3.17** L'architettura in virgola mobile del RISC-V introdotta finora. Potete trovare queste informazioni anche nella scheda tecnica riassuntiva del RISC-V in fondo al libro.

### Interfaccia hardware/software

Un problema che deve essere affrontato dai progettisti di calcolatori per supportare le operazioni in virgola mobile è decidere se utilizzare gli stessi registri impiegati dalle istruzioni intere o se aggiungere un insieme apposito di registri per i numeri in virgola mobile. Poiché i programmi normalmente eseguono le operazioni intere e quelle in virgola mobile su dati diversi, la separazione dei registri incrementa solo leggermente il numero di istruzioni necessarie per l'esecuzione di un programma. L'impatto maggiore è legato alla creazione di un insieme separato di istruzioni di trasferimento dati per spostarli tra i registri e la memoria.

I benefici che derivano dall'avere registri separati sono la disponibilità di un numero doppio di registri senza dover utilizzare un numero maggiore di bit all'interno delle istruzioni, la disponibilità di una banda di trasferimento doppia con i registri (in quanto vi sono insiemi separati di registri per le operazioni intere e per quelle in virgola mobile) e la possibilità di specializzare i registri per la virgola mobile; alcuni calcolatori, per esempio, convertono gli operandi di tutte le dimensioni esistenti nei registri in un unico formato interno.

### Traduzione di un programma C in virgola mobile nel codice assembler RISC-V

#### ESEMPIO

Convertiamo una data temperatura da gradi Fahrenheit a gradi Celsius utilizzando il seguente frammento di codice C:

```
float f2c (float fahr)
{
    return ((5,0f/9,0f) * (fahr - 32,0f));
}
```

(continua)

(continua)

Si supponga che l'argomento in virgola mobile fahr sia passato nel registro f10 e che il risultato debba andare in f10. Qual è il corrispondente codice assembler RISC-V?

Si supponga che il compilatore metta le tre costanti in virgola mobile in memoria, in maniera facilmente raggiungibile tramite il registro x3. Le prime due istruzioni caricano le costanti 5,0 e 9,0 nei registri in virgola mobile:

f2c:

```
flw f0, const5(x3) // f0 = 5,0f
flw f1, const9(x3) // f1 = 9,0f
```

Le due costanti vengono poi divise tra loro per ottenere la frazione 5,0/9,0:

```
fdiv.s f0, f0, f1 // f0 = 5,0f / 9,0f
```

(Molti compilatori dividerebbero 5,0 per 9,0 al momento della compilazione e salverebbero in memoria la sola costante 5,0/9,0, evitando così la divisione del numero durante l'esecuzione.) In seguito viene caricata la costante 32,0, poi sottratta da fahr (f10):

```
flw f1, const32(x3) // f1 = 32,0f
fsub.s f10, f10, f1 // f10 = fahr - 32,0f
```

Infine si moltiplicano i due risultati intermedi e si inserisce il prodotto in f10, che contiene così il risultato da restituire; infine si esce dalla procedura:

```
fmul.s f10, f0, f10 // f10 = (5,0f / 9,0f) * (fahr - 32,0f)
jalr x0, 0(x1) // ritorna
```

Esamineremo ora le operazioni in virgola mobile sulle matrici che si possono incontrare frequentemente nei programmi scientifici.

### Traduzione in assembler RISC-V di una procedura C in virgola mobile che utilizza matrici bidimensionali

La maggior parte dei calcoli in virgola mobile sono eseguiti in doppia precisione. Consideriamo ora la seguente operazione tra matrici  $C = C + A * B$ . Questa operazione viene chiamata in gergo *DGEMM*, che sta per *Double precision GEneral Matrix Multiply* (Moltiplicazione generale tra matrici, in doppia precisione). Vedremo altre versioni di DGEMM nel paragrafo 3.8 e quindi nei Capitoli 4, 5 e 6. Si assuma che C, A e B siano tutte matrici quadrate con 32 elementi per ciascuna dimensione.

```
void mm (double c[][], double a[][], double b[][])
{
    int i, j, k;
    for (i = 0; i < 32; i = i + 1)
        for (j = 0; j < 32; j = j + 1)
            for (k = 0; k < 32; k = k + 1)
                c[i][j] = c[i][j] + a[i][k] * b[k][j];
}
```

ESEMPIO

(continua)

(continua)

Gli indirizzi di partenza delle matrici sono parametri e quindi saranno contenuti in  $x10$ ,  $x11$  e  $x12$ . Si assuma che le variabili intere si trovino rispettivamente in  $x5$ ,  $x6$  e  $x7$ . Qual è il codice assembler RISC-V corrispondente?

**SOLUZIONE**

Si noti che  $c[i][j]$  viene utilizzato nel ciclo più interno. Poiché l'indice del ciclo è  $k$ , tale indice non influenza  $c[i][i]$  e quindi si può evitare di caricare e scaricare  $c[i][j]$  a ogni iterazione. Piuttosto, il compilatore carica  $c[i][j]$  in un registro, al di fuori del ciclo, accumula la somma dei prodotti di  $a[i][k]$  con  $b[k][j]$  in questo stesso registro e poi, al termine del ciclo più interno, memorizza la somma in  $c[i][j]$ . Si può rendere il codice più semplice utilizzando la pseudoistruzione del linguaggio assembler `li` (che carica una costante in un registro).

Il corpo della procedura inizia con il salvataggio del valore di fine ciclo, 32, in un registro temporaneo, e procede poi con l'inizializzazione delle tre variabili dei cicli `for`:

```
mm: ...
    li x28,32 // x28 = 32 (dimensione della riga/fine
                // del ciclo)
    li x5,0   // i = 0; inizializzazione primo ciclo
    ll:  li x6,0 // j = 0; punto di inizio secondo ciclo
    L2:  li x7,0 // k = 0; punto di inizio terzo ciclo
```

Per calcolare l'indirizzo di  $c[i][j]$ , occorre sapere come viene scritta in memoria una matrice  $32 \times 32$ . Come ci si potrebbe aspettare, la disposizione è la stessa di quella che avrebbero 32 vettori monodimensionali, ciascuno composto da 32 elementi. Quindi, la prima fase consiste nel saltare i primi "vettori monodimensionali", o righe, per arrivare alla posizione desiderata. Per fare questo si moltiplica l'indice della prima dimensione per l'ampiezza della riga, ossia 32. Poiché 32 è una potenza di 2, si può usare uno scorriamento:

```
slli x30,x5,5 // x30 = i * 25 (ampiezza della riga di c)
```

Poi si somma il secondo indice per selezionare l'elemento  $j$ -esimo della riga desiderata:

```
add x30,x30,x6 // x30 = i * ampiezza(riga) + j
```

Per trasformare il risultato in un indice riferito al byte, lo si moltiplica per la dimensione in byte degli elementi della matrice. Poiché in doppia precisione ogni elemento è costituito da 8 byte, possiamo scrivere il prodotto come uno scorriamento a sinistra di 3 posizioni:

```
slli x30,x30,3 // x30 = spiazzamento in byte di [i][j]
```

Si somma poi il valore così ottenuto all'indirizzo di partenza di  $c$ , ottenendo l'indirizzo di  $c[i][j]$ , e si carica quindi il numero in doppia precisione  $c[i][j]$  in  $f0$ :

```
add x30,x10,x30 // x30 = indirizzo in byte di c[i][j]
fld f0,0(x30)    // f0 = 8 byte di c[i][j]
```

Le cinque istruzioni successive sono virtualmente identiche alle ultime cinque: calcolano l'indirizzo di  $b[k][j]$  e caricano il corrispondente numero in doppia precisione.

(continua)

```
L3: slli x29,x7,5    // x29 = k * 25 (ampiezza riga di b)
    add x29,x29,x6   // x29 = k * ampiezza(riga) + j
    slli x29,x29,3    // x29 = spiazzamento in byte di [k][j]
    add x29,x12,x29  // x29 = indirizzo in byte di b[k][j]
    fld f1,0(x29)    // f1 = 8 byte di b[k][j]
```

Analogamente, le cinque istruzioni che seguono sono simili alle ultime cinque: calcolano l'indirizzo di  $a[i][k]$  e caricano il numero in doppia precisione corrispondente.

```
slli x29,x5,5    // x29 = i * 25 (ampiezza riga di a)
add x29,x29,x7   // x29 = i * ampiezza(riga) + j
slli x29,x29,3    // x29 = spiazzamento in byte di [i][k]
add x29,x11,x29  // x29 = indirizzo in byte di a[i][k]
fld f2,0(x29)    // f2 = a[i][k]
```

Ora che tutti i dati sono stati caricati, possiamo finalmente eseguire qualche operazione in virgola mobile. Possiamo moltiplicare gli elementi di  $a$  e  $b$  contenuti nei registri  $f2$  e  $f1$ , e accumulare la somma in  $f0$ .

```
fmul.d f1,f2,f1  // f1 = a[i][k] * b[k][j]
fadd.d f0,f0,f1  // f0 = c[i][j] + a[i][k] * b[k][j]
```

Il blocco finale del codice incrementa l'indice  $k$  e rimane all'interno del ciclo fino a che l'indice non raggiunge il valore 32. Quando lo raggiunge siamo arrivati alla fine del ciclo più interno e occorre memorizzare in  $c[i][j]$  la somma che abbiamo accumulato in  $f0$ .

```
addi x7,x7,1      // k = k + 1
bltu x7,x28,L3   // se (k < 32) vai a L3
fsd f0,0(x30)    // c[i][j] = f0
```

Analogamente, le quattro istruzioni seguenti incrementano la variabile indice del ciclo intermedio e del ciclo esterno, eseguendo una nuova iterazione se l'indice è inferiore a 32 e uscendo dal ciclo quando l'indice è uguale a 32.

```
addi x6,x6,1      // j = j+1
bltu x6,x28,L2   // se (j < 32) vai a L2
addi x5,x5,1      // i = i + 1
bltu x5,x28,L1   // se (i < 32) vai a L1
...
```

La Figura 3.20 mostra il codice assembler dell'x86 che implementa la versione leggermente differente dell'istruzione DGEMM riportata in Figura 3.19.

**Approfondimento.** Il C e molti altri linguaggi di programmazione utilizzano la disposizione degli elementi nelle matrici presentata nell'esempio. Questa viene chiamata *disposizione per righe*. Al contrario, il Fortran utilizza la *disposizione per colonne*, in cui la matrice viene memorizzata colonna per colonna.

**Approfondimento.** Un'altra ragione per avere registri separati per numeri interi e in virgola mobile è che i microprocessori negli anni '80 non avevano a disposizione abbastanza transistor per poter mettere l'unità in virgola mobile sullo stesso chip utilizzato per l'unità dell'aritmetica intera. Quindi le unità aritmetiche in virgola mobile, compresi i registri in virgola mobile, erano resi disponibili eventualmente su un secondo chip. Questi chip opzionali di accelerazione vengono chiamati *co-processori*. A partire dai primi anni '90, i microprocessori hanno integrato l'aritme-

tica in virgola mobile (e quasi tutto il resto) su un unico chip, e quindi il termine "coprocessore" oggi rientra tra i termini obsoleti, come "accumulatore" e "nucleo di memoria".

**Approfondimento.** Come accennato nel paragrafo 3.4, è più difficile accelerare la divisione che la moltiplicazione. Oltre alla tecnica SRT, un'altra tecnica per realizzare un divisore veloce è la *tecnica iterativa di Newton*, nella quale la divisione viene ricondotta alla ricerca dello 0 di una funzione che calcola il reciproco del divisore,  $1/c$ , che viene poi moltiplicato per il dividendo. Le tecniche iterative non consentono un adeguato arrotondamento del risultato, a meno che non venga calcolato un numero elevato di bit aggiuntivi. Questo problema veniva risolto da un chip della Texas Instrument, calcolando un reciproco estremamente preciso.

**Approfondimento.** Java abbraccia alla lettera l'IEEE 754 nella definizione dei tipi di dati e delle operazioni in virgola mobile. Quindi, il codice del primo esempio potrebbe appartenere a buon diritto anche a un metodo di una classe Java che converte Fahrenheit in Celsius. Il secondo esempio riportato sopra utilizza i vettori multidimensionali (matrici), che non vengono supportati esplicitamente in Java. Java consente di definire vettori di vettori, dove ogni vettore può avere una lunghezza diversa, a differenza dei vettori multidimensionali del C. Come per gli esempi del Capitolo 2, una versione Java di questo secondo esempio richiederebbe diverse istruzioni per controllare che vengano rispettati i limiti dei vettori e altre istruzioni per il calcolo della nuova lunghezza di una riga, dopo ogni accesso alla riga stessa. Inoltre, il codice dovrebbe controllare anche che il riferimento all'oggetto non sia nullo.

## Accuratezza dell'aritmetica

**Guardia:** il primo dei due bit aggiunti alla destra dell'ultimo bit durante i calcoli intermedi nelle operazioni in virgola mobile; è utilizzato per migliorare l'accuratezza dell'arrotondamento.

**Arrotondamento:** metodo che fa sì che il risultato in virgola mobile intermedio sia contenuto all'interno del formato in virgola mobile prescelto. Scopo del metodo è trovare il numero più vicino che possa essere rappresentato in quel formato.

A differenza dei numeri interi, per i quali è possibile rappresentare con esattezza ogni numero compreso tra il valore più piccolo e quello più grande, i numeri in virgola mobile costituiscono in genere un'approssimazione di numeri che in realtà non possono essere rappresentati. La ragione è che, per esempio, esistono infiniti numeri reali compresi tra 0 e 1 ma non più di 253 di questi possono essere rappresentati esattamente nella codifica in doppia precisione; la cosa migliore che si possa fare è trovare la rappresentazione più vicina al numero considerato. Per questo motivo, lo standard IEEE 754 offre diverse modalità di arrotondamento che permettono al programmatore di scegliere l'approssimazione desiderata.

L'arrotondamento sembra un'operazione abbastanza semplice, ma per eseguirlo in modo accurato occorre che l'hardware possa utilizzare dei bit aggiuntivi nei calcoli. Negli esempi precedenti siamo rimasti piuttosto vaghi sul numero di bit delle rappresentazioni intermedie dei numeri, ma è chiaro che se ogni risultato intermedio dovesse essere troncato al numero esatto di bit previsto per il risultato, non ci sarebbe alcuna possibilità di arrotondare in modo intelligente. Per questo lo standard IEEE 754 mantiene sempre 2 bit aggiuntivi alla destra della mantissa durante i calcoli intermedi, chiamati rispettivamente **guardia** e **arrotondamento**. Vediamo ora un esempio che illustra la loro importanza.

### Arrotondamento con cifre di guardia

#### ESEMPIO

Sommare  $2,56_{\text{dec}} \times 10^0$  e  $2,34_{\text{dec}} \times 10^2$  supponendo di disporre di tre cifre significative. Si arrotondi al numero decimale più vicino contenente tre cifre, prima con le cifre di guardia e arrotondamento, poi senza.

#### SOLUZIONE

(continua)

Innanzitutto si deve scalare il numero più piccolo a destra per allineare gli esponenti, così  $2,56_{\text{dec}} \times 10^0$  diventa  $0,0256_{\text{dec}} \times 10^2$ . Poiché ci sono le

(continua)

due cifre di guardia e di arrotondamento, si possono rappresentare anche le due cifre meno significative del primo operando quando si allineano gli esponenti: la cifra di guardia vale 5 e quella di arrotondamento 6. La somma è:

$$\begin{array}{r} 2,3400_{\text{dec}} \\ + 0,0256_{\text{dec}} \\ \hline 2,3656_{\text{dec}} \end{array}$$

La somma vale quindi  $2,3656_{\text{dec}} \times 10^2$ . Poiché le cifre da arrotondare sono due, i valori da 0 a 49 saranno arrotondati verso il basso e quelli da 51 a 99 saranno arrotondati verso l'alto, con il numero 50 che funge da separatore. Arrotondando la somma a tre cifre significative si ottiene  $2,37_{\text{dec}} \times 10^2$ .

Ripetendo i calcoli *senza* le cifre di guardia e arrotondamento, nei calcoli si perdono due cifre decimali. La nuova somma sarà quindi:

$$\begin{array}{r} 2,34_{\text{dec}} \\ + 0,02_{\text{dec}} \\ \hline 2,36_{\text{dec}} \end{array}$$

Il risultato finale sarà quindi  $2,36_{\text{dec}} \times 10^2$ , la cui ultima cifra è minore di una unità rispetto alla somma corretta, ottenuta in precedenza.

Poiché il caso peggiore per l'arrotondamento si ha quando il numero in virgola mobile da rappresentare si trova a metà strada tra due numeri rappresentabili, l'accuratezza della rappresentazione in virgola mobile si misura normalmente in termini di numero di bit sbagliati tra i bit meno significativi della mantissa; tale misura è nota come numero di **unità in ultima posizione o ulp** (*units in the last place*). Se un numero è diverso da quello da rappresentare per i 2 bit meno significativi, allora diciamo che è "diverso di 2 ulp". A meno che non si verifichi un'eccezione di overflow, underflow o di operazione non valida, lo standard IEEE 754 garantisce che il calcolatore utilizzi sempre numeri approssimati a meno di mezzo ulp.

**Approfondimento.** Benché l'esempio precedente richiedesse soltanto un bit aggiuntivo, la moltiplicazione potrebbe richiederne due. Se un prodotto binario ha il bit iniziale a 0, occorre normalizzarlo facendo scorrere a sinistra il prodotto di una posizione. Questo fa sì che il bit di guardia si sposti nella posizione del bit meno significativo del prodotto, lasciando al solo bit di arrotondamento il compito di assistere nell'arrotondamento accurato del prodotto.

Lo standard IEEE 754 prevede quattro modalità di arrotondamento: arrotondamento sempre al valore superiore (verso  $+\infty$ ), arrotondamento sempre al valore inferiore (verso  $-\infty$ ), troncamento e arrotondamento al numero pari più vicino. Quest'ultima modalità definisce che cosa fare se il numero si trova esattamente a metà tra due numeri rappresentabili. Per esempio, nel modello 730 della dichiarazione dei redditi la somma di 0,50 Euro viene arrotondata per eccesso all'Euro superiore, a beneficio di chi effettua la dichiarazione se si tratta di imposte già versate. Una soluzione più equa consisterebbe nell'arrotondare tale valore per eccesso in metà dei casi e per difetto nell'altra metà. Per i numeri a metà strada tra due numeri rappresentabili, lo standard IEEE 754 prevede di arrotondare il numero al numero successivo se il bit meno significativo è dispari e di troncare il numero se è pari. Questo metodo produce sempre uno 0 nel bit meno significativo, da cui l'origine del nome. Questa modalità è quella più utilizzata ed è l'unica supportata da Java.

Lo scopo dei bit aggiuntivi di arrotondamento è di consentire al calcolatore di ottenere lo stesso risultato che si otterebbe se i risultati intermedi fossero calcolati

**Unità in ultima posizione (ulp):** il numero di bit meno significativi del significando che possono essere diversi nel numero effettivamente rappresentato rispetto al numero vero.

**Bit di presenza (sticky bit):** un bit utilizzato nell'arrotondamento in aggiunta a quelli di guardia e arrotondamento. Esso è posto a 1 quando ci sono bit diversi da 0 alla destra del bit di arrotondamento.

**Moltiplicazione e somma integrate:** una singola istruzione in virgola mobile che esegue una moltiplicazione e poi una somma, ed effettua l'arrotondamento una volta sola, dopo l'addizione.

con precisione infinita e poi arrotondati. Per fare ciò, lo standard prevede un terzo bit in aggiunta a quelli di guardia e arrotondamento, che viene impostato a 1 ogni volta che ci sono dei bit non nulli alla destra del bit di arrotondamento. Questo bit (detto **bit di presenza**, o **sticky bit**) permette al calcolatore di apprezzare la differenza tra  $0,50 \dots 00_{dec}$  e  $0,50 \dots 01_{dec}$  quando viene eseguito l'arrotondamento.

Può succedere che il bit di presenza venga impostato a 1, per esempio, durante la somma, quando il numero più piccolo viene fatto scorrere verso destra. Si supponga di sommare  $5,01_{dec} \times 10^{-1}$  e  $2,34_{dec} \times 10^2$ . Anche utilizzando le cifre di guardia e arrotondamento finiremmo per sommare 0,0050 a 2,34, e la somma darebbe 2,3450. In questo caso il bit di presenza verrebbe impostato a 1, dato che ci sono bit diversi da 0 alla destra dell'ultima cifra. Senza il bit di presenza che ricorda che degli 1 della mantissa sono stati portati fuori, concluderemmo che la somma è uguale a 2,345000 ... 00 e arrotonderemmo al più vicino intero, cioè 2,34. Con il bit di presenza, invece, ricordiamo che il numero è più grande di 2,345000 ... 00, e quindi il numero verrà arrotondato a 2,35.

**Approfondimento.** Le architetture RISC-V, MIPS-64, PowerPC, SPARC64, AMD SSE5 e Intel AVX forniscono una singola istruzione per eseguire una moltiplicazione seguita da un'addizione su tre registri:  $a = a + (b \times c)$ . È chiaro che questa istruzione fornisce potenzialmente prestazioni superiori sulle comuni operazioni in virgola mobile. È altrettanto importante che, invece di eseguire due arrotondamenti successivi, dopo la moltiplicazione e dopo l'addizione, come succederebbe se fossero eseguite due istruzioni in sequenza, l'istruzione esegua una sola operazione di arrotondamento alla fine: un solo passaggio implica più precisione nel risultato. Operazioni di questo tipo, che richiedono un unico passaggio di arrotondamento, sono chiamate **moltiplicazione e somma integrate** (*fused multiply add*). Queste istruzioni sono state aggiunte nello standard IEEE 754-2008 (par. 3.11).

## Riepilogo

La sezione *Quadro di insieme* che segue rafforza il concetto di programma memorizzato introdotto nel Capitolo 2; il significato dell'informazione non può essere determinato semplicemente dall'analisi dei bit, dal momento che gli stessi bit possono rappresentare una pluralità di oggetti. Questo paragrafo mostra che l'aritmetica dei calcolatori ha precisione finita e può quindi non coincidere con l'aritmetica naturale. Per esempio, la rappresentazione in virgola mobile secondo lo standard IEEE 754:

$$(-1)^s \times (1 + \text{frazione}) \times 2^{(\text{Esponente} - \text{Polarizzazione})}$$

contiene quasi sempre un'approssimazione del numero reale. I sistemi di elaborazione devono cercare di minimizzare la distanza tra l'aritmetica del calcolatore e l'aritmetica del mondo reale, e i programmatore a volte devono essere consci delle implicazioni derivanti da tale approssimazione.

## QUADRO D'INSIEME

Le combinazioni di bit non hanno un significato intrinseco. Esse possono rappresentare numeri interi con segno o senza segno, numeri in virgola mobile, istruzioni ecc. Ciò che rappresentano dipende dall'istruzione che opera sui bit all'interno della parola. La differenza principale tra i numeri dei calcolatori e i numeri del mondo reale sta nel fatto che i numeri dei calcolatori hanno una dimensione limitata: attraverso il calcolo è possibile trovare numeri che sono troppo grandi o troppo piccoli per essere rappresentati in una parola. I programmatore devono ricordarsi di questi limiti e scrivere i loro programmi di conseguenza. ■

Tipi del C	Tipi di Java	Trasferimento dati	Operazioni
long long int	long	ld, sd	add, sub, addi, mul, mulh, mulhu, mulhsu, div, divu, rem, remu, and, andi, or, ori, xor, xori
unsigned long long int	-	ld, sd	add, sub, addi, mul, mulh, mulhu, mulhsu, div, divu, rem, remu, and, andi, or, ori, xor, xori
char	-	lb, sb	add, sub, addi, mul, div, divu, rem, remu, and, andi, or, ori, xor, xori
short	char	lh, sh	add, sub, addi, mul, div, divu, rem, remu, and, andi, or, ori, xor, xori
float	float	flw, fsw	fadd.s, fsub.s, fmul.s, fdiv.s, feq.s, flt.s, fle.s
double	double	fld, fsd	fadd.d, fsub.d, fmul.d, fdiv.d, feq.d, flt.d, fle.d

Nel capitolo precedente avevamo presentato i tipi di variabili del linguaggio C (vedi la sezione **Interfaccia hardware/software** del paragrafo 2.7). La tabella precedente mostra alcuni tipi di dati del C e di Java, assieme alle istruzioni di trasferimento dati e alle istruzioni che operano sui diversi tipi di dati introdotti nel Capitolo 2. Si noti che Java non considera gli interi senza segno.

## Interfaccia hardware/software

### Autovalutazione

Si supponga che esista un formato di codifica in virgola mobile su 16 bit dello standard IEEE 754-2008 con cinque bit per l'esponente. Qual è l'intervallo dei numeri rappresentabili?

- Da  $1,0000\ 00 \times 2^0$  a  $1,1111\ 1111\ 11 \times 2^{31}, 0$
- Da  $\pm 1,0000\ 0000\ 0 \times 2^{-14}$  a  $\pm 1,1111\ 1111\ 1 \times 2^{15}, \pm 0, \pm \infty, \text{NaN}$
- Da  $\pm 1,0000\ 0000\ 00 \times 2^{-14}$  a  $\pm 1,1111\ 1111\ 11 \times 2^{15}, \pm 0, \pm \infty, \text{NaN}$
- Da  $\pm 1,0000\ 0000\ 00 \times 2^{-15}$  a  $\pm 1,1111\ 1111\ 11 \times 2^{14}, \pm 0, \pm \infty, \text{NaN}$

**Approfondimento.** Per consentire confronti che includano NaN, lo standard prevede le due opzioni di comparazione *ordinata* e *non ordinata*. Il RISC-V non fornisce istruzioni per la comparazione non ordinata, ma un attento ordinamento dei confronti produce lo stesso effetto (Java non supporta i confronti non in ordine.)

Nel tentativo di utilizzare tutti i bit di precisione in ciascuna operazione in virgola mobile, lo standard permette di rappresentare alcuni numeri in forma non normalizzata. Piuttosto che accettare un buco tra lo 0 e il più piccolo numero normalizzato, lo standard IEEE ammette dei *numeri denormalizzati* (detti anche *denorm* o *subnormal*). Questi hanno lo stesso esponente dello 0, ma una mantissa non nulla. Questo meccanismo permette a un numero di diminuire gradualmente fino a 0, realizzando un meccanismo chiamato *underflow graduale*. Per esempio, il numero normalizzato più piccolo in singola precisione è:

$$1,000\ 0000\ 0000\ 0000\ 000_{\text{due}} \times 2^{-126}$$

ma il numero denormalizzato più piccolo in singola precisione è:

$$0,0000\ 0000\ 0000\ 0000\ 0001_{\text{due}} \times 2^{-126}, \text{ o } 1,0_{\text{due}} \times 2^{-149}$$

Per la doppia precisione, l'intervallo dei numeri denormalizzati va da  $1,0 \times 2^{-1022}$  a  $1,0 \times 2^{-1074}$ .

La gestione degli operandi denormalizzati ha procurato molti mal di testa ai progettisti dei circuiti per le operazioni in virgola mobile, sempre alla ricerca di nuovi modi per costruire unità più veloci. Per questo motivo molti calcolatori generano un'eccezione se un operando è denormalizzato, lasciando al software il compito di completare l'operazione.

Benché le implementazioni software siano perfettamente accettabili, le loro minori prestazioni hanno diminuito la diffusione dei numeri denormalizzati nel software in virgola mobile portabile. D'altro canto, se i programmatori non tenessero conto dei numeri denormalizzati, i loro programmi potrebbero generare errori inaccettabili.

### 3.6 Parallelismo e aritmetica dei calcolatori: parallelismo a livello di parola

Dato che ogni microprocessore contenuto in un cellulare, calcolatore portatile o tablet possiede, per definizione, un proprio monitor grafico, era inevitabile che all'aumentare del numero di transistor le architetture offrissero anche il supporto per le operazioni della grafica.

Molti sistemi grafici inizialmente utilizzavano 8 bit per rappresentare ciascuno dei 3 colori primari e 8 bit per definire la posizione di un pixel. Quando sono stati introdotti gli altoparlanti e i microfoni per le teleconferenze e i videogiochi, è stato aggiunto il supporto anche per l'audio. I campioni di un segnale audio richiedono più di 8 bit di precisione: 16 bit sono sufficienti.

I microprocessori possiedono un supporto speciale per i byte e le mezze parole, così che occupino meno spazio quando vengono salvati in memoria (vedi par. 2.9). Tuttavia, dato che le operazioni aritmetiche sono rare su questi tipi di dati nei programmi tipici che fanno uso degli interi, il supporto era limitato al trasferimento dei dati da e verso la memoria. I progettisti hanno capito che molte delle operazioni sulla grafica e sull'audio vengono eseguite su vettori di dati e non su singoli dati. Spezzando opportunamente la catena dei riporti degli addizionatori a 128 bit, un processore può sfruttare il **parallelismo** per eseguire simultaneamente operazioni su vettori di sedici numeri da 8 bit, otto da 16 bit, quattro da 32 bit o due da 64 bit.

Il costo di questo partizionamento degli addizionatori è piccolo mentre l'accelerazione può essere molto grande.

Dato che il parallelismo in questo caso viene applicato a una parola ampia, queste estensioni delle architetture vengono classificate come *parallelismo a livello di parola* (*subword parallelism*). Vengono classificate anche nell'insieme più ampio del *parallelismo a livello di dati*. Le architetture con queste estensioni vengono anche chiamate vettoriali o SIMD (*Single Instruction Multiple Data*, singola istruzione per dati multipli; vedi par. 6.6). Il rapido aumento delle applicazioni multimediali ha guidato l'introduzione di istruzioni aritmetiche che supportino operazioni su dati più piccoli di una parola che possono essere eseguite in parallelo. Al momento della scrittura di questo libro il RISC-V non possiede istruzioni aggiuntive per sfruttare il parallelismo a livello di parola, ma il prossimo paragrafo presenta degli esempi reali di architetture di questo tipo.



### 3.7 Un caso reale: le estensioni SIMD per lo streaming e le estensioni avanzate dell'x86 per il calcolo vettoriale

Le istruzioni originali dell'estensione MMX (*MultiMedia eXtension*) dell'x86 comprendevano istruzioni che operavano su piccoli vettori di interi. Successivamente l'SSE (*Streaming SIMD Extension*) ha messo a disposizione istruzioni che operavano su piccoli vettori di numeri in virgola mobile a precisione singola. Nel Capitolo 2 abbiamo visto che nel 2001 Intel ha aggiunto 144 istruzioni alla sua architettura, all'interno dell'estensione SSE2, con i relativi registri e operazioni in virgola mobile a doppia precisione. Fanno parte di questa estensione otto registri a 64 bit che possono essere utilizzati per gli operandi in virgola mobile. AMD ha espanso il numero di questi registri a 16 nell'AMD64, chiamato XMM, che Intel ha rinominato EM64T nelle sue estensioni. La Figura 3.18 riassume le istruzioni delle estensioni SSE e SSE2.

Oltre a permettere di inserire in un registro numeri in singola o doppia precisione, Intel consente di compattare all'interno dello stesso registro SSE2 a 128 bit più operandi in virgola mobile: quattro in singola precisione o due in doppia precisione. Quindi, i 16 registri in virgola mobile del SSE2 sono di fatto ampi 128 bit. Se gli operandi possono essere scritti in memoria allineati come dati su 128 bit, allora le istruzioni di trasferimento a 128 bit consentono di caricare e memorizzare più operandi per ogni singola istruzione. Questo formato compatto degli operandi in virgola mobile è supportato da operazioni aritmetiche, che possono operare contemporaneamente su 4 operandi in singola precisione (PS) o su 2 in doppia precisione (PD).

Nel 2011 Intel ha raddoppiato ancora l'ampiezza dei registri, ora denominati YMM, nella nuova estensione *Advanced Vector Extension* (AVX). Perciò oggi una singola istruzione può lanciare otto operazioni in virgola mobile su 32 bit o quattro su 64 bit. Le vecchie istruzioni SSE e SSE2 continuano a operare sui

Trasferimento dati	Aritmetica	Comparazione
MOV{AU}{SS PS SD PD} xmm, {mem xmm}	ADD{SS PS SD PD} xmm, {mem xmm}	CMP{SS PS SD PD}
	SUB{SS PS SD PD} xmm, {mem xmm}	
MOV{HL}{PS PD} xmm, {mem xmm}	MUL{SS PS SD PD} xmm, {mem xmm}	
	DIV{SS PS SD PD} xmm, {mem xmm}	
	SQRT{SS PS SD PD} {mem xmm}	
	MAX{SS PS SD PD} {mem xmm}	
	MIN{SS PS SD PD} {mem xmm}	

**Figura 3.18** Le istruzioni in virgola mobile delle estensioni SSE/SSE2 dell'x86. xmm indica che un operando è un registro di 128 bit SSE2, mentre [mem|xmm] indica che il secondo operando si trova in memoria o in uno dei registri SSE2. La tabella utilizza delle espressioni per mostrare le varianti delle istruzioni. Quindi MOV [AU] {SS|PS|SD|PD} rappresenta le otto istruzioni MOVASS, MOVAPS, MOVASD, MOVAPD, MOVSUSS, MOVUPS, MOVUSD e MOVUPD. Utilizziamo le parentesi quadre [] per mostrare le alternative espresse da lettere singole; A significa che l'operando di 128 bit è allineato in memoria; U significa che l'operando di 128 bit non è allineato in memoria; H richiede il trasferimento della metà superiore di un operando di 128 bit e L richiede il trasferimento della metà inferiore. Varianti delle operazioni di base espresse da più di una lettera sono indicate con le parentesi graffe: { SS } sta per singola precisione scalare in virgola mobile (*Scalar Single*), ovvero un operando di 32 bit in un registro a 128 bit; { PS } sta per singola precisione impacchettato in virgola mobile (*Packed Single*), ovvero quattro operandi di 32 bit in un registro a 128 bit; { SD } sta per doppia precisione scalare (*Scalar Double*) impacchettata in virgola mobile, ovvero un operando di 64 bit in un registro a 128 bit. { PD } sta per doppia precisione impacchettata in virgola mobile (*Packed Double*), ovvero due operandi di 64 bit in un registro a 128 bit.

128 bit meno significativi dei registri YMM. Per passare dalle operazioni su 128 bit a quelle su 256 bit è sufficiente aggiungere il prefisso "v" (vettore) alle istruzioni SSE2 dell'assembler Intel e utilizzare i registri YMM al posto dei registri XMM. Per esempio, l'istruzione SSE2 che esegue la moltiplicazione di due coppie di numeri a 64 bit in virgola mobile:

`addpd %xmm0, %xmm4`

diventa

`vaddpd %ymm0, %ymm4`

che comanda la moltiplicazione di quattro coppie di numeri in virgola mobile su 64 bit. Intel ha annunciato di avere pianificato di ampliare i registri AVX prima a 512 bit e successivamente a 1024 bit nelle versioni più recenti dell'architettura x86.

**Approfondimento.** L'estensione AVX ha aggiunto anche istruzioni su tre registri. Per esempio si può richiedere il salvataggio del risultato della somma su un terzo registro attraverso l'istruzione vaddpd:

`vaddpd %ymm0, %ymm1, %ymm4 // %ymm4 = %ymm0 + %ymm1`

invece di utilizzare solo due registri:

`addpd %xmm0, %xmm4 // %xmm4 = %xmm4 + %xmm0`

(A differenza dei RISC-V, negli x86 il registro di destinazione è scritto a destra.) L'utilizzo di tre registri può ridurre il numero di registri e di istruzioni richieste per i calcoli.

### 3.8 Come andare più veloci: il parallelismo a livello di parola applicato alla moltiplicazione di matrici

Per mostrare l'impatto del parallelismo a livello di parola sulle prestazioni, eseguiamo lo stesso codice su un Core i7 di Intel prima senza utilizzare l'estensione AVX e poi utilizzandola. La Figura 3.19 mostra la versione C non ottimizzata della procedura che segue il prodotto di due matrici. Come abbiamo visto nel paragrafo 3.5, questa procedura viene chiamata comunemente *DGEMM*, che sta per *Double precision GEneral Matrix Multiply*. A partire da questa edizione, abbiamo aggiunto un nuovo paragrafo intitolato *Come andare più veloci*, che mostra i miglioramenti delle prestazioni che si ottengono adattando il software

```

1. void dgemm (int n, double* A, double* B, double* C)
2. {
3.     for (int i = 0; i < n; ++i)
4.         for (int j = 0; j < n; ++j)
5.     {
6.         double cij = C[i+j*n]; /* cij = C[i][j] */
7.         for(int k = 0; k < n; k++)
8.             cij += A[i+k*n] * B[k+j*n]; /* cij += A[i][k]*B[k][j] */
9.         C[i+j*n] = cij; /* C[i][j] = cij */
10.    }
11. }
```

**Figura 3.19** Procedura C non ottimizzata che implementa la moltiplicazione di matrici in doppia precisione, nota anche come DGEMM (*Double precision GEneral Matrix Multiply*). A differenza delle matrici bidimensionali utilizzate nel paragrafo 3.5, la procedura lavora qui su versioni monodimensionali (vettorizzate) delle matrici C, A e B, per cui viene passata alla procedura una sola dimensione, n. Inoltre viene utilizzata l'aritmetica basata sugli indirizzi (e non sui valori) per ottenere prestazioni migliori. I commenti si riferiscono alla versione che utilizza le matrici per rendere il codice più leggibile.

all'hardware, in questo caso la versione Sandy Bridge del microprocessore Core i7 di Intel. In questo paragrafo, contenuto nei Capitoli 3, 4, 5 e 6, descriveremo come migliorare le prestazioni della procedura DGEMM utilizzando le idee descritte nei capitoli.

La Figura 3.20 mostra il codice assembler dell'x86 generato dalla compilazione del ciclo interno di Figura 3.19. Le cinque istruzioni in virgola mobile iniziano con la lettera **v** come le istruzioni AVX, ma utilizzano i registri XMM invece dei registri YMM, e contengono **sd** nel loro nome, che sta per scalare in doppia precisione. Definiremo fra poco le corrispondenti istruzioni che sfruttano il parallelismo a livello di parola.

Probabilmente i progettisti dei compilatori riusciranno in futuro a generare codice assembler di alta qualità che sfrutti le istruzioni AVX dell'x86; per ora dobbiamo "barare" utilizzando le funzioni intrinseche del C per dire al compilatore più o meno esattamente come produrre un buon codice compilato. La Figura 3.21 mostra la versione migliorata del codice AVX mostrato in Figura 3.19 prodotto dal compilatore Gnu. La Figura 3.22 mostra il corrispondente codice x86 commentato, prodotto dal compilatore gcc utilizzando l'ottimizzazione di livello -O3.

La dichiarazione alla linea 6 di Figura 3.21 utilizza il tipo di dati **\_mm256d**, che informa il compilatore che la variabile conterrà 4 numeri in virgola mobile in doppia precisione. La funzione intrinseca del C **\_mm256\_load\_pd()** alla linea 6 utilizza un'istruzione AVX per caricare 4 numeri in virgola mobile in

```

1. vmovsd (%r10),%xmm0          // Caricamento di 1 elemento di C in %xmm0
2. mov    %rsi,%rcx             // Registro %rcx = %rsi
3. xor    %eax,%eax            // Registro %eax = 0
4. vmovsd (%rcx),%xmm1          // Caricamento di 1 elemento di B in %xmm1
5. add    $0x1,%rax             // Registro %rcx = %rcx + %r9
6. vmulsd (%r8,%rax,8),%xmm1,%xmm1 // Moltiplicare %xmm1 con elemento di A
7. add    $0x1,%rax             // Registro %rax = %rax + 1
8. cmp    %eax,%edi             // Comparazione di %eax e %edi
9. vaddsd %xmm1,%xmm0,%xmm0     // Somma %xmm1, %xmm0
10. jg   30 <dgemm+0x30>        // Salta se %eax > %edi
11. add    $0x1,%r11             // Registro %r11 = %r11 + 1
12. vmovsd %xmm0,(%r10)          // Scrittura di %xmm0 nell'elemento di C

```

**Figura 3.20** Codice assembler x86 generato dal compilatore per il ciclo interno dalla procedura in Figura 3.19. Anche se i dati sono su 64 bit, il compilatore utilizza la versione AVX delle istruzioni invece della versione SSE2, presumibilmente per potere utilizzare le istruzioni di trasferimento dati su tre registri invece che su due (vedi la sezione *Approfondimento* del paragrafo 3.7).

```

1. //include <x86intrin.h>
2. void dgemm (size_t n, double* A, double* B, double* C)
3. {
4.   for (int i = 0; i < n; i+=4)
5.     for (int j = 0; j < n; j++){
6.       _mm256d c0 = _mm256_load_pd(C+i+j*n); /* c0 = C[i][j] */
7.       for(int k = 0; k < n; k++)
8.         c0 = _mm256_add_pd(c0, /* c0 += A[i][k]*B[k][j] */
9.                           _mm256_mul_pd(_mm256_load_pd(A+i+k*n),
10.                           _mm256_broadcast_sd(B+k+j*n)));
11.       _mm256_store_pd(C+i+j*n, c0); /* C[i][j] = c0 */
12.     }
13. }

```

**Figura 3.21** Versione C della procedura DGEMM ottimizzata utilizzando le funzioni intrinseche del C per generare le istruzioni assembler AVX che implementano il parallelismo a livello di parola nell'x86. Il codice assembler corrispondente al ciclo interno di questa procedura C è riportato in Figura 3.22.

```

1. vmovapd    (%r11),%ymm0      // Caricamento di 4 elementi di C in %ymm0
2. mov        %rbx,%rcx       // Registro %rcx = %rbx
3. xor        %eax,%eax      // Registro %eax = 0
4. vbroadcastsd (%rax,%r8.1),%ymml // Creazione di 4 copie dell'elemento di B
5. add        $0x8,%rax       // Registro %rax = %rax + 8
6. vmulpd    (%rcx),%ymml,%ymml // Moltiplicaz parallela %ymml con 4 elem di A
7. add        %r9,%rcx       // Registro %rcx = %rcx + %r9
8. cmp        %r10,%rax       // Comparazione di %r10 e %rax
9. vaddpd    %ymml,%ymm0,%ymm0 // Somma parallela %ymml, %ymm0
10. jne       50 <dgemm+0x50>   // Salta se %r10 non è uguale a %rax
11. add        $0x1,%esi      // Registro %esi = %esi + 1
12. vmovapd    %ymm0,(%r11)    // Scrittura di %ymm0 in 4 elementi di C

```

**Figura 3.22** Codice assembler x86 generato dal compilatore per il corpo del ciclo interno della procedura ottimizzata di Figura 3.21. Si noti la somiglianza con il codice in Figura 3.20. Le differenze principali sono che le operazioni in virgola mobile utilizzano ora i registri YMM e che le istruzioni sono nella versione “pd”, parallele in doppia precisione, invece che “sd”, scalari in doppia precisione.

doppia precisione in parallelo (`_pd`) dalla matrice C in `c0`. L'indirizzo `C+i+j*n` calcolato alla linea 6 corrisponde all'elemento `C[i+j*n]`. Simmetricamente, l'ultimo passo alla linea 11 utilizza la funzione intrinseca `_mm256_store_pd()` per trasferire 4 numeri in virgola mobile in doppia precisione da `c0` alla matrice C. Dato che a ogni iterazione ci spostiamo di 4 elementi, il ciclo `for` esterno alla linea 4 incrementa i di 4 invece che di 1 come alla linea 3 di Figura 3.19.

All'interno dei cicli, alla linea 9, prima carichiamo 4 elementi di A utilizzando `_mm256_load_pd()`. Per moltiplicare questi 4 elementi per un elemento di B, utilizziamo alla linea 10 dapprima la funzione intrinseca `_mm256_broadcast_sd()`, che fa 4 copie identiche di un numero scalare in doppia precisione, un elemento di B in questo caso, in uno dei registri YMM. Utilizziamo poi la funzione intrinseca `_mm256_mul_pd()` alla linea 9 per effettuare le moltiplicazioni dei 4 numeri in parallelo. Infine utilizziamo `_mm256_add_pd()`, alla linea 8, che somma i 4 prodotti alle 4 somme contenute in `c0`.

La Figura 3.22 mostra il codice x86 prodotto dal compilatore per il corpo interno del ciclo. Potete vedere le cinque istruzioni AVX (iniziano tutte con v e quattro su cinque utilizzano pd che sta per doppia precisione), che corrispondono alle cinque funzioni intrinseche menzionate sopra. Il codice è molto simile a quello in Figura 3.20; entrambi utilizzano 12 istruzioni, le istruzioni intere sono praticamente identiche (ma lavorano su registri diversi) e la differenza tra le istruzioni in virgola mobile è generalmente soltanto nel passaggio dalla *forma scalare in doppia precisione (scalar double, sd)* utilizzata con i registri XMM alla forma impaccata in doppia precisione (*packed double, pd*) utilizzata con i registri YMM. L'unica eccezione è alla linea 4 della Figura 3.22. Ciascun elemento di A deve essere moltiplicato con un elemento di B. Una soluzione è inserire 4 copie identiche degli elementi di B su 64 bit, fianco a fianco in un registro YMM a 256 bit; questa operazione viene eseguita dall'istruzione `vbroadcastsd`.

Per matrici di dimensioni 32 per 32, la versione non ottimizzata di DGEMM di Figura 3.19, viene eseguita a 1,7 GFLOPS (*Giga Floating point Operations Per Second*) su uno dei core del Core i7 di Intel (Sandy Bridge) a 2,6 GHz. La versione ottimizzata del codice di Figura 3.21 viene eseguita a 6,4 GFLOPS. La versione AVX è perciò 3,85 volte più veloce, molto vicino all'incremento di un fattore 4 che ci si può aspettare eseguendo il quadruplo delle operazioni nello stesso tempo sfruttando il **parallelismo** a livello di parola.



**Approfondimento.** Abbiamo già accennato nell'Approfondimento del paragrafo 1.6 che Intel offre la modalità Turbo per fare funzionare il microprocessore a una frequenza più elevata fino a che non si sia scaldato troppo. L'Intel Core i7 Sandy

Bridge può aumentare la sua frequenza da 2,6 GHz a 3,3 GHz in modalità Turbo. I risultati riportati sopra sono stati ottenuti con la modalità Turbo non attivata; se la attiviamo, otteniamo un miglioramento delle prestazioni pari al rapporto tra le frequenze di clock:  $3,3/2,6 = 1,27$ , ovvero 2,1 GFLOPS per il codice DGEMM non ottimizzato e 8,1 GFLOPS per il codice ottimizzato con l'estensione AVX. La modalità Turbo funziona particolarmente bene quando si utilizza uno solo degli otto core di cui è dotato un microprocessore quale l'i7, perché in questo caso consentiamo al singolo core di utilizzare molta più potenza della sua quota nominale dato che gli altri core sono inattivi.

## 3.9 Errori e trabocchetti

Gli errori e i trabocchetti riguardanti l'aritmetica hanno generalmente origine dalla differenza che esiste tra la precisione limitata dell'aritmetica dei calcolatori e la precisione illimitata dell'aritmetica naturale.

*Errore: come un'istruzione di scorrimento a sinistra può sostituire una moltiplicazione intera per una potenza di 2, così uno scorrimento a destra coincide con una divisione intera per una potenza di 2.*

Si ricordi che un numero binario  $x$ , per il quale  $x_i$  indica l' $i$ -esimo bit, rappresenta il numero:

$$\dots + (x^3 \times 2^3) + (x^2 \times 2^2) + (x^1 \times 2^1) + (x^0 \times 2^0)$$

Lo scorrimento a destra dei bit del numero  $c$  di  $n$  posizioni sembra coincidere con la divisione per  $2^n$ . E questo è vero per i numeri interi senza segno. Il problema nasce con i numeri interi dotati di segno. Per esempio, si supponga di voler dividere  $-5_{dec}$  per  $4_{dec}$ ; il quoziente dovrebbe essere  $-1_{dec}$ . La rappresentazione in complemento a 2 di  $-5_{dec}$  è:

$$11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 111111011_{due}$$

Eseguendo lo scorrimento a destra di due posizioni, si dovrebbe ottenere la divisione per  $4_{dec}$  ( $2^2$ ):

$$00111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111110_{due}$$

ma il risultato è chiaramente errato, dal momento che troviamo uno 0 nel bit di segno: il numero così ottenuto mediante lo scorrimento a destra è in realtà  $4\ 611\ 686\ 018\ 427\ 387\ 902_{dec}$  e non  $-1_{dec}$ .

Una soluzione potrebbe consistere in un'operazione di scorrimento a destra aritmetica che estende il bit di segno, anziché inserire a sinistra della parola degli 0. Un'operazione di scorrimento aritmetico a destra di 2 bit eseguita su  $-5_{dec}$  produrrebbe:

$$11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111111\ 11111110_{due}$$

Il risultato sarebbe  $-2_{dec}$  anziché  $-1_{dec}$ ; vicino al risultato desiderato, ma ancora non esatto.

*La matematica può quindi essere definita come quella materia a proposito della quale non sappiamo mai di che cosa stiamo parlando, né se ciò che stiamo dicendo sia vero.*

Bertrand Russel, *Recent Words on the Principles of Mathematics*, 1901

**Trabocchetto:** l'addizione in virgola mobile non è associativa.

L'associatività vale per una sequenza di addizioni di numeri interi in complemento a 2, anche se si verifica un overflow. Tuttavia, non vale per le addizioni di numeri in virgola mobile perché sono un'approssimazione dei numeri reali e perché l'aritmetica sui numeri in virgola mobile ha una precisione limitata. Dato che i numeri in virgola mobile possono rappresentare un insieme molto

ampio di numeri, si possono verificare dei problemi quando vengono sommati due numeri molto grandi, di segno opposto, e viene poi sommato al risultato un numero piccolo. Per esempio, vediamo se  $c + (a + b) = (c + a) + b$ . Si supponga che  $c = -1,5_{\text{dec}} \times 10^{38}$ ,  $a = 1,5_{\text{dec}} \times 10^{38}$  e  $b = 1,0$ , e che siano codificati tutti in singola precisione.

$$\begin{aligned} c + (a + b) &= -1,5_{\text{dec}} \times 10^{38} + (1,5_{\text{dec}} \times 10^{38} + 1,0) \\ &= -1,5_{\text{dec}} \times 10^{38} + (1,5_{\text{dec}} \times 10^{38}) \\ &= 0,0 \end{aligned}$$

$$\begin{aligned} c + (a + b) &= (-1,5_{\text{dec}} \times 10^{38} + 1,5_{\text{dec}} \times 10^{38}) + 1,0 \\ &= (0,0_{\text{dec}}) + 1,0 \\ &= 1,0 \end{aligned}$$

Poiché i numeri in virgola mobile hanno una precisione limitata e che il risultato di un'operazione può essere un'approssimazione di un numero reale, il numero  $1,5_{\text{dec}} \times 10^{38}$  è tanto più grande di  $1,0_{\text{dec}}$  che  $1,5_{\text{dec}} \times 10^{38} + 1,0$  è ancora uguale a  $1,5_{\text{dec}} \times 10^{38}$ . Questo è il motivo per cui la somma di  $c$ ,  $a$ , e  $b$  risulta 0,0 o 1,0 a seconda dell'ordine in cui vengono eseguite le somme in virgola mobile, e quindi l'addizione in virgola mobile *non* gode della proprietà associativa.

*Errore: le strategie per l'esecuzione parallela che funzionano per numeri interi funzionano anche per numeri in virgola mobile.*

I programmi vengono tipicamente scritti per essere eseguiti in modo sequenziale prima di essere riscritti per essere eseguiti in modo concorrente, perciò è naturale chiedersi: "le due versioni danno lo stesso risultato?". Se la risposta fosse no, dovremmo supporre che nella versione parallela c'è un baco che deve essere identificato.

Questo modo di ragionare presuppone che l'aritmetica dei calcolatori non influenzi il risultato quando l'esecuzione passa da sequenziale a parallela; cioè, si deve ottenere lo stesso risultato se si utilizza un processore o 1000 processori per sommare un milione di numeri tra loro. Tale assunzione è vera per l'aritmetica dei numeri interi in complemento a 2, perché in questo caso la somma gode della proprietà associativa. Non è vera invece per la somma di numeri in virgola mobile, perché in questo caso la proprietà associativa non vale.

Questo errore è ancora più delicato sui calcolatori paralleli nei quali lo scheduler del sistema operativo può decidere di utilizzare un numero diverso di processori a seconda degli altri programmi già in esecuzione. Utilizzando un numero diverso di processori, gli addendi possono essere sommati in ordine diverso e il risultato della somma può risultare leggermente diverso nonostante il codice e i dati in virgola mobile siano identici, cosa che può sconcertare i programmati di codice parallelo che non conoscono questo problema.

Perciò i programmati di codice parallelo devono verificare se i risultati di una somma di numeri in virgola mobile è attendibile anche se può non essere lo stesso risultato della somma sequenziale. La disciplina che analizza questi problemi si chiama *analisi numerica* e viene trattata in libri di testo dedicati. Questa preoccupazione è uno dei motivi della diffusione delle librerie di calcolo numerico quali LAPACK e ScaLAPACK, che sono state validate per il calcolo sia sequenziale sia parallelo.

*Errore: solo i matematici teorici si preoccupano dell'accuratezza dei calcoli in virgola mobile.*

I titoli dei giornali del novembre 1994 dimostrano che questa affermazione è sbagliata (Figura 3.23). Riportiamo di seguito la storia che sta dietro a quei titoli.



**Figura 3.23** Un campione degli articoli di giornali e riviste del novembre 1994, tra cui figurano il *New York Times*, il *San Jose Mercury News*, il *San Francisco Chronicle* e *Infoworld*. L'errore nella divisione in virgola mobile del Pentium entrò anche nella "Top 10 List" del *David Letterman Late Show* in televisione. L'Intel spese circa 300 milioni di dollari per rimpiazzare i chip che contenevano il baco.

Il Pentium utilizzava un algoritmo standard di divisione in virgola mobile che, a ogni passo, genera diversi bit del quoziente, utilizzando i bit più significativi del divisore e del dividendo per prevedere i 2 bit successivi del quoziente. I valori previsti derivano da una tabella predefinita che contiene -2, -1, 0, +1 o +2. Il valore previsto viene moltiplicato per il divisore e sottratto al resto attuale per generare un nuovo valore del resto.

Come nella divisione senza ripristino, se il valore generato al passo precedente produce un resto troppo grande, il resto parziale viene ridotto in un passo successivo.

A quanto pare, c'erano cinque elementi della tabella predefinita dell'80486 ai quali gli ingegneri di Intel pensavano che nessuno avrebbe mai tentato di accedere, e la PLA che implementava la tabella predefinita venne quindi ottimizzata nel Pentium in modo da restituire 0 anziché 2 in queste situazioni. Ma l'Intel si sbagliava: mentre i primi 11 bit erano sempre corretti, si presentavano occasionalmente degli errori nei bit 12-52, ossia nelle cifre decimali dalla quarta alla quindicesima.

Un professore di matematica del Lynchburg College in Virginia, Thomas Nicely, scoprì il baco nel settembre 1994. Dopo aver chiamato il supporto tecnico di Intel e non aver ricevuto alcuna risposta ufficiale, annunciò su Internet la sua scoperta; all'annuncio seguì un articolo su una rivista che forzò Intel a rilasciare un comunicato stampa in cui definì il problema un "anomalia" che avrebbe potuto riguardare solamente i matematici teorici, mentre per l'utente medio di un foglio elettronico questo baco avrebbe potuto generare un errore ogni 27 000 anni.

Tuttavia, poco dopo, la divisione ricerca di IBM contestò i calcoli di Intel e ipotizzò che il baco avrebbe potuto produrre un errore ogni 24 giorni per l'uten-

te medio di un foglio elettronico. Intel fu costretta a capitolare e il 21 dicembre 1994 diffuse il seguente comunicato:

*Noi di Intel desideriamo scusarci sinceramente per il modo in cui abbiamo recentemente reso pubblico il malfunzionamento del processore Pentium. Il simbolo "Intel Inside" significa che il vostro calcolatore possiede un microprocessore che non è secondo a nessuno in termini di qualità e prestazioni. Migliaia di dipendenti di Intel lavorano con grande impegno per far sì che questo sia vero. Ma nessun microprocessore è perfetto. Ciò che Intel continua a credere è che un problema, dal punto di vista tecnico considerato di second'ordine, abbia iniziato ad avere una propria vicenda autonoma. Benché sosteniamo fermamente la qualità dell'attuale versione del processore Pentium, riconosciamo che molti utenti sono preoccupati. Vogliamo fugare queste preoccupazioni. Intel sostituirà l'attuale versione del processore Pentium con la versione aggiornata, in cui il difetto nella divisione in virgola mobile è stato eliminato, a chiunque lo richieda, senza alcun costo, per tutta la durata della vita del suo calcolatore.*

Gli analisti stimano che le sostituzioni siano costate all'Intel 500 milioni di dollari.

Questa vicenda evidenzia alcuni aspetti su cui vale la pena riflettere. Quanto più economico sarebbe stato correggere il baco fin dall'inizio? Qual è stato il costo per riparare il danno di immagine di Intel? E quali sono le responsabilità di una società nel rendere pubblico un baco in un prodotto così ampiamente usato, e su cui si fa così tanto affidamento, come un microprocessore?

### 3.10 Note conclusive

Nel corso degli anni l'aritmetica dei calcolatori è stata oggetto di una diffusa opera di standardizzazione che ha fortemente migliorato la portabilità dei programmi. L'aritmetica sui numeri interi in complemento a 2 si trova in tutti i calcolatori venduti oggi; se contiene anche il supporto per i numeri in virgola mobile, questi sono codificati secondo lo standard IEEE754.

L'aritmetica dei calcolatori si distingue dall'aritmetica svolta a mano per i vincoli dettati dalla precisione limitata. Questi limiti possono portare a operazioni non valide quando si calcolano numeri più grandi o più piccoli dei limiti prefissati. Queste anomalie, chiamate overflow e underflow, possono produrre eccezioni o interrupt, eventi eccezionali simili a chiamate non previste a procedure. Le eccezioni saranno trattate in maggior dettaglio nei Capitoli 4 e 5.

L'aritmetica in virgola mobile presenta in più il problema di essere basata su un'approssimazione dei numeri reali, e occorre molta attenzione per garantire che il numero calcolato sia la rappresentazione più vicina al numero vero. I problemi della mancanza di precisione e della limitatezza della rappresentazione sono oggetto di studio nel campo dell'analisi numerica. Il recente passaggio al parallelismo manterrà i riflettori puntati sull'analisi numerica, poiché soluzioni da lungo tempo considerate sicure sui calcolatori sequenziali devono essere tuttavia riconsiderate per i calcolatori paralleli nella progettazione di algoritmi più veloci che devono comunque garantire risultati corretti.

Il parallelismo a livello dei dati e, in particolare, il parallelismo a livello di parola offrono una soluzione semplice per ottenere prestazioni più elevate dei programmi che utilizzano in modo intensivo operazioni su interi o su numeri in virgola mobile. Abbiamo visto che si può velocizzare la moltiplicazione tra matrici quasi di un fattore quattro utilizzando istruzioni che possono eseguire quattro operazioni in virgola mobile in parallelo.



Istruzione RISC-V	Nome	Frequenza	Cumulativa
Add immediato	addi	14,36%	14,36%
Load doppia parola	ld	8,27%	22,63%
Load doppia parola in VM	fld	6,83%	29,46%
Add registri	add	6,23%	35,69%
Load parola	lw	4,38%	40,07%
Store parola doppia	sd	4,29%	44,36%
Salta se non uguale	bne	4,14%	48,50%
Shift a sinistra immediato	slli	3,65%	52,15%
Mult-add fuse assieme, parole doppie	fmadd.d	3,49%	55,64%
Salta se uguale	beq	3,27%	58,91%
Add immediato, parola	addiw	2,86%	61,77%
Store doppia parola in VM	fsd	2,24%	64,00%
Mult doppia parola in VM	fmul.d	2,02%	66,02%
Load bit più significativi, immediata	lui	1,56%	67,59%
Store di parola	sw	1,52%	69,10%
Jump and link	jal	1,38%	70,49%
Salta se minore	blt	1,37%	71,86%
Add di parola	addw	1,34%	73,19%
Sottrazione di parole doppie in VM	fsub.d	1,28%	74,47%
Salta se maggiore uguale	bge	1,27%	75,75%

**Figura 3.24** Frequenza delle istruzioni RISC-V negli SPEC CPU2006. Le 20 istruzioni più popolari, che collettivamente costituiscono il 76% di tutte le istruzioni eseguite, vengono riportate nella tabella. Le pseudoistruzioni vengono convertite in istruzioni RISC-V prima della loro esecuzione, e quindi non appaiono.

Assieme alla spiegazione dell'aritmetica dei calcolatori contenuta in questo capitolo, abbiamo anche descritto un'altra parte dell'insieme di istruzioni del RISC-V.

La Figura 3.24 riporta la popolarità delle venti istruzioni più comuni RISC-V per i benchmark SPEC CPU2006 interi e in virgola mobile. Come si può vedere, un numero relativamente piccolo di istruzioni domina questa classifica. Questa osservazione ha delle implicazioni importanti per la progettazione dei processori, come vedremo nel Capitolo 4.

Occorre avere presente che qualsiasi sia l'insieme delle istruzioni e la sua dimensione (RISC-V, MIPS, x86), le stringhe di bit non hanno un significato proprio perché possono rappresentare un intero con segno, senza segno, un numero in virgola mobile, una stringa, un'istruzione e altro ancora. Nei calcolatori a programma memorizzato, la stringa di bit acquista un significato in funzione delle operazioni che sono consentite.

## 3.11 Inquadramento storico e approfondimenti

Questo paragrafo, disponibile online , contiene una rassegna storica della virgola mobile risalendo nel tempo fino a von Neumann; vengono in esso descritti i sorprendenti sforzi che hanno consentito di superare le controversie sullo standard IEEE e viene inoltre fornita la motivazione che ha portato alla realizzazione dell'architettura a stack per la virgola mobile di 80 bit dell'x86.

*La legge di Gresham ("La cattiva moneta fa scappare quella buona") per i calcolatori suonerebbe: "Il Veloce fa scappare il Lento, anche se il Veloce è sbagliato".*  
William Kahan, 1992

## 3.12 Esercizi

*Non arrendersi mai, non arrendersi mai, mai, mai e poi mai – in niente, grande o piccolo, importante o insignificante – non arrendersi mai.*

Winston Churchill, discorso alla Harrow School, 1941

**3.1** [5] <3.2> Qual è il risultato della sottrazione 5ED4 – 07A4 se questi numeri sono rappresentati come numeri esadecimali senza segno su 16 bit? Scrivere il risultato in esadecimale e mostrare come è stata eseguita l'operazione.

**3.2** [5] <3.2> Qual è il risultato della sottrazione 5ED4 – 07A4 se questi numeri sono rappresentati come numeri esadecimali dotati di segno su 16 bit, codificati in modulo e segno? Scrivere il risultato in esadecimale e mostrare come è stata eseguita l'operazione.

**3.3** [10] <3.2> Convertire 5ED4 in binario. Che cosa rende la base 16 (esadecimale) un sistema di numerazione così attraente per rappresentare i numeri in un calcolatore?

**3.4** [5] <3.2> Qual è il risultato di 4365 – 3412 se essi sono rappresentati come numeri ottali senza segno su 12 bit? Scrivere il risultato in ottale e mostrare come è stata eseguita l'operazione.

**3.5** [5] <3.2> Qual è il risultato di 4365 – 3412 se essi sono rappresentati come numeri dotati di segno su 12 bit, codificati in modulo e segno? Scrivere il risultato in ottale e mostrare come è stata eseguita l'operazione.

**3.6** [5] <3.2> Si supponga che 185 e 122 siano numeri decimali interi senza segno su 8 bit. Calcolare il risultato della sottrazione 185 – 122 e dichiarare se si verifica overflow, underflow o nessuno dei due.

**3.7** [5] <3.2> Si supponga che 185 e 122 siano numeri decimali interi dotati di segno su 8 bit, codificati in modulo e segno. Calcolare 185 + 122 e dichiarare se si verifica overflow, underflow o nessuno dei due.

**3.8** [5] <3.2> Si supponga che 185 e 122 siano numeri decimali interi dotati di segno su 8 bit, codificati in modulo e segno. Calcolare 185 – 122 e dichiarare se si verifica overflow, underflow o nessuno dei due.

**3.9** [10] <3.2> Si supponga che 151 e 214 siano numeri decimali interi dotati di segno su 8 bit, codificati in complemento a 2. Calcolare 151 + 214 utilizzando l'aritmetica con saturazione e scrivere il risultato in

notazione decimale. Mostrare come è stato eseguito il calcolo.

**3.10** [10] <3.2> Si supponga che 151 e 214 siano numeri decimali interi dotati di segno su 8 bit, codificati in complemento a 2. Calcolare 151 – 214 utilizzando l'aritmetica con saturazione e scrivere il risultato in notazione decimale. Mostrare come è stato eseguito il calcolo.

**3.11** [10] <3.2> Si supponga che 151 e 214 siano numeri decimali interi senza segno su 8 bit. Calcolare 151 + 214 utilizzando l'aritmetica con saturazione e scrivere il risultato in notazione decimale. Mostrare come è stato eseguito il calcolo.

**3.12** [20] <3.3> Utilizzando una tabella simile a quella mostrata in Figura 3.6, calcolare il prodotto dei numeri interi senza segno 62 e 12, codificati in base ottale su 6 bit, servendosi dell'hardware descritto in Figura 3.3. Mostrare il contenuto di tutti i registri per ogni passo.

**3.13** [20] <3.3> Utilizzando una tabella simile a quella mostrata in Figura 3.6, calcolare il prodotto dei numeri interi senza segno 62 e 12, codificati in base esadecimale su 8 bit, servendosi dell'hardware descritto in Figura 3.5. Mostrare il contenuto di tutti i registri per ogni passo.

**3.14** [10] <3.3> Calcolare il tempo necessario per eseguire una moltiplicazione utilizzando gli approcci mostrati nelle Figure 3.3 e 3.4, supponendo che un intero sia codificato su 8 bit e che ciascun passo dell'operazione richieda 4 unità di tempo. Si supponga che al passo la venga sempre eseguita un'addizione: o si somma il moltiplicando, oppure il valore zero. Si assuma, inoltre, che i registri siano già stati inizializzati: occorre quindi misurare solamente quanto tempo impiega il ciclo della moltiplicazione vero e proprio. Se la moltiplicazione viene eseguita in hardware, gli scorrimenti del moltiplicando e del moltiplicatore possono essere eseguiti contemporaneamente; se invece è eseguita in software, i due scorrimenti devono essere effettuati in sequenza, uno dopo l'altro. Risolvere il problema sia per la moltiplicazione hardware sia per quella software.

**3.15** [10] <3.3> Calcolare il tempo necessario per eseguire una moltiplicazione utilizzando l'approccio descritto nel testo (31 sommatori impilati verticalmente), supponendo che gli interi siano codificati su 8 bit e che l'addizionatore richieda 4 unità di tempo.

**3.16** [20] <3.3> Calcolare il tempo necessario per eseguire una moltiplicazione utilizzando l'approccio mostrato in Figura 3.7, supponendo che gli interi siano codificati su 8 bit e che l'addizionatore richieda 4 unità di tempo.

**3.17** [20] <3.3> Come discusso nel testo, un possibile modo per migliorare le prestazioni è eseguire uno scorrimento seguito da una somma al posto di una moltiplicazione. Per esempio, dato che  $9 \times 6$  può essere scritto come  $(2 \times 2 \times 2 + 1) \times 6$ , possiamo calcolare  $9 \times 6$  facendo scorrere il numero 6 verso sinistra di tre posizioni, per poi sommare 6 al risultato. Mostrare il modo migliore per calcolare  $0 \times 33 \times 0 \times 55$  utilizzando le operazioni di scorrimento e di somma/sottrazione. Si supponga che moltiplicando e moltiplicatore siano interi senza segno su 8 bit.

**3.18** [20] <3.4> Utilizzando una tabella simile a quella mostrata in Figura 3.10, calcolare 74 diviso 21 servendosi dell'hardware descritto in Figura 3.8. Mostrare il contenuto di tutti i registri per ogni passo. Si supponga che entrambi gli input siano interi senza segno su 6 bit.

**3.19** [30] <3.4> Utilizzando una tabella simile a quella mostrata in Figura 3.10, calcolare 74 diviso 21 servendosi dell'hardware descritto in Figura 3.11. Mostrare il contenuto di tutti i registri per ogni passo. Si supponga che dividendo e divisore siano interi senza segno su 6 bit. Questo algoritmo richiede un approccio leggermente differente da quello mostrato in Figura 3.9. Per capire come eseguire correttamente la divisione, occorre fare un'analisi approfondita del problema ed effettuare alcuni tentativi, oppure cercare informazioni sul web. (Suggerimento: una possibile soluzione è basata sul fatto che il registro Resto può essere fatto scorrere in entrambe le direzioni, come mostrato in Figura 3.11.)

**3.20** [5] <3.5> Quale numero decimale è rappresentato dalla combinazione di bit 0x0C000000 nel caso in cui si tratti di un numero intero in complemento a 2? E se si trattasse di un intero senza segno?

**3.21** [10] <3.5> Se la combinazione di bit 0x0C000000 fosse scritta nel registro istruzioni, quale istruzione MIPS verrebbe eseguita?

**3.22** [10] <3.5> Quale numero in base dieci è rappresentato dalla combinazione di bit 0x0C000000 nel caso in cui si tratti di un numero in virgola mobile? Si utilizzi lo standard IEEE 754.

**3.23** [10] <3.5> Scrivere la rappresentazione binaria del numero decimale 63,25 supponendo che il numero sia codificato nel formato IEEE 754 in singola precisione.

**3.24** [10] <3.5> Scrivere la rappresentazione binaria del numero decimale 63,25 supponendo che il numero sia codificato nel formato IEEE 754 in doppia precisione.

**3.25** [10] <3.5> Scrivere la rappresentazione binaria del numero decimale 63,25 supponendo che il numero sia stato memorizzato utilizzando il formato IBM in singola precisione (base 16, invece che base 2, con 7 bit di esponente).

**3.26** [20] <3.5> Scrivere la rappresentazione binaria del numero  $-1,5625 \times 10^{-1}$  se il numero è codificato in un formato simile a quello adottato dal PDP-8 della Digital Equipment (i 12 bit più significativi contengono l'esponente memorizzato in complemento a 2, e gli altri 24 bit più a destra contengono la mantissa memorizzata anch'essa come numero in complemento a 2). Non viene omesso l'1 prima della virgola. Confrontare l'intervallo di rappresentazione e la risoluzione di questo formato su 36 bit con lo standard IEEE 754 in singola e doppia precisione.

**3.27** [20] <3.5> Il formato IEEE 754-2008 prevede la mezza precisione che codifica numeri su soli 16 bit. Il bit più a sinistra è sempre il bit di segno, l'esponente è su 5 bit e viene memorizzato in notazione polarizzata, con polarizzazione pari a 15, e la mantissa è lunga 10 bit. L'1 prima della virgola non viene scritto. Scrivere la combinazione di cifre binarie che rappresenta il numero  $-1,5625 \times 10^{-1}$  supponendo di avere a disposizione una versione di questo formato che utilizza 16 bit in più per rappresentare l'esponente. Confrontare l'intervallo di rappresentazione e la risoluzione di questo formato su 16 bit con lo standard IEEE 754 in singola precisione.

**3.28** [20] <3.5> I modelli 2114, 2115 e 2116 della Hewlett-Packard utilizzavano un formato che prevedeva i 16 bit più a sinistra per la mantissa, memorizzata in complemento a 2, seguita da un campo di altri 16 bit che conteneva negli 8 bit più a sinistra un'estensione della mantissa (che portava la mantissa a un totale di 24 bit), e negli 8 bit più a destra l'esponente. Tuttavia, con un curioso ribaltamento rispetto all'IEEE 754, l'esponente era memorizzato nel formato

modulo e segno, con il bit di segno contenuto nel bit più a destra! Scrivere la rappresentazione binaria del numero  $-1,5625 \times 10^{-1}$  secondo questo formato. Non viene omesso l'1 prima della virgola. Confrontare l'intervallo di rappresentazione e la risoluzione di questo formato su 32 bit con lo standard IEEE 754 in singola precisione.

**3.29** [20] <3.5> Calcolare la somma di  $2,6125 \times 10^1$  e  $4,150390625 \times 10^{-1}$  a mano, supponendo che i due addendi siano memorizzati nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27. Si supponga un bit di guardia, un bit di arrotondamento e uno di presenza. Arrotondare il risultato al numero pari più vicino e mostrare tutti i passi.

**3.30** [30] <3.5> Calcolare il prodotto di  $-8,0546875 \times 10^0$  e  $-1,79931640625 \times 10^{-1}$  a mano, supponendo che i due numeri siano memorizzati nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27. Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi della moltiplicazione. Si può eseguire la moltiplicazione in una forma comprensibile, come nell'esempio riportato nel testo, invece di utilizzare le tecniche adottate negli Esercizi 3.12-3.14. Indicare se si verifica un overflow o un underflow e scrivere il risultato sia in binario nel formato su 16 bit descritto nell'Esercizio 3.27 sia in decimale. Qual è l'accuratezza del risultato? Qual è l'accuratezza rispetto a quella che si ottiene con una calcolatrice?

**3.31** [30] <3.5> Calcolare a mano la divisione di  $8,625 \times 10^1$  per  $-4,875 \times 10^0$ , mostrando tutti i passaggi necessari a effettuare l'operazione. Si supponga che siano definiti i bit di guardia, di arrotondamento e di presenza, da utilizzarsi dove necessario. Scrivere il risultato nel formato in virgola mobile su 16 bit descritto nell'Esercizio 3.27 e nel formato decimale, e confrontare quest'ultimo con quello fornito dalla calcolatrice.

**3.32** [20] <3.10> Calcolare  $(3,984375 \times 10^{-1} + 3,4375 \times 10^{-1}) + 1,771 \times 10^3$  a mano, supponendo che ciascuno dei tre addendi sia codificato nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e descritto anche nel testo). Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.33** [20] <3.10> Calcolare  $3,984375 \times 10^{-1} + (3,4375 \times 10^{-1} + 1,771 \times 10^3)$  a mano, supponendo che ciascuno dei tre addendi sia codificato nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e de-

scritto anche nel testo). Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.34** [10] <3.10> Analizzare la soluzione degli Esercizi 3.32 e 3.33: si può affermare che  $(3,984375 \times 10^{-1} + 3,4375 \times 10^{-1}) + 1,771 \times 10^3 = 3,984375 \times 10^{-1} + (3,4375 \times 10^{-1} + 1,771 \times 10^3)$ ?

**3.35** [30] <3.10> Calcolare  $(3,41796875 \times 10^{-3} \times 6,34765625 \times 10^{-3}) \times 1,05625 \times 10^2$  a mano, supponendo che ciascun numero sia memorizzato nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e descritto anche nel testo). Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.36** [30] <3.10> Calcolare  $3,41796875 \times 10^{-3} \times (6,34765625 \times 10^{-3} \times 1,05625 \times 10^2)$  a mano, supponendo che ciascun numero sia memorizzato nel formato mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e descritto anche nel testo). Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.37** [10] <3.10> Analizzare la soluzione degli Esercizi 3.35 e 3.36: si può affermare che  $(3,41796875 \times 10^{-3} \times 6,34765625 \times 10^{-3}) \times 1,05625 \times 10^2 = 3,41796875 \times 10^{-3} \times (6,34765625 \times 10^{-3} \times 1,05625 \times 10^2)$ ?

**3.38** [30] <3.10> Calcolare  $1,666015625 \times 10^0 \times (1,9760 \times 10^4 + -1,9744 \times 10^4)$  a mano, supponendo che ciascun numero sia memorizzato nel formato a mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e descritto anche nel testo). Arrotondare al numero pari più vicino supponendo la presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli, e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.39** [30] <3.10> Calcolare  $(1,666015625 \times 10^0 \times 1,9760 \times 10^4) + (1,666015625 \times 10^0 \times -1,9744 \times 10^4)$  a mano, supponendo che ciascun numero sia memorizzato nel formato a mezza precisione su 16 bit descritto nell'Esercizio 3.27 (e descritto anche nel testo). Arrotondare al numero pari più vicino supponendo la

presenza di un bit di guardia, uno di arrotondamento e uno di presenza. Mostrare tutti i passi dei calcoli e scrivere il risultato sia in formato in virgola mobile su 16 bit sia in decimale.

**3.40** [10] <3.10> Analizzare la soluzione degli Esercizi 3.38 e 3.39: si può affermare che  $(1,666015625 \times 10^0 \times 1,9760 \times 10^4) + (1,666015625 \times 10^0 \times -1,9744 \times 10^4) = 1,666015625 \times 10^0 \times (1,9760 \times 10^4 + -1,9744 \times 10^4)$ ?

**3.41** [10] <3.5> Scrivere la sequenza di bit che rappresenta il numero  $-1/4$  nel formato in virgola mobile IEEE 754. Si può rappresentare  $-1/4$  in modo esatto?

**3.42** [10] <3.5> Quale numero si ottiene se si somma  $-1/4$  a se stesso 4 volte? Quanto vale  $-1/4 \times 4$ ? I due risultati sono uguali? Quali dovrebbero essere?

**3.43** [10] <3.5> Scrivere la sequenza di bit che rappresenta il numero  $1/3$  nel formato in virgola mobile IEEE 754 che utilizza numeri binari nella mantissa. Si supponga di avere a disposizione 24 bit e che non ci sia bisogno di normalizzare. Si può rappresentare  $1/3$  in modo esatto?

**3.44** [10] <3.5> Scrivere i bit che costituiscono la mantissa per un formato in virgola mobile che utilizza numeri BCD (Binari in Codifica Decimale – base 10) per la mantissa invece che numeri in base 2. Si supponga che ci siano 24 bit a disposizione e che non sia richiesta la normalizzazione. È esatta la rappresentazione dei numeri?

**3.45** [10] <3.5> Scrivere i bit che costituiscono la mantissa per un formato in virgola mobile che utilizza la base 15 invece della base 2. (Le cifre della base 16 sono 0-9 e A-F; quelle della base 15 saranno quindi 0-9

e A-E). Si supponga che ci siano 24 bit a disposizione e che non sia richiesta la normalizzazione. È esatta la rappresentazione dei numeri?

**3.46** [20] <3.5> Scrivere i bit che costituiscono la mantissa per un formato in virgola mobile che utilizza la base 30 invece della base 2. (Le cifre della base 16 sono 0-9 e A-F; quelle della base 30 saranno quindi 0-9 e A-T). Si supponga che ci siano 20 bit a disposizione e che non sia richiesta la normalizzazione. È esatta la rappresentazione dei numeri?

**3.47** [45] <3.6, 3.7> Il seguente frammento di codice C implementa un filtro FIR a quattro campioni sul vettore di input sig\_in. Si supponga che tutti i vettori siano costituiti da valori su 16 bit a virgola fissa.

```
for (i=3;i<128;i++)
    sig_out[i]=sig_in[i-3]*f[0]+sig_in[i-2]*
        f[1]+sig_in[i-1]*f[2]+sig_in[i]*f[3];
```

Si supponga di dover scrivere una versione ottimizzata del filtro in assembler per un processore dotato di istruzioni SIMD e registri a 128 bit. Senza conoscere i dettagli dell'insieme delle istruzioni, descrivere brevemente come implementereste questo frammento di codice in assembler, massimizzando le operazioni su sottoinsiemi della parola e minimizzando la quantità di dati che devono essere trasferiti tra i registri e la memoria. Dichiarate tutte le ipotesi che fate sulle istruzioni da utilizzare.

## Risposte alle domande di autovalutazione

Paragrafo 3.2, pagina 154 – 2.

Paragrafo 3.5, pagina 191 – 3.

# 4

## Il processore

*In un problema importante, nessun dettaglio è insignificante.*

Proverbio francese

### 4.1 Introduzione

Nel Capitolo 1 abbiamo visto che le prestazioni di un calcolatore sono determinate da tre fattori chiave: numero di istruzioni, durata del ciclo di clock e *numero di cicli di clock per istruzione* (CPI). Il compilatore e l'architettura dell'insieme delle istruzioni, che abbiamo esaminato nel Capitolo 2, determinano il numero di istruzioni richieste da un dato programma. Tuttavia, sia il periodo di clock sia il numero di cicli di clock per istruzione dipendono dalla particolare implementazione del processore. In questo capitolo, costruiremo un'unità di elaborazione dati (*datapath*) e l'unità di controllo (*control unit*) associata a due differenti implementazioni hardware dell'insieme di istruzioni RISC-V.

In particolare, il capitolo fornisce una spiegazione dei principi e delle tecniche adottate nella costruzione di un processore, partendo, in questo paragrafo, da un livello di astrazione elevato e da una visione semplificata. Passeremo poi a descrivere passo dopo passo come ideare un'unità di elaborazione e una versione semplificata di processore sufficienti a implementare un insieme di istruzioni come quello del MIPS. Nella parte centrale del capitolo descriveremo un'implementazione più realistica del RISC-V, basata sulla *pipeline*. Infine, seguirà un paragrafo nel quale saranno sviluppati i concetti necessari a implementare insiemi di istruzioni più complessi, come l'x86.

Per i lettori interessati all'interpretazione ad alto livello delle istruzioni e del loro impatto sulle prestazioni di un programma, i concetti che stanno alla base della pipeline sono presentati in questo paragrafo iniziale e nel paragrafo 4.5.



PIPELINE

Gli sviluppi più recenti sono descritti nel paragrafo 4.10, mentre il paragrafo successivo è dedicato alla descrizione delle architetture Core i7 Intel e Cortex-A8 ARM. Nel paragrafo 4.12 viene mostrato come sfruttare il parallelismo a livello di istruzioni per raddoppiare le prestazioni della funzione che effettua la moltiplicazione di due matrici, discussa nel paragrafo 3.8. Questi paragrafi consentono di comprendere ad alto livello il funzionamento di una pipeline.

I paragrafi 4.3, 4.4 e 4.6 saranno particolarmente utili ai lettori più interessati al funzionamento del processore e alle sue prestazioni, mentre coloro che vogliono comprendere come si costruisce un processore dovrebbero leggere anche i paragrafi 4.2, 4.7, 4.8 e 4.9. Per i lettori interessati alla progettazione moderna dell'hardware, il paragrafo 4.13  descrive i linguaggi di progettazione dell'hardware e gli strumenti CAD impiegati per implementare l'hardware e per descrivere l'implementazione di una pipeline. Questo paragrafo, inoltre, propone altri esempi di come siano eseguite le istruzioni dall'hardware basato su pipeline.

## Un'implementazione di base del RISC-V

Esamineremo ora un'implementazione che comprende il seguente sottoinsieme di istruzioni di base del RISC-V:

- le istruzioni di riferimento alla memoria *load doubleword* (ld) e *store doubleword* (sd);
- le istruzioni aritmetico-logiche add, sub, and e or;
- le istruzioni di salto condizionato dal risultato di un test di uguaglianza, *branch equal* (beq).

Un tale sottoinsieme non comprende tutte le istruzioni sugli interi (per es., sono assenti le istruzioni di scorrimento, moltiplicazione e divisione), e non include neppure alcuna istruzione sui numeri in virgola mobile, ma consente di illustrare i principi chiave utilizzati per costruire un'unità di elaborazione e per progettare l'unità di controllo relativa. L'implementazione delle restanti istruzioni è molto simile.

Nell'esaminare l'implementazione, avremo l'opportunità di comprendere come l'architettura dell'insieme delle istruzioni determini molti aspetti dell'implementazione stessa, e come la scelta tra le varie strategie di implementazione influenzi sia la frequenza di clock sia il CPI del calcolatore. Molti principi chiave introdotti nel Capitolo 2, quale *La semplicità favorisce la regolarità*, trovano diretta applicazione nell'implementazione delle istruzioni. Inoltre, molti concetti introdotti in questo capitolo per implementare il sottoinsieme di istruzioni RISC-V sono gli stessi su cui si basa la realizzazione della maggior parte dei calcolatori, dai server ad alte prestazioni ai microprocessori universali, ai processori embedded.

## Una panoramica sull'implementazione

Nel Capitolo 2 abbiamo visto l'insieme di base delle istruzioni RISC-V, comprese le istruzioni aritmetico-logiche sugli interi, le istruzioni di accesso alla memoria e le istruzioni di salto condizionato. Molti componenti necessari per realizzare queste istruzioni sono gli stessi, indipendentemente dal tipo di istruzione. Per ogni istruzione i primi due passi sono identici:

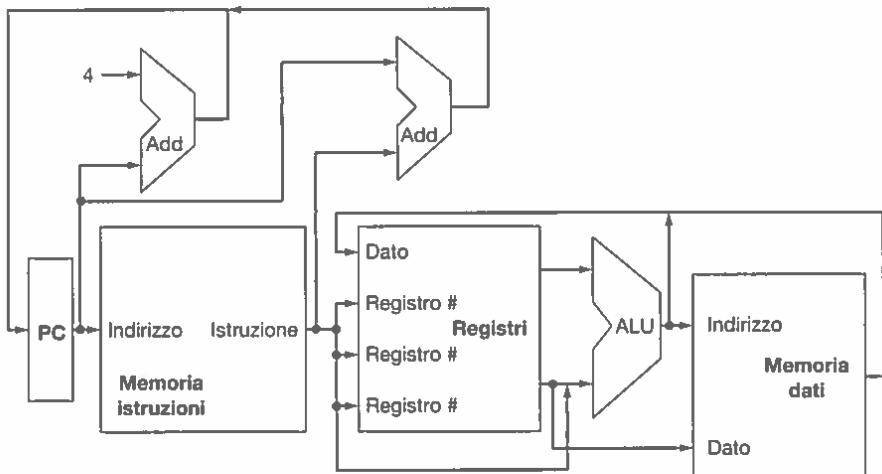
1. Inviare il contenuto del *program counter* (PC) alla memoria che contiene il programma e prelevare l'istruzione dalla memoria; questa operazione viene anche indicata con il termine inglese di *fetch*, letteralmente "raggiungere".
2. Leggere il contenuto di uno o due registri utilizzando i campi dell'istruzione per selezionare i registri. Per l'istruzione *load word* (ld) abbiamo bisogno di leggere un solo registro, ma la maggior parte delle altre istruzioni richiede la lettura di due registri.

Dopo questi due passi, le azioni richieste per completare l'esecuzione delle istruzioni dipendono dalla loro tipologia. Fortunatamente, per ognuna delle tre classi (istruzioni di accesso alla memoria, aritmetico-logiche e di salto), le azioni da compiere sono, a grandi linee, le stesse, indipendentemente dal particolare codice operativo dell'istruzione. La semplicità e la regolarità dell'insieme di istruzioni RISC-V semplificano l'implementazione rendendo simile l'esecuzione di istruzioni di classi diverse.

Per esempio, tutti i tipi di istruzioni, eccetto i salti incondizionati, utilizzano l'unità aritmetico-logica (ALU) dopo aver letto i registri; le istruzioni di accesso alla memoria utilizzano la ALU per il calcolo dell'indirizzo, le istruzioni aritmetico-logiche per l'esecuzione della loro operazione e i salti condizionati per effettuare i confronti. Dopo avere utilizzato la ALU, le azioni richieste per completare l'esecuzione variano a seconda del tipo di istruzione. Un'istruzione di accesso alla memoria avrà bisogno di accedere alla memoria sia per scrivere i dati (con una store) sia per leggerli (mediante una load). Un'istruzione aritmetico-logica, invece, dovrà scrivere i dati prodotti dalla ALU nel registro destinazione. Infine, in un'istruzione di salto condizionato l'indirizzo dell'istruzione successiva potrà cambiare in base al risultato del confronto; in assenza di salto, il PC verrà incrementato di 4 per puntare all'istruzione successiva.

La Figura 4.1 mostra uno schema ad alto livello di astrazione dell'implementazione di un RISC-V, focalizzato sulle varie unità funzionali e sulle loro interconnessioni. Sebbene questa figura mostri che la maggior parte del flusso dei dati "scorre" attraverso il processore, omette due importanti aspetti dell'esecuzione di un'istruzione.

Per prima cosa, in diversi punti, la Figura 4.1 mostra dati che provengono da due diverse sorgenti e arrivano alla stessa unità funzionale. Per esempio, il valore scritto nel PC può provenire da uno dei due sommatori ("Add" in figura);



**Figura 4.1** Uno schema astratto dell'implementazione di un sottoinsieme delle istruzioni RISC-V che mostra le principali unità funzionali e le principali connessioni tra loro. Tutte le istruzioni iniziano utilizzando il PC per fornire il loro indirizzo alla memoria istruzioni. Dopo che l'istruzione è stata caricata (*fetch*), il numero d'ordine dei registri contenenti gli operandi può essere letto nei campi opportuni dell'istruzione stessa. Una volta caricati gli operandi, si può elaborare il loro contenuto per determinare un indirizzo di memoria (nel caso di load o store), per eseguire un calcolo effettivo (operazioni aritmetico-logiche su interi) o per eseguire un confronto (salto condizionato). Se l'istruzione è di tipo aritmetico-logico, il risultato dalla ALU deve venire scritto in un registro. Se l'operazione è una load o una store, il risultato della ALU viene utilizzato come indirizzo, rispettivamente, per leggere dalla memoria il dato da scrivere nel register file o per scrivere il contenuto di un registro in memoria. I salti condizionati richiedono l'utilizzo dell'uscita della ALU per determinare l'indirizzo dell'istruzione successiva; tale indirizzo può provenire dalla ALU, nella quale PC e offset vengono sommati, oppure da un sommatore apposito che incrementa il PC di 4. Le linee che collegano tra loro le unità funzionali rappresentano i bus, formati da gruppi di segnali di controllo. Le frecce indicano la direzione in cui l'informazione scorre. Dato che le linee dei segnali possono incrociarsi senza necessariamente essere connesse, la connessione tra due linee viene esplicitamente segnalata con un punto.

il dato scritto nel register file può provenire sia dalla ALU sia dalla memoria. Inoltre, il secondo input della ALU può provenire da un registro o dal campo immediato dell'istruzione. In pratica, queste linee multiple per i dati non possono semplicemente essere connesse tra loro, ma dobbiamo introdurre un circuito logico che selezioni da quale, fra le differenti sorgenti collegate, debba essere preso il dato, connettendo la sorgente opportuna con la sua destinazione. Questa selezione viene comunemente effettuata con un dispositivo chiamato *multiplexer*, che funge da “selettore dei dati” (*data selector*).

Il circuito logico del multiplexer viene descritto in dettaglio nell'Appendice A ; esso porta in uscita il contenuto di uno degli ingressi a seconda del valore presente sulle sue linee di controllo. Le linee di controllo sono impostate principalmente in base alle informazioni prelevate dall'istruzione stessa che viene eseguita.

Il secondo aspetto omesso nella figura è che i segnali di controllo di molte unità funzionali presenti nello schema dipendono dal tipo di istruzione. Per esempio, la memoria dati deve essere letta da una load, mentre deve essere scritta da una store. Il register file deve essere scritto solamente nel caso di un'istruzione di load o di un'istruzione aritmetico-logica. Naturalmente, la ALU deve poi poter scegliere l'operazione corretta tra le tante che ha a disposizione, come abbiamo visto nel Capitolo 2 (l'Appendice A descrive in dettaglio la progettazione dei circuiti della ALU). Come nei multiplexer, le linee di controllo impostate sulla base del contenuto dei vari campi dell'istruzione dirigono l'elaborazione.

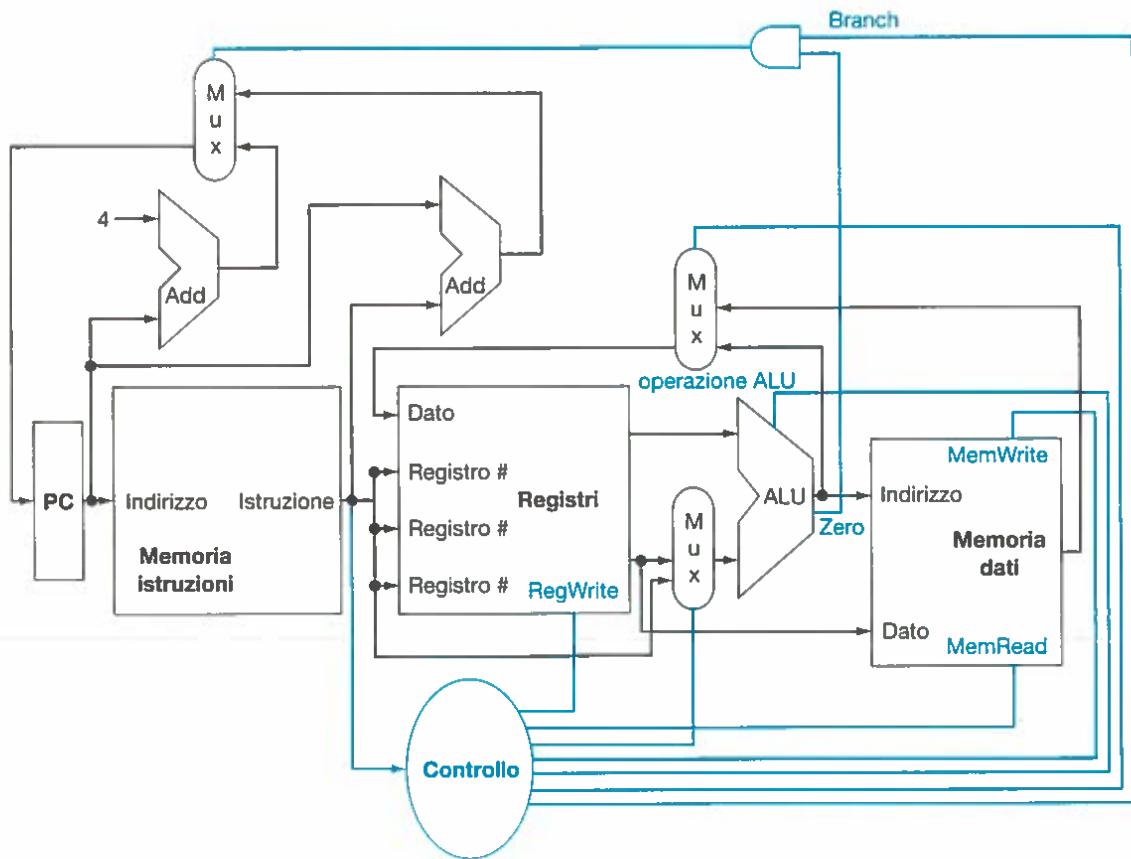
La Figura 4.2 mostra l'unità di elaborazione di Figura 4.1 integrata con i tre multiplexer richiesti e con le linee di controllo per le principali unità funzionali. L'*unità di controllo*, che ha come ingresso l'istruzione, viene utilizzata per determinare come impostare le linee di controllo per le unità funzionali e per due dei tre multiplexer. Il terzo multiplexer, che determina se nel PC debba essere scritto il valore  $PC + 4$  o l'indirizzo di destinazione del salto, viene controllato dall'uscita “Zero” della ALU; questa uscita è impostata a 1 quando l'uguaglianza tra i due operandi è verificata, e quindi riporta il risultato del confronto richiesto dall'istruzione beq. La regolarità e la semplicità dell'insieme delle istruzioni RISC-V consentono di utilizzare un processo di decodifica semplice per impostare le linee di controllo.

Nel resto del capitolo raffineremo questo schema introducendo dei dettagli che richiederanno l'aggiunta di altre unità funzionali, l'aumento del numero di connessioni tra le unità, e, naturalmente, il miglioramento dell'unità di controllo affinché essa sia in grado di stabilire quali azioni debbano essere compiute per ogni diverso tipo di istruzione. I paragrafi 4.3 e 4.4 descrivono una semplice implementazione che utilizza un singolo ciclo di clock, molto lungo, per ogni istruzione ed è basata sulla forma generale del cammino di elaborazione descritto nelle Figure 4.1 e 4.2. In questa CPU ogni istruzione inizia l'esecuzione su un fronte del segnale di clock e completa l'esecuzione sul fronte di clock successivo.

Sebbene sia facile da comprendere, questo approccio non è pratico, dal momento che il ciclo di clock deve essere allungato di molto per dare il tempo alle istruzioni più lente di terminare l'esecuzione. Dopo aver progettato l'unità di controllo di questo semplice calcolatore, esamineremo un'implementazione mediante pipeline con tutte le sue complessità, compresa la gestione delle eccezioni.

## Autovalutazione

Quanti dei cinque componenti classici di un calcolatore (unità di controllo, processore, memoria, input e output) sono mostrati nelle Figure 4.1 e 4.2?



**Figura 4.2 Implementazione base del sottoinsieme di istruzioni RISC-V, comprendente anche i multiplexer necessari e le linee di controllo.**

I multiplexer "Mux" in alto controlla ciò che viene scritto nel PC: PC + 4 oppure l'indirizzo di destinazione del salto condizionato; questo multiplexer viene controllato da una porta logica che esegue l'AND tra il segnale di Zero in uscita dalla ALU e un segnale di controllo che indica che l'istruzione è un'istruzione di salto condizionato. Il multiplexer posizionato centralmente nella parte alta della figura, la cui uscita ritorna ai register file, viene utilizzato per portare l'uscita della ALU (in caso di istruzione aritmetico-logica) oppure l'uscita della memoria dati (in caso di load) al register file per la scrittura. Infine l'ultimo multiplexer, posizionato più in basso, viene utilizzato per determinare se il secondo ingresso della ALU provenga dai registri (per un'istruzione aritmetico-logica oppure per un'istruzione di salto condizionato) o dal campo costante dell'istruzione stessa (per un'operazione di load o store). Le linee di controllo aggiunte hanno un significato evidente e determinano le operazioni che devono essere eseguite dalla ALU, oppure se la memoria debba essere letta o scritta, o ancora se si debba eseguire un'operazione di scrittura sul register file. Le linee di controllo sono mostrate in blu in modo da risultare più visibili.

## 4.2 Convenzioni del progetto logico

Nell'analizzare il progetto di un calcolatore, dobbiamo decidere come funzionerà la logica hardware che implementa il calcolatore e come esso dovrà essere sincronizzato. Questo paragrafo riprende alcuni concetti fondamentali di logica digitale che verranno utilizzati largamente in tutto il capitolo. Il lettore che non abbia una buona preparazione nel campo della progettazione dei circuiti digitali, troverà utile leggere l'Appendice A prima di proseguire oltre.

**Elemento combinatorio:** un elemento in grado di eseguire delle operazioni, per esempio una porta AND o una ALU.

Gli elementi funzionali che costituiscono l'unità di elaborazione del RISC-V sono costituiti da due diverse classi di elementi logici: elementi che operano sui dati ed elementi che contengono lo stato. Gli elementi che operano sui dati sono di tipo **combinatorio**, il che significa che in ogni istante i loro output dipendono solamente dagli input ricevuti nello stesso istante. Dati gli stessi valori di ingresso, un elemento combinatorio produrrà sempre il medesimo valore in uscita. La ALU raffigurata in Figura 4.1, e descritta nell'Appendice A , costituisce un esempio di elemento combinatorio: dato un certo insieme di valori in ingresso, essa produce sempre la stessa uscita perché non possiede alcun elemento interno di memoria.

Gli elementi dell'unità di elaborazione che non sono combinatori contengono lo *stato*. Un elemento può contenere lo stato solo se ha al suo interno elementi di memoria. Chiameremo questi **elementi di stato**, perché se spegniamo e riaccendiamo il calcolatore, possiamo far ripartire il flusso di lavoro del calcolatore esattamente dallo stesso punto caricando gli elementi di stato con i valori che essi contenevano prima di avere spento. Inoltre, salvando e ripristinando gli elementi di stato, il calcolatore si comporterebbe come se non fosse mai stato spento. Di conseguenza, gli elementi di stato caratterizzano completamente il funzionamento del calcolatore. Nella Figura 4.1 la memoria istruzioni, la memoria dati e i registri sono esempi di elementi di stato.

Un elemento di stato possiede almeno due ingressi e un'uscita. Gli ingressi richiesti sono il valore da scrivere nell'elemento e il *clock*, che determina quando scrivere. L'uscita di un elemento di stato è il valore contenuto al suo interno, scritto in un ciclo di *clock* passato. Per esempio, uno degli elementi di stato più semplici, dal punto di vista logico, è il *flip-flop* di tipo D (*vedi Appendice A*), che ha esattamente questi due ingressi (il valore del dato e il *clock*) e un'uscita. Oltre ai *flip-flop*, l'architettura RISC-V utilizza altri due tipi di elementi di stato: le memorie e i registri, che compaiono anch'essi in Figura 4.1. Il *clock* viene utilizzato per determinare quando l'elemento di stato deve essere scritto, mentre la lettura può avvenire in qualsiasi istante.

I componenti logici che contengono lo stato sono anche detti *sequenziali*, in quanto le loro uscite dipendono sia dagli ingressi sia dal valore del loro stato interno. Per esempio, l'uscita dell'unità funzionale che contiene i registri dipende sia dal numero del registro specificato sugli ingressi sia dal valore contenuto da quel particolare registro, che era stato scritto in precedenza. L'*Appendice A* descrive in modo dettagliato il funzionamento e la struttura degli elementi combinatori e di quelli sequenziali.

## Metodologia di temporizzazione

La **metodologia di temporizzazione** definisce quando i segnali possono essere scritti e quando possono essere letti. È importante temporizzare le operazioni di lettura e scrittura, perché se un segnale venisse letto e contemporaneamente scritto, il valore letto potrebbe corrispondere al valore precedente la scrittura, a quello successivo, o addirittura a una combinazione dei due. Ovviamente i progettisti dei calcolatori non possono accettare questa imprevedibilità nei risultati, che viene evitata proprio grazie alla metodologia di temporizzazione.

Per semplicità, supporremo di utilizzare una **metodologia di temporizzazione sensibile ai fronti** (*edge-triggered*); essa garantisce che il valore memorizzato all'interno di un elemento logico sequenziale venga aggiornato solamente in corrispondenza di un fronte del segnale di *clock*, che è una veloce transizione dal livello basso al livello alto e viceversa (Figura 4.3). Poiché solo gli elementi di stato possono memorizzare i dati, un qualsiasi circuito combinatorio deve ricevere l'input da un insieme di elementi di stato e può scrivere il risultato dell'elaborazione in un secondo insieme di elementi di stato. Gli ingressi sono i valori che erano stati scritti in un ciclo di *clock* precedente, mentre gli output del circuito combinatorio sono valori che potranno essere utilizzati in un ciclo di *clock* successivo. La Figura 4.3 mostra un blocco di logica combinatoria posto tra due elementi di stato. Esso lavora all'interno di un singolo ciclo di *clock*: tutti i segnali si devono propagare a partire dall'elemento di stato 1 e devono attraversare la logica combinatoria fino a raggiungere l'elemento di stato 2 nel tempo di un periodo di *clock*. Il tempo richiesto perché i segnali raggiungano l'elemento di stato 2 definisce quindi la durata minima del periodo di *clock*.

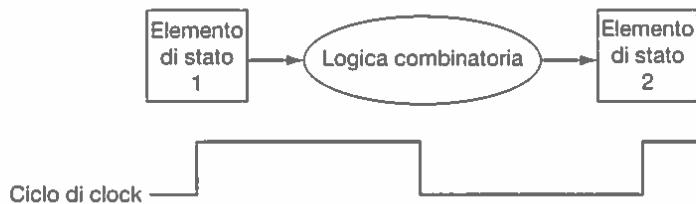
Per semplicità, non si riporterà esplicitamente il **segnale di controllo** della scrittura per quegli elementi di stato che vengono scritti a ogni fronte attivo del

**Elemento di stato:** un elemento di memoria, ad esempio un registro o una memoria.

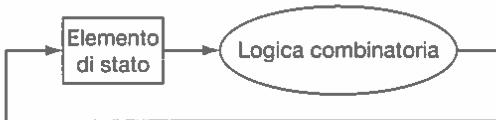
**Metodologia di temporizzazione:** l'approccio utilizzato per determinare quando un dato è valido e stabile in relazione al segnale di *clock*.

**Temporizzazione sensibile ai fronti:** una tecnica di temporizzazione nella quale tutti i cambiamenti di stato avvengono su un fronte del segnale di *clock*.

**Segnale di controllo:** un segnale utilizzato per la selezione in un multiplexer o per dirigere il funzionamento di un'unità funzionale; si distingue dai segnali che portano i dati, che contengono l'informazione che viene elaborata dalle unità funzionali.



**Figura 4.3** Logica combinatoria, elementi di stato e clock sono strettamente correlati. In un sistema digitale sincrono, il segnale di clock determina quando gli elementi di stato scrivono la loro memoria interna. Gli ingressi a un elemento di stato devono raggiungere un valore stabile (cioè devono aver raggiunto un valore che non cambierà almeno fino a dopo il fronte di clock successivo) prima che il fronte attivo del clock provochi l'aggiornamento dello stato. In questo capitolo, supporremo che tutti gli elementi di stato, comprese le memorie, siano sensibili al fronte di salita, cioè cambino contenuto quando il clock passa dallo stato basso allo stato alto.



**Figura 4.4** Una metodologia sensibile ai fronti permette di leggere e scrivere un elemento di stato nello stesso ciclo di clock, senza creare una competizione che potrebbe causare indeterminazione nel contenuto dell'elemento. Ovviamente, il ciclo di clock deve avere una durata sufficiente a garantire che gli ingressi siano stabili quando arriva il fronte attivo del clock. Non si può verificare alcuna retroazione all'interno del singolo ciclo di clock, grazie all'aggiornamento sul fronte del clock dell'elemento di stato. Se fosse possibile aggiornare l'elemento di stato prima del fronte di clock, il circuito non potrebbe funzionare correttamente. I circuiti presentati in questo capitolo e nel prossimo si basano sulla metodologia di temporizzazione sensibile ai fronti e su strutture simili a quele mostrata in questa figura.

clock. Viceversa, se un elemento di stato non viene aggiornato a ogni ciclo di clock, allora si renderà necessario un segnale esplicito di controllo della scrittura. Sia il segnale di clock sia il segnale di controllo della scrittura sono ingressi, e il contenuto dell'elemento di stato può essere modificato solamente sul fronte attivo del clock, quando cioè il segnale di controllo della scrittura è asserito.

Utilizzeremo il termine **asserito** per indicare un segnale che si trova nello stato logico alto e il termine **asserrire** per indicare che un segnale deve essere portato al livello logico alto; utilizzeremo il termine **non asserito** per indicare un segnale che si trova nello stato logico basso e il termine **non asserrire** per indicare che un segnale deve essere portato al livello logico basso. Abbiamo introdotto i termini “asserrire” e “non asserrire” perché un 1 può rappresentare a volte un valore logico alto e altre volte un valore logico basso.

La metodologia sensibile ai fronti permette di leggere il contenuto di un registro, inviare il valore attraverso uno o più blocchi di logica combinatoria e scrivere lo stesso registro nello stesso ciclo di clock, come mostrato in Figura 4.4. Non occorre specificare se la scrittura avvenga sul fronte di salita (da basso ad alto) o di discesa (da alto a basso) del segnale di clock, dal momento che gli ingressi a un circuito logico combinatorio possono cambiare solo sul fronte attivo del clock. In questo testo utilizzeremo il fronte di salita del clock. Con la metodologia di temporizzazione sensibile ai fronti non si corre il rischio di innescare una retroazione all'interno dello stesso ciclo di clock, e il circuito logico illustrato in Figura 4.4 funziona correttamente. Nell'Appendice A sono discussi brevemente alcuni vincoli temporali aggiuntivi (come i tempi di set-up e di hold) e vengono presentate altre metodologie di temporizzazione.

Quasi tutti gli elementi di stato e combinatori dell'architettura RISC-V a 64 bit hanno ingressi e uscite di ampiezza pari a 64 bit, essendo questa l'am-

**Asserito:** un segnale associato a un livello logico alto, o vero.

**Non asserito:** un segnale associato a un livello logico basso, o falso.

piezza della maggior parte dei dati elaborati dal processore. Quando incontreremo unità con ingressi o uscite di ampiezza diversa, la loro presenza sarà indicata esplicitamente. Nelle figure, i bus, ossia le linee che trasportano segnali di ampiezza maggiore di 1 bit, saranno indicati con linee più spesse. In alcuni casi due o più bus verranno uniti per formare un bus di ampiezza maggiore; per esempio, possiamo ottenere un bus da 64 bit unendo due bus da 32 bit. In questi casi, etichette sulle linee indicheranno chiaramente che due bus sono stati concatenati per ottenere un unico bus. Verranno inoltre aggiunte delle frecce per indicare la direzione del flusso dei dati tra gli elementi. Infine, un diverso **colore** consentirà di distinguere i segnali di controllo dai dati, distinzione che diverrà più chiara procedendo con la lettura del capitolo.

### Autovalutazione

Vero o falso? Poiché il register file può essere letto e scritto nello stesso ciclo di clock, le unità di elaborazione RISC-V che utilizzano la metodologia di temporizzazione sensibile ai fronti devono possedere più di una copia del register file.

**Approfondimento.** Esiste anche una versione dell'architettura RISC-V a 32 bit. Naturalmente, la maggior parte dei cammini interni di questa architettura è a 32 bit.

## 4.3 | Realizzazione di un'unità di elaborazione

Un modo ragionevole per iniziare la costruzione di un'unità di elaborazione consiste nell'esaminare i componenti principali che servono per eseguire i diversi tipi di istruzioni RISC-V. Procediamo dall'inizio, analizzando quali **elementi dell'unità di elaborazione** (*datapath elements*) siano richiesti da ciascuna classe di istruzioni per poi approfondire i diversi elementi procedendo attraverso i diversi livelli di **astrazione**. Quando mostreremo i diversi elementi dell'unità di elaborazione, mostreremo anche i relativi segnali di controllo associati a quel livello di astrazione.

La Figura 4.5a mostra il primo elemento di cui abbiamo bisogno: un'unità di memoria in cui salvare le istruzioni del programma che sia in grado di fornire in uscita l'istruzione associata all'indirizzo dato in ingresso. La Figura 4.5b mostra anche il **program counter (PC)** che, come abbiamo già visto nel Capitolo 2, è un registro utilizzato per memorizzare l'indirizzo dell'istruzione corrente. Infine, è necessario un sommatore per incrementare di 4 il PC e ottenere l'indirizzo dell'istruzione successiva. Tale sommatore è un elemento combinatorio e può essere costruito a partire dalla ALU descritta in dettaglio nell'Appendice A; è sufficiente impostare i segnali di controllo della ALU in modo che venga specificata sempre un'operazione di somma. Identificheremo negli schemi circuituali una ALU di questo tipo, configurata permanentemente come sommatore, con l'etichetta *Add*, come mostrato in Figura 4.5c; essa non può eseguire le altre operazioni tipiche di una ALU.

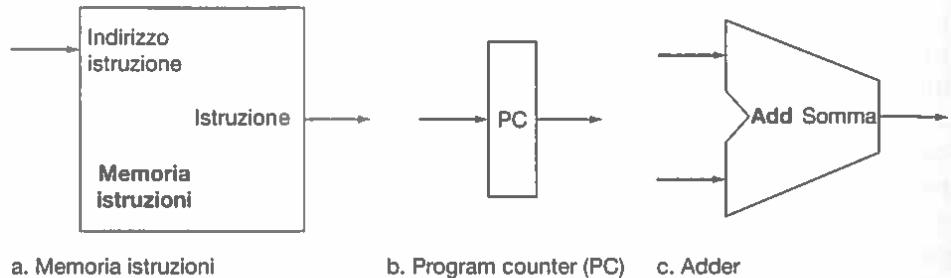
Per eseguire una qualsiasi istruzione occorre anzitutto prelevare l'istruzione stessa dalla memoria. Per prepararsi a eseguire l'istruzione successiva, bisogna incrementare il program counter in modo che punti a quell'istruzione, cioè 4 byte dopo. La Figura 4.6 mostra come si possono combinare i tre elementi della Figura 4.5 per formare un'unità di elaborazione che prelevi le istruzioni e incrementi il PC per ottenere l'indirizzo dell'istruzione successiva del programma.



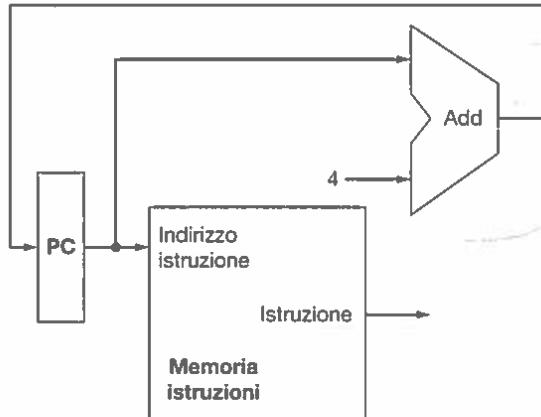
ASTRAZIONE

**Elemento dell'unità di elaborazione:** un'unità funzionale utilizzata per eseguire operazioni o memorizzare i dati all'interno del processore. Nell'implementazione RISC-V, gli elementi dell'unità di elaborazione sono la memoria dati, la memoria istruzioni, il register file, l'unità aritmetico-logica (ALU) e i sommatori.

**Program counter (PC):** il registro che contiene l'indirizzo dell'istruzione in esecuzione di un programma.



**Figura 4.5** Sono necessari due elementi di stato per memorizzare e accedere alle istruzioni, e serve un sommatore per calcolare l'indirizzo dell'istruzione successiva. Gli elementi di stato sono la memoria delle istruzioni e il program counter (PC). La memoria delle istruzioni fornisce solamente un accesso in lettura, perché essa non viene scritta dall'unità di elaborazione. Per questo tratteremo la memoria istruzioni come un circuito logico combinatorio: l'uscita in ogni istante corrisponde al contenuto della locazione specificata dall'indirizzo in ingresso, e non è necessario alcun segnale di controllo esplicito. La memoria delle istruzioni dovrà essere scritta solamente nella fase di caricamento del programma; non è difficile aggiungere questa funzionalità e quindi, per semplicità, in questo libro non la considereremo. Il program counter è un registro a 64 bit che viene scritto al termine di ogni ciclo di clock e non richiede quindi un segnale di controllo esplicito per la scrittura. Il sommatore è una ALU configurata in modo da eseguire sempre la somma dei suoi due ingressi a 64 bit e fornire in uscita il risultato della somma.



**Figura 4.6** Parte dell'unità di elaborazione utilizzata per prelevare le istruzioni e incrementare il program counter. L'istruzione prelevata viene utilizzata dalle altre parti dell'unità di elaborazione.

Consideriamo ora le istruzioni in formato R (vedi Figura 2.19). Tutte le istruzioni di questo tipo leggono due registri, eseguono un'operazione con la ALU sul contenuto di questi due registri e, alla fine, scrivono il risultato in un registro. Chiameremo queste *istruzioni di tipo R* o *istruzioni aritmetico-logiche*, dato che eseguono operazioni aritmetiche o logiche. Questa classe di istruzioni comprende add, sub, and e or, introdotte nel Capitolo 2. Un tipico esempio di istruzione di questo tipo è add x<sub>1</sub>, x<sub>2</sub>, x<sub>3</sub>, che legge il contenuto di x<sub>2</sub> e x<sub>3</sub> e scrive il risultato dell'addizione nel registro x<sub>1</sub>.

I registri universali a 32 bit del processore sono raccolti in una struttura detta **register file** (banco di registri), un insieme di registri in cui ciascun registro può essere letto o scritto specificando il numero ad esso associato all'interno dell'insieme. Il register file contiene lo stato dei registri del calcolatore. Avremo inoltre bisogno di una ALU per operare sui valori letti dai registri.

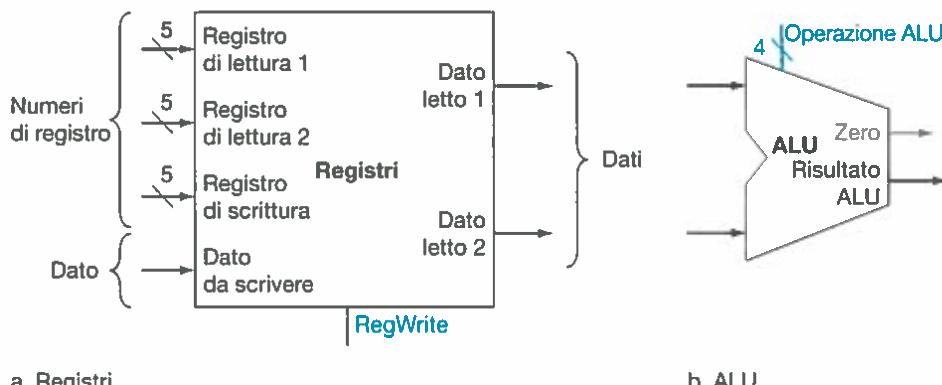
Poiché le istruzioni di tipo R hanno tre registri come operandi, per ciascuna istruzione dovremo leggere due dati di una parola ciascuno dal register file e poi scrivere il risultato, sempre nel register file. Per ogni parola letta dai regi-

**Register file:** un elemento dello stato che contiene un insieme di registri che possono essere letti o scritti fornendo il numero del registro sul quale si vuole operare.

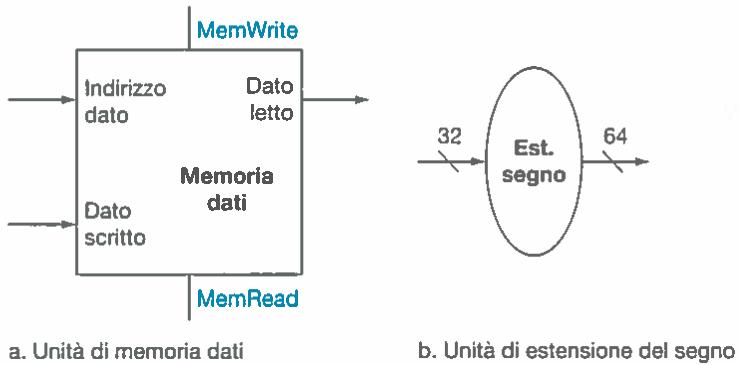
stri è necessario prevedere un ingresso al register file che specifichi il *numero del registro da leggere* e un'uscita dal register file sulla quale viaggerà il valore letto dal registro. Per scrivere un dato di una parola serviranno due ingressi: il primo ingresso deve specificare il *numero del registro di scrittura* e il secondo deve fornire il *dato* da scrivere nel registro. Il register file fornisce in qualunque momento in uscita il contenuto del registro il cui numero d'ordine è specificato sull'ingresso relativo al *numero del registro di lettura*. La scrittura, invece, viene controllata da un segnale di controllo esplicito, "RegWrite", che deve essere asserito perché il registro possa essere scritto in corrispondenza del fronte del clock. Quindi serviranno complessivamente quattro ingressi (tre per i numeri del registro e uno per il dato da scrivere) e due uscite (entrambe per i dati), come mostrato in Figura 4.7a. Gli ingressi che specificano il numero dei registri hanno ampiezza di 5 bit, in modo da poter specificare uno dei 32 registri ( $32 = 2^5$ ), mentre i bus dei dati in ingresso e in uscita sono di 64 bit ciascuno.

La Figura 4.7b mostra una ALU che riceve due input di 64 bit e produce un risultato su 64 bit; inoltre produce un segnale a un bit che vale 1 se il risultato dell'operazione è 0. I segnali di controllo della ALU, su 4 bit, vengono descritti in dettaglio nell'Appendice A ; li rivedremo brevemente quando dovremo definire il circuito che ne imposta il valore opportuno.

Consideriamo ora le istruzioni RISC-V di caricamento di un registro (load) e di trasferimento alla memoria (store). Esse hanno la forma generale  $ld \ x_1, offset(x_2)$ ,  $sd \ x_1, offset(x_2)$ , e calcolano un indirizzo di memoria sommando il contenuto del registro base, in questo caso  $x_2$ , al campo offset di 12 bit, contenuto nell'istruzione stessa, considerato dotato di segno. Se l'istruzione è una store il dato da memorizzare deve essere letto dal register file dove risiede in  $x_1$ , mentre se è una load il valore letto dalla memoria deve essere scritto all'interno del register file nel registro specificato, in questo caso  $x_1$ . Quindi, per eseguire queste istruzioni occorreranno sia il register file sia la ALU, come riportato in Figura 4.7.



**Figura 4.7** I due elementi richiesti per implementare le operazioni di tipo R sono il register file e la ALU. Il register file contiene tutti i registri e dispone di due porte di lettura e una di scrittura. Lo schema circuitale di un register file a più porte è discusso nel paragrafo A.8 dell'Appendice A . In ogni istante il register file fornisce in uscita il contenuto dei registri corrispondenti agli ingressi "Registro di lettura", senza richiedere alcun segnale di controllo; viceversa, la scrittura di un registro deve essere indicata esplicitamente asserendo il segnale di controllo della scrittura, "RegWrite". Si ricordi che la scrittura avviene sul fronte attivo del clock, per cui tutti gli ingressi interessati (cioè il dato da scrivere, il numero di registro e il segnale di controllo) devono essere validi sul fronte attivo del clock. Essendo la scrittura attiva sui fronti, l'architettura implementata nel nostro progetto può tranquillamente leggere e scrivere lo stesso registro nello stesso ciclo di clock: la lettura fornirà il dato contenuto nel registro scritto in un ciclo di clock precedente, mentre il valore scritto potrà essere letto a partire dal ciclo di clock successivo. Gli ingressi che specificano al register file il numero d'ordine dei registri hanno ampiezza di 5 bit, mentre i bus che trasportano i dati sono ampi 64 bit. L'operazione che la ALU deve eseguire viene specificata dal segnale "Operazione ALU", che è ampio 4 bit per la ALU descritta nell'Appendice A . L'uscita "Zero" della ALU, verrà utilizzata tra breve per implementare i salti condizionati.



**Figura 4.8** Le due unità necessarie per l'implementazione delle istruzioni di load e di store: oltre al register file e alla ALU mostrate in Figura 4.7, sono richieste anche un'unità di memoria dati e un'unità di estensione del segno. L'unità di memoria è un elemento di stato che ha come ingressi l'indirizzo del dato e il dato da scrivere, e possiede una sola uscita per il dato letto. I segnali di controllo per la lettura e la scrittura sono separati, anche se all'interno di un ciclo di clock solo uno dei due può essere asserito. L'unità di memoria ha bisogno di un segnale di lettura esplicito, dal momento che, diversamente da quanto avviene per il register file, leggere il contenuto di un indirizzo non valido può causare dei problemi, come vedremo nel Capitolo 5. L'unità di estensione del segno riceve in ingresso un numero a 12 bit per una load, una store e un salto condizionato all'uguaglianza, e lo fornisce in uscita esteso a 64 bit dotato di segno (vedi Cap. 2). Si suppone che la scrittura della memoria sia attiva su un fronte del segnale di clock. I chip di memoria standard, in realtà, possiedono un segnale di abilitazione alla scrittura che permette, quando asserito, l'esecuzione dell'operazione di scrittura. Sebbene il segnale di abilitazione alla scrittura nei chip di memoria reali non sia sensibile ai fronti, l'unità di elaborazione attiva sui fronti presentata in questo capitolo può essere facilmente adattata per funzionare anche con i chip di memoria reali. Per una discussione più approfondita sul funzionamento dei chip di memoria reali vedi il paragrafo A.8 dell'Appendice A.

**Estensione del segno:** aumento della dimensione di un dato ottenuta replicando il bit più significativo, che rappresenta il segno del dato originale, nei bit più significativi del dato di destinazione.

**Indirizzo di destinazione del salto:** l'indirizzo che viene specificato in un salto condizionato; esso diventerà il nuovo indirizzo contenuto nel program counter (PC) se il salto condizionato viene eseguito. Nell'architettura RISC-V, l'indirizzo di destinazione di un salto è ottenuto mediante la somma del campo offset contenuto nell'istruzione stessa con l'indirizzo dell'istruzione di salto condizionato.

Inoltre, saranno necessarie anche un'unità per l'**estensione del segno** del campo offset su 12 bit, contenuto nell'istruzione, da estendere a 64 bit, con segno, e un'unità di memoria dati in cui scrivere o da cui leggere il dato. La memoria dati viene scritta dalle istruzioni di store, e sarà quindi dotata dei segnali di controllo sia di lettura sia di scrittura; riceverà in ingresso, inoltre, l'indirizzo e il dato che deve essere scritto. La Figura 4.8 mostra questi due nuovi elementi.

L'istruzione `beq` ha tre operandi: due registri il cui contenuto viene confrontato per determinare se è uguale, e un offset di 16 bit utilizzato per calcolare l'**indirizzo di destinazione del salto** a partire dall'indirizzo dell'istruzione stessa. La forma di questa istruzione è `beq x1, x2, offset`; per eseguirla, occorre calcolare l'indirizzo di destinazione del salto, sommando il campo offset dell'istruzione (dopo averlo esteso a 32 bit con segno) al PC. Nella definizione delle istruzioni di salto condizionato ci sono due dettagli (vedi Cap. 2) ai quali è bene prestare attenzione.

- L'architettura dell'insieme delle istruzioni specifica che l'indirizzo di base per il calcolo dell'indirizzo di salto è quello dell'istruzione di salto stessa.
- L'architettura stabilisce anche che il campo offset sia spostato di 1 bit a sinistra, in modo tale che l'offset non codifichi lo spiazzamento in numero di byte, ma in numero di mezze parole. Di fatto, tale spostamento aumenta lo spazio di indirizzamento dell'offset di un fattore 2 rispetto alla codifica dello spiazzamento in byte.

Per poter gestire quest'ultima complicazione è necessario far scorrere il campo offset di 1 bit a sinistra.

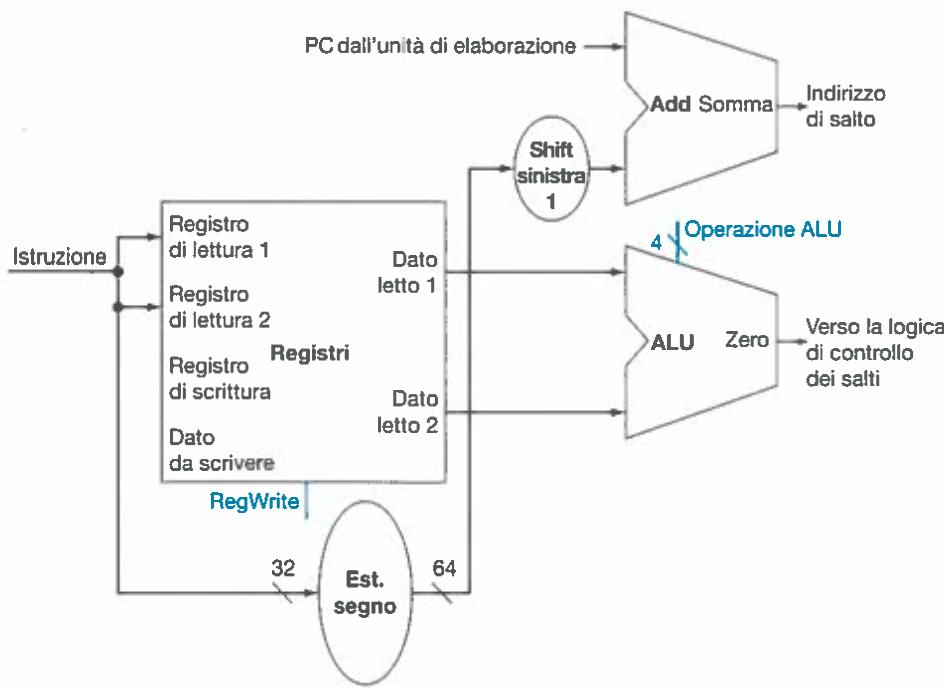
Oltre a calcolare l'indirizzo di destinazione del salto, bisogna determinare se l'istruzione da eseguire dopo sia quella presente nella posizione di memo-

ria successiva oppure quella contenuta all'indirizzo di destinazione del salto. Quando la condizione di salto è vera (cioè quando i due operandi sono uguali), l'indirizzo di destinazione del salto diventa il nuovo valore del PC e si parla di **salto condizionato eseguito** (*branch taken*, salto preso); se il contenuto degli operandi è diverso, il valore del PC incrementato di 4 sostituisce il valore corrente nel PC, proprio come in tutte le altre istruzioni; in questo caso si parla di **salto non eseguito** (*branch not taken*, salto non preso).

Per i salti condizionati, quindi, l'unità di elaborazione deve eseguire due operazioni: calcolare l'indirizzo di destinazione del salto e confrontare il contenuto di due registri. Inoltre, come si vedrà tra poco, i salti hanno anche un effetto sulla parte dell'unità di elaborazione incaricata di prelevare le istruzioni. La parte dell'unità di elaborazione che gestisce i salti è visualizzata in Figura 4.9. Per il calcolo dell'indirizzo di destinazione, si utilizzano un'unità di estensione del segno, analoga a quella mostrata in Figura 4.8, e un sommatore. Per eseguire il confronto possiamo utilizzare il register file mostrato in Figura 4.7a, che ci permette di leggere il contenuto dei due registri contenenti gli operandi; osserviamo che, in questo caso, non è richiesto di scrivere il risultato nel register file. Il confronto tra i due operandi viene fatto dalla ALU descritta nell'Appendice A. Dato che la ALU fornisce un segnale in uscita che indica se il risultato è pari a 0, si possono inviare i due operandi alla ALU specificando tramite i segnali di controllo di effettuare una sottrazione: se l'uscita Zero della ALU viene asserita, i due operandi sono uguali. Segnaliamo esplicitamente

**Salto condizionato eseguito:** un salto condizionato nel quale la condizione di salto è soddisfatta; in questo caso nel program counter (PC) viene scritto l'indirizzo di destinazione del salto. Le istruzioni di salto incondizionato possono essere assimilate a dei salti condizionati che vengono sempre eseguiti.

**Salto condizionato non eseguito:** un salto condizionato in cui la condizione di salto non è soddisfatta; in questo caso nel program counter (PC) viene scritto l'indirizzo dell'istruzione successiva all'istruzione di salto condizionato.



**Figura 4.9** Per i salti condizionati l'unità di elaborazione impiega la ALU per valutare la condizione di salto e un sommatore a parte per determinare l'indirizzo di destinazione del salto, calcolato come somma del PC e dei 12 bit dell'istruzione, considerati dotati di segno (lo spiazzamento del salto), fatti scorrere a sinistra di 1 bit ed estesi a 32 bit con segno. L'unità etichettata "Shift a sinistra 1" corrisponde semplicemente a un particolare instradamento dei segnali che rappresentano lo spiazzamento; l'instradamento è effettuato in modo tale che, a valle dell'estensione del segno, venga aggiunto un segnale uguale a  $0_{\text{segno}}$  all'estremità inferiore del campo offset: poiché l'entità dello spostamento è costante, non è richiesto alcun componente dedicato. Sapendo che lo spiazzamento viene esteso a partire dal 12° bit, questa operazione di scorrimento scatterà solamente dei "bit di segno". La logica di controllo deciderà, in funzione dell'uscita Zero della ALU, se nel PC debba essere caricato il contenuto del PC incrementato di 4, oppure l'indirizzo di destinazione del salto.

che, anche se l'uscita Zero indica sempre che il risultato di un'operazione è 0, noi la utilizzeremo solamente per il test di uguaglianza dei salti condizionati. In seguito, vedremo precisamente come debbano essere collegati i segnali di controllo della ALU all'interno dell'unità di elaborazione.

L'istruzione di salto condizionato opera sommando al PC i 12 bit dell'istruzione fatti scorrere a sinistra di 1 posizione. Questo scorrimento si ottiene semplicemente concatenando degli 0 all'offset dell'istruzione di salto condizionato, come già descritto nel Capitolo 2.

### Progettazione di un'unità di elaborazione unificata

Ora che abbiamo esaminato i componenti dell'unità di elaborazione necessari a eseguire i diversi tipi di istruzioni, possiamo combinarli in una singola unità di elaborazione aggiungendo i segnali di controllo opportuni. L'unità di elaborazione più semplice cercherà di eseguire tutte le istruzioni in un singolo ciclo di clock. Questo comporta che nessuna risorsa dell'unità di elaborazione possa essere utilizzata più di una volta per ogni istruzione e che ogni elemento funzionale utilizzato più di una volta debba essere duplicato: servirà quindi una memoria delle istruzioni separata da quella dei dati. Anche se dovremo duplicare alcune unità funzionali, molti elementi potranno essere condivisi dal flusso di esecuzione di istruzioni diverse.

La condivisione da parte di istruzioni di tipo diverso di uno stesso elemento potrebbe richiedere la presenza di collegamenti multipli agli ingressi di quell'elemento; in questo caso occorrerà introdurre un multiplexer, con i suoi segnali di controllo, per poter selezionare l'input corretto all'unità funzionale tra tutti gli input possibili.

### Costruzione di un'unità di elaborazione

#### ESEMPIO

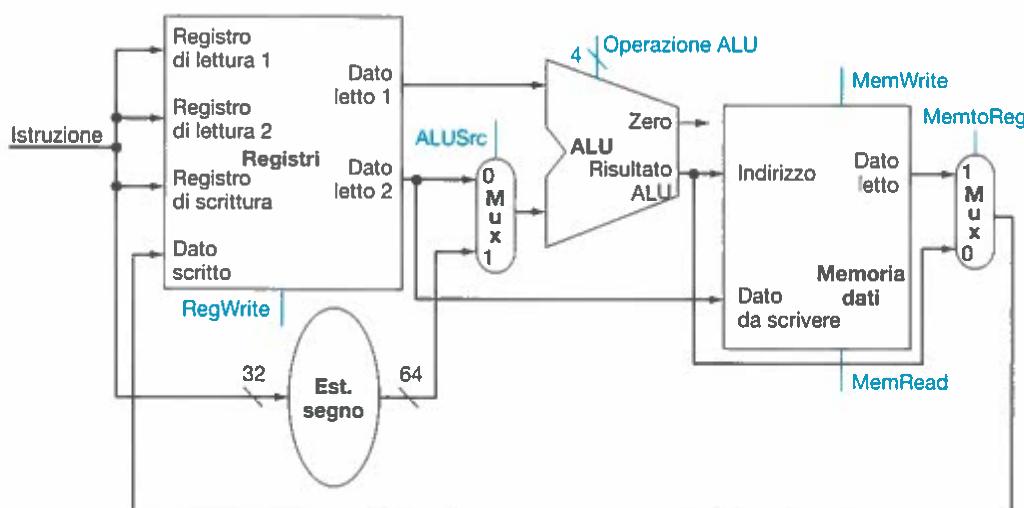
L'unità di elaborazione per le istruzioni aritmetico-logiche (o di tipo R) e quella per le istruzioni di accesso alla memoria sono molto simili. Le differenze principali sono le seguenti.

- Le istruzioni aritmetico-logiche utilizzano la ALU con gli ingressi provenienti da due registri. Le istruzioni di memoria utilizzano la ALU per calcolare l'indirizzo della memoria dati, ma prendono il secondo ingresso alla ALU dal campo offset a 12 bit dell'istruzione stessa, con l'estensione del segno.
- Il valore da scrivere nel registro destinazione proviene dalla ALU (per le istruzioni di tipo R) o dalla memoria (per l'istruzione load).

Si mostri come costruire un'unità di elaborazione relativa alla parte di calcolo per le istruzioni di accesso a memoria e per le istruzioni aritmetico-logiche, utilizzando un solo register file e una sola ALU per entrambi i tipi di istruzioni; si introducano i multiplexer necessari.

#### SOLUZIONE

Per creare un'unità di elaborazione che utilizzi un solo register file e una sola ALU, dobbiamo implementare un meccanismo che possa prelevare il secondo operando della ALU da due sorgenti diverse; lo stesso dobbiamo fare per il dato da scrivere nel register file. Occorre quindi inserire un primo multiplexer all'ingresso della ALU e un secondo multiplexer all'ingresso del register file. La Figura 4.10 mostra la parte di esecuzione dell'unità di elaborazione che supporta questi due tipi di istruzioni.

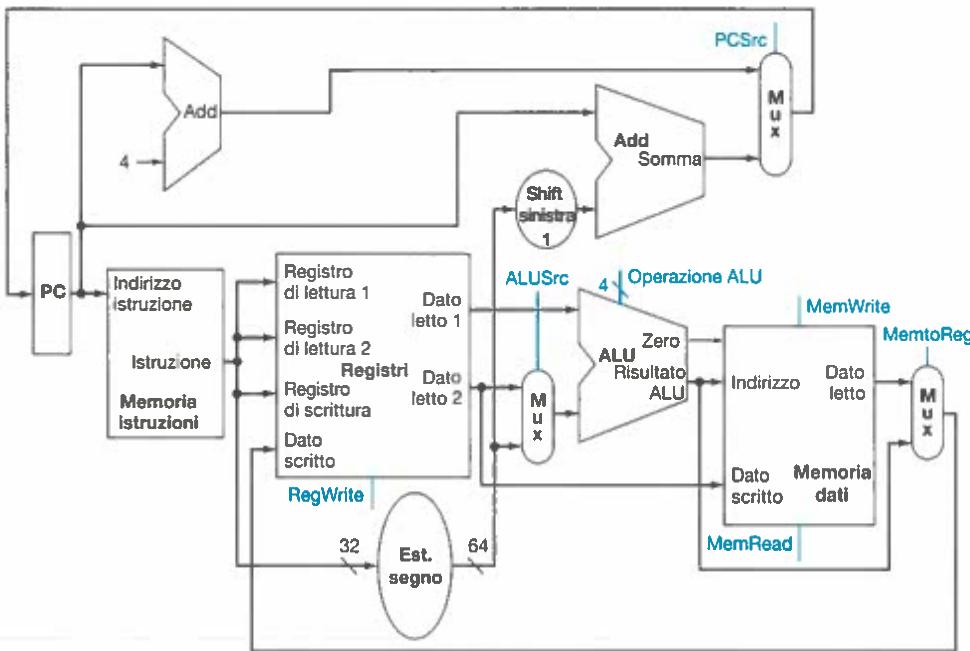


**Figura 4.10** Unità di elaborazione per le istruzioni di accesso alla memoria e per le istruzioni di tipo R. Questo esempio mostra come si possa ottenere un'unica unità di elaborazione utilizzando i componenti delle Figure 4.7 e 4.8 con l'aggiunta di multiplexer. In particolare, sono necessari due multiplexer.

A questo punto, possiamo combinare tra loro tutti i pezzi introdotti finora per costruire una semplice unità di elaborazione per il nucleo dell'architettura RISC-V: possiamo aggiungere la parte relativa al prelevamento delle istruzioni (fetch, Figura 4.6), la parte di esecuzione delle istruzioni di tipo R e delle istruzioni di accesso alla memoria (Figura 4.10) e, infine, la parte dell'unità di elaborazione per i salti condizionati (Figura 4.9). La Figura 4.11 mostra l'unità di elaborazione che otteniamo dalla composizione delle diverse parti. L'istruzione di salto condizionato utilizza la ALU principale per il confronto degli operandi, contenuti in due registri; occorre perciò continuare a utilizzare il sommatore di Figura 4.9 per calcolare l'indirizzo di destinazione del salto. Un altro multiplexer consente di selezionare quale tra l'indirizzo dell'istruzione successiva ( $PC + 4$ ) e l'indirizzo di destinazione del salto debba essere scritto nel PC.

## Autovalutazione

- I. Quale delle seguenti affermazioni è corretta per un'istruzione di load (Figura 4.10)?
  - a. Il segnale **MemtoReg** deve essere asserito per fare in modo che il dato proveniente dalla memoria sia inviato al register file.
  - b. Il segnale **MemtoReg** deve essere asserito per fare in modo che il registro di destinazione corretto sia inviato al register file.
  - c. Non è importante il modo in cui viene impostato il segnale **MemtoReg** per le istruzioni di load.
- II. L'unità di elaborazione a singolo ciclo, descritta a grandi linee in questo paragrafo, deve avere due memorie separate per le istruzioni e i dati, perché:
  - a. il formato dei dati e delle istruzioni nel RISC-V è diverso e quindi occorrono due memorie distinte;
  - b. avere due memorie separate risulta meno costoso;
  - c. il processore opera in un singolo ciclo di clock e quindi non può utilizzare una memoria con una sola porta di uscita per due accessi diversi all'interno dello stesso ciclo.



**Figura 4.11** Una semplice unità di elaborazione per il nucleo dell'architettura RISC-V che combina tra loro gli elementi richiesti dalle istruzioni appartenenti ai diversi tipi. I componenti sono mostrati nelle Figure 4.6, 4.9 e 4.10. Questa unità di elaborazione può eseguire le istruzioni di base (load/store di un registro, operazioni con la ALU e salti condizionati) in un singolo ciclo di clock. Per integrare i salti condizionati si è resa necessaria solamente l'aggiunta di un altro multiplexer.

Ora che abbiamo completato quest'unità di elaborazione dati elementare, possiamo aggiungere l'unità di controllo: essa dovrà poter accettare dei valori di input e generare un segnale di scrittura per ciascun elemento di stato, un segnale di selezione per ciascun multiplexer e i segnali di controllo per la ALU. Il controllo della ALU è, per molti versi, particolare e conviene perciò progettarlo prima delle altre parti dell'unità di controllo.

**Approfondimento.** L'unità logica che genera il campo immediato deve scegliere se estendere il segno del campo a 12 bit contenuto nei bit 31:20 dell'istruzione per le load, nei bit 31:25 e 11:7 per le istruzioni di store o i bit 31, 7, 30:25 e 11:8 per le istruzioni di salto condizionato. Dato che l'input è sempre di 32 bit per tutte le istruzioni, l'unità logica può utilizzare i bit del codice operativo dell'istruzione per selezionare il campo appropriato. Nei RISC-V il bit 6 del codice operativo è pari a 0 per le istruzioni di trasferimento dati e 1 per le istruzioni di salto condizionato, il bit del codice operativo è pari a 0 per le istruzioni di load e 1 per le istruzioni di store. Quindi i bit 5 e 6 controllano un multiplexer 3:1, contenuto all'interno dell'unità logica di generazione del campo immediato che seleziona il gruppo di 12 bit appropriato alle istruzioni di load, store e salto condizionato.

## 4.4 Uno schema semplice di implementazione

In questo paragrafo esamineremo quella che può essere considerata una semplice implementazione hardware possibile del sottoinsieme RISC-V che abbiamo scelto; utilizzeremo l'unità di elaborazione dati costruita nel paragrafo precedente, a cui aggiungeremo la relativa semplice unità di controllo. Questa semplice implementazione consente di eseguire le istruzioni *load doubleword* (ld), *store doubleword* (sd), *branch if equal* (beq) e le istruzioni aritmetico-logiche add, sub, and e or.

## Unità di controllo della ALU

La ALU del RISC-V descritta nell'Appendice A  prevede le 6 seguenti combinazioni dei suoi quattro input di controllo:

Linee di controllo della ALU	Operazione
0000	AND
0001	OR
0010	somma
0110	sottrazione

A seconda del tipo di istruzione, la ALU eseguirà una di queste quattro operazioni. Per le istruzioni load e store la ALU deve eseguire una somma per calcolare l'indirizzo di memoria; per le istruzioni di tipo R deve invece eseguire una delle quattro operazioni (AND, OR, somma o sottrazione) in funzione del valore dei 7 bit del campo *funz7* (bit 31:25) e dei 3 bit del campo *funz3* (bit 14:12) dell'istruzione (*vedi Cap. 2*). Per l'istruzione di salto condizionato, la ALU deve eseguire una sottrazione tra i due operandi e controllare se il risultato è 0.

I 4 bit di controllo della ALU possono essere generati utilizzando una piccola unità di controllo che riceve in ingresso i campi *funz7* e *funz3* dell'istruzione e un campo di controllo su 2 bit, chiamato *ALUOp*; *ALUOp* indica se l'operazione da eseguire è una somma (*ALUOp* = 00) per le istruzioni di load e store, una sottrazione (01) per le *beq*, oppure indica se l'operazione viene determinata dal contenuto dei campi *funz7* e *funz3* (10). L'uscita dell'unità di controllo della ALU è un segnale di 4 bit che controlla direttamente la ALU, generando una delle combinazioni di 4 bit mostrate nella tabella precedente.

Nella tabella di Figura 4.12 sono riportati i valori da assegnare agli ingressi di controllo della ALU in funzione del segnale di controllo *ALUOp* di 2 bit e del contenuto dei campi *funct7* e *funct3*; più avanti vedremo come l'unità di controllo principale imposta i bit di *ALUOp*.

L'utilizzo di livelli multipli di decodifica (in questo caso l'unità di controllo principale imposta i bit *ALUOp*, poi utilizzati come ingressi dell'unità di controllo della ALU che, a sua volta, genera i segnali effettivi di controllo della ALU) è una tecnica di progettazione utilizzata di frequente. Più livelli di controllo possono ridurre le dimensioni dell'unità di controllo principale;

Codice operativo istruzione	ALUOp	Operazione eseguita dall'istruzione	Campo <i>funz7</i>	Campo <i>funz3</i>	Operazione dell'ALU	Ingresso di controllo alla ALU
ld	00	load di 1 parola doppia	XXXXXX	XXX	somma	0010
sd	00	store di 1 parola doppia	XXXXXX	XXX	somma	0010
beq	01	salto condizionato all'uguaglianza	XXXXXX	XXX	sottrazione	0110
Tipo R	10	add	0000000	000	somma	0010
Tipo R	10	sub	0100000	000	sottrazione	0110
Tipo R	10	and	0000000	111	AND	0000
Tipo R	10	or	0000000	110	OR	0001

**Figura 4.12** Il valore dei bit di controllo della ALU dipende, per le istruzioni di tipo R, dai bit di *ALUOp* e dai diversi codici funzione. Il codice operativo, riportato nella prima colonna, determina il valore dei bit di *ALUOp*. Si noti che quando il codice *ALUOp* vale 00 o 01, l'operazione che la ALU deve eseguire non dipende dal contenuto dei campi *funz7* e *funz3*; in questo caso si dice che il contenuto del campo funzione è *indifferente*, e viene indicato con delle X. Invece, quando *ALUOp* vale 10 vengono utilizzati i campi *funz7* e *funz3* per determinare il valore degli ingressi di controllo alla ALU. Vedi l'Appendice A  per maggiori dettagli.

ALUOp		Campo funz7									Campo funz3			Operazione
ALUOp1	ALUOp2	I[31]	I[30]	I[29]	I[28]	I[27]	I[26]	I[25]	I[14]	I[13]	I[12]			
0	0	X	X	X	X	X	X	X	X	X	X	0010		
X	1	X	X	X	X	X	X	X	X	X	X	0110		
1	X	0	0	0	0	0	0	0	0	0	0	0010		
1	X	0	1	0	0	0	0	0	0	0	0	0110		
1	X	0	0	0	0	0	0	0	1	1	1	0000		
1	X	0	0	0	0	0	0	0	1	1	0	0001		

**Figura 4.13** Tabella della verità dei 4 bit di controllo della ALU (detti “operazione”). Gli ingressi della tabella della verità sono i campi ALUOp e i campi funz. Vengono riportate solamente le righe della tabella per cui i segnali di controllo della ALU devono essere asseriti. Sono mostrate anche alcune righe associate a valori indifferenti: per esempio, ALUOp non utilizza il codice 11, quindi la tabella della verità conterrà le combinazioni 1X e X1 di ALUOp anziché 10 e 01. Anche se mostriamo tutti e dieci i bit dei campi funz, notate che gli unici bit che assumono valori diversi per le quattro istruzioni di formato R sono i bit 30, 14, 13 e 12. Quindi, occorrono solamente questi quattro bit dei campi funzione come input all’unità di controllo della ALU.

inoltre, l’uso di più unità di controllo di dimensioni ridotte può ridurre la latenza dell’unità di controllo. Queste ottimizzazioni sono importanti, poiché spesso l’unità di controllo è l’elemento critico per la definizione della durata del ciclo di clock.

Ci sono diversi modi per implementare la corrispondenza tra i 2 bit del campo ALUOp e i bit dei campi funz con i 4 bit di controllo della ALU che selezionano l’operazione. Dato che interessa solamente un piccolo sottoinsieme dei possibili valori che possono essere assunti dai campi funz, e che questi ultimi vengono utilizzati solamente quando i bit di ALUOp sono pari a 10, possiamo ricorrere a un piccolo circuito logico che riconosca il sottoinsieme dei valori possibili e imposti di conseguenza i bit di controllo della ALU in modo corretto.

Come primo passo nella progettazione di questo circuito, è utile costruire una *tabella della verità* che contenga le combinazioni di interesse dei campi funz e del segnale di ALUOp, come mostrato in Figura 4.13. Tale **tabella della verità** mostra i 4 bit di controllo della ALU in funzione di questi tre campi in ingresso. Dato che la tabella della verità completa è molto lunga e che non siamo interessati al valore che assumono i segnali di controllo della ALU per le combinazioni che non servono (cioè, sono indifferenti), vengono riportate solamente le righe in corrispondenza delle quali deve essere specificato un preciso valore dei segnali di controllo. Nel resto di questo capitolo adotteremo la convenzione di riportare nelle tabelle della verità solamente gli elementi che devono essere asseriti, tralasciando quelli che sono non asseriti o che sono indifferenti. Questa scelta ha uno svantaggio, che sarà discusso nel paragrafo B.2 dell’Appendice B.

Poiché in molti casi non serve conoscere i valori assunti da alcuni input, essi sono considerati **termini indifferenti**, il che consente di compattare le tabelle. Un termine indifferente di una tabella della verità, rappresentato da una X in una delle colonne associate alle variabili in ingresso, indica che l’uscita non dipende dal valore dell’ingresso associato a quella colonna. Per esempio, quando i bit di ALUOp valgono 00, come nella prima riga della tabella di Figura 4.13, i segnali di controllo della ALU vengono sempre impostati a 0010, indipendentemente dal contenuto dei campi funz; in tal caso gli ingressi dei campi funzione sono indifferenti. Incontreremo in seguito altri tipi di termini indifferenti; il lettore che non abbia familiarità con il concetto di termine indifferente è invitato a consultare l’Appendice A per maggiori dettagli.

Una volta costruita la tabella della verità, questa può essere ottimizzata e convertita in un circuito costituito da porte logiche connesse tra loro: tale processo

**Tabella della verità:** utilizzata nella logica booleana, è la rappresentazione di una funzione logica. In una tabella della verità vengono elencate tutte le possibili combinazioni dei valori che possono assumere le variabili in ingresso e, per ogni combinazione, viene associata l’uscita relativa.

**Termine indifferente:** una variabile di ingresso di una funzione logica che non ha effetto sull’uscita della funzione. I termini indifferenti possono assumere un valore arbitrario.

è di natura completamente meccanica. Perciò, invece di mostrare qui questi passi, abbiamo preferito inserire la descrizione del processo di ottimizzazione e il risultato finale nel paragrafo B.2 dell'Appendice B .

## Progettazione dell'unità di controllo principale

Ora che abbiamo terminato di descrivere la progettazione di una ALU che utilizza come ingressi di controllo il campo funzione e un segnale su 2 bit, possiamo tornare ad analizzare il resto dell'unità di controllo. Iniziamo con l'identificare i diversi campi delle istruzioni e i segnali di controllo richiesti dall'unità di elaborazione mostrata in Figura 4.11. Per capire come connettere i campi di un'istruzione all'unità di elaborazione, conviene rivedere prima il formato delle tre classi di istruzioni: tipo R, salto condizionato e load/store. Questi formati sono mostrati in Figura 4.14.

Ci sono alcune considerazioni importanti da fare in merito a questo formato delle istruzioni, su cui ci baseremo in fase di progettazione:

- Il campo **codice operativo** o **codop** (vedi Cap. 2) è sempre contenuto nei bit 6:0. A seconda del codice operativo, il campo funz3 (bit 14:12) e funz7 (31:25) servono come campi di estensione del codice operativo.
- Il primo registro operando si trova sempre nei bit in posizione 19:15 (rs1) per le istruzioni di tipo R e le istruzioni di salto condizionato. Questo campo specifica anche il registro base per le istruzioni di load e di store.
- Il secondo registro operando si trova sempre nei bit in posizione 24:20 (rs2) per le istruzioni di tipo R e le istruzioni di salto condizionato. Questo campo specifica anche il registro il cui contenuto viene copiato in memoria nelle istruzioni di store.
- Il secondo operando può anche essere costituito dai 12 bit di offset delle istruzioni di branch o di load/store.
- Il registro destinazione si trova sempre nei bit in posizione 11:7 (rd) per le istruzioni di tipo R e di load.

**Codice operativo (codop):** il campo che denota il tipo di operazione e il formato di un'istruzione.

Il primo principio di progettazione introdotto nel Capitolo 2 – *La semplicità favorisce la regolarità* – risulta particolarmente utile nella progettazione dell'unità di controllo.

Nome (posizione dei bit)	Campi					
	31:25	24:20	19:15	14:12	11:7	6:0
(a) Tipo R	funz7	rs2	rs1	funz3	rd	codop
(b) Tipo I	immediate[11:0]		rs1	funz3	rd	codop
(c) Tipo S	immed[11:5]	rs2	rs1	funz3	immed[4:0]	codop
(d) Tipo SB	immed[12,10:5]	rs2	rs1	funz3	immed[4:1,11]	codop

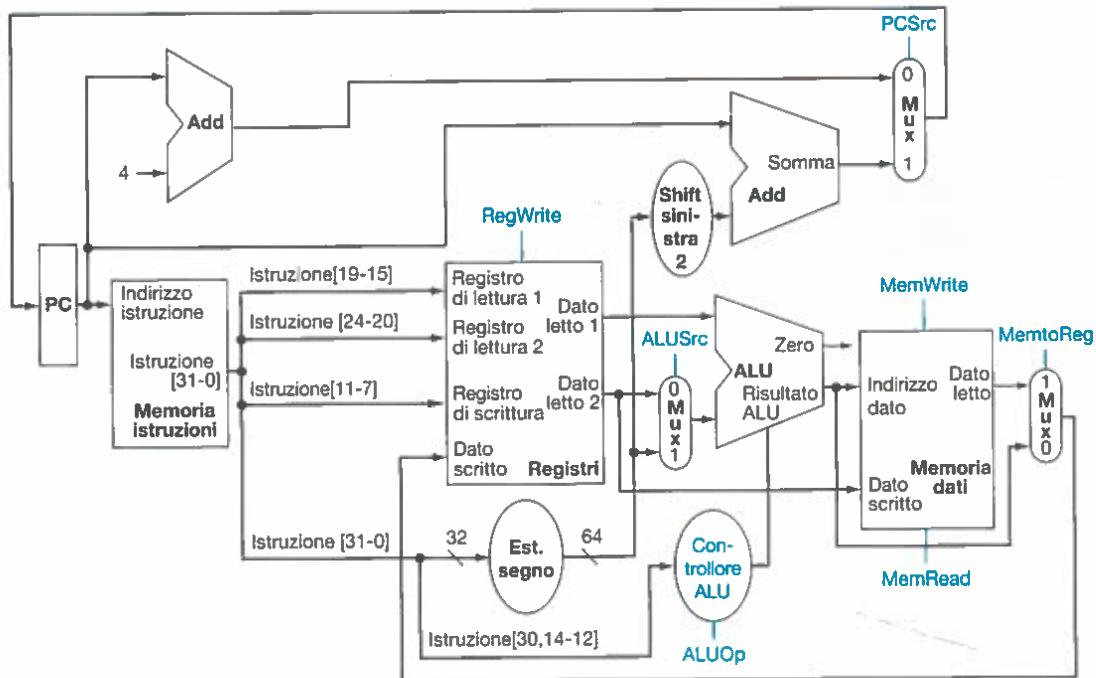
**Figura 4.14** Le quattro classi di istruzioni (tipo R, load, store e salto condizionato) utilizzano quattro diversi formati di istruzione. (a) Formato delle istruzioni di tipo R (codice operativo = 51<sub>dec</sub>). Queste hanno tre operandi contenuti nei registri: rs1, rs2 e rd; i campi rs1 e rs2 contengono il numero dei due registri sorgenti e rd contiene il numero del registro destinazione. L'operazione che la ALU deve eseguire è codificata nei campi funz3 e funz7 e verrà decodificata dall'unità di controllo della ALU, progettata nel paragrafo precedente. Le istruzioni di tipo R implementate sono add, sub, and e or. (b) Formato di tipo I delle istruzioni di load (codice operativo = 3<sub>dec</sub>) di tipo I. Il registro rs1 è il registro base il cui contenuto viene sommato al campo immediato di 12 bit per ottenere l'indirizzo del dato in memoria. Il campo rd è il registro destinazione per il valore letto dalla memoria. (c) Formato di tipo S delle istruzioni di store (codice operativo = 35<sub>dec</sub>) di tipo S. Il registro rs1 è il registro base il cui contenuto viene sommato al campo immediato di 12 bit per ottenere l'indirizzo del dato in memoria (il campo immediato è suddiviso in un gruppo di 7 e in un gruppo di 5 bit). Il campo rs2 è il registro sorgente il cui valore viene copiato nella memoria. (d) Formato di tipo SB delle istruzioni branch if equal (codice operativo = 99<sub>dec</sub>). I registri rs1 e rs2 vengono confrontati. Il campo indirizzo immediato di 12 bit viene preso, il suo bit di segno esteso, fatto scorrere a sinistra di una posizione e sommato al PC per calcolare l'indirizzo di destinazione del salto.

Utilizzando le osservazioni riportate sopra, possiamo ora aggiungere all'unità di elaborazione elementare che abbiamo costruito le etichette sui campi dell'istruzione e un multiplexer aggiuntivo per l'ingresso "registro di scrittura" del register file (la Figura 4.15 mostra queste aggiunte). Sono stati inoltre inseriti: il blocco che rappresenta l'unità di controllo della ALU, i segnali di scrittura per gli elementi di stato, il segnale di lettura della memoria dati e i segnali di controllo dei multiplexer. Dato che tutti i multiplexer hanno due ingressi, ciascuno di essi richiede un solo segnale di controllo.

La Figura 4.15 mostra i sei segnali di controllo di 1 bit ciascuno e il segnale di controllo ALUOp di 2 bit. Abbiamo già spiegato come funziona questo segnale, ed è utile definire informalmente come debbano operare gli altri segnali di controllo prima di determinare il modo in cui devono essere impostati durante l'esecuzione delle istruzioni. La Figura 4.16 descrive il funzionamento di questi sei segnali di controllo.

Ora che abbiamo visto il funzionamento di ciascuno dei segnali di controllo, possiamo passare a esaminare come debbano essere impostati. L'unità di controllo può impostare tutti i segnali tranne uno, basandosi esclusivamente sul codice operativo dell'istruzione stessa. L'eccezione è costituita dal segnale di controllo PCSrc: esso deve essere asserito se l'istruzione è una branch if equal, decisione che l'unità di controllo è in grado di prendere, *ma anche* se l'uscita Zero della ALU, utilizzata per il confronto di uguaglianza, è vera. Per generare il segnale PCSrc occorre dunque collegare in AND un segnale proveniente dall'unità di controllo, che verrà chiamato *Branch*, con l'uscita Zero della ALU.

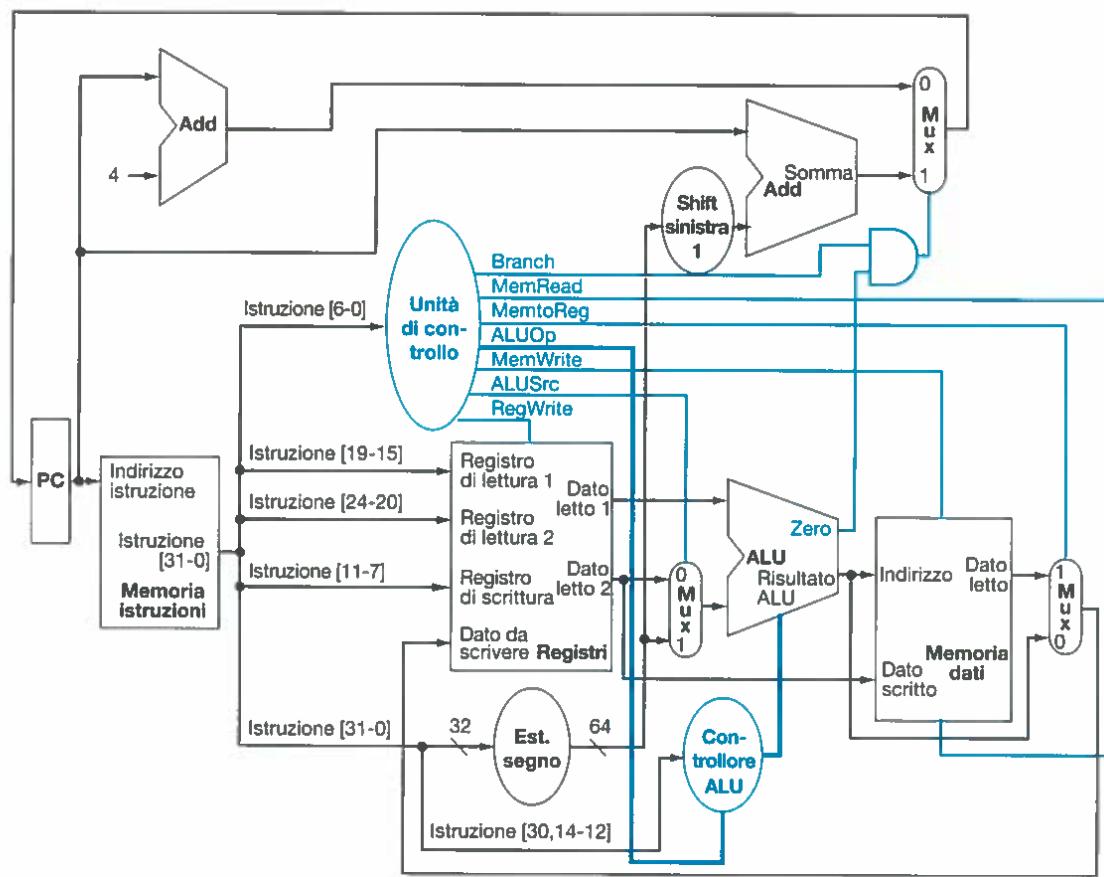
Questi otto segnali di controllo, i sei di Figura 4.16 e i due di ALUOp, possono ora essere impostati sulla base dei segnali di ingresso all'unità di controllo, che sono i sette bit del codice operativo (bit 6:0). La Figura 4.17 mostra l'unità di elaborazione completa, con l'unità di controllo e i relativi segnali.



**Figura 4.15** L'unità di elaborazione di Figura 4.11 con tutti i multiplexer e i segnali di controllo necessari. I segnali di controllo sono mostrati in blu. È stato introdotto anche il blocco che rappresenta il controllore della ALU. Il PC non richiede un segnale di controllo della scrittura esplicito, dal momento che viene scritto al termine di ogni ciclo di clock; la logica di controllo dei salti determina se nel PC debba essere scritto il valore precedente del PC incrementato di 4 oppure l'indirizzo di destinazione del salto.

Nome del segnale	Effetto quando non asserito	Effetto quando asserito
RegWrite	Nullo	Il dato viene scritto nel register file nel registro individuato dal numero del registro di scrittura
ALUSrc	Il secondo operando della ALU proviene dalla seconda uscita del register file (Dato letto 2)	Il secondo operando della ALU proviene dall'estensione del segno dei 12 bit del campo immediato dell'istruzione
PCSrc	Nel PC viene scritta l'uscita del sommatore che calcola il valore di $PC + 4$	Nel PC viene scritta l'uscita del sommatore che calcola l'indirizzo di salto
MemRead	Nullo	Il dato della memoria nella posizione puntata dall'indirizzo viene inviato in uscita sulla linea "Dato letto"
MemWrite	Nullo	Il contenuto della memoria nella posizione puntata dall'indirizzo viene sostituito con il dato presente sulla linea "Dato scritto"
MemtoReg	Il dato inviato al register file per la scrittura proviene dalla ALU	Il dato inviato al register file per la scrittura proviene dalla Memoria Dati

**Figura 4.16 Effetto di ciascuno dei sette segnali di controllo.** Quando il segnale di controllo a 1 bit di un multiplexer a due vie è asserito, il multiplexer seleziona l'ingresso etichettato con 1; in caso contrario, cioè se il segnale di controllo non è asserito, il multiplexer seleziona l'ingresso etichettato con 0. Si ricordi che tutti gli elementi di stato ricevono il clock come ingresso implicito e che il clock controlla le operazioni di scrittura. Far commutare il clock mediante un segnale esterno a un elemento di stato potrebbe causare problemi di temporizzazione (vedi l'Appendice A  per la discussione di questo problema).



**Figura 4.17 L'unità di elaborazione costruita finora, con la sua unità di controllo.** L'input dell'unità di controllo è il campo a 7 bit dell'istruzione che costituisce il codice operativo. L'output dell'unità di controllo è formata da due segnali a 1 bit utilizzati per controllare i multiplexer (ALUSrc e MemtoReg), tre segnali a 1 bit per controllare lettura e scrittura del register file e della memoria dati (RegWrite, MemRead e MemWrite), un segnale a 1 bit utilizzato come segnale di controllo per i salti condizionati (Branch) e un segnale di controllo a 2 bit per la ALU (ALUOp). Una porta AND viene utilizzata per combinare il segnale di Branch e l'uscita Zero della ALU; l'uscita di questa porta determina da dove prendere il valore successivo del PC e genera il segnale di controllo PCSrc per il multiplexer posizionato in alto (Figura 4.15). Si noti che PCSrc viene derivato e non proviene direttamente dall'unità di controllo; per questo motivo, ometteremo di indicarlo nelle figure successive.

Istruzione	ALUSrc	MemtoReg	RegWrite	MemRead	MemWrite	Branch	ALUOp1	ALUOp0
Tipo R	0	0	1	0	0	0	1	0
ld	1	1	1	1	0	0	0	0
sd	1	X	0	0	1	0	0	0
beq	0	X	0	0	0	1	0	1

**Figura 4.18** Il valore delle linee di controllo è completamente determinato dal campo del codice operativo dell'istruzione. La prima riga della tabella corrisponde alle istruzioni in formato R (add, sub, and e or), per le quali i campi registro sorgente sono rs1 e rs2 e il campo registro destinazione è rd; questo stabilisce come impostare il segnale ALUSrc. Inoltre, le istruzioni di tipo R scrivono un registro (RegWrite = 1), ma non leggono né scrivono nella memoria dati. Quando il segnale di controllo Branch è impostato a 0, il PC viene aggiornato incondizionatamente a PC + 4; altrimenti, il contenuto del PC viene sostituito con l'indirizzo della destinazione del salto, ma solo quando anche l'uscita Zero della ALU è alta. Il campo ALUOp per le istruzioni di tipo R è impostato a 10 per indicare che i segnali di controllo della ALU dovranno essere generati a partire dai campi funz. La seconda e la terza riga di questa tabella riportano l'impostazione dei segnali di controllo per ld e sd, per le quali ALUSrc e ALUOp assumono il valore richiesto perché venga calcolato l'indirizzo in memoria, mentre MemRead e MemWrite determinano il tipo di accesso alla memoria. Infine, in una load, RegWrite viene impostato in modo tale che il dato letto dalla memoria venga scritto nel registro rd. Nell'istruzione di salto condizionato il campo ALUOp indica la sottrazione (01), utilizzata per controllare l'uguaglianza degli operandi. Si noti che il segnale MemtoReg risulta irrilevante quando il segnale RegWrite è uguale a 0: dal momento che l'istruzione non scrive nel register file, il valore del dato presente sulla porta di scrittura non viene utilizzato; quindi il valore di MemtoReg nelle ultime due righe viene sostituito con una X per indicare che è indifferente. La scelta di impostare questi segnali di controllo con un valore indifferente deve essere fatta dal progettista, poiché dipende dalla conoscenza del funzionamento dell'unità di elaborazione.

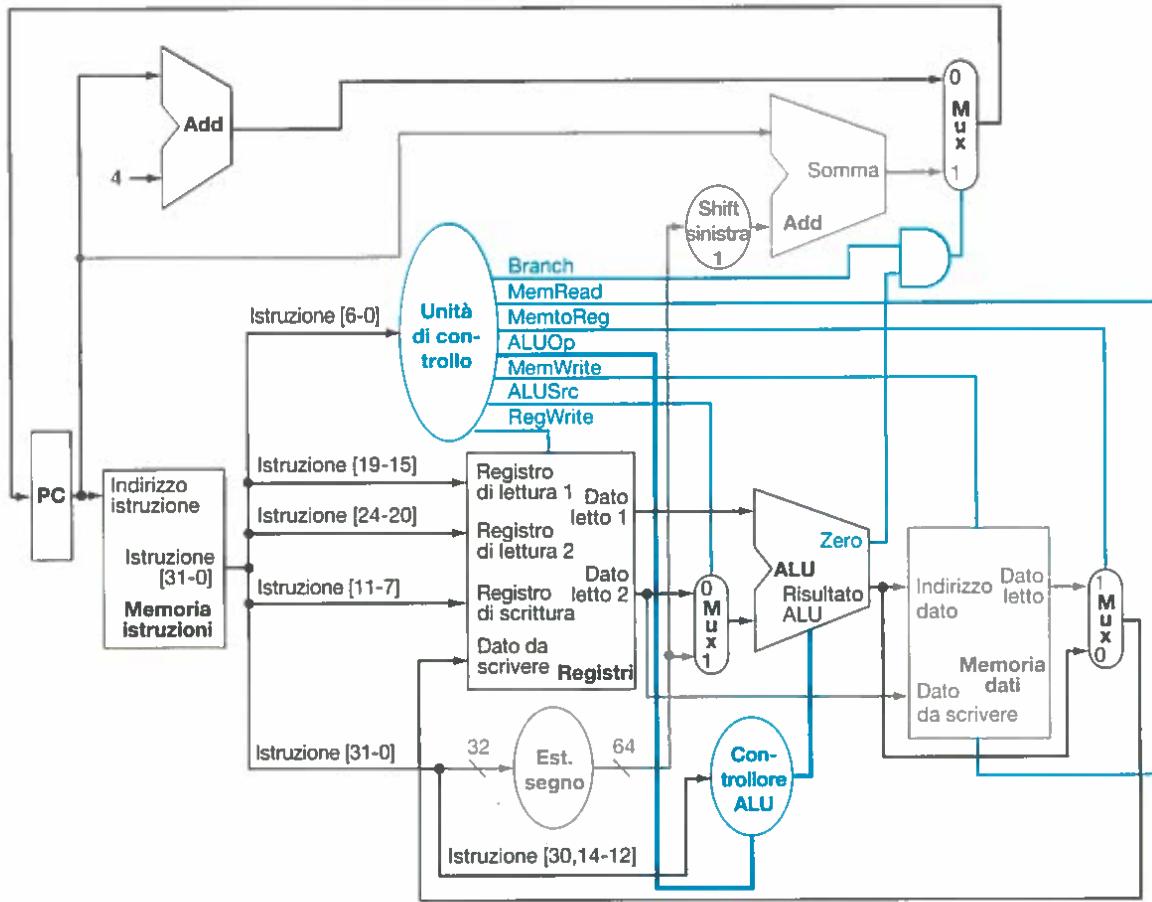
Prima di iniziare a scrivere un insieme di equazioni o una tabella della verità per specificare il modo di operare dell'unità di controllo, è utile cercare di definire informalmente il suo funzionamento. Dato che il valore di ciascuna linea di controllo dipende solamente dal codice operativo, possiamo dire se ciascun segnale debba valere 0, 1 o sia indifferente (X) per i diversi valori che il codice operativo può assumere. La Figura 4.18 riporta il valore di ciascun segnale di controllo in funzione del valore del codice operativo, ed è ottenuta direttamente dalle informazioni contenute nelle Figure 4.12, 4.16 e 4.17.

## Funzionamento dell'unità di elaborazione

Con le informazioni contenute nelle Figure 4.16 e 4.18 potremmo progettare la logica dell'unità di controllo; tuttavia, prima di procedere esaminiamo l'utilizzo che ciascuna istruzione fa dell'unità di elaborazione. Nelle prossime figure sarà riportato il flusso di elaborazione attraverso l'unità di elaborazione di tre istruzioni appartenenti alle tre diverse classi. I segnali di controllo assegnati e gli elementi attivi dell'unità di elaborazione saranno evidenziati. Si noti che un multiplexer il cui segnale di controllo valga 0 compie un'azione ben definita, anche se il suo segnale di controllo non viene evidenziato. I segnali a più bit saranno evidenziati se risultano assegnati almeno uno dei segnali associati ai singoli bit.

La Figura 4.19 mostra il modo in cui l'unità di elaborazione esegue un'istruzione di tipo R, per esempio add x1, x2, x3. Sebbene l'esecuzione avvenga in un singolo ciclo di clock, possiamo pensare che l'esecuzione proceda attraverso una sequenza di quattro passi, che possiamo identificare analizzando il flusso dei dati:

- l'istruzione viene prelevata dalla memoria e si incrementa il PC;
- i due registri x2 e x3 vengono letti dal register file, mentre l'unità di controllo principale calcola il valore da attribuire alle linee di controllo;
- la ALU elabora i dati letti dal register file, utilizzando alcuni bit del codice operativo per selezionare l'operazione della ALU;
- il risultato calcolato dalla ALU viene scritto nel registro destinazione (x1) del register file.



**Figura 4.19** L'unità di elaborazione durante l'esecuzione di un'istruzione di tipo R, come add x1, x2, x3. Le linee di controllo, le unità funzionali dell'unità di elaborazione e le connessioni attive sono evidenziate in blu.

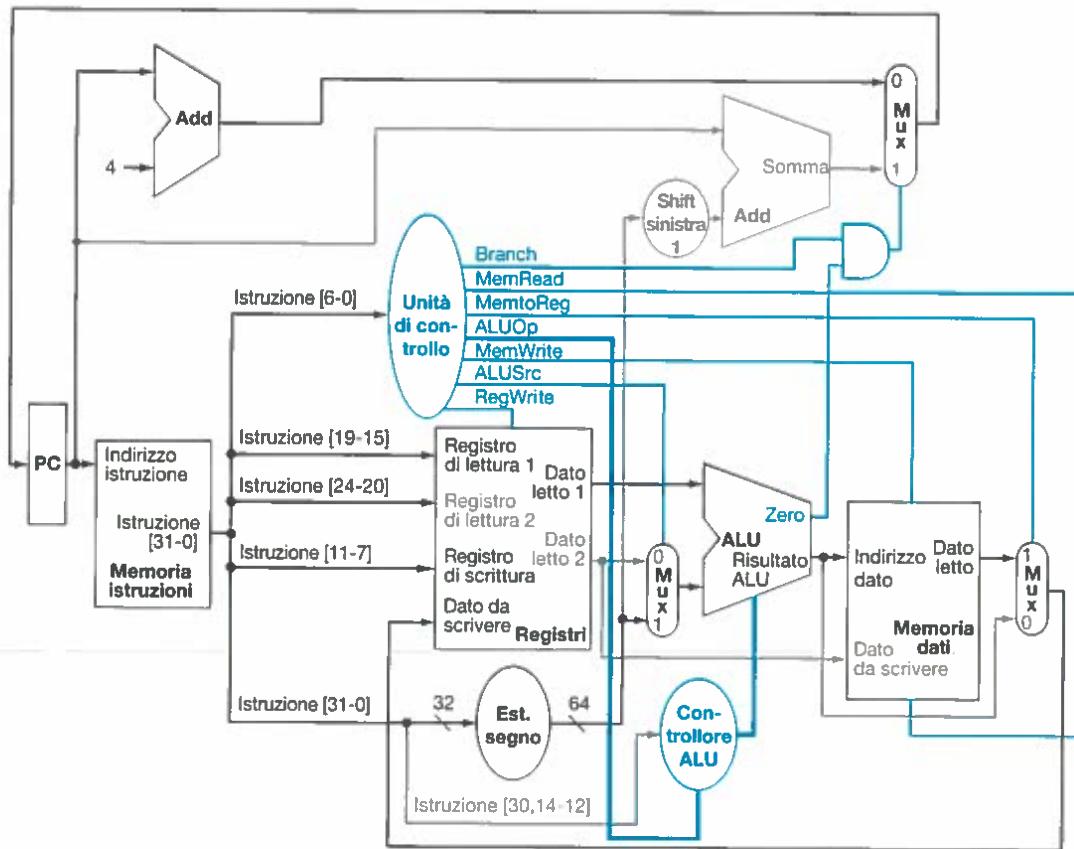
In modo analogo, si può illustrare l'esecuzione di un'istruzione di caricamento di un registro, per esempio

ld x1, offset(x2)

utilizzando lo stesso approccio di Figura 4.19. Le unità funzionali attive e i segnali di controllo asseriti per l'elaborazione di un'istruzione di load sono mostrati in Figura 4.20. Così come avevamo scomposto l'istruzione di tipo R in quattro passi sequenziali, possiamo scomporre l'istruzione di load in cinque passi sequenziali:

1. l'istruzione viene prelevata dalla memoria istruzioni e si incrementa il PC;
2. viene letto il contenuto di un registro (x2) dal register file;
3. la ALU somma il valore letto dal register file ai 12 bit del campo offset dell'istruzione, dotati di segno ed estesi a 32 bit;
4. la somma calcolata dalla ALU viene utilizzata come indirizzo per la memoria dati;
5. il dato proveniente dall'unità di memoria dati viene scritto nel register file nel registro x1.

Infine, possiamo utilizzare la stessa tecnica per esaminare l'esecuzione di un'istruzione di salto condizionato, per esempio beq x1, x2, offset. Essa si comporta in modo molto simile a un'istruzione di tipo R, con la differenza che l'uscita della ALU viene utilizzata per scegliere se scrivere nel PC il valore



**Figura 4.20** L'unità di elaborazione durante l'esecuzione di un'istruzione di load. Le linee di controllo, le unità funzionali dell'unità di elaborazione e le connessioni attive sono evidenziate in blu. Un'istruzione di store viene eseguita in modo molto simile; le differenze principali sono tre: i segnali di controllo della memoria devono specificare una scrittura (del contenuto del secondo registro letto), il secondo registro letto contiene il dato da scrivere in memoria e non è presente l'operazione di scrittura del dato nel register file.

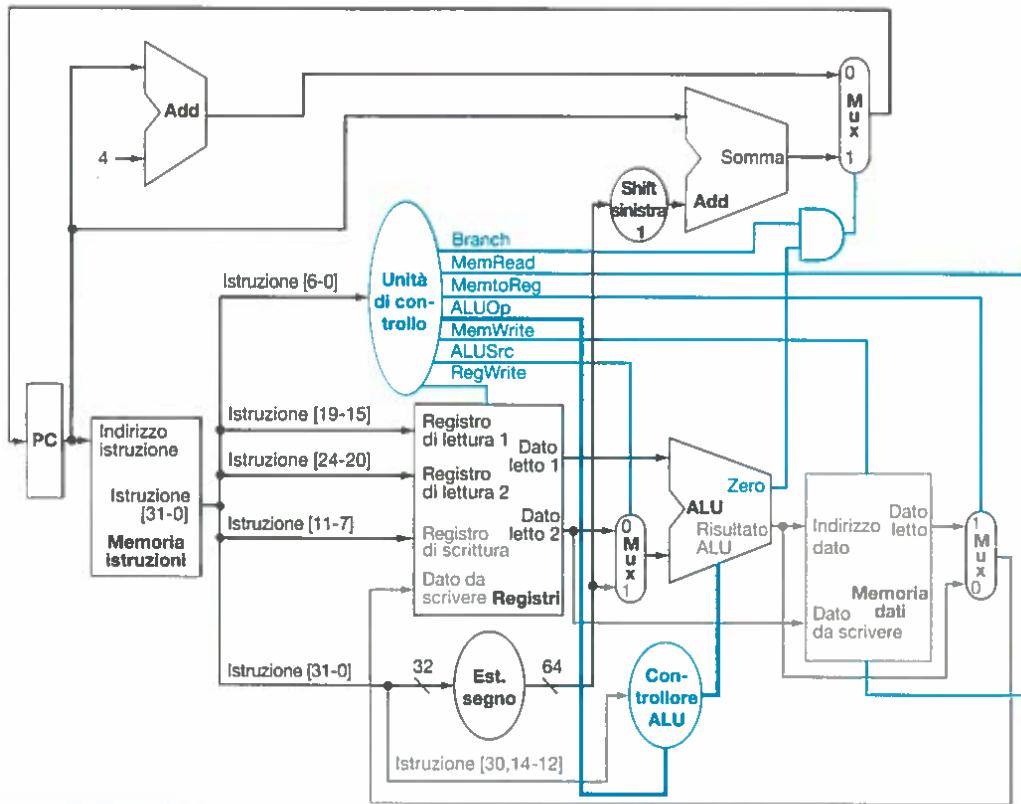
PC + 4 oppure l'indirizzo di destinazione del salto. I quattro passi sequenziali dell'esecuzione di una `beq` sono mostrati in Figura 4.21:

1. l'istruzione viene prelevata dalla memoria istruzioni e si incrementa il PC;
2. vengono letti dal register file i due registri  $x_1$  e  $x_2$ .
3. la ALU esegue la sottrazione del contenuto dei due registri letti dal register file; il valore del PC viene sommato ai 12 bit del campo offset dell'istruzione, dotati di segno, estesi a 32 bit e fatti scorrere di una posizione a sinistra; il risultato costituisce l'indirizzo di destinazione del salto;
4. la linea Zero in uscita dalla ALU viene utilizzata per determinare da quale sommatore prendere l'indirizzo successivo da scrivere nel PC.

### Compleramento dell'unità di controllo

Dopo aver esaminato i diversi passi di esecuzione delle istruzioni, possiamo procedere all'implementazione dell'unità di controllo. La funzione di controllo può essere determinata in modo esatto a partire dal contenuto di Figura 4.18. Le uscite dell'unità di controllo sono i segnali di controllo, mentre i suoi input sono costituiti dai 7 bit del codice operativo. Possiamo quindi scrivere una tabella della verità per ciascuna delle uscite, basata sulla codifica binaria dei codici operativi.

La Figura 4.22 mostra la funzione logica implementata nell'unità di controllo sotto forma di una grande tabella della verità che contiene tutte le uscite e i



**Figura 4.21** L'unità di elaborazione durante l'esecuzione di un'istruzione branch equal. Le linee di controllo, le unità funzionali dell'unità di elaborazione e le connessioni attive sono evidenziate in blu. Dopo aver utilizzato il register file e la ALU per eseguire il confronto, l'uscita Zero della ALU viene impiegata per scegliere il valore successivo del program counter tra i due possibili candidati.

Input o output	Nome del segnale	Formato R	Id	sd	beq
Input	I[6]	0	0	0	1
	I[5]	1	0	1	1
	I[4]	1	0	0	0
	I[3]	0	0	0	0
	I[2]	0	0	0	0
	I[1]	1	1	1	1
	I[0]	1	1	1	1
Output	ALUSrc	0	1	1	0
	MemtoReg	0	1	X	X
	RegWrite	1	1	0	0
	MemRead	0	1	0	0
	MemWrite	0	0	1	0
	Branch	0	0	0	1
	ALUOp01	1	0	0	0
	ALUOp02	0	0	0	1

**Figura 4.22** La funzione logica dell'unità di controllo per l'implementazione a singolo ciclo di clock è completamente specificata da questa tabella della verità. La metà superiore della tabella riporta le combinazioni del codice operativo associate alle quattro classi di istruzioni esaminate, una per ogni colonna; i segnali di controllo vengono impostati in funzione di questi ingressi. La parte inferiore della tabella fornisce invece le uscite in funzione del codice operativo. Quindi l'uscita RegWrite viene assegnata in corrispondenza di due diverse combinazioni degli input. Se si considerano solamente i quattro codici operativi riportati, la tabella della verità si può semplificare utilizzando la condizione di indifferenza per alcuni ingressi: per esempio, si può identificare un'istruzione in formato R con l'espressione  $Op4 \cdot Op5$ , dato che questa è sufficiente a distinguere le istruzioni in formato R dalle 1d, sd e beq. Tuttavia non sfrutteremo queste semplificazioni, dato che in un'implementazione completa dell'insieme delle istruzioni anche i restanti codici operativi vengono utilizzati nell'implementazione completa del RISC-V.

bit del codice operativo come ingressi. Tale tabella specifica completamente la funzione di controllo e può essere implementata direttamente mediante porte logiche in modo automatico. Quest'ultimo passaggio viene descritto nel paragrafo B.2 dell'Appendice B .

### Perché non si utilizzano più implementazioni a singolo ciclo?

Sebbene un processore a singolo ciclo funzioni correttamente, esso non viene più utilizzato nelle implementazioni moderne perché troppo inefficiente. Per capire il perché, si noti che, in questa implementazione, il periodo di clock deve avere la stessa durata per tutte le istruzioni; quindi, il periodo di clock è determinato dal cammino di elaborazione più lungo all'interno del processore (cammino critico). Tale cammino sarà molto probabilmente associato a un'istruzione di load, dato che essa utilizza cinque unità funzionali in sequenza: la memoria istruzioni, il register file, la ALU, la memoria dati e, di nuovo, il register file. Anche se il CPI è pari a 1 (vedi Cap. 1), le prestazioni complessive di un'implementazione a singolo ciclo saranno probabilmente scarse, poiché la durata del periodo di clock risulterà troppo lunga.

Il prezzo da pagare per un'implementazione a singolo ciclo con una durata del clock fissa è significativo, ma potrebbe comunque essere considerato accettabile per un insieme limitato di istruzioni, come quello considerato. I primi calcolatori, dotati di un insieme di istruzioni semplice, utilizzavano questa implementazione. Tuttavia, se si volesse implementare un'unità di elaborazione in virgola mobile o un insieme di istruzioni comprendente istruzioni più complesse, un progetto a singolo clock non funzionerebbe altrettanto bene.

Dato che dobbiamo assicurare che la durata del ciclo di clock sia pari alla durata dell'istruzione più lunga da eseguire, non si possono sfruttare quelle tecniche che consentono di ridurre il tempo di esecuzione delle istruzioni più comuni ma che non riducono la durata dell'esecuzione dell'istruzione più lenta: un'implementazione a un solo ciclo è in contrasto quindi con una delle grandi idee introdotte nel Capitolo 2, che richiede di **rendere veloci le situazioni più comuni**.

Nel prossimo paragrafo illustreremo una tecnica di implementazione diversa, chiamata pipeline. Questa implementazione utilizza un'unità di elaborazione simile a quella a singolo ciclo ma molto più efficace, perché consente un throughput sensibilmente maggiore. Le pipeline sono più efficienti grazie alla sovrapposizione temporale dell'esecuzione di diverse istruzioni.



SITUAZIONI COMUNI

### Autovalutazione

Analizzate i segnali di controllo presenti in Figura 4.22. Si possono combinare fra loro? Esiste qualche segnale di controllo che può essere sostituito con l'inverso di un altro? (Suggerimento: si prendano in considerazione i segnali indifferenti.) Se sì, è possibile sostituire un segnale di controllo con un altro senza utilizzare un invertitore?

## 4.5 | Introduzione alla pipeline

*Il tempo non va sprecato.*  
Proverbo americano

La **pipeline** è una tecnica di implementazione hardware di un insieme di istruzioni che prevede la sovrapposizione temporale dell'esecuzione delle diverse istruzioni; attualmente la **pipeline** viene utilizzata quasi dalla totalità dei calcolatori.

In questo paragrafo utilizzeremo spesso delle analogie per introdurre i termini e le problematiche legate alla pipeline. I lettori non interessati a una trattazione approfondita possono limitarsi a leggere questo paragrafo e poi passare ai paragrafi 4.10 e 4.11, che contengono un'introduzione alle tecniche di pipeline avanzate adottate da alcuni fra i processori più recenti, quali il Core i7 Intel e il Cortex-A53 ARM. Chi, invece, è interessato a esaminare a fondo la struttura dei calcolatori dotati di pipeline, potrà considerare questo paragrafo una buona introduzione ai concetti che verranno trattati nei paragrafi 4.6-4.9.

Chiunque abbia dovuto fare il bucato molte volte ha sicuramente già utilizzato, in un certo senso, la tecnica della pipeline. Fare il bucato, secondo un approccio *non basato sulla pipeline*, richiederebbe i passi seguenti:

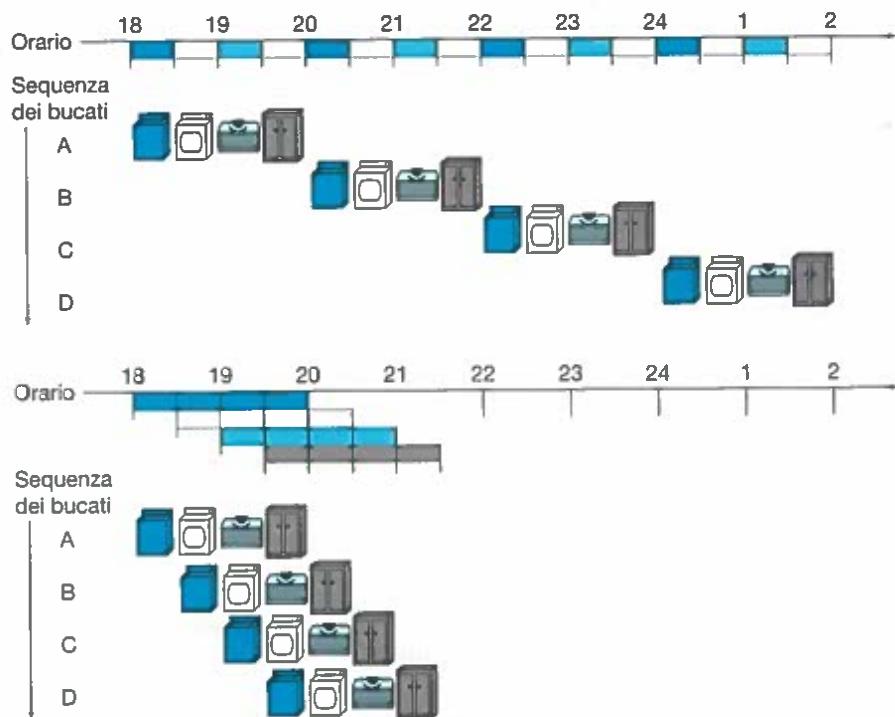
1. mettere un carico di biancheria sporca nella lavatrice;
2. una volta terminato il ciclo di lavaggio, mettere la biancheria ancora bagnata nell'asciugatrice;
3. dopo che l'asciugatrice ha terminato, mettere la biancheria asciugata sul tavolo da stiro e stirarla;
4. quando l'operazione di stiratura è finita, chiedere al proprio coinquilino di riporre la biancheria nell'armadio.

Non appena il coinquilino avrà terminato, si potrà cominciare con il carico successivo di biancheria sporca.

Un approccio basato sul concetto di *pipeline* richiede molto meno tempo. Come mostrato in Figura 4.23, quando la lavatrice finisce il primo ciclo di lavaggio e la biancheria viene trasferita all'asciugatrice, si può riempire subito la lavatrice con il secondo carico di biancheria sporca. Dopo che il primo carico di biancheria è stato asciugato, si può iniziare l'operazione di stiratura, e si



**Pipeline:** una tecnica di implementazione hardware grazie alla quale l'esecuzione di più istruzioni viene sovrapposta nel tempo, in maniera simile a quanto avviene in una catena di montaggio.



**Figura 4.23 Analogia della pipeline con il bucato.** Anna, Bruno, Carla e Dario hanno, ciascuno, dei panni sporchi da lavare, asciugare, stirare e riporre negli armadi. Ogni operazione richiede 30 minuti. Fare i bucati in sequenza richiede quindi 8 ore per completare i quattro cicli di lavaggio, asciugatura, stiratura e sistemazione della biancheria, mentre il bucato con pipeline richiede solamente 3,5 ore. Lo stadio di avanzamento dei diversi carichi all'interno della pipeline, al variare del tempo, è indicato dalla posizione sull'asse dei tempi delle quattro unità funzionali; queste vengono indicate per ogni bucato per chiarezza, ma in realtà c'è una sola unità funzionale di ogni tipo.

può spostare il carico bagnato nell'asciugatrice e mettere il successivo carico di biancheria sporca in lavatrice. Una volta stirata la biancheria del primo carico, il vostro coinquilino potrà prendere gli abiti dal tavolo da stiro e riporli nell'armadio; nel frattempo inizierà la stiratura del secondo carico, il terzo carico si troverà nell'asciugatrice e potrete mettere il quarto carico nella lavatrice. A questo punto tutti i passaggi, detti *stadi* nel linguaggio delle pipeline, sono attivi contemporaneamente: purché si abbiano risorse separate per ciascuno stadio, il bucato può essere messo in pipeline.

Il paradosso delle pipeline è che il tempo necessario per lavare, asciugare, stirare e mettere nell'armadio un singolo calzino non viene affatto ridotto; il motivo per cui le pipeline sono più veloci quando si hanno molti carichi di biancheria è che tutte le unità lavorano in parallelo e quindi si possono fare più bucati nello stesso intervallo di tempo: in sostanza la pipeline aumenta il *throughput* (vedi Cap. 1) della nostra lavanderia. Perciò, la pipeline non diminuisce il tempo per completare il ciclo di lavoro di un singolo bucato, ma quando ci sono più bucati da fare in sequenza, il miglioramento del throughput fa sì che il tempo complessivo diminuisca.

Se tutti gli stadi richiedessero circa lo stesso tempo e ci fosse abbastanza lavoro da fare, l'incremento di velocità dovuto alla pipeline sarebbe pari al numero dei suoi stadi, in questo caso quattro: lavaggio, asciugatura, stiratura e sistemazione della biancheria. Quindi, fare il bucato sfruttando la pipeline è potenzialmente quattro volte più veloce rispetto a farlo senza pipeline: 20 bucati richiederebbero circa 5 volte il tempo richiesto da un singolo bucato, mentre 20 bucati fatti in sequenza richiedono 20 volte il tempo di un singolo bucato. In Figura 4.23 la pipeline risulta solo 2,3 volte più veloce, perché ci sono solamente 4 bucati in tutto. Si noti che all'inizio e alla fine dell'attività la pipeline di Figura 4.23 non è completamente piena, e le fasi di avvio e di fine possono diminuire significativamente le prestazioni se il numero di bucati non è abbastanza grande rispetto al numero degli stadi della pipeline. Se ci fossero molti più di 4 bucati, allora tutti gli stadi sarebbero occupati per la maggior parte del tempo e il miglioramento del throughput sarebbe molto vicino a 4.

Gli stessi principi si applicano ai processori, dove ciò che si vuole porre in pipeline è l'esecuzione delle istruzioni. Le istruzioni RISC-V richiedono tipicamente cinque passi:

1. fetch dell'istruzione (cioè caricamento dell'istruzione dalla memoria);
2. lettura dei registri e decodifica dell'istruzione;
3. esecuzione di un'operazione o calcolo di un indirizzo;
4. accesso a un operando nella memoria dati (se richiesto);
5. scrittura del risultato in un registro (se richiesto).

Quindi la pipeline RISC-V, che esploreremo in questo capitolo, ha cinque stadi. L'esempio seguente mostra come la pipeline acceleri l'esecuzione delle istruzioni proprio come avviene per il bucato.

### Confronto fra le prestazioni di un processore a singolo ciclo e quelle di un processore con pipeline

#### ESEMPIO

Per rendere concreto il discorso, consideriamo una pipeline ben precisa. In questo esempio, e nel resto del capitolo, limiteremo la nostra attenzione a otto istruzioni: load di una parola doppia (`ld`), store di una parola doppia (`sw`), somma (`add`), sottrazione (`sub`), AND (`and`), OR (`or`) e branch if equal (`beq`).

(continua)

(continua)

Confrontiamo il tempo medio necessario per l'esecuzione di queste istruzioni in un'implementazione a singolo ciclo, in cui tutte le istruzioni richiedono un ciclo di clock, con il tempo medio richiesto per un'implementazione con pipeline. La durata delle operazioni delle unità funzionali più importanti è, per questo esempio, di 200 ps per l'accesso alla memoria, 200 ps per le operazioni della ALU e 100 ps per la lettura e la scrittura del register file. Nel modello a singolo ciclo tutte le istruzioni impiegano esattamente un ciclo di clock, quindi la durata del ciclo di clock va dilatata in modo tale da poter contenere l'esecuzione dell'istruzione più lenta.

Il tempo richiesto da ciascuna delle sette istruzioni è riportato in Figura 4.24. Il progetto a singolo ciclo deve permettere la completa esecuzione dell'istruzione più lenta, la ld in Figura 4.24, per cui il tempo richiesto da ciascuna istruzione sarà di 800 ps. In modo simile a quanto mostrato in Figura 4.23, la Figura 4.25 confronta l'esecuzione con e senza pipeline di tre istruzioni di load: il tempo trascorso tra la prima e la quarta istruzione nell'architettura senza pipeline è pari a  $3 \times 800$  ps, ossia 2400 ps.

Tutti gli stadi della pipeline richiedono un ciclo di clock, per cui il periodo di clock deve essere abbastanza lungo da contenere l'operazione più lenta. Come nel progetto a singolo ciclo si deve considerare la durata dell'istruzione più lenta (che è di 800 ps), anche se alcune istruzioni potrebbero richiedere solo 500 ps, nel caso della pipeline si deve considerare la durata dello stadio più lento (che è di 200 ps), anche se alcuni stadi richiedono solo 100 ps. La pipeline offre comunque un miglioramento delle prestazioni di quattro volte: l'intervallo di tempo tra la prima e la quarta istruzione è infatti di  $3 \times 200$  ps, ossia 600 ps.

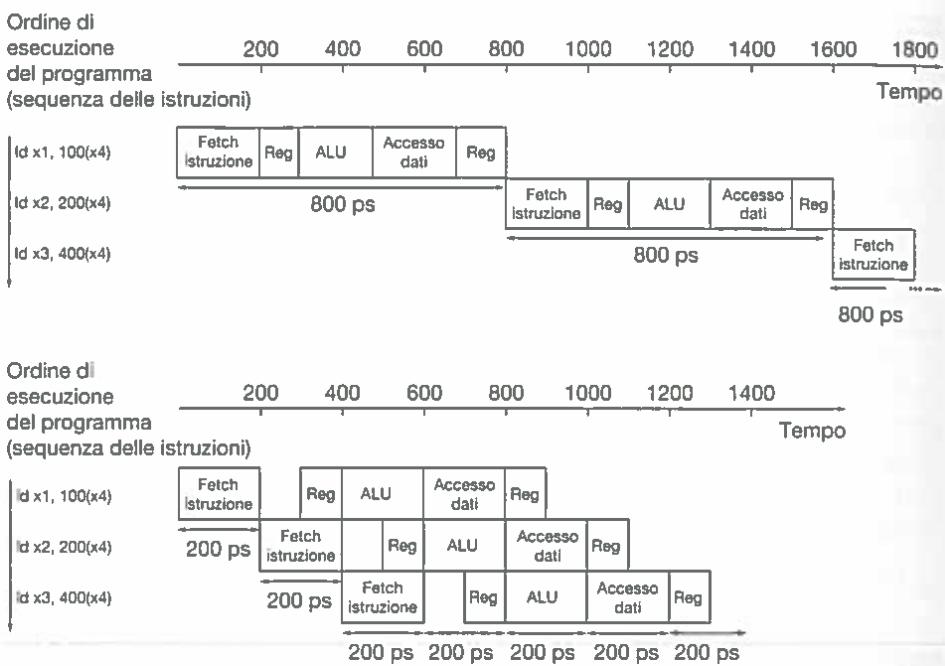
**SOLUZIONE**

Tipo di istruzione	Lettura dell'istruzione	Lettura dei registri	Operazione con la ALU	Accesso ai dati in memoria	Scrittura del register file	Tempo totale
Load doubleword (ld)	200 ps	100 ps	200 ps	200 ps	100 ps	800 ps
Store doubleword (sd)	200 ps	100 ps	200 ps	200 ps		700 ps
Formato R (add, sub e or)	200 ps	100 ps	200 ps		100 ps	600 ps
Salto condizionato (beq)	200 ps	100 ps	200 ps			500 ps

**Figura 4.24** Tempo di esecuzione totale per le diverse istruzioni calcolato a partire dal tempo di elaborazione richiesto dalle singole unità funzionali. Il calcolo assume che i multiplexer, l'unità di controllo, l'accesso al PC e l'unità di estensione del segno non introducano ritardi.

Dalla trattazione precedente è possibile ricavare una formula che fornisce l'incremento della velocità prodotto dalla pipeline: se gli stadi sono perfettamente bilanciati, il tempo di esecuzione delle istruzioni in un calcolatore con pipeline (in condizioni ideali) è pari a:

$$\text{Tempo tra due istruzioni}_{\text{con pipeline}} = \frac{\text{Tempo tra due istruzioni}_{\text{senza pipeline}}}{\text{Numero di stadi della pipeline}}$$



**Figura 4.25** Confronto tra un'esecuzione a singolo ciclo di clock senza pipeline (in alto) e un'esecuzione con pipeline (in basso). In entrambi i casi si utilizzano gli stessi componenti hardware, il cui tempo di elaborazione è riportato in Figura 4.24. Si trova una diminuzione del tempo medio di esecuzione delle istruzioni di un fattore 4: esso scende da 800 ps a 200 ps. Si confronti questa figura con la Figura 4.23: per il bucato avevamo ipotizzato che tutti gli stadi avessero la stessa durata. Se l'asciugatrice fosse stata più lenta, avrebbe determinato la durata degli altri stadi. Anche nel calcolatore la durata dei diversi stadi di elaborazione è determinata dalla risorsa funzionale più lenta: dal calcolo nella ALU o dall'accesso alla memoria. Assumiamo che la scrittura nel register file avvenga nella prima metà del ciclo di clock e che la lettura dal register file abbia luogo nella seconda metà; tale assunzione sarà valida per tutto questo capitolo.

In condizioni ideali, con un numero molto grande di istruzioni, l'incremento di velocità dovuto alla pipeline è pari al numero dei suoi stadi: una pipeline a cinque stadi sarebbe quindi quasi cinque volte più veloce.

La formula suggerisce che una pipeline a cinque stadi dovrebbe diminuire di quasi cinque volte il tempo di esecuzione rispetto alla versione senza pipeline, che nel nostro esempio era di 800 ps; cioè, dovrebbe avere un clock di 160 ps. Tuttavia, l'esempio mostra anche che gli stadi possono non essere perfettamente bilanciati; la pipeline introduce inoltre dei tempi (la cui origine sarà più chiara tra poco) che si aggiungono al tempo di esecuzione. In definitiva, il tempo di esecuzione per la singola istruzione in un calcolatore con pipeline sarà maggiore del minimo teorico, e quindi l'incremento di velocità sarà minore del numero di stadi.

È bene notare che nell'esempio appena discusso non siamo riusciti a ottenere neppure il miglioramento di un fattore 4 che avevamo dichiarato. Se osserviamo il tempo totale di esecuzione delle tre istruzioni, infatti, in presenza di pipeline ci sono voluti 1400 ps rispetto ai 2400 ps dell'implementazione a singolo ciclo. Naturalmente, questo è dovuto al fatto che il numero di istruzioni non è grande. Che cosa accadrebbe se aumentassimo il numero di istruzioni?

Supponiamo di estendere l'esempio a 1 000 003 istruzioni (ne aggiungiamo 1 000 000 alle tre considerate prima). In presenza di pipeline, ogni istruzione contribuirebbe a incrementare di 200 ps il tempo di esecuzione totale, che passerebbe quindi a  $1\ 000\ 000 \times 200\text{ ps} + 1400\text{ ps}$ , ovvero a 200 001 400 ps. In assenza di pipeline ogni istruzione in più richiederebbe 800 ps, dando un tempo di esecuzione totale di  $1\ 000\ 000 \times 800\text{ ps} + 2400\text{ ps}$ , cioè di 800 002 400 ps. In queste condizioni, quindi, il rapporto tra il tempo di esecuzione del processore

senza pipeline e quello del processore dotato di pipeline si avvicina al rapporto tra la durata delle istruzioni nei due tipi di calcolatore:

$$\frac{800\ 002\ 400 \text{ ps}}{200\ 001\ 400 \text{ ps}} \approx \frac{800 \text{ ps}}{200 \text{ ps}} \approx 4,00$$

La tecnica della pipeline migliora le prestazioni *aumentando il throughput delle istruzioni*, ma *non riduce il tempo di esecuzione della singola istruzione*. Tuttavia il throughput delle istruzioni rimane la metrica più importante, dal momento che i programmi reali eseguono miliardi di istruzioni.

## Progettazione dell'insieme di istruzioni per architetture dotate di pipeline

Anche da questa semplice introduzione alla pipeline è possibile evidenziare alcuni elementi chiave dell'insieme delle istruzioni del RISC-V, che è stato progettato espressamente per l'esecuzione in pipeline.

Anzitutto, tutte le istruzioni RISC-V hanno la stessa lunghezza: tale proprietà semplifica notevolmente la fase di fetch delle istruzioni nel primo stadio della pipeline e la loro decodifica nel secondo stadio. In un insieme di istruzioni come quello dell'x86, in cui la lunghezza delle istruzioni varia da 1 byte a 15 byte, la progettazione della pipeline sarebbe notevolmente più complicata. Le implementazioni più recenti dell'architettura x86, in realtà, traducono le istruzioni x86 in microistruzioni più semplici che somigliano a quelle del RISC-V e inviano quindi alla pipeline queste semplici microistruzioni invece delle istruzioni originali x86 (par. 4.10).

In secondo luogo, le istruzioni RISC-V adottano un numero ridotto di formati diversi, e in tutti i formati i registri sorgente sono sempre specificati nella stessa posizione, qualunque sia l'istruzione. Tale regolarità permette al secondo stadio di iniziare a leggere il register file mentre l'unità di controllo sta determinando il tipo dell'istruzione caricata dalla memoria: se il formato delle istruzioni RISC-V non fosse regolare, sarebbe necessario spezzare in due questo secondo stadio, portando la pipeline a un totale di sei stadi. Vedremo tra poco lo svantaggio di avere pipeline più lunghe.

In terzo luogo, nel RISC-V gli operandi residenti in memoria possono comparire solo nelle istruzioni di load e di store: questo vincolo permette di utilizzare lo stadio di esecuzione per calcolare l'indirizzo di memoria, per accedere poi alla memoria nello stadio successivo. Se si potessero effettuare delle operazioni sugli operandi residenti in memoria, come avviene nell'x86, gli stadi 3 e 4 dovrebbero essere espansi in tre stadi: uno stadio per il calcolo dell'indirizzo, uno per l'accesso alla memoria dati e un terzo per l'esecuzione delle operazioni. Vedremo tra poco quali siano gli inconvenienti di pipeline più lunghe.

## Hazard nelle pipeline

In una pipeline si possono verificare situazioni in cui l'istruzione successiva non può essere eseguita nel ciclo di clock immediatamente successivo. Tali eventi sono detti *hazard* (o *criticità*) e possono essere di tre tipi.

### Hazard strutturali

Il primo tipo di hazard è l'**hazard strutturale**, che si verifica quando le risorse hardware presenti non sono in grado di supportare la combinazione di istruzioni che si vorrebbe eseguire nello stesso ciclo di clock. Nel caso del bucato, si verificherebbe un hazard strutturale se si utilizzasse un'unica macchina per

**Hazard strutturale:** una situazione in cui un'istruzione non può essere eseguita nel ciclo di clock previsto perché l'hardware non supporta la combinazione di istruzioni da eseguire.

lavare e asciugare anziché due macchine separate, oppure se il nostro coinquino fosse impegnato in qualche altro compito e non potesse riporre gli abiti nell'armadio. In questi casi, la nostra attenta programmazione basata sulla pipeline diventerebbe inutile.

Come già osservato, l'insieme di istruzioni RISC-V è stato progettato appositamente per essere eseguito in pipeline, facendo in modo che per i progettisti fosse semplice evitare gli hazard strutturali nel progetto dell'unità di elaborazione. Supponiamo, tuttavia, di disporre di una singola memoria anziché di due: se nella pipeline di Figura 4.25 ci fosse una quarta istruzione, la prima istruzione dovrebbe accedere ai dati in memoria, mentre la quarta istruzione dovrebbe essere prelevata dalla stessa memoria nel medesimo ciclo di clock. Senza due memorie distinte, l'architettura con pipeline presenterebbe in questo caso un hazard strutturale.

### Hazard sui dati

**Hazard sui dati:** si verifica quando un'istruzione non può essere eseguita in un certo ciclo di clock perché i dati di cui ha bisogno l'istruzione non sono ancora disponibili.

Un **hazard sui dati** si verifica quando la pipeline deve essere messa in stallo perché uno stadio deve attendere che termini l'elaborazione in un altro stadio della pipeline. Ritornando all'esempio della lavanderia, supponiamo di trovare un calzino spaiato sul tavolo da stirto. Una possibilità è di cercare il calzino mancante nell'armadio della vostra stanza; naturalmente, il bucato che ha terminato l'asciugatura (ed è pronto per essere stirato) dovrà attendere che abbiate terminato la ricerca, e così pure il bucato che è uscito dalla fase di lavaggio ed è pronto per essere asciugato.

In un calcolatore dotato di pipeline, gli hazard sui dati nascono quando un'istruzione dipende dal risultato di un'istruzione precedente che si trova ancora all'interno della pipeline (questa è una dipendenza che non esiste realmente quando si fa il bucato). Supponiamo, per esempio, di avere un'istruzione di somma seguita da una sottrazione che utilizza il risultato della somma (x19):

```
add x19, x0, x1
sub x2, x19, x3
```

Se non si prendessero delle contromisure, gli hazard sui dati potrebbero causare tanti stalli della pipeline: l'istruzione `add`, infatti, non scrive il proprio risultato prima del quinto stadio, e ciò significa che si dovrebbero perdere tre cicli di clock.

Anche se possiamo affidare ai compilatori il compito di evitare questo tipo di hazard, il risultato non sarebbe soddisfacente: queste dipendenze sono troppo frequenti e il ritardo che viene introdotto è troppo lungo per poter sperare che il compilatore possa risolvere il problema.

La soluzione più utilizzata è basata sull'osservazione che non è necessario che termini l'istruzione per cercare di risolvere un hazard sui dati: nella sequenza di istruzioni appena vista, non appena il risultato della somma viene generato dalla ALU, questo potrebbe essere già utilizzato come ingresso alla ALU per l'operazione di sottrazione. La tecnica che prevede l'aggiunta di un circuito che fornisce in un punto dell'esecuzione il dato mancante, prendendolo da una risorsa interna, viene chiamata **propagazione** o **bypassing**.

### Propagazione nel caso di due istruzioni

#### ESEMPIO

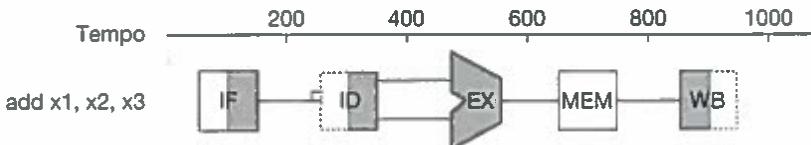
(continua)

Nel caso delle due istruzioni sopra riportate, mostrare quali stadi della pipeline dovranno essere collegati con la tecnica della propagazione. Si utilizzi lo schema riportato in Figura 4.26 per rappresentare l'unità

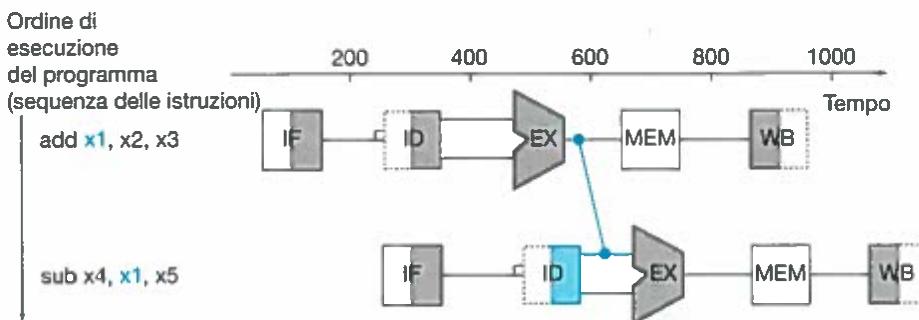
(continua)

di elaborazione durante l'elaborazione nei cinque stadi della pipeline. Si faccia una copia dell'unità di elaborazione per ciascuna istruzione e la si allinei sull'asse dei tempi in modo analogo all'esempio del bucato di Figura 4.25.

La Figura 4.27 mostra i collegamenti necessari a propagare il valore di  $x_1$ , preso al termine della fase di esecuzione dell'istruzione add, come input nella fase di esecuzione dell'istruzione sub.

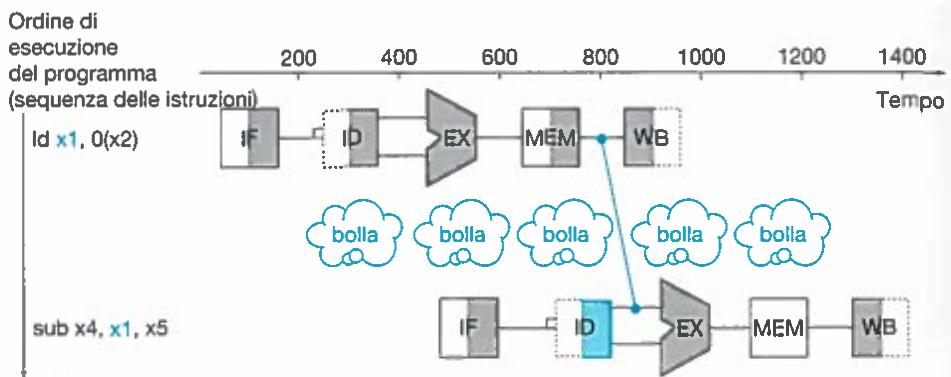
**SOLUZIONE**

**Figura 4.26** Rappresentazione grafica di un'istruzione nella pipeline ispirata alla pipeline del bucato di Figura 4.23. In tutto il capitolo utilizzeremo dei simboli grafici per rappresentare le risorse fisiche della pipeline; all'interno dei simboli verrà riportata un'abbreviazione, corrispondente al nome dello stadio associato. Le sigle dei cinque stadi sono le seguenti: *IF* per lo stadio di prelevamento dell'istruzione (da *instruction fetch*), con il quadrato associato che rappresenta la memoria istruzioni; *ID* per lo stadio di decodifica delle istruzioni e lettura dal register file (da *instruction decode*), con il simbolo associato che rappresenta il register file (che in questo stadio viene letto); *EX* per lo stadio di esecuzione (da *execution*), con il simbolo associato che rappresenta la ALU; *MEM* per lo stadio di accesso alla memoria dati, con il quadrato associato che rappresenta la memoria dati; *WB* per lo stadio di scrittura del risultato nel register file (da *write-back*), con il simbolo associato che rappresenta il register file (che in questo stadio viene scritto). L'ombreggiatura indica che l'elemento funzionale viene utilizzato dall'istruzione; per esempio, *MEM* viene rappresentato senza ombreggiatura perché l'istruzione *add* non accede alla memoria dati. L'ombreggiatura della metà destra del register file o della memoria indica che l'elemento viene letto in quello stadio, mentre l'ombreggiatura della metà sinistra indica che l'elemento viene scritto; per esempio, la metà destra di *ID* risulta ombreggiata nel secondo stadio perché il register file viene letto, mentre nel quinto stadio la metà sinistra di *WB* è ombreggiata per indicare che il register file viene scritto.



**Figura 4.27** Rappresentazione grafica della propagazione. Il segmento in blu che connette le due copie della pipeline mostra il cammino che segue la propagazione del dato dall'uscita dello stadio EX della add all'input dello stadio EX della sub, dove sostituisce il contenuto del registro  $x_1$  letto dal register file nel secondo stadio della sub.

Da questa rappresentazione grafica dell'esecuzione nella pipeline si evince che la propagazione funziona solamente se lo stadio a cui il dato viene propagato è successivo nel tempo allo stadio dal quale viene prelevato. Per esempio, non ci potrebbe essere una propagazione valida tra l'uscita della fase di accesso alla memoria di un'istruzione e l'ingresso della fase di esecuzione dell'istruzione successiva, poiché ciò vorrebbe dire propagare un dato a ritroso nel tempo.



**Figura 4.28** Anche se si fa uso della propagazione, è necessario uno stallo quando un'istruzione *in* formato R, che segue un'istruzione *load*, deve utilizzare il dato trasferito dalla memoria dalla *load*. Senza lo stallo, la propagazione del dato dall'uscita dello stadio di accesso alla memoria all'ingresso dello stadio di esecuzione procederebbe a ritroso nel tempo, il che è impossibile. Questa figura è in realtà una semplificazione, non potendo sapere se lo stallo sarà necessario o meno fino a quando l'istruzione di sottrazione non verrà letta e decodificata. Nel paragrafo 4.7 mostreremo in dettaglio che cosa accade realmente in questo tipo di hazard.

**Hazard sui dati di una load:** è un tipo particolare di hazard sui dati, originato dal fatto che il dato letto dalla memoria da un'istruzione *load* non è ancora disponibile quando servirebbe a un'altra istruzione.

**Stallo della pipeline:** detto anche **bolla**, è uno stallo richiesto per risolvere un hazard.

La tecnica della propagazione funziona molto bene e verrà descritta in dettaglio nel paragrafo 4.7, ma non è in grado di evitare tutti gli stalli di una pipeline. Per esempio, supponiamo che la prima istruzione di Figura 4.27 sia un'istruzione di caricamento dalla memoria in *x1* (*load*), anziché una *add*. Come si può facilmente vedere esaminando la Figura 4.27, il dato desiderato sarebbe disponibile solamente *dopo* il quarto stadio della prima istruzione, ossia troppo tardi per essere utilizzato come *input* del terzo stadio della sottrazione. Di conseguenza, anche disponendo del meccanismo della propagazione, dovremmo imporre uno stallo della durata di un ciclo di clock a causa di un **hazard sui dati di una load** (*load-use data hazard*) come mostrato in Figura 4.28. Questa figura mostra una situazione particolarmente importante della pipeline, chiamata formalmente **stallo della pipeline** ma spesso indicata con il termine **bolla** (*bubble*)<sup>1</sup>.

In seguito vedremo altre situazioni di stallo della pipeline. Nel paragrafo 4.7 mostreremo come possono essere gestiti i casi più difficili (come quello appena illustrato), utilizzando dei circuiti dedicati al riconoscimento degli hazard e alla messa in stallo della pipeline, oppure utilizzando un software che sia in grado di riordinare il codice per evitare gli hazard dovuti all'utilizzo dei dati delle *load*, come mostrato dal prossimo esempio.

### Riordinare il codice per evitare stalli nella pipeline

#### ESEMPIO

Si consideri il seguente frammento di codice C:

```
a = b + e;
c = b + f;
```

Di seguito è riportato il codice RISC-V corrispondente, assumendo che tutte le variabili siano contenute in memoria e possano essere indirizzate

(continua)

<sup>1</sup> Più propriamente, stallo è la condizione in cui si vuole mettere la pipeline, la bolla è il modo per ottenere questa condizione, come sarà chiaro tra poco.

come spiazzamento rispetto a x31:

```
ld  x1,0(x31)  // load b
ld  x2,8(x31)  // load e
add x3,x1,x2  // b + e
sd  x3,24(x31) // store a
ld  x4,16(x31) // load f
add x5,x1,x4  // b + f
sd  x5,32(x31) // store c
```

Identificare gli hazard nel codice riportato sopra e riordinare le istruzioni in modo da evitare stalli nella pipeline.

Entrambe le istruzioni add presentano un hazard, perché uno dei loro operandi dipende dall'istruzione `ld` immediatamente precedente. Si noti che la propagazione elimina alcuni potenziali hazard, tra cui la dipendenza della prima add dalla prima `ld` e tutti gli hazard sulle istruzioni di store. Spostando di due posizioni verso l'alto la terza istruzione `ld`, si possono eliminare entrambi gli hazard:

```
ld  x1,0(x31)
ld  x2,8(x31)
ld  x4,16(x31)
add x3,x1,x2
sd  x3,24(x31)
add x5,x1,x4
sd  x5,32(x31)
```

Su un processore con pipeline dotato di propagazione, la sequenza riordinata richiederà due cicli di clock in meno della versione originale.

(continua)

### SOLUZIONE

Il meccanismo della propagazione illustra un altro elemento chiave dell'architettura RISC-V, oltre alle tre caratteristiche descritte in precedenza in questo paragrafo. Un'istruzione RISC-V scrive un solo risultato e lo fa sempre nell'ultimo stadio della pipeline. La propagazione sarebbe molto più difficile se ogni istruzione producesse più risultati da propagare o se questi venissero scritti prima del termine dell'esecuzione dell'istruzione.

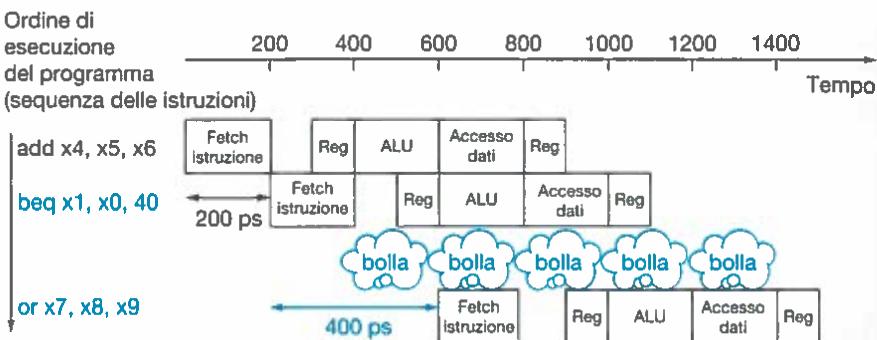
**Approfondimento.** Il termine “propagazione” deriva dall’idea che il risultato è passato da un’istruzione verso quella successiva. Il termine inglese *bypassing* (scavalcamiento) sottolinea il passaggio diretto di un risultato all’unità che ne ha bisogno, senza passare per il register file.

## Hazard sul controllo

Il terzo tipo di hazard è rappresentato dagli **hazard sul controllo** o **hazard sui salti condizionati**. Questi hazard avvengono quando bisogna prendere una decisione in funzione del risultato dell'esecuzione di un'istruzione e altre istruzioni sono già state avviate all'esecuzione.

Riferendoci all'esempio del bucato, supponiamo che ai nostri amici sia affidato l'infelice compito di lavare le divise di una squadra di calcio: a seconda di quanto sono sporche, occorre determinare la quantità di detersivo e la temperatura dell'acqua in modo da effettuare un lavaggio abbastanza energico da pulire le divise, ma non troppo energico per evitare che le divise si usurino in fretta. Nella pipeline per il bucato sarebbe necessario attendere il completamento del

**Hazard sul controllo:** detti anche **hazard sui salti condizionati**, si verificano quando un'istruzione non può essere eseguita nel ciclo di clock corretto della pipeline perché un'altra istruzione è stata caricata al suo posto; ossia, il flusso degli indirizzi delle istruzioni non è quello che la pipeline si aspetta.



**Figura 4.29** Esempio di pipeline messa in stallo a ogni istruzione di salto condizionato come possibile soluzione agli hazard sul controllo. In questo esempio si suppone che il salto venga eseguito e che l'istruzione di destinazione del salto sia la or. Dopo l'istruzione di salto condizionato si verifica uno stallo della pipeline, o bolla, della durata di uno stadio. In realtà il processo di generazione di uno stallo è leggermente più complicato, come vedremo nel paragrafo 4.8; l'effetto sulle prestazioni è comunque lo stesso dell'inserimento di una bolla.

secondo stadio, quando si possono esaminare le divise asciutte, prima di poter decidere se sia necessario modificare le impostazioni della lavatrice. Che cosa si può fare in questo caso?

Nel seguito esamineremo due possibili soluzioni agli hazard sul controllo applicate al caso del bucato e forniremo la soluzione equivalente per i calcolatori.

**Stallo:** lavare il bucato in modo sequenziale fino a quando il primo carico è asciutto; se il risultato ottenuto non è quello desiderato, ripetere i cicli di lavaggio e di asciugatura finché non si ottengono le impostazioni corrette.

Questa scelta prudente è certamente corretta, ma è lenta.

In un calcolatore l'equivalente di un'operazione di scelta è l'istruzione di salto condizionato. Si noti che la lettura dalla memoria istruzioni dell'istruzione che segue la branch viene fatta nel ciclo di clock immediatamente successivo al prelevamento della branch; tuttavia la pipeline non ha alcun modo di sapere quale sarà l'istruzione successiva perché *ha appena ricevuto* l'istruzione di branch dalla memoria! Esattamente come nel caso del bucato, una soluzione potrebbe essere: mettere in stallo la pipeline immediatamente dopo aver caricato dalla memoria un'istruzione di branch e attendere fino a che la pipeline non abbia determinato il risultato del confronto prima di calcolare l'indirizzo dell'istruzione successiva.

Ipotizziamo di avere inserito sufficienti risorse hardware addizionali da poter confrontare il contenuto di due registri, calcolare l'indirizzo di destinazione del salto e aggiornare il PC già nel secondo stadio (per i dettagli *vedi* par. 4.8). Anche disponendo di tali circuiti aggiuntivi, in presenza di salti condizionati la pipeline si comporterebbe come quella di Figura 4.29: l'istruzione ld, che deve essere eseguita subito dopo la beq nel caso in cui i due operandi di quest'ultima risultino diversi, richiede uno stallo pari a un ciclo di clock, cioè 200 ps, prima di poter essere caricata dalla memoria istruzioni.

### Impatto sulle prestazioni di uno stallo dovuto ai salti condizionati

#### ESEMPIO

Stimare l'impatto sul numero di *cicli di clock per istruzione* (CPI) degli stalli dovuti alle istruzioni di salto condizionato. Si assuma che tutte le altre istruzioni abbiano un CPI pari a 1.

(continua)

La Figura 3.28 mostra che i salti condizionati costituiscono circa il 17% delle istruzioni eseguite dagli SPECint2006. Dato che le altre istruzioni eseguite hanno un CPI pari a 1 e i salti condizionati richiedono un ciclo di clock aggiuntivo dovuto allo stallo, si otterrebbe un CPI totale pari a 1,17, e conseguentemente un rallentamento di 1,17 rispetto al caso ideale.

(continua)

**SOLUZIONE**

Se non si può decidere nel secondo stadio l'esito del confronto associato ai salti condizionati, come spesso accade nelle pipeline più lunghe, si genera un rallentamento ancora maggiore della pipeline quando occorre metterla in stallo per una branch. Il costo di questa soluzione sarebbe troppo elevato (per quasi tutti i calcolatori) per poterla prendere in considerazione. Si ricorre, quindi, a una seconda soluzione per risolvere gli hazard sul controllo, che è basata su una delle grandi idee del Capitolo 1:

**Predizione:** se siamo ragionevolmente sicuri di conoscere le impostazioni ottimali della lavatrice per lavare le divise, possiamo *supporre* che queste impostazioni funzioneranno e iniziare a lavare il secondo carico di divise mentre il primo si sta ancora asciugando.

In tal modo, se le impostazioni funzionano, non si rallenta la pipeline. In caso contrario, occorrerà ripetere il lavaggio del primo carico di divise che era stato lavato utilizzando le impostazioni che avevamo supposto ottimali.

In effetti, i calcolatori utilizzano estensivamente la **predizione** nella gestione dei salti condizionati. L'approccio più semplice consiste nel predire sempre che il salto venga non eseguito (*untaken branch*). Se la predizione si rivela corretta, la pipeline procede al massimo della velocità, e solo quando il salto doveva essere eseguito si verifica uno stallo. I due casi possibili sono esemplificati in Figura 4.30.

Una versione più sofisticata di **predizione dei salti** consente di prevedere se un salto debba essere eseguito o meno. Tornando all'analogia con il bucato, le uniformi scure richiederebbero certe impostazioni di lavaggio e quelle chiare altre. Nei programmi software, consideriamo i salti che si trovano al termine dei cicli e che sono utilizzati per ritornare all'inizio del ciclo: poiché questi salti vanno all'indietro ed è verosimile che siano eseguiti la maggior parte delle volte, il processore potrebbe decidere che i salti verso indirizzi minori vadano sempre eseguiti.

Questi approcci rigidi alla predizione dei salti ipotizzano un comportamento che è sempre lo stesso e non tengono conto delle individualità specifiche di ciascuna istruzione di salto. I predittori hardware *dinamici*, invece, determinano la predizione per ciascun salto in funzione del comportamento precedente di quell'istruzione di salto, e possono modificare la predizione di ogni salto durante l'esecuzione del programma. Tornando all'esempio del bucato, nella predizione dinamica ci sarebbe una persona che esamina quanto erano sporche le divise e predice di conseguenza le impostazioni del lavaggio, modificando la predizione per la volta successiva sulla base del successo delle predizioni fatte in precedenza.

Un'implementazione diffusa della predizione dinamica dei salti consiste nel memorizzare la storia di ciascun salto, ricordando quando è stato preso o no, e nell'utilizzare il comportamento recente per predire il futuro. Come vedremo più avanti, la quantità e il tipo di informazione che viene memorizzata a questo scopo sono cresciuti molto negli anni, con il risultato che i circuiti

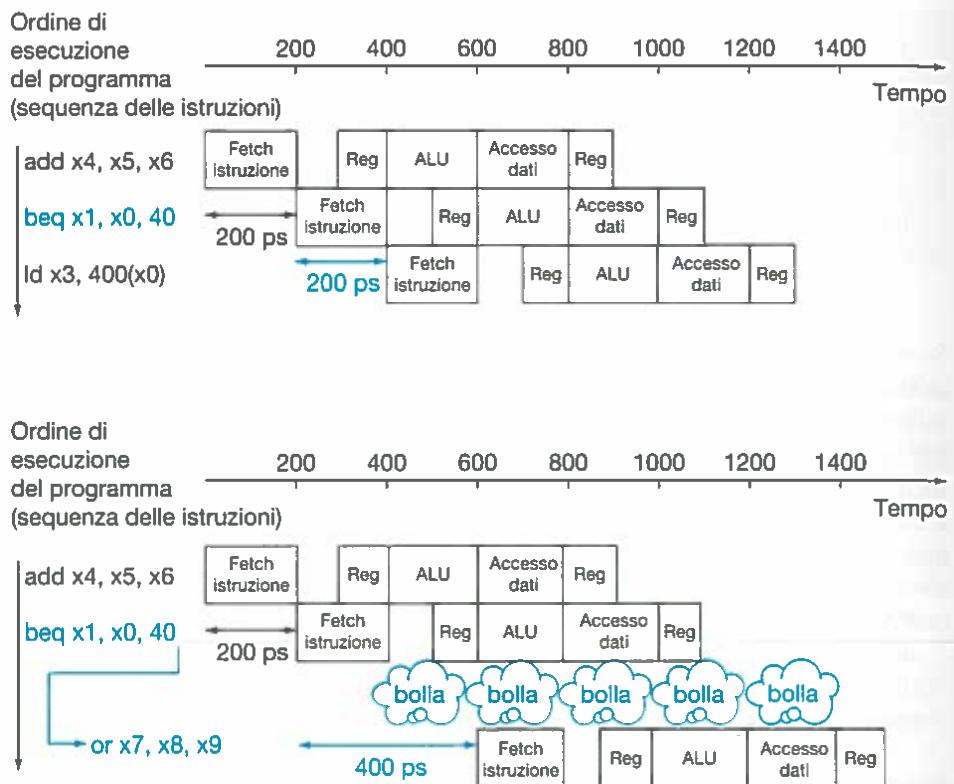


PREDIZIONE

**Predizione dei salti:** un metodo per risolvere gli hazard sul controllo basato sull'ipotesi che il confronto associato a un salto condizionato dia un certo risultato; l'esecuzione procede basandosi su questa ipotesi, invece di aspettare che sia verificato il risultato effettivo del confronto.



PREDIZIONE



**Figura 4.30** Pipeline in cui gli hazard sul controllo vengono risolti predicendo che i salti condizionati non vengano eseguiti. Lo schema mostra la pipeline nel caso in cui il salto non viene eseguito (sopra) e nel caso in cui il salto viene eseguito (sotto). Come sottolineato in Figura 4.29, l'inserimento di una bolla nel modo mostrato in questo schema costituisce una semplificazione di ciò che accade realmente, almeno durante il primo ciclo di clock successivo alla beq; i dettagli verranno mostrati nel paragrafo 4.8.

predittori dei salti possono raggiungere attualmente una precisione superiore al 90% (par. 4.8). Quando la predizione si rivela errata, l'unità di controllo della pipeline deve garantire che le istruzioni che seguivano quella di salto e sono già state caricate nella pipeline non abbiano alcun effetto, e deve far ripartire la pipeline dall'indirizzo corretto. Nell'analogia con la lavanderia, dobbiamo smettere di accettare nuovi bucato in modo da poter far ripartire la pipeline dal bucato lavato in maniera non corretta.

Come avviene in tutte le altre soluzioni degli hazard sul controllo, il problema si aggrava al crescere della lunghezza della pipeline, poiché aumenta il costo delle predizioni errate. Le diverse soluzioni agli hazard sul controllo sono descritte in maggior dettaglio nel paragrafo 4.8.

**Approfondimento.** C'è un terzo approccio possibile alla soluzione degli hazard sul controllo, detto *decisione ritardata*. Proseguendo con l'analogia con il bucato, ogni volta che occorre prendere una decisione possiamo mettere nella lavatrice un carico di abiti diverso dalle divise dei calciatori, nell'attesa che le divise del carico precedente si asciughino. A patto di avere a disposizione una sufficiente quantità di abiti sporchi il cui lavaggio non dipende da quello delle divise, questa soluzione funziona correttamente.

Tale soluzione, a cui abbiamo accennato in precedenza, nel mondo dei calcolatori è detta *salto ritardato* ed è la soluzione effettivamente adottata nell'architettura MIPS. In un salto ritardato, l'istruzione successiva all'istruzione di salto viene

sempre eseguita e il salto viene eseguito, eventualmente, dopo questa istruzione. Questo comportamento è invisibile al programmatore assembler MIPS, poiché l'assemblatore è in grado di riordinare automaticamente la sequenza delle istruzioni in modo da garantire il comportamento dei salti previsto dal programmatore. I compilatori MIPS inseriscono sempre un'istruzione non dipendente dal risultato del salto immediatamente dopo l'istruzione di salto ritardato; quando il salto viene preso, viene modificato l'indirizzo dell'istruzione che segue questa istruzione "sicura". Nell'esempio di Figura 4.29, l'istruzione `add` prima della branch non ha influenza sull'istruzione di salto, può quindi essere spostata subito dopo la branch per coprire completamente il ritardo nel salto. Dal momento che i salti ritardati sono utili quando i ritardi sono brevi, nessun processore utilizza questo tipo di salti con un ritardo maggiore di un ciclo di clock. Per ritardi più lunghi solitamente si utilizza la predizione hardware sui salti.

## Riepilogo sulla pipeline

La pipeline è una tecnica che sfrutta il **parallelismo** tra le istruzioni di un programma software scritto in maniera sequenziale. Rispetto alla programmazione dei multiprocessori (vedi Cap. 6), ha il grosso vantaggio di essere sostanzialmente invisibile al programmatore.

Nei prossimi paragrafi di questo capitolo svilupperemo il concetto di pipeline basandoci sul sottoinsieme di istruzioni RISC-V utilizzato nel paragrafo 4.4 per l'implementazione a singolo ciclo e mostreremo una versione semplificata della pipeline RISC-V. Esamineremo, quindi, i problemi introdotti dall'uso delle pipeline e le prestazioni che possono essere ottenute in condizioni operative tipiche.

I lettori maggiormente interessati al software e all'impatto della pipeline sulle prestazioni hanno già acquisito le conoscenze fondamentali e possono passare al paragrafo 4.10, che introduce alcuni concetti avanzati su questa tecnica (come la pipeline superscalare e il riordinamento dinamico del codice). Nel paragrafo 4.11, invece, viene esaminata la pipeline di alcuni dei microprocessori più recenti.

Se invece siete interessati a comprendere come si possa implementare una pipeline e alle sfide poste dalla gestione degli hazard, potete passare a esaminare il progetto di un'unità di elaborazione con pipeline e l'unità di controllo di base associata, trattati nel paragrafo 4.6. Utilizzando i concetti appresi in questo paragrafo, potete poi passare a esaminare i meccanismi di implementazione della propagazione e dello stallo, trattati nel paragrafo 4.7. Infine, potete leggere il paragrafo 4.8 per vedere alcune soluzioni agli hazard sul controllo e il paragrafo 4.9 per la gestione delle eccezioni.



## Autovalutazione

Per ciascuna sequenza di istruzioni riportate nella tabella sottostante, stabilire se occorre mettere in stallo la pipeline, se si possono evitare gli stalli utilizzando la propagazione, o se si può eseguire il codice senza stalli né propagazione.

Sequenza 1	Sequenza 2	Sequenza 3
<code>ld x10, 0(x10)</code>	<code>add x11, x10, x10</code>	<code>addi x11, x10, 1</code>
<code>add x11, x10, x10</code>	<code>addi x12, x10, 5</code>	<code>addi x12, x10, 2</code>
	<code>addi x14, x11, 5</code>	<code>addi x13, x10, 3</code>
		<code>addi x14, x10, 4</code>
		<code>addi x15, x10, 5</code>

## Capire le prestazioni dei programmi



Al di fuori del sottosistema di memoria, la pipeline è di solito l'elemento più importante per determinare il CPI di un processore e quindi le sue prestazioni. Come vedremo nel paragrafo 4.10, la comprensione delle prestazioni di un moderno processore con pipeline a esecuzione parallela è difficile, perché i problemi che si incontrano sono più complessi di quelli di un semplice processore a singola pipeline. In ogni caso, gli hazard strutturali, sui dati e sul controllo rimangono rilevanti sia per le pipeline semplici sia per quelle più sofisticate.

Per le pipeline moderne, gli hazard strutturali spesso riguardano le unità a virgola mobile che possono operare non completamente in pipeline, mentre gli hazard sul controllo sono di solito un problema che riguarda maggiormente i programmi che lavorano su numeri interi (questi tendono ad avere più salti condizionati sui quali è difficile fare previsioni). Gli hazard sui dati possono rappresentare un vero e proprio collo di bottiglia per le prestazioni sia dei programmi che lavorano su numeri interi sia di quelli che lavorano su numeri in virgola mobile. Spesso è più facile trattare gli hazard sui dati nei programmi in virgola mobile: in questi casi, infatti, per la ridotta frequenza dei salti condizionati e per l'accesso più regolare ai dati in memoria, diventa più facile per il compilatore riordinare la sequenza delle istruzioni in modo da evitare gli hazard. È invece più difficile effettuare questa ottimizzazione nei programmi che lavorano su interi, i quali accedono ai dati in memoria in modo più irregolare, dato che utilizzano maggiormente i puntatori. Come vedremo nel paragrafo 4.10, esistono tecniche più ambiziose, basate sia sui compilatori sia sull'hardware, per ridurre le dipendenze tra i dati attraverso la riorganizzazione dell'ordine di esecuzione delle istruzioni.

### QUADRO D'INSIEME



**Latenza:** numero di stadi di una pipeline o numero di stadi tra due istruzioni successive in esecuzione.



La pipeline aumenta il numero di istruzioni che vengono eseguite contemporaneamente e la frequenza con cui viene iniziata e terminata l'esecuzione delle istruzioni. Questa tecnica non riduce il tempo richiesto per completare l'esecuzione della singola istruzione, chiamato anche **latenza**. Per esempio, una pipeline a cinque stadi richiede sempre 5 cicli di clock per completare l'esecuzione di un'istruzione. Utilizzando la terminologia introdotta nel Capitolo 1, la pipeline migliora il *throughput* delle istruzioni ma *non il tempo di esecuzione*, o *latenza*, della singola istruzione. Gli insiemi di istruzioni possono semplificare o rendere particolarmente difficile il compito ai progettisti delle pipeline, i quali devono comunque gestire gli hazard strutturali, sul controllo e sui dati. I meccanismi di predizione dei salti condizionati e la propagazione aiutano a rendere veloce l'elaborazione mantenendo un funzionamento corretto del calcolatore. ■

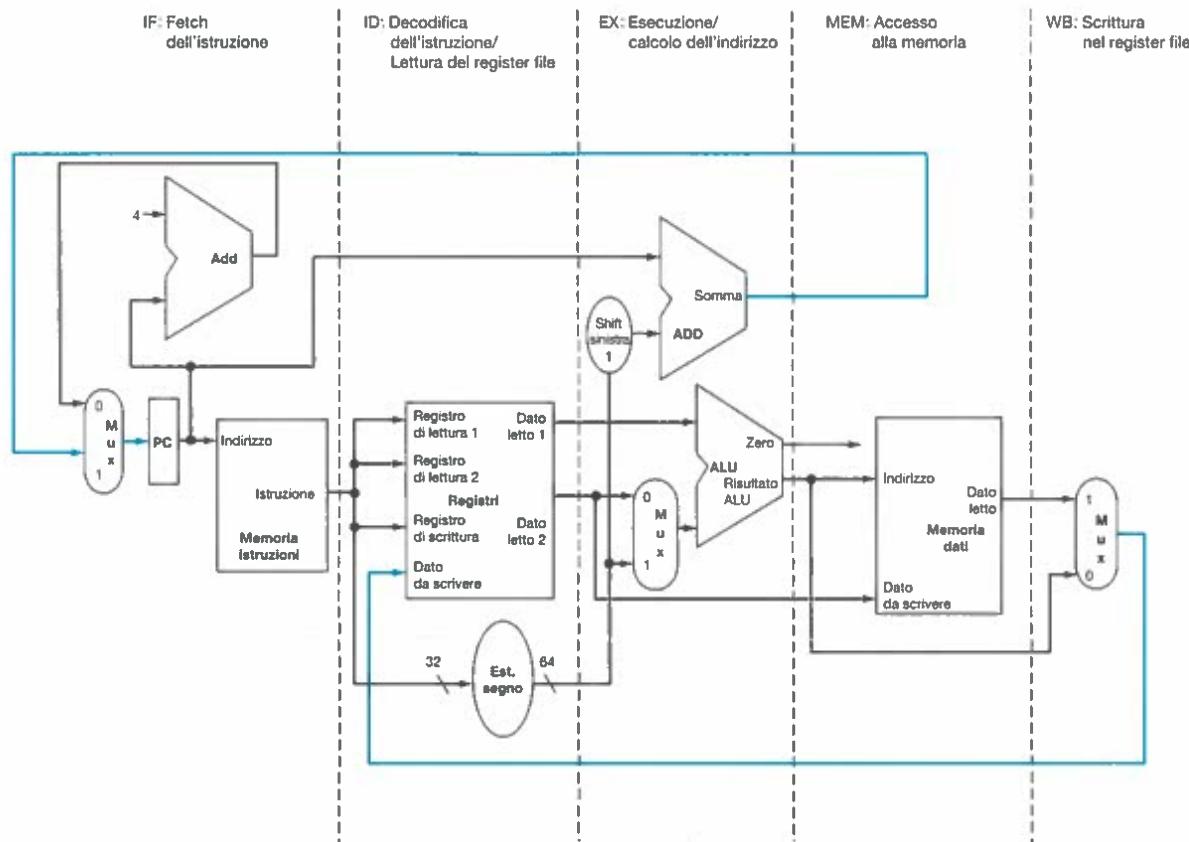
## 4.6 Unità di elaborazione con pipeline e unità di controllo associata

In questo, c'è meno di quanto vedano gli occhi.

Tallulah Bankhead, riferendosi ad Alexander Woollcott, 1922

La Figura 4.31 mostra l'unità di elaborazione a singolo ciclo progettata nel paragrafo 4.4 con i cinque stadi della pipeline ben identificati. La suddivisione di un'istruzione in cinque fasi implica una pipeline a cinque stadi, il che significa che in un singolo ciclo di clock ci potranno essere fino a cinque istruzioni in esecuzione contemporaneamente. Occorre quindi separare l'unità di elaborazione in cinque parti, ciascuna delle quali prende il nome dalla fase di esecuzione dell'istruzione corrispondente:

1. IF: fetch dell'istruzione (*instruction fetch*);
2. ID: decodifica dell'istruzione (*instruction decode*) e lettura del register file;



**Figura 4.31** L'unità di elaborazione a singolo ciclo del paragrafo 4.4 (simile a quella di Figura 4.17). Ciascun passaggio dell'esecuzione di un'istruzione può essere messo in corrispondenza con una parte dell'unità di elaborazione, seguendo il flusso di elaborazione da sinistra verso destra. Le sole eccezioni, mostrate in blu, sono costituite dall'aggiornamento del PC e dalla scrittura del risultato: nella fase di WB il risultato calcolato dalla ALU o il dato letto dalla memoria dati viene portato alla sinistra dello schema per scriverlo nel register file (di solito utilizziamo il colore blu per evidenziare i segnali di controllo, ma in questa figura vogliamo evidenziare linee dati).

3. EX: esecuzione (*execution*) o calcolo dell'indirizzo;
4. MEM: accesso alla memoria dati (*memory*);
5. WB: scrittura del risultato nel register file (*write-back*).

Queste cinque parti corrispondono approssimativamente al modo in cui è stata disegnata l'unità di elaborazione dei dati di Figura 4.31; durante l'esecuzione, istruzioni e dati si spostano generalmente da sinistra a destra, attraversando le cinque fasi. Tornando all'analogia con il bucato, i panni diventano più puliti, più asciutti e più ordinati mano a mano che si prosegue la lavorazione, e non devono mai ritornare alle fasi precedenti.

Vi sono tuttavia due eccezioni al flusso di esecuzione da sinistra verso destra:

- lo stadio di scrittura del risultato, nel quale viene scritto il risultato dell'esecuzione nel register file al centro dell'unità di elaborazione;
- la selezione del valore successivo del PC, che viene scelto tra il valore del PC incrementato di 4 e l'indirizzo di destinazione del salto proveniente dallo stadio MEM.

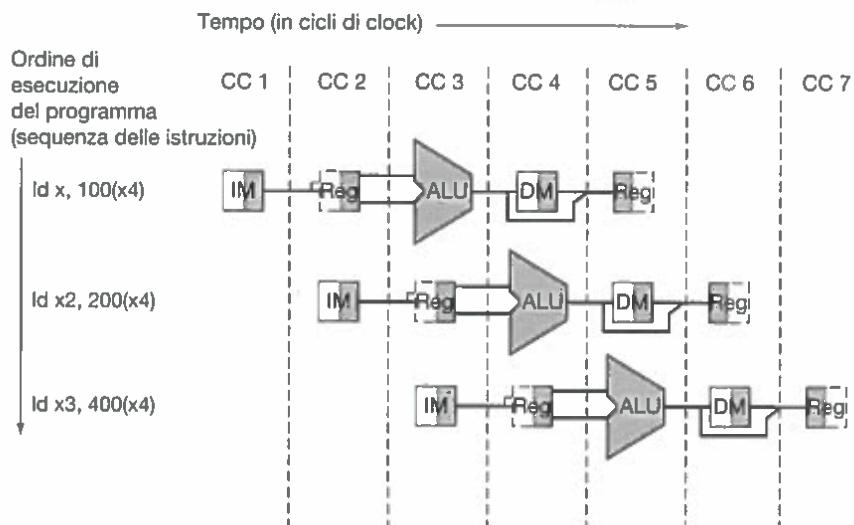
Il flusso dei dati da destra a sinistra non ha effetto sull'istruzione corrente: questo spostamento all'indietro dei dati influenzera' solo le istruzioni successive presenti nella pipeline. Si noti che la freccia da destra verso sinistra posta nella parte inferiore dello schema può causare hazard sui dati, mentre la seconda freccia, nella parte superiore dello schema, può provocare hazard sul controllo.

Un modo per descrivere ciò che succede nell'esecuzione con pipeline è quello di ipotizzare che ciascuna istruzione disponga di una propria unità di elaborazione e di inserire tale unità nella corretta posizione lungo l'asse dei tempi, visualizzando le relazioni tra le diverse istruzioni. La Figura 4.32 mostra l'esecuzione delle stesse istruzioni di Figura 4.25, riportando per ciascuna di esse l'unità di elaborazione posizionata correttamente su un asse unico dei tempi. In questa figura è stata utilizzata la versione schematizzata dell'unità di elaborazione di Figura 4.31.

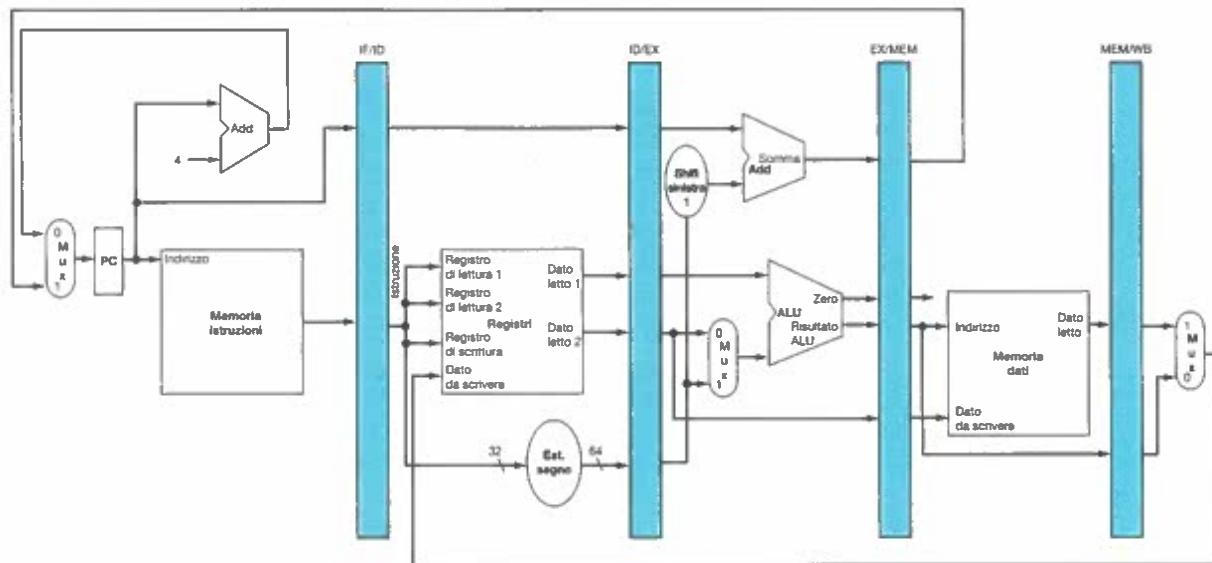
La Figura 4.32 sembra suggerire che tre istruzioni richiedano tre unità di elaborazione. Invece, possiamo aggiungere dei registri per salvare i dati parziali, in modo tale che le diverse parti di una singola unità di elaborazione possano essere condivise dalle diverse istruzioni durante la loro esecuzione.

Per esempio, come mostrato in Figura 4.32, la memoria istruzioni viene utilizzata in uno solo dei cinque stadi di un'istruzione, per cui potrà essere impiegata dalle istruzioni successive quando l'istruzione passa attraverso gli altri quattro stadi. Per preservare l'istruzione prelevata dalla memoria a beneficio degli altri quattro stadi, essa deve essere salvata in un registro. Considerazioni analoghe valgono per tutti gli altri stadi della pipeline, per cui occorrerà aggiungere dei registri all'interfaccia tra le coppie contigue di stadi della pipeline di Figura 4.31. Ricordando l'analogia con la lavandaia, ci potrebbe essere un cestino tra due stadi di lavorazione dove riporre i panni destinati alla fase successiva.

La Figura 4.33 mostra l'unità di elaborazione con pipeline in cui sono evidenziati i registri della pipeline. In ogni ciclo di clock, tutte le istruzioni procedono da un registro di pipeline a quello successivo. I registri prendono il nome dai due stadi che separano: per esempio, il registro di pipeline inserito tra gli stadi IF e ID è chiamato IF/ID.



**Figura 4.32** Esecuzione in pipeline delle stesse istruzioni eseguite dall'unità di elaborazione a singolo ciclo di Figura 4.31. In modo analogo alle Figure 4.26-4.28, si suppone che ciascuna istruzione disponga della propria unità di elaborazione e si evidenzia la parte utilizzata di ciascuna unità funzionale. A differenza delle figure precedenti, qui i diversi stadi sono etichettati con il nome delle risorse fisiche utilizzate, che corrisponde a una parte dell'unità di elaborazione di Figura 4.31. IM rappresenta la memoria istruzioni, Reg sta per register file, che assieme all'unità di estensione del segno viene utilizzato nello stadio di decodifica dell'istruzione e lettura del register file (ID) ecc. Per garantire una corretta sequenza temporale, questa unità di elaborazione schematizzata divide il register file in due parti logiche: i registri letti durante la fase di decodifica (ID) e i registri scritti nella fase finale (WB). Questo duplice impiego del register file viene rappresentato tracciando con una linea tratteggiata, nello stadio ID, il bordo della parte sinistra del register file che non è evidenziato (poiché non avviene alcuna operazione di scrittura in questo stadio); nello stadio WB, invece, è tratteggiato il bordo della parte destra non evidenziata (poiché non avvengono operazioni di lettura). Come già fatto in precedenza, si assume che il register file venga letto nella prima metà del periodo di clock e scritto nella seconda.



**Figura 4.33** Versione con pipeline dell'unità di elaborazione di Figura 4.31. I registri di pipeline, evidenziati in blu, separano i diversi stadi della pipeline e per etichettarli si usa il nome dei relativi stadi che separano; per esempio, il primo registro viene etichettato con *IF/D*, poiché separa lo stadio *IF*, di prelevamento dell'istruzione, dallo stadio *ID*, di decodifica dell'istruzione. I registri devono avere un'ampiezza sufficiente a memorizzare tutti i dati corrispondenti alle linee che li devono attraversare. Il registro *IF/D*, per esempio, deve essere ampio 96 bit, poiché deve contenere sia l'istruzione prelevata dalla memoria su 32 bit sia il contenuto del PC incrementato di 4, su 64 bit. Gli altri registri verranno esaminati più in dettaglio nel seguito del capitolo; anticipiamo solamente che gli altri tre registri di pipeline contengono rispettivamente 256, 193 e 128 bit.

Si noti che non c'è alcun registro di pipeline alla fine dello stadio di write-back. Tutte le istruzioni devono modificare qualche elemento dello stato del processore: il register file, la memoria o il PC, per cui un registro di pipeline separato sarebbe ridondante. Per esempio, un'istruzione di load scrive il proprio risultato in uno dei 32 registri del register file, e se l'istruzione successiva richiedesse tale valore potrebbe caricarlo direttamente da lì.

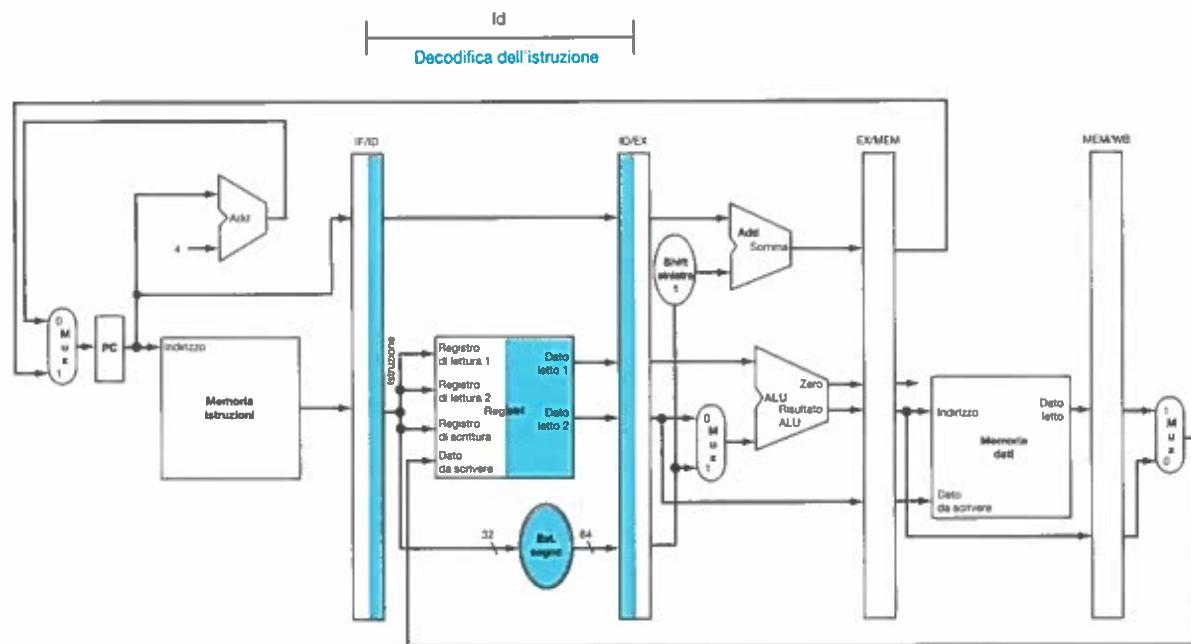
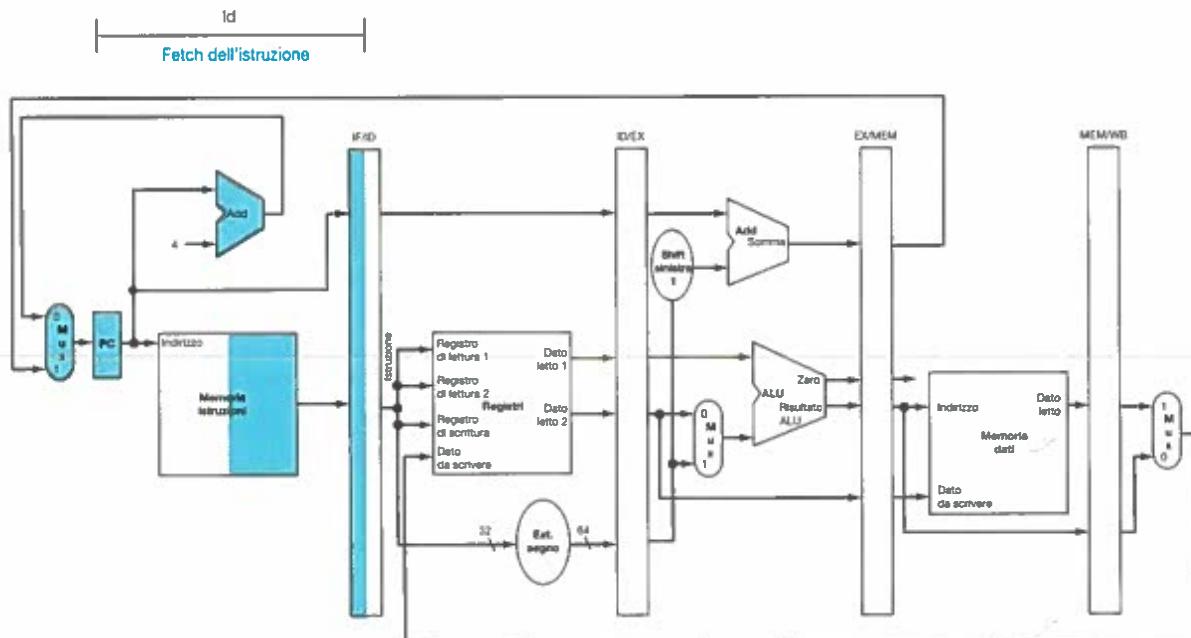
Ovviamente ogni istruzione aggiorna il PC, incrementandolo di 4 o scrivendo l'indirizzo di destinazione del salto. Il PC può essere visto come il registro di pipeline che alimenta lo stadio IF. Tuttavia, a differenza dei registri di pipeline evidenziati in Figura 4.33, il PC fa parte dello stato dell'architettura, visibile all'esterno: il suo contenuto deve essere salvato quando si verifica un'eccezione, mentre il contenuto dei registri di pipeline veri e propri viene perso. Nell'analogia con la lavanderia, il PC rappresenta il cestino che contiene i panni sporchi prima dell'inizio del bucato.

Per mostrare il funzionamento della pipeline, in questo capitolo riporteremo delle sequenze di schemi che rappresentano le operazioni effettuate dalla pipeline in istanti di tempo successivi. Potreste pensare che occorra molto tempo per comprendere il contenuto di tutte queste pagine, ma non è così: le sequenze di schemi si analizzano in maniera assai veloce, poiché possono essere semplicemente confrontate tra loro per vedere che cosa cambia in ciascun ciclo di clock. Il paragrafo 4.7 descrive che cosa succede quando si verificano degli hazard sui dati nelle istruzioni presenti nella pipeline, ma per ora ignoriamo queste complicazioni.

Le Figure dalla 4.34 alla 4.36 costituiscono la prima sequenza di schemi e mostrano (evidenziata in blu) la parte dell'unità di elaborazione attiva nei diversi stadi di esecuzione di un'istruzione di load. Questa istruzione viene esaminata per prima, essendo attiva in tutti e cinque gli stadi. Come nelle Figure 4.28-4.30, verrà evidenziata la *metà destra* dei registri o della memoria quando queste unità vengono *lette*, e la *metà sinistra* quando vengono *scritte*.

In ciascuna figura si riporterà il nome dell'istruzione, 1d, insieme al nome dello stadio attivo della pipeline. I cinque stadi sono i seguenti.

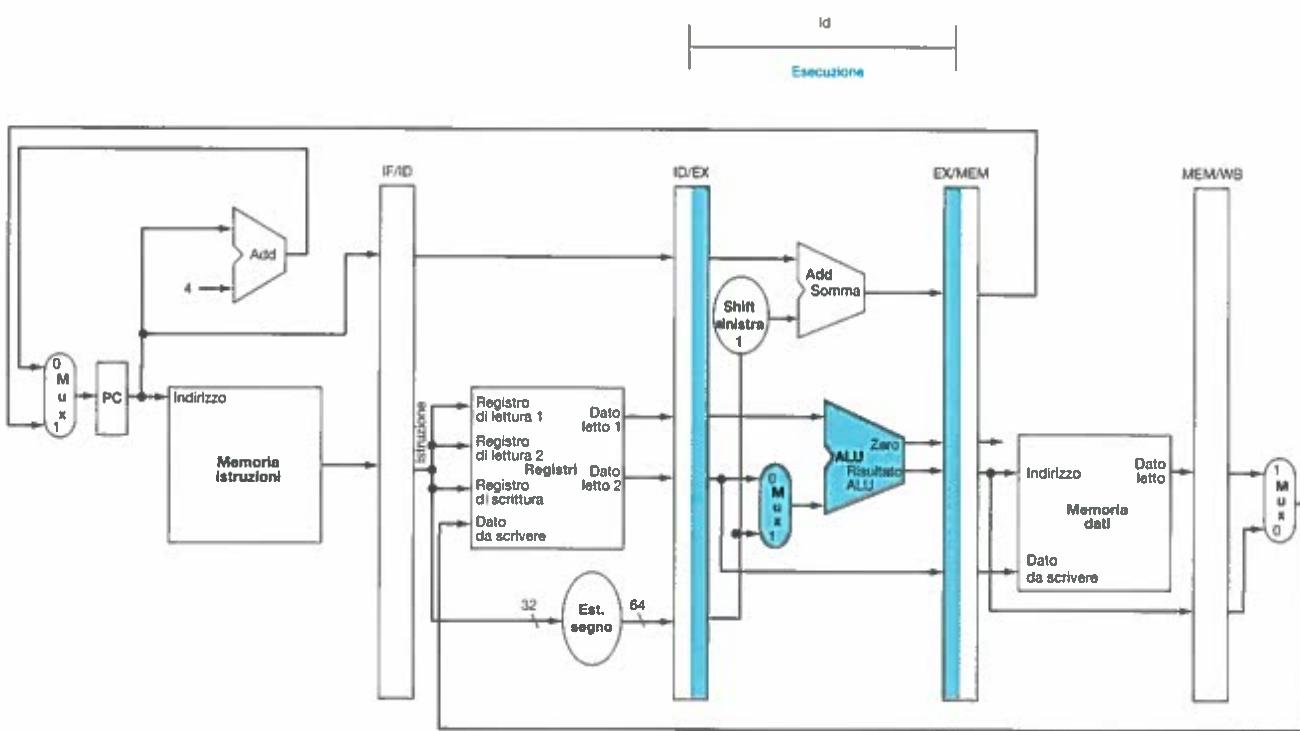
1. *Fetch dell'istruzione*. La parte superiore della Figura 4.34 mostra la fase in cui l'istruzione viene prelevata dalla memoria istruzioni (utilizzando l'indirizzo contenuto nel PC) e scritta nel registro di pipeline IF/ID. L'indirizzo contenuto nel PC viene incrementato di 4 e riscritto nel PC stesso, in modo



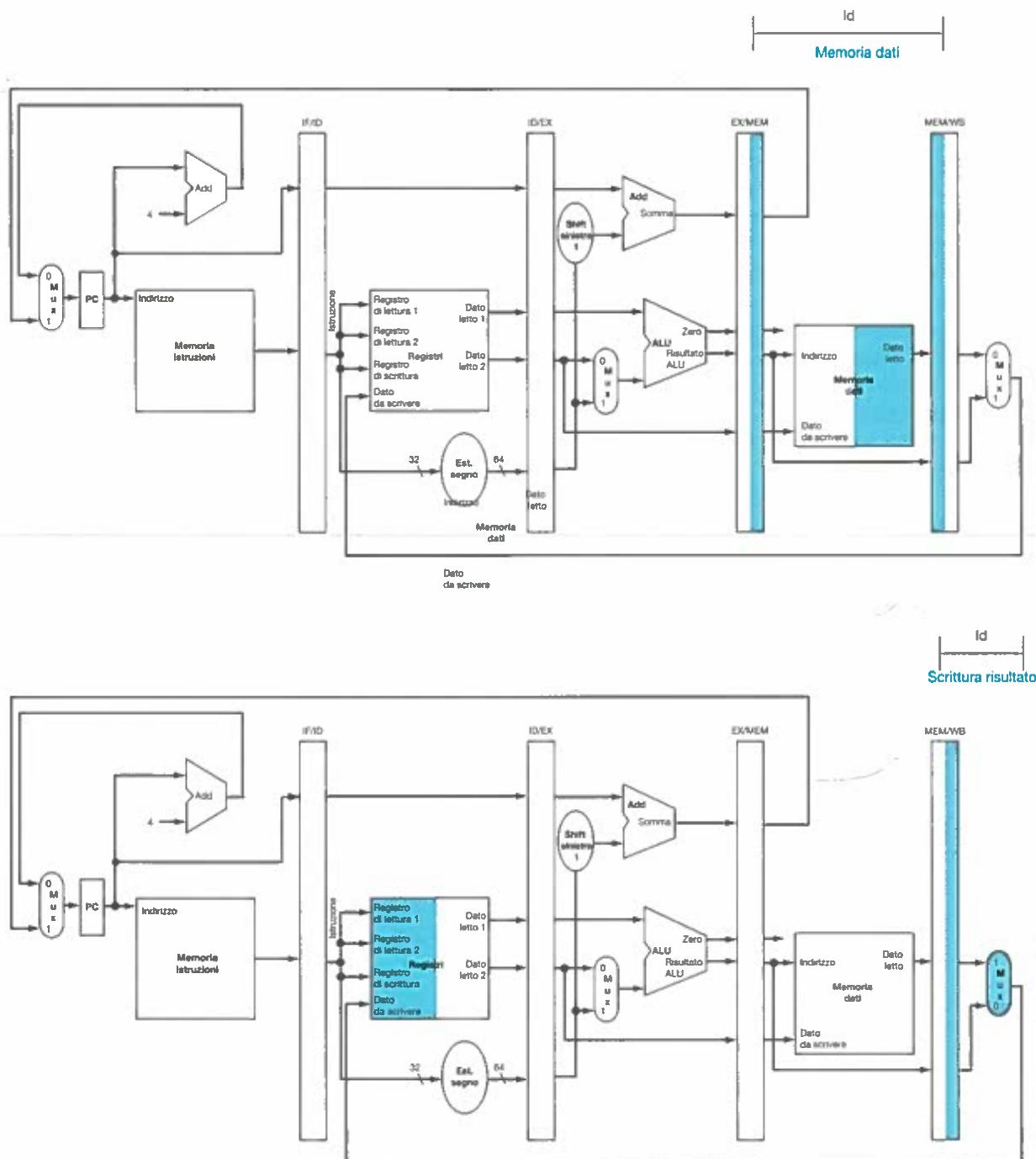
**Figura 4.34 IF e ID:** il primo e il secondo stadio della pipeline durante l'esecuzione di un'istruzione di load; le parti attive dell'unità di elaborazione mostrata in Figura 4.33 sono evidenziate in blu. La convenzione per le parti evidenziate è la stessa utilizzata in Figura 4.26. Come visto nel paragrafo 4.2, non si può confondere la lettura con la scrittura dei registri, in quanto il loro contenuto varia solamente in corrispondenza del fronte del clock. Sebbene la load nello stadio 2 richieda solo il contenuto del primo registro, il processore non può sapere quale istruzione sta decodificando, per cui estende il segnale del contenuto del campo immediato e legge il contenuto di entrambi i registri scrivendolo nel registro di pipeline ID/EX. Tre operandi non servono, ma la realizzazione dell'unità di controllo risulta più semplice se tutti e tre gli operandi vengono comunque scritti nel registro ID/EX.

da essere pronto per il ciclo di clock successivo. Il PC viene anche salvato nel registro di pipeline IF/ID, perché potrebbe essere richiesto da un'istruzione successiva, per esempio da una `beq`. Il calcolatore non può prevedere quale istruzione verrà prelevata dalla memoria, per cui deve preparare i dati per l'esecuzione di tutte le possibili istruzioni, trasportando lungo la pipeline tutte le informazioni potenzialmente necessarie.

2. *Decodifica dell'istruzione e lettura del register file*. La parte inferiore di Figura 4.34 mostra che la parte contenente l'istruzione del registro di pipeline IF/ID fornisce il campo immediato su 12 bit, che viene convertito a 64 bit mediante estensione del segno. IF/ID fornisce anche il numero dei due registri da leggere. Il contenuto dei due registri e il campo immediato esteso a 64 bit vengono memorizzati nel registro di pipeline ID/EX insieme al valore del PC. Anche in questo stadio, infatti, viene trasferito tutto ciò che potrebbe essere necessario per l'esecuzione delle fasi successive di tutte le possibili istruzioni.
3. *Esecuzione o calcolo dell'indirizzo*. La Figura 4.35 mostra che l'istruzione di load legge dal registro di pipeline ID/EX il contenuto del registro 1 e il contenuto del campo immediato, esteso a 64 bit e dotato di segno, li somma utilizzando la ALU e scrive il risultato della somma nel registro di pipeline EX/MEM.
4. *Accesso alla memoria dati*. La parte superiore di Figura 4.36 mostra l'istruzione di load mentre legge la memoria dati utilizzando come indirizzo il valore letto dal registro di pipeline EX/MEM. Il dato viene quindi scritto nel registro di pipeline MEM/WB.
5. *Scrittura del risultato*. La parte inferiore di Figura 4.36 riporta la fase finale dell'esecuzione: il dato viene letto dal registro di pipeline MEM/WB e scritto nel file posto al centro della figura.



**Figura 4.35 EX:** Il terzo stadio della pipeline durante l'esecuzione di un'istruzione di load; le parti attive dell'unità di elaborazione mostrata in Figura 4.33 sono evidenziate in blu. Il contenuto del registro viene sommato al contenuto del campo immediato dotato di segno, esteso a 64 bit, e il risultato della somma viene scritto nel registro di pipeline EX/MEM.

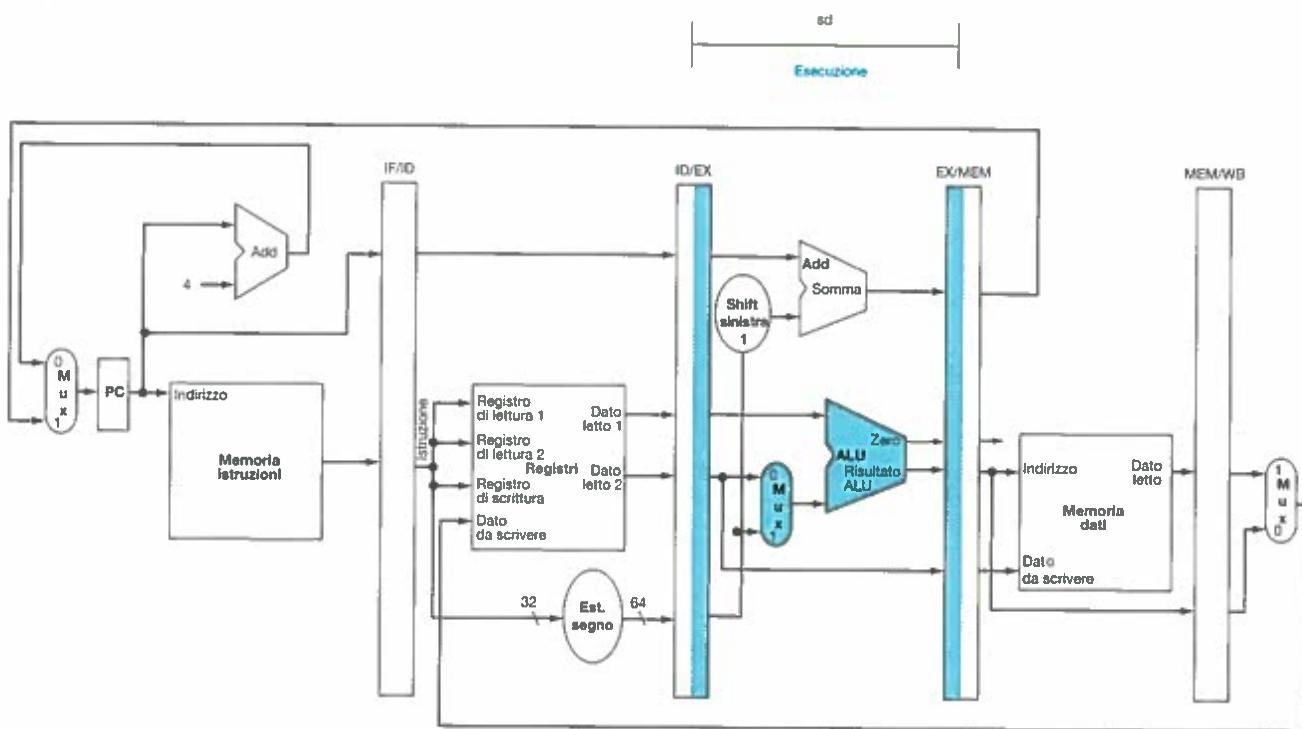


**Figura 4.36 MEM e WB:** il quarto e il quinto stadio della pipeline durante l'esecuzione di un'istruzione di load; le parti attive dell'unità di elaborazione mostrata in Figura 4.33 sono evidenziate in blu. La memoria dati viene letta utilizzando l'indirizzo contenuto nel registro di pipeline EX/MEM e il dato letto viene scritto nel registro di pipeline MEM/WB. Poi, nello stadio WB, il dato viene letto dal registro di pipeline MEM/WB e scritto nel register file, al centro della figura. Si noti che c'è un baco in questo schema che verrà corretto in Figura 4.39.

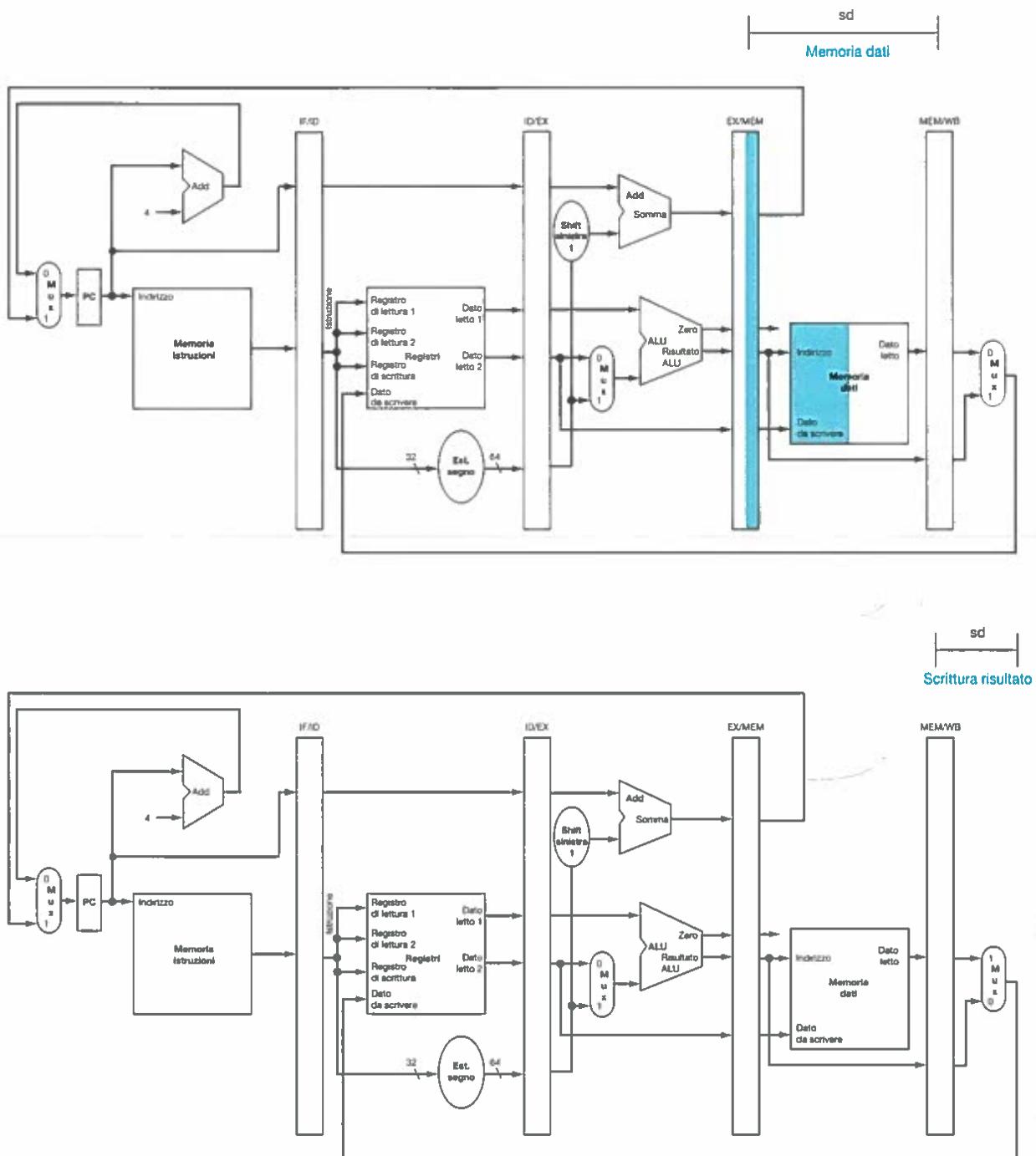
Questa carrellata sull'esecuzione dell'istruzione di load mostra che tutte le informazioni necessarie a uno stadio successivo della pipeline devono essere trasportate attraverso i registri di pipeline. La stessa descrizione dell'esecuzione di un'istruzione di store evidenzia le analogie con la ld nell'esecuzione e nel passaggio delle informazioni agli stadi successivi. Ecco i cinque stadi di pipeline per un'istruzione di store.



1. *Fetch dell'istruzione.* L'istruzione viene prelevata dalla memoria utilizzando l'indirizzo contenuto nel PC e viene quindi scritta nel registro di pipeline IF/ID. Questo stadio viene eseguito prima che l'istruzione possa essere identificata, per cui la parte superiore di Figura 4.34 è valida non solo per la load ma anche per la store.
2. *Decodifica dell'istruzione e lettura del register file.* L'istruzione presente nel registro di pipeline IF/ID fornisce il numero dei due registri per leggere il campo immediato dotato di segno ed estenderlo a 64 bit. Il contenuto di questi tre campi a 64 bit, il contenuto dei due registri e il campo immediato esteso, vengono memorizzati nel registro di pipeline ID/EX, insieme al valore del PC. La parte inferiore della Figura 4.34 relativa all'istruzione di load descrive il funzionamento del secondo stadio anche per le istruzioni di store. I primi due stadi funzionano nello stesso modo per tutte le istruzioni, perché è troppo presto per capire quale sia il tipo di istruzione. (L'istruzione di store utilizza il campo rs2 per leggere il secondo registro in questo stadio della pipeline. Questo dettaglio non è riportato in questo diagramma di pipeline, per cui possiamo utilizzare lo stesso schema per entrambi i casi.)
3. *Esecuzione e calcolo dell'indirizzo.* La Figura 4.37 mostra il terzo passo nel quale l'indirizzo della memoria viene scritto nel registro di pipeline EX/MEM.
4. *Accesso alla memoria dati.* La parte superiore di Figura 4.38 mostra in che modo avviene la scrittura del dato in memoria. Si noti che il registro contenente il dato da scrivere era stato letto nel precedente stadio di decodifica e il suo contenuto era stato memorizzato nel registro ID/EX. L'unico modo per rendere questo dato disponibile anche nello stadio MEM è trasportarlo dal registro ID/EX al registro EX/MEM durante lo stadio EX, scrivendolo in EX/MEM esattamente come viene scritto l'indirizzo della memoria.



**Figura 4.37 EX:** il terzo stadio della pipeline durante l'esecuzione di un'istruzione di store. A differenza del terzo stadio di esecuzione dell'istruzione di load mostrato in Figura 4.35, il contenuto del secondo registro viene scritto nel registro di pipeline EX/MEM per essere utilizzato nello stadio successivo. Sebbene non sia sbagliato scrivere sempre il contenuto del secondo registro nel registro di pipeline EX/MEM, effettueremo questa operazione solo per le istruzioni di store, in modo da rendere più facile la comprensione della pipeline.



**Figura 4.38 MEM e WB: il quarto e il quinto stadio della pipeline durante l'esecuzione di un'istruzione di store.** Nel quarto stadio il dato viene scritto nella memoria dati; si noti che esso proviene dal registro di pipeline EX/MEM e che nulla viene modificato nel registro di pipeline MEM/WB. Una volta scritto il dato in memoria, per l'istruzione di store non rimane più nulla da fare, per cui nello stadio 5 non accade niente.

5. **Scrittura del risultato.** La parte inferiore di Figura 4.38 riporta la fase finale dell'esecuzione dell'istruzione: si vede che nel caso della store non succede nulla. Poiché l'esecuzione di tutte le istruzioni che seguono la store deve ancora essere completata, non c'è alcun modo per accelerarle; quindi un'istruzione deve sempre passare attraverso tutti gli stadi, anche se a volte non ci sono operazioni da compiere, poiché le istruzioni successive stanno già procedendo alla massima velocità possibile.

Il funzionamento dell'istruzione di store sottolinea ancora una volta come, per trasportare un'informazione da uno stadio della pipeline a quello successivo, quell'informazione debba essere scritta in un registro di pipeline; se così non fosse, l'informazione verrebbe persa nel momento in cui l'istruzione successiva raggiungesse quella fase. Per l'istruzione di store abbiamo dovuto trasportare il contenuto del secondo dei due registri letti dallo stadio ID fino allo stadio MEM, dove viene scritto in memoria. Il dato proveniente dal register file viene prima scritto nel registro di pipeline ID/EX e poi copiato nel successivo registro EX/MEM.

Le istruzioni di load e di store permettono di fare una seconda considerazione: ciascun componente funzionale dell'unità di elaborazione, come la memoria istruzioni, le porte di lettura del register file, la ALU, la memoria dati e la porta di scrittura del register file, può essere utilizzato in *un solo* stadio della pipeline; in caso contrario, si incorrerebbe in un *hazard strutturale* (par. 4.5). Di conseguenza, questi componenti funzionali e i loro segnali di controllo possono essere associati a un unico stadio della pipeline.

A questo punto, possiamo svelare un baco nell'implementazione hardware dell'istruzione di load che alcuni lettori avranno già notato. Quale registro viene modificato nello stadio finale della load? In particolare, quale istruzione fornisce il numero del registro da scrivere? In questo schema il numero del registro da scrivere è contenuto nel relativo campo dell'istruzione, presente nel registro di pipeline IF/ID. Tuttavia tale numero di registro si riferisce a un'istruzione che è *successiva* alla load!

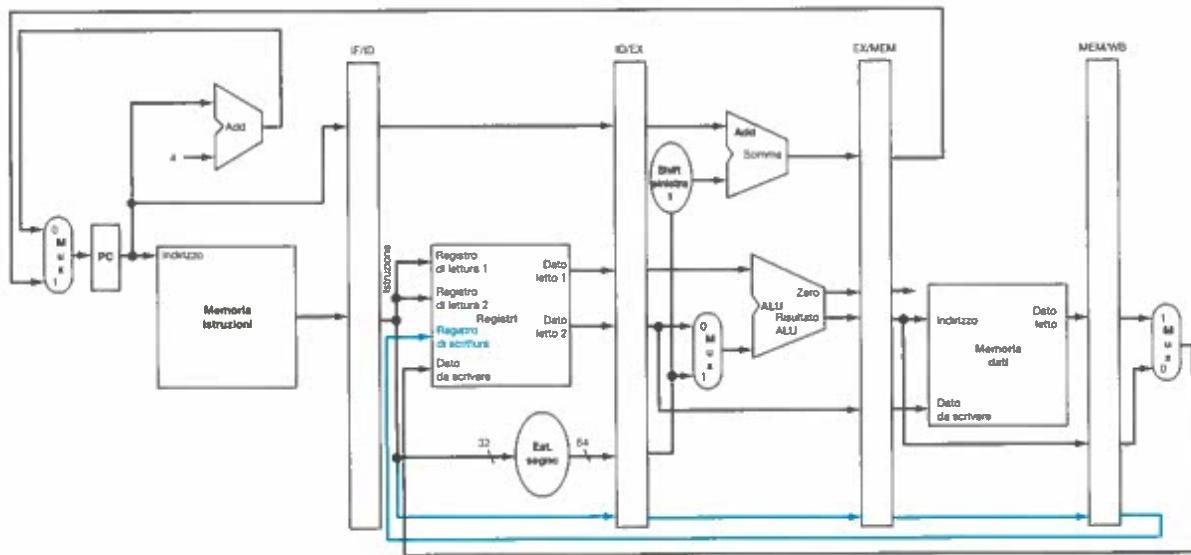
Quindi, è necessario preservare il numero del registro di destinazione dell'istruzione di load: così come per l'istruzione di store avevamo trasportato il *contenuto* del registro rs2 dal registro di pipeline ID/EX al registro EX/MEM (per poterlo utilizzare nello stadio MEM), così per l'istruzione di load dobbiamo trasportare il *numero* del registro da scrivere da ID/EX, attraverso EX/MEM, fino a MEM/WB per poterlo poi utilizzare nello stadio di WB. Un modo alternativo di vedere questo passaggio del numero del registro da scrivere è pensare che affinché un'unità di elaborazione dotata di pipeline possa essere condivisa da più istruzioni, occorre preservare alcuni campi dell'istruzione prelevata durante lo stadio IF, in modo tale che ciascun registro di pipeline contenga la parte di istruzione che deve essere utilizzata da quello stadio e dagli stadi successivi.

La Figura 4.39 mostra la versione corretta dell'unità di elaborazione: il numero del registro da scrivere viene dapprima scritto nel registro di pipeline ID/EX, poi nel registro EX/MEM e infine nel registro MEM/WB; viene infine utilizzato nello stadio WB per scrivere il dato caricato dalla memoria.

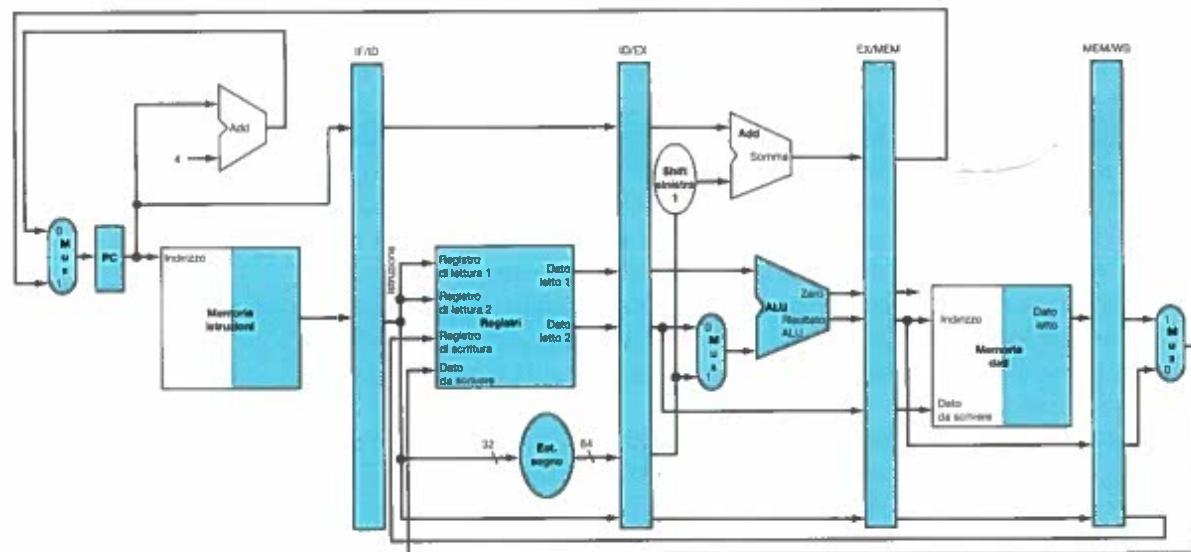
La Figura 4.40 mostra in un unico schema l'unità di elaborazione corretta ed evidenzia i componenti funzionali utilizzati nei cinque stadi di esecuzione dell'istruzione load word, riportati nelle Figure 4.34-4.36. Il paragrafo 4.8 mostra l'implementazione che serve a eseguire correttamente l'istruzione di salto condizionato.

## Rappresentazione grafica delle pipeline

Può essere difficile comprendere il funzionamento delle pipeline, perché in ogni ciclo di clock ci sono molte istruzioni che sono eseguite contemporaneamente nella stessa unità di elaborazione. Per facilitare la comprensione, ci sono due modi principali per disegnare i diagrammi che rappresentano le pipeline: i *diagrammi delle pipeline a più cicli di clock*, come quello di Figura 4.32, e i *diagrammi delle pipeline a singolo ciclo di clock*, come quelli mostrati nelle Figure 4.34-4.38. I diagrammi della pipeline a più cicli di clock sono più sem-



**Figura 4.39** L'unità di elaborazione con pipeline modificata per elaborare correttamente l'istruzione load. Il numero del registro da scrivere ora proviene dal registro di pipeline MEM/WB, così come il dato da scrivere. Il numero del registro è stato trasportato dallo stadio ID fino al registro di pipeline MEM/WB; abbiamo quindi aggiunto cinque bit agli ultimi tre registri di pipeline. Il nuovo cammino è mostrato in blu.



**Figura 4.40** La parte di unità di elaborazione di Figura 4.39 utilizzata per l'esecuzione dei cinque stadi di un'istruzione load.

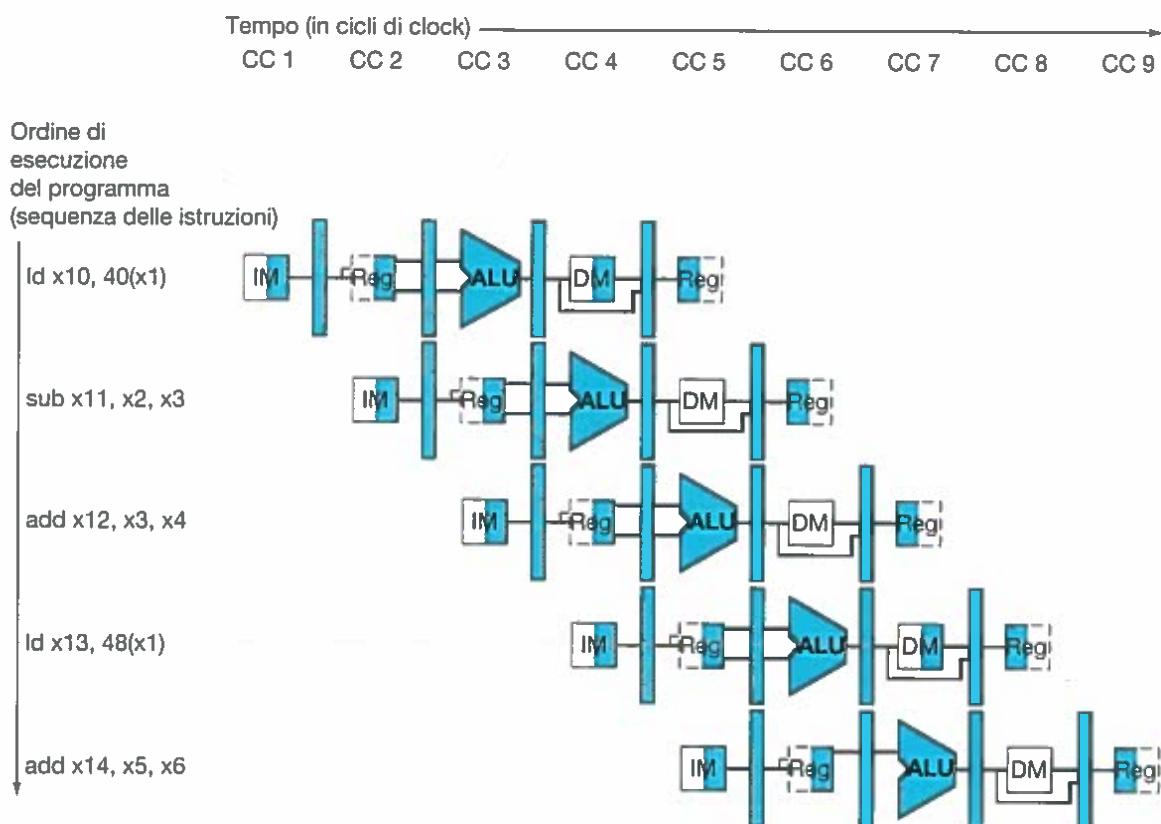
plici ma meno dettagliati. Per esempio, consideriamo la seguente sequenza di cinque istruzioni:

```

ld    x10, 40(x1)
sub   x11, x2, x3
add   x12, x3,x4
lw    x13, 48(x1)
add   x14, x5, x6

```

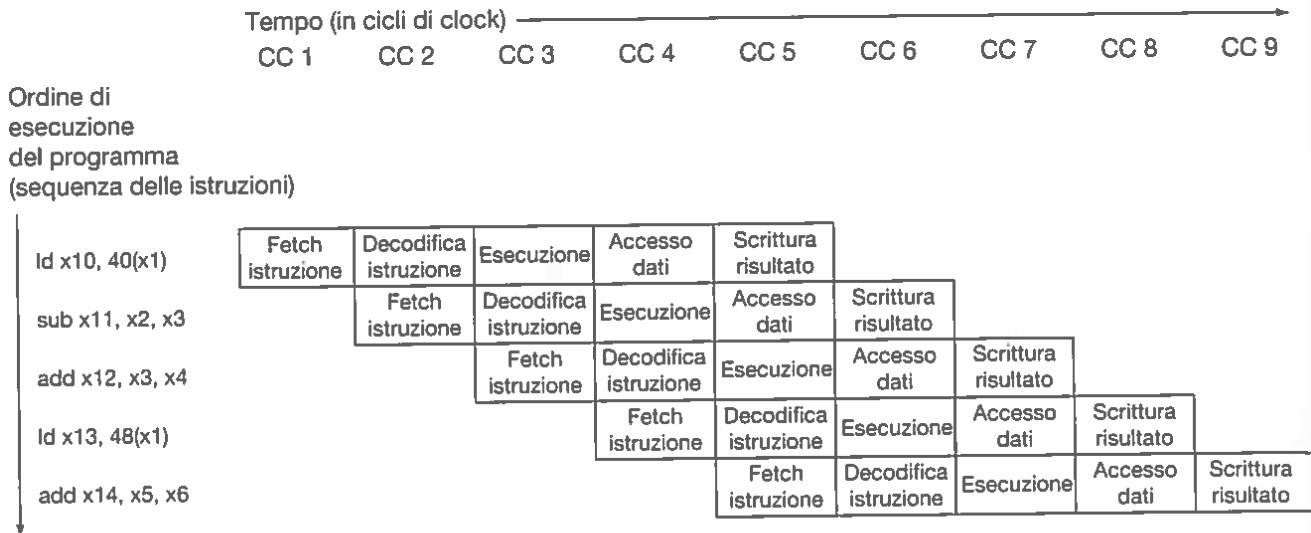
La Figura 4.41 mostra il diagramma della pipeline a più cicli di clock relativo a tali istruzioni: in questo tipo di diagrammi il tempo scorre da sinistra verso



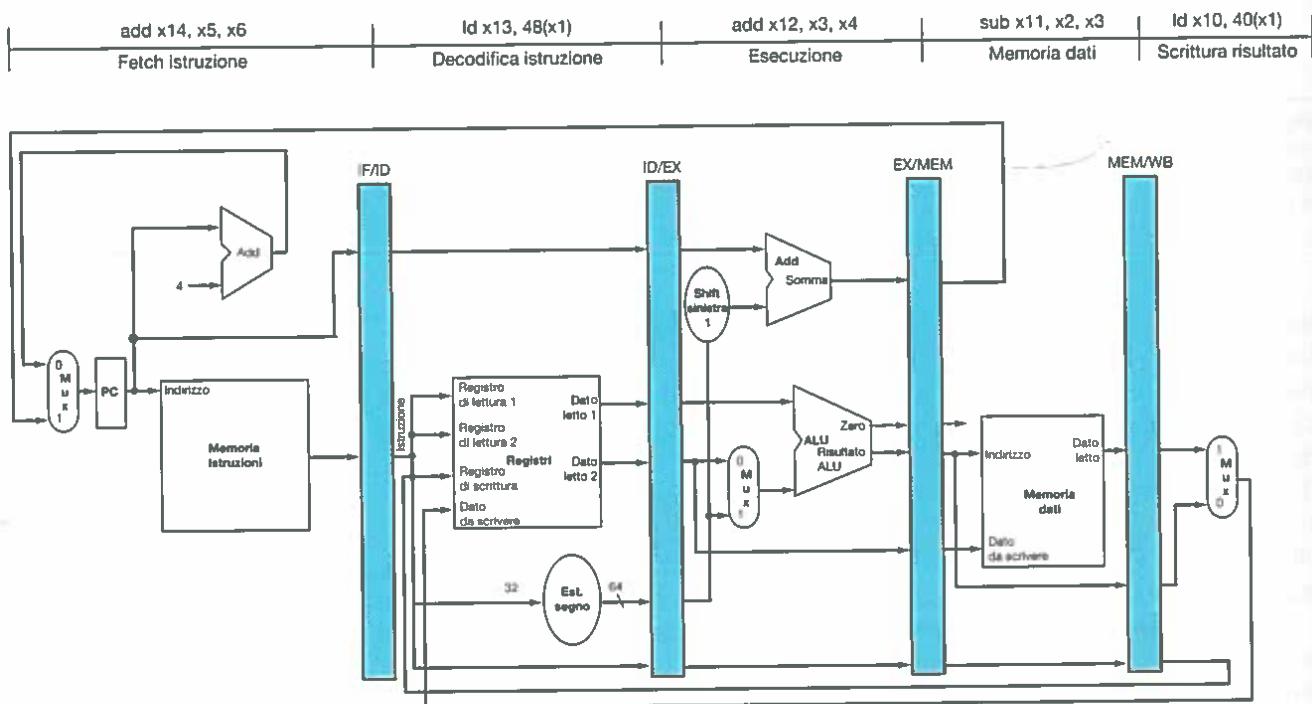
**Figura 4.41** Diagramma della pipeline a più cicli di clock durante l'esecuzione di cinque istruzioni. Questo stile di rappresentazione delle pipeline mostra in una sola figura l'esecuzione completa delle istruzioni di un frammento di codice. Le istruzioni sono riportate sull'asse verticale in ordine di esecuzione (dall'alto verso il basso), mentre i cicli di clock procedono da sinistra verso destra. A differenza di Figura 4.26, qui sono indicati i registri di pipeline inseriti tra le coppie contigue di stadi della pipeline. La Figura 4.42 mostra lo schema tradizionale di questo diagramma.

destra, mentre le istruzioni procedono dall'alto verso il basso (come nella pipeline del bucato di Figura 4.23). Per ogni istruzione, lungo l'asse verticale viene rappresentata una pipeline specificando i diversi stadi posizionati nel ciclo di clock corretto: le cinque fasi vengono rappresentate da una versione stilizzata dell'unità funzionale associata a quello stadio, ma un semplice rettangolo con il nome dello stadio andrebbe bene ugualmente. La Figura 4.42 mostra la versione più tradizionale del diagramma di una pipeline a più cicli di clock; si noti che mentre in Figura 4.41 si visualizzano le risorse fisiche utilizzate in ciascuno stadio, in Figura 4.42 viene utilizzato il *nome* di ciascuno stadio.

I diagrammi delle pipeline a singolo ciclo di clock mostrano lo stato dell'intera unità di elaborazione durante un ciclo di clock e, solitamente, le cinque istruzioni caricate nella pipeline vengono identificate da etichette poste sopra lo stadio in cui si trovano. Utilizzeremo questo tipo di diagramma per mostrare i dettagli di ciò che avviene nella pipeline in ciascun periodo di clock, cosicché saranno mostrati gruppi di diagrammi per illustrare il funzionamento della pipeline su più cicli di clock consecutivi. I diagrammi a più cicli di clock serviranno a fornire un quadro della situazione della pipeline durante l'esecuzione (se siete interessati a vedere ulteriori dettagli sulla Figura 4.41, nel paragrafo 4.13 potete trovare altri gruppi di diagrammi a singolo ciclo). Un diagramma a singolo ciclo rappresenta una sezione verticale, in un certo ciclo di clock, di un insieme di diagrammi a più cicli di clock e mostra come viene utilizzata l'unità di elaborazione da ciascuna delle istruzioni caricate nella pipeline in quel determinato ciclo. Per esempio, la Figura 4.43 mostra il diagramma



**Figura 4.42** Diagramma della pipeline a più cicli di clock in forma tradizionale per le stesse istruzioni di Figura 4.41.



**Figura 4.43** Diagramma della pipeline a singolo ciclo di clock corrispondente al quinto ciclo di clock della pipeline delle Figure 4.41 e 4.42. Come si può vedere, un diagramma a singolo ciclo di clock è una sezione verticale di un diagramma a più cicli di clock.

a singolo ciclo di clock corrispondente al quinto ciclo di clock delle Figure 4.41 e 4.42. Ovviamente i diagrammi a singolo ciclo contengono più dettagli, ma richiedono anche uno spazio molto maggiore per mostrare lo stesso numero di cicli di clock. Negli esercizi verrà richiesto di disegnare tali diagrammi per altre sequenze di istruzioni.

## Autovalutazione

Un gruppo di studenti sta discutendo l'efficienza della pipeline a cinque stadi, quando uno di loro fa notare che non tutte le istruzioni sono attive in ogni stadio della pipeline. Dopo aver deciso di ignorare gli effetti degli hazard, gli studenti fanno le seguenti affermazioni. Quali di queste sono corrette?

1. Permettere ai salti incondizionati e condizionati e alle istruzioni che operano sulla ALU di utilizzare meno stadi dei cinque richiesti dall'istruzione load aumenterà le prestazioni della pipeline in tutte le situazioni.
2. Consentire ad alcune istruzioni di utilizzare meno cicli di clock non aiuta, poiché il throughput è determinato dalla durata del ciclo di clock; il numero degli stadi della pipeline per ciascuna istruzione influenza la latenza, non il throughput.
3. Non si può ridurre il numero di cicli di clock utilizzati dalle istruzioni che operano sulla ALU, per via della scrittura del risultato nella fase di write-back; i salti condizionati e incondizionati possono però utilizzare meno cicli di clock, dunque c'è spazio per un miglioramento.
4. Invece di tentare di fare in modo che le istruzioni durino meno cicli di clock, dovremmo esplorare la possibilità di rendere la pipeline più lunga; le istruzioni, quindi, dovrebbero utilizzare più cicli di clock, ma ciascun ciclo avrebbe una durata inferiore, e ciò produrrebbe un miglioramento delle prestazioni.

## Unità di controllo della pipeline

Aggiungeremo ora l'unità di controllo dell'unità di elaborazione con pipeline analogamente a quanto avevamo fatto nel paragrafo 4.4 per l'unità di elaborazione a singolo ciclo. Inizieremo con una versione semplificata, in cui non si considerano gli elementi più complessi dell'architettura.

Il primo passo consiste nell'etichettare i segnali di controllo dell'unità di elaborazione esistente, come mostrato in Figura 4.44. Cercheremo di riutilizzare il più possibile l'unità di controllo della CPU a singolo ciclo di Figura 4.17; in particolare, utilizzeremo la stessa unità di controllo della ALU, la stessa logica di controllo dei salti e le stesse linee di controllo. Questi componenti sono descritti nelle Figure 4.12, 4.16 e 4.18 e gli elementi chiave sono riprodotti nelle Figure 4.45-4.47.

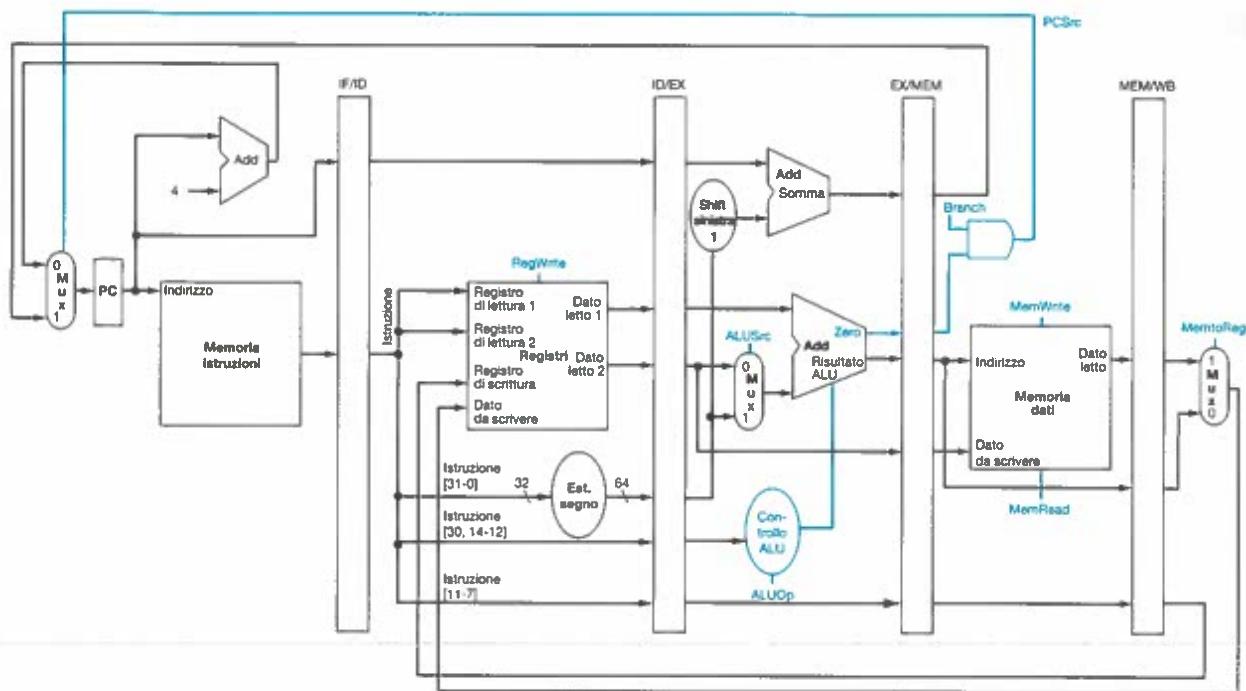
Come nell'implementazione a singolo ciclo, supporremo che il registro PC venga scritto a ogni ciclo di clock, per cui non ci sarà nessun segnale di controllo esplicito per questa operazione. Per lo stesso motivo non ci saranno segnali di scrittura esplicativi per i registri di pipeline (IF/ID, ID/EX, EX/MEM e MEM/WB), dato che anch'essi vengono scritti a ogni ciclo di clock.

Per definire l'unità di controllo della pipeline è sufficiente impostare il valore corretto dei segnali di controllo in ciascuno stadio. Dato che ogni segnale è associato a un componente e che ciascun componente è attivo in uno solo degli stadi della pipeline, si possono suddividere i segnali di controllo in cinque gruppi in base alla fase in cui i segnali sono attivi: vediamoli.

1. *Fetch dell'istruzione.* I segnali di controllo per leggere la memoria istruzioni e per scrivere il PC sono sempre asseriti, quindi non ci sono unità da controllare in questo stadio.
2. *Decodifica dell'istruzione/lettura del register file.* Come nello stadio precedente, a ogni ciclo di clock si verificano sempre le stesse azioni, quindi non ci sono segnali di controllo esplicativi associati a questo stadio.

*Nel Calcolatore 6600, forse più che in ogni altro calcolatore precedente, era l'unità di controllo a fare la differenza.*

*James Thornton, Design of a Computer: The Control Data 6600, 1970*



**Figura 4.44** L'unità di elaborazione con pipeline di Figura 4.39 completa dei segnali di controllo. L'unità di elaborazione adotta la logica di controllo descritta nel paragrafo 4.4 per stabilire l'input al PC, il numero del registro destinazione e i segnali di controllo della ALU. Si noti che ora servono i campi funzionali dell'istruzione nello stadio EX come input del controllore della ALU; occorre quindi salvare anche questi bit nel registro di pipeline ID/EX.

3. *Esecuzione/calcolo dell'indirizzo.* I segnali da impostare sono ALUOp e ALUSrc (Figure 4.45 e 4.46); essi selezionano il registro sul quale scrivere il risultato, l'operazione che deve essere eseguita dalla ALU e scelgono se inviare alla ALU "Dato letto 2" oppure il campo immediato con segno, esteso a 64 bit.
4. *Accesso alla memoria.* Le linee di controllo da impostare per questo stadio sono Branch, MemRead e MemWrite. Questi segnali devono essere asseriti rispettivamente per le istruzioni branch if equal, load e store. Si ricordi che PCSrc (Figura 4.46) seleziona di default l'indirizzo dell'istruzione seguente, a meno che il segnale di Branch non sia asserito e il risultato della ALU valga 0.
5. *Scrittura del risultato.* Le due linee di controllo sono MemtoReg, che sceglie se inviare al register file l'uscita della ALU oppure il valore proveniente dalla memoria, e RegWrite, che comanda la scrittura nel register file del dato selezionato.

Codice operativo istruzione	ALUOp	Operazione eseguita dall'istruzione	Campo funz7	Campo funz3	Operazione dell'ALU	Ingresso di controllo alla ALU
ld	00	load di 1 parola doppia	XXXXXX	XXX	somma	0010
sd	00	store di 1 parola doppia	XXXXXX	XXX	somma	0010
beq	01	salto condizionato all'uguaglianza	XXXXXX	XXX	sottrazione	0110
Tipo R	10	add	0000000	000	somma	0010
Tipo R	10	sub	0100000	000	sottrazione	0110
Tipo R	10	and	0000000	111	AND	0000
Tipo R	10	or	0000000	110	OR	0001

**Figura 4.45** Copia di Figura 4.12. La figura mostra come determinare il valore dei bit di controllo della ALU a partire dai bit di controllo di ALUOp e, per le istruzioni di tipo R, a partire anche dal contenuto del campo funzione.

Nome del segnale	Effetto quando non asserito	Effetto quando asserito
RegWrite	Nullo	Il dato viene scritto nel register file nel registro individuato dal numero del registro di scrittura
ALUSrc	Il secondo operando della ALU proviene dalla seconda uscita del register file (Dato letto 2)	Il secondo operando della ALU proviene dall'estensione del segno dei 12 bit del campo immediato dell'istruzione
PCSrc	Nel PC viene scritta l'uscita del sommatore che calcola il valore di PC + 4	Nel PC viene scritta l'uscita del sommatore che calcola l'indirizzo di salto
MemRead	Nullo	Il dato della memoria nella posizione puntata dall'indirizzo viene inviato in uscita sulla linea "Dato letto"
MemWrite	Nullo	Il contenuto della memoria nella posizione puntata dall'indirizzo viene sostituito con il dato presente sulla linea "Dato scritto"
MemtoReg	Il dato inviato al register file per la scrittura proviene dalla ALU	Il dato inviato al register file per la scrittura proviene dalla Memoria Dati

**Figura 4.46 Copia di Figura 4.16.** La tabella riporta il funzionamento dei sette segnali di controllo. Quelli della ALU (ALUOp) sono stati definiti nella seconda colonna di Figura 4.45. Quando il segnale di controllo a 1 bit di un multiplexer a due vie viene asserito, il multiplexer seleziona l'ingresso etichettato con 1; in caso contrario, cioè se il segnale di controllo è non asserito, il multiplexer seleziona l'ingresso 0. Si noti che in Figura 4.44 PCSrc viene controllato dall'uscita di una porta AND. Se il segnale di Branch e il segnale di Zero in uscita dalla ALU assumono entrambi il valore 1, PCSrc viene impostato a 1, altrimenti viene impostato a 0. La logica di controllo impone il segnale di Branch a 1 solamente in corrispondenza delle istruzioni beq, in tutti gli altri casi, quindi, PCSrc viene impostato a 0.

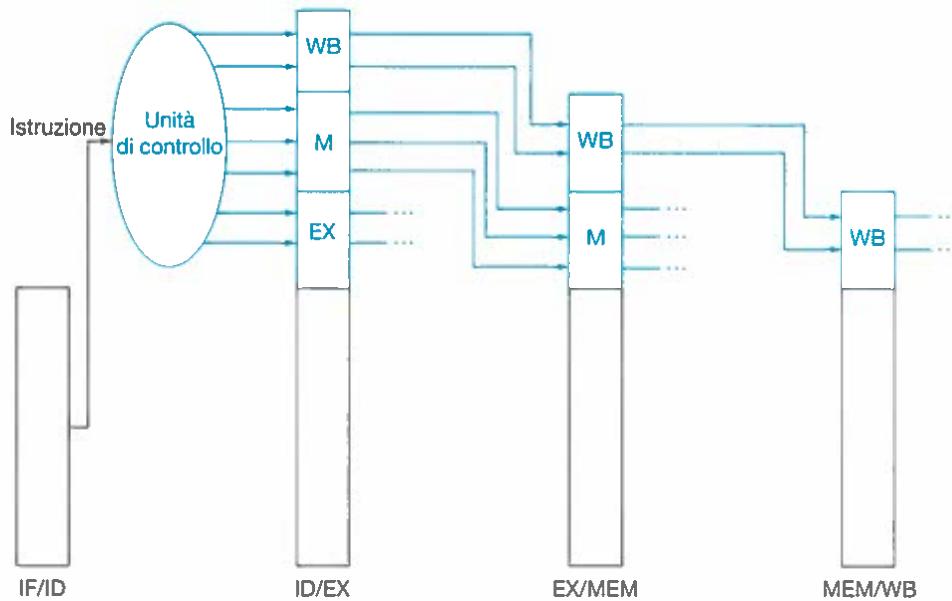
Istruzione	Segnali di controllo dello stadio di esecuzione/calcolo dell'indirizzo		Segnali di controllo dello stadio di accesso alla memoria dati			Segnali di controllo dello stadio di scrittura	
	ALUOp	ALUSrc	Branch	MemRead	MemWrite	RegWrite	MemtoReg
Tipo R	10	0	0	0	0	1	0
ld	00	1	0	1	0	1	1
sd	00	1	0	0	1	0	X
beq	01	0	1	0	0	0	X

**Figura 4.47** I segnali di controllo sono gli stessi di Figura 4.18, ma qui sono suddivisi in tre gruppi, ciascuno attivo in uno degli ultimi tre stadi della pipeline.

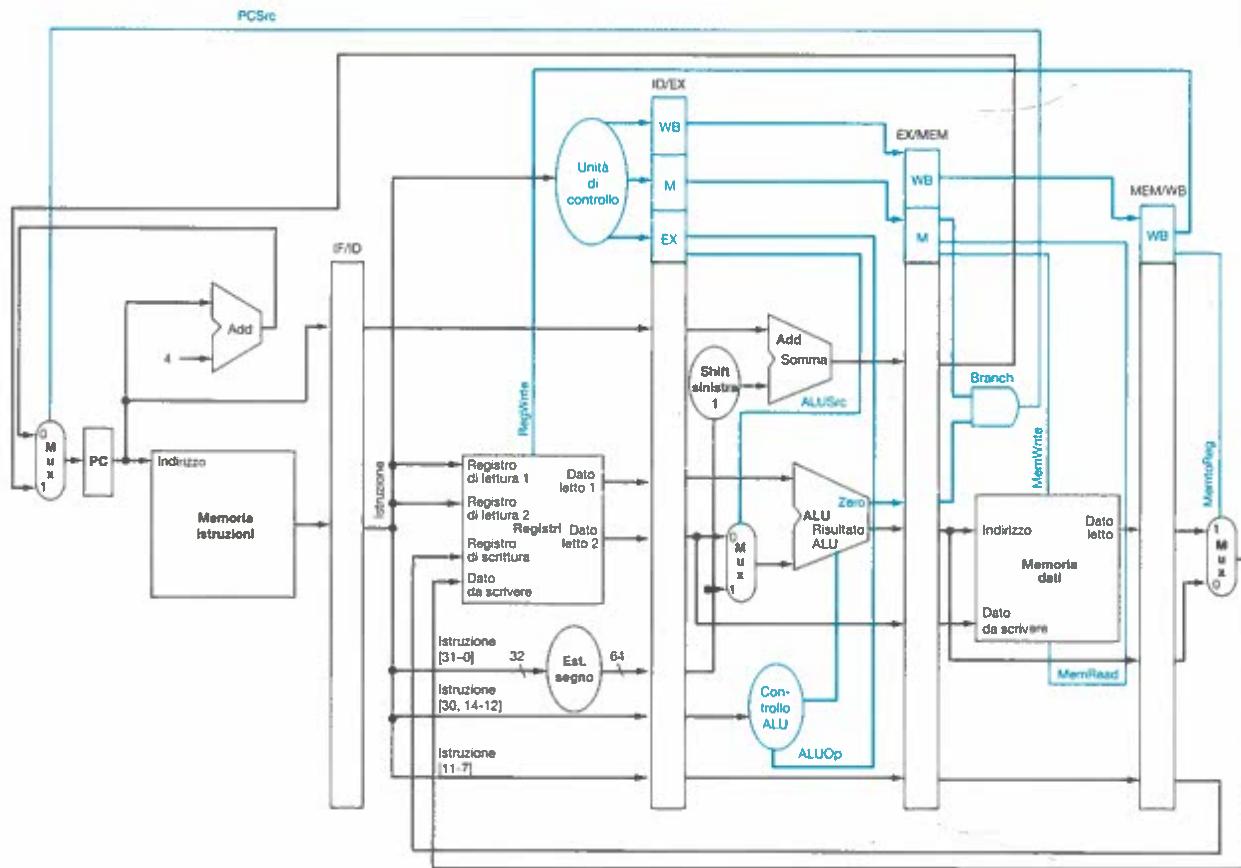
Dato che la versione con pipeline dell'unità di elaborazione non modifica il significato dei segnali di controllo, si possono utilizzare gli stessi nomi adottati in precedenza. La Figura 4.47 riporta gli stessi segnali di controllo introdotti nel paragrafo 4.4, con la differenza che ora le sette linee di controllo sono raggruppate secondo lo stadio della pipeline in cui sono attive.

Implementare l'unità di controllo significa impostare correttamente il valore dei sette segnali di controllo per ciascuno stadio della pipeline e per ciascuna istruzione. Il modo più semplice per farlo è estendere i registri di pipeline per poter salvare anche questi segnali.

Dato che i segnali di controllo esplicativi vengono utilizzati a partire dallo stadio EX, essi possono essere determinati durante la decodifica dell'istruzione per poi venire utilizzati negli stadi successivi. Il modo più semplice per propagare questi segnali di controllo è estendere i registri di pipeline per includere l'informazione sul controllo. La Figura 4.48 mostra che questi segnali di controllo vengono poi utilizzati nello stadio appropriato della pipeline man mano che l'istruzione scorre lungo la pipeline stessa, analogamente a quanto accade per il registro destinazione della load (Figura 4.39). La Figura 4.49 mostra l'unità di elaborazione completa dei registri di pipeline estesi e dei segnali di controllo associati allo stadio opportuno. Per i lettori interessati ad approfondire l'ar-



**Figura 4.48** I sette segnali di controllo per gli ultimi tre stadi. Si noti che due dei sette segnali di controllo sono utilizzati nello stadio EX, mentre i rimanenti cinque vengono trasportati al registro di pipeline EX/MEM, esteso in modo da poterli contenere. Di questi cinque segnali di controllo, tre vengono impiegati nello stadio MEM, mentre gli altri due vengono trasportati al registro MEM/WB per essere poi utilizzati nello stadio WB.



**Figura 4.49** L'unità di elaborazione con pipeline di Figura 4.44 completa dei segnali di controllo uscenti dalla porzione dei registri di pipeline dedicata a memorizzare le informazioni sul controllo. I valori dei segnali di controllo degli ultimi tre stadi sono generati durante lo stadio di decodifica dell'istruzione e vengono salvati nel registro di pipeline ID/EX. In ciascuno degli stadi successivi sono impiegati alcuni di questi segnali di controllo, mentre i rimanenti vengono inviati allo stadio successivo della pipeline.

gomento, nel paragrafo 4.13 sono riportati altri diagrammi a singolo ciclo che illustrano l'esecuzione di frammenti di programmi RISC-V su hardware dotato di pipeline.

## 4.7 Hazard sui dati: propagazione o stallo

Gli esempi visti nel paragrafo precedente mostrano le potenzialità dell'esecuzione in pipeline e illustrano come la pipeline possa essere implementata in hardware. A questo punto, però, è necessario abbandonare la situazione ideale e rendersi conto di quello che può accadere durante l'esecuzione di un programma reale. Le istruzioni RISC-V nelle Figure 4.41-4.43 erano indipendenti: nessuna di loro utilizzava i risultati prodotti dalle istruzioni precedenti. Tuttavia abbiamo già visto nel paragrafo 4.5 che gli hazard sui dati costituiscono un ostacolo all'esecuzione in pipeline.

Esaminiamo una sequenza di istruzioni contenenti molte dipendenze, evidenziate in blu:

```
sub x2, x1, x3 // sub scrive nel registro x2
and x12, x2, x5 // Il primo operando (x2) dipende da sub
or x13, x6, x2 // Il secondo operando (x2) dipende da sub
add x14, x2, x2 // Il primo (x2) e il secondo (x2)
                  // operando dipendono da sub
sw x15, 100(x2) // Il registro base (x2) dipende da sub
```

Le ultime quattro istruzioni dipendono tutte dal risultato della prima istruzione, scritto nel registro x2. Supponiamo che il registro x2 contenga il valore 10 prima dell'istruzione di sottrazione e -20 dopo: il programmatore vorrebbe che fosse utilizzato il valore -20 in tutte le istruzioni che seguono la sub e che fanno riferimento al registro x2.

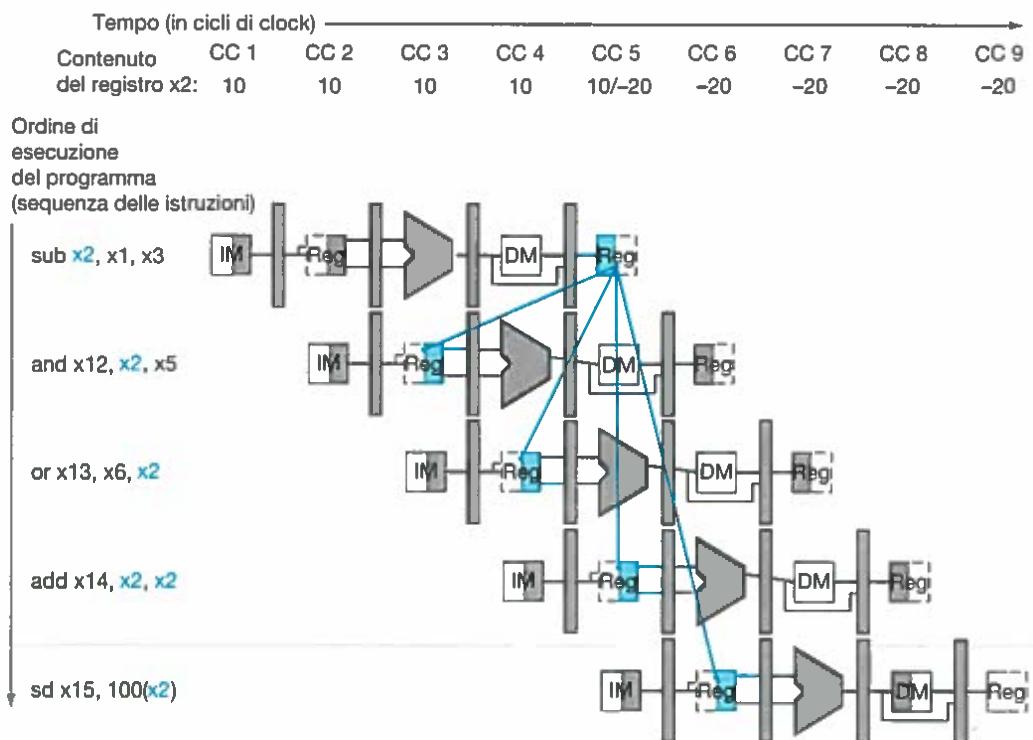
Ma come si comporterebbe la sequenza di istruzioni se venisse eseguita sulla pipeline vista in precedenza? La Figura 4.50 mostra l'esecuzione di queste istruzioni utilizzando un diagramma a più cicli di clock. Per illustrare ciò che accade all'interno della pipeline, la parte in alto di Figura 4.50 riporta il contenuto del registro x2, che cambia solo a metà del quinto ciclo di clock, quando l'istruzione sub scrive il risultato della sottrazione nel register file.

L'ultimo hazard potenziale può essere risolto progettando opportunamente l'hardware che implementa il register file: che cosa succede quando un registro viene letto e scritto nello stesso ciclo di clock? Si può assumere che la scrittura avvenga nella prima metà del ciclo di clock e la lettura nella seconda metà: in questo modo si può leggere nello stesso ciclo di clock il valore che viene scritto. Questa situazione è tipica di molte implementazioni hardware del register file, e in questo caso non si hanno hazard sui dati. La Figura 4.50 mostra che il dato letto dal registro x2 *non* può essere quello prodotto dall'istruzione sub, a meno che l'operazione di lettura non avvenga durante il quinto ciclo di clock o nei cicli successivi, per cui le sole istruzioni che leggerebbero il valore corretto (-20) sarebbero add e sd; and e or leggerebbero invece il valore 10, che è sbagliato! Osservando il diagramma, possiamo affermare che questo tipo di problema si verifica quando c'è una linea di dipendenza che va a ritroso nel tempo.

Come si è detto nel paragrafo 4.5, il risultato dell'istruzione sub è già disponibile al termine dello stadio EX, cioè nel ciclo di clock 3. Ma quando le istruzioni and e or hanno realmente bisogno di questo dato? All'inizio dello stadio EX, cioè nei cicli di clock 4 e 5, rispettivamente. Di conseguenza, è possibile eseguire questo frammento di codice senza stalli semplicemente *propagando* il dato non appena è disponibile a qualsiasi unità che lo richieda, ancor prima che esso sia disponibile per essere letto dal register file.

*Cosa vuoi dire, perché dovrebbe essere costruito?  
È una scorciatoia. Tu devi costruire scorciatoie.*

Douglas Adams,  
*The Hitchhiker's Guide to the Galaxy*, 1979



**Figura 4.50** Dipendenze in una pipeline che contiene una sequenza di cinque istruzioni: viene utilizzata l'unità di elaborazione semplificata vista finora e le dipendenze sono evidenziate. Tutte le funzioni dipendenti sono disegnate in blu; nella parte alta "CC 1" indica il primo ciclo di clock. La prima istruzione scrive il risultato in  $x_2$  e tutte le istruzioni successive devono leggerne il contenuto. L'operazione di scrittura di tale registro avviene solo nel quinto ciclo di clock, quindi il valore corretto non è disponibile prima di questo ciclo. Si ricordi che la lettura di un registro in un certo ciclo di clock restituisce il valore scritto nel registro nella prima metà del ciclo di clock stesso. Le linee blu che partono dalla prima unità di elaborazione e sono dirette verso le unità di elaborazione disegnate sotto mostrano le dipendenze: le linee che vanno a ritroso nel tempo costituiscono degli hazard sui dati.

Come funziona la propagazione? Per semplicità, nel seguito di questo paragrafo considereremo solo la propagazione di un dato per un'operazione che deve essere eseguita nello stadio EX, la quale può essere una generica operazione con la ALU oppure il calcolo di un indirizzo. Questo significa che quando un'istruzione nel suo stadio EX deve utilizzare un dato di un registro, che però deve ancora essere scritto da un'istruzione precedente nella sua fase di WB, occorre che il dato venga portato all'ingresso della ALU.

Una notazione che assegna dei nomi ai campi dei registri di pipeline permette di formalizzare le dipendenze con maggior precisione. Per esempio, "ID/EX.RegistroRs1" si riferisce al numero di uno dei registri che si possono trovare nel registro di pipeline ID/EX, cioè il registro associato alla prima porta di lettura del register file. La prima parte del nome, a sinistra del punto, è il nome del registro di pipeline, mentre la seconda parte corrisponde al nome di un campo all'interno del registro. Utilizzando questa notazione, le due coppie di condizioni che generano un hazard sui dati si possono esprimere come segue:

- 1a. EX/MEM.RegistroRd = ID/EX.RegistroRs1
- 1b. EX/MEM.RegistroRd = ID/EX.RegistroRs2
- 2a. MEM/WB.RegistroRd = ID/EX.RegistroRs1
- 2b. MEM/WB.RegistroRd = ID/EX.RegistroRs2

Il primo hazard nella sequenza di istruzioni dell'esempio che stiamo considerando è sul registro  $x_2$  e si verifica tra la scrittura del risultato dell'istruzione `sub x2, x1, x3` e l'operazione di lettura del primo operando dell'istruzione

and x12, x2, x5. Tale hazard può essere rilevato quando l'istruzione and si trova nello stadio EX e l'istruzione precedente si trova nello stadio MEM; si tratta quindi di un hazard di tipo 1a:

$$\text{EX/MEM.RegistroRd} = \text{ID/EX.RegistroRs1} = x2$$

## Rilevamento delle dipendenze

Classificare le dipendenze all'interno della seguente sequenza di istruzioni (identica a quella riportata all'inizio del paragrafo):

```
sub x2, x1, x3 // sub scrive nel registro x2
and x12, x2, x5 // Il primo operando (x2)
                  // dipende da sub
or   x13, x6, x2 // Il secondo operando (x2)
                  // dipende da sub
add  x14, x2, x2 // Il primo (x2) e il secondo (x2)
                  // operando dipendono da sub
sd   x15, 100(x2) // Il registro base (x2) dipende da sub
```

Come già detto, l'hazard tra sub e and è di tipo 1a. Gli altri hazard sono classificabili nel seguente modo:

- tra sub e or c'è un hazard di tipo 2b:

$$\text{MEM/WB.RegistroRd} = \text{ID/EX.RegistroRs2} = x2$$

- le due dipendenze tra sub e add non costituiscono hazard, dal momento che il register file fornisce i dati corretti già durante lo stadio ID dell'istruzione add;
- non c'è hazard sui dati tra sub e sd, poiché sd legge x2 nel ciclo di clock successivo a quello in cui sub scrive in x2.

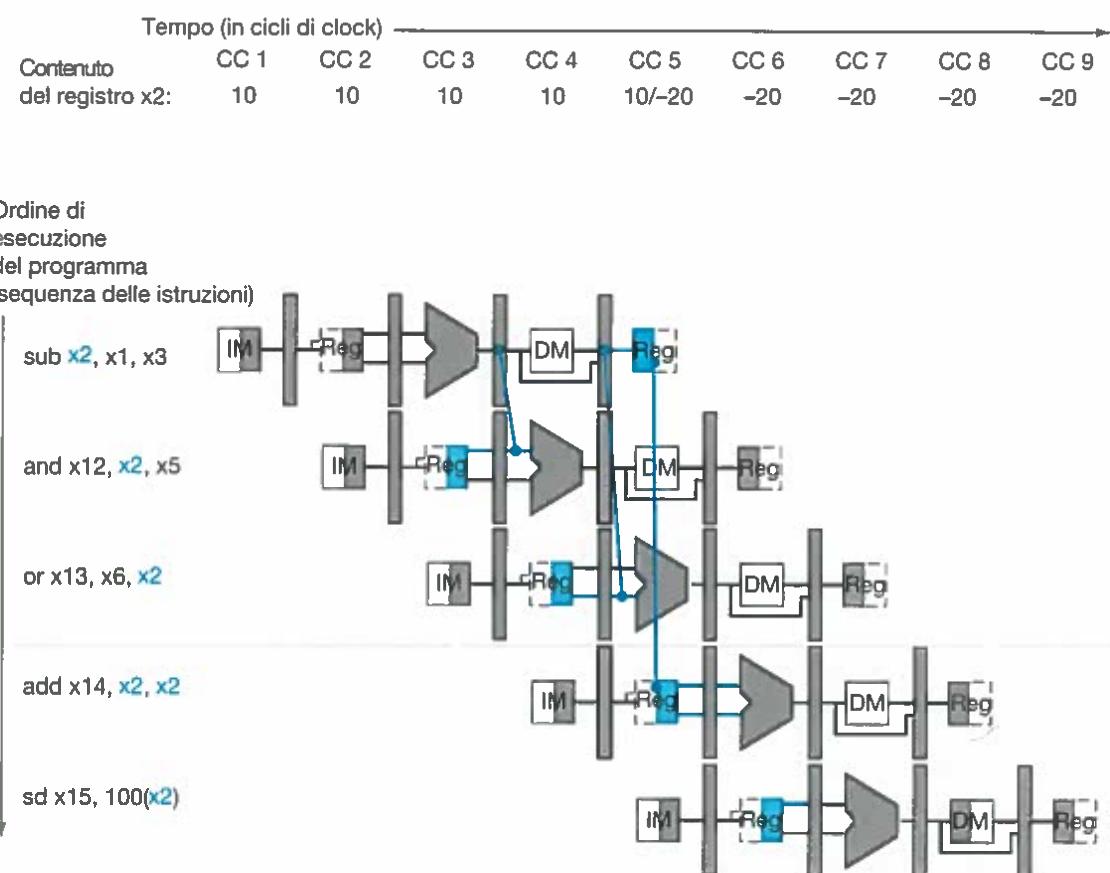
## ESEMPIO

## SOLUZIONE

Dato che non tutte le istruzioni scrivono il risultato nel register file, questa strategia non è precisa: in alcuni casi si propagherebbero dei dati anche quando non si deve. Una possibile soluzione consiste nel verificare se il segnale RegWrite è attivo; a tal fine, si può controllare se sia asserito il segnale RegWrite contenuto nella porzione dei registri di pipeline EX/MEM e MEM/WB riservata ai segnali di controllo attivi nello stadio di WB. Ricordiamo, inoltre, che il RISC-V richiede che il registro x0 contenga sempre il valore 0 quando viene utilizzato come operando. Nel caso in cui un'istruzione nella pipeline abbia x0 come registro destinazione (per es. addi x0, x1, 2), si cerca di evitare che il risultato dell'operazione, eventualmente non nullo, sia propagato; in questo modo si evita che il programmatore assembler e il compilatore non possano considerare il registro x0 come registro destinazione. Le condizioni riportate sopra funzionano quindi correttamente se si aggiunge la condizione EX/MEM.RegistroRd ≠ 0 alla prima condizione di hazard e MEM/WB.RegistroRd ≠ 0 alla seconda.

Ora che siamo in grado di rilevare gli hazard abbiamo risolto metà del problema, ma dobbiamo ancora propagare i dati corretti.

La Figura 4.51 mostra le dipendenze tra i registri di pipeline e gli ingressi della ALU per la stessa sequenza di istruzioni di Figura 4.50. La differenza consiste nel fatto che le linee di dipendenza partono dai *registri di pipeline*, anziché dallo stadio WB nel quale viene scritto il register file. Queste linee di dipendenza evidenziano che i dati necessari, che devono essere propagati, sono



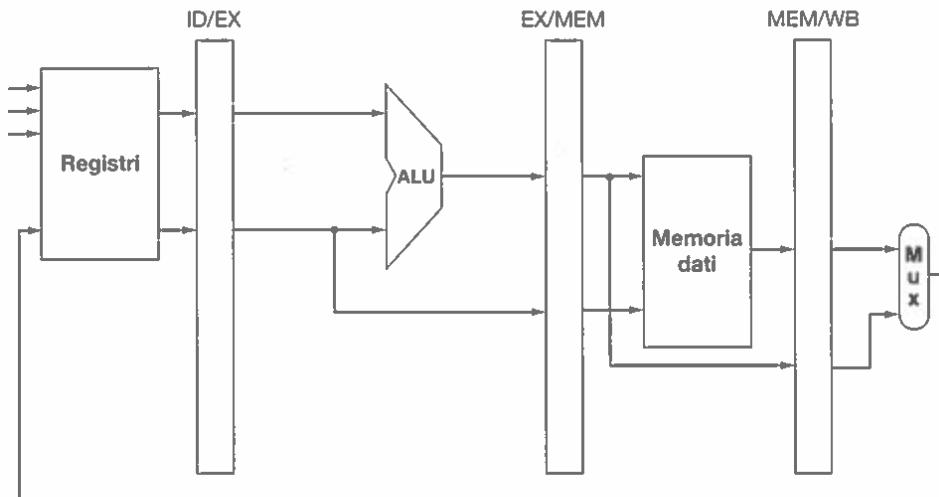
**Figura 4.51** Le dipendenze tra i registri di pipeline si muovono in avanti nel tempo; è così possibile fornire alla ALU gli ingressi corretti per l'esecuzione dell'istruzione `and` e dell'istruzione `or`, propagando il risultato della `sub` presente nei registri di pipeline. Il contenuto del registo di pipeline mostra che il dato desiderato è disponibile prima che esso venga scritto nel register file. Si suppone che il register file propaghi i valori scritti al suo interno nello stesso ciclo di clock, quindi la `add` non deve essere messa in stallo perché riceve regolarmente i dati dal register file e non serve la propagazione dai registri di pipeline. La propagazione all'interno del register file avviene perché è possibile leggere il dato che viene scritto nello stesso ciclo di clock; questo è il motivo per cui nel quinto ciclo di clock il registro `x2` contiene il valore 10 all'inizio e -20 alla fine.

contenuti nei registri di pipeline e sono quindi disponibili in tempo utile per le istruzioni successive.

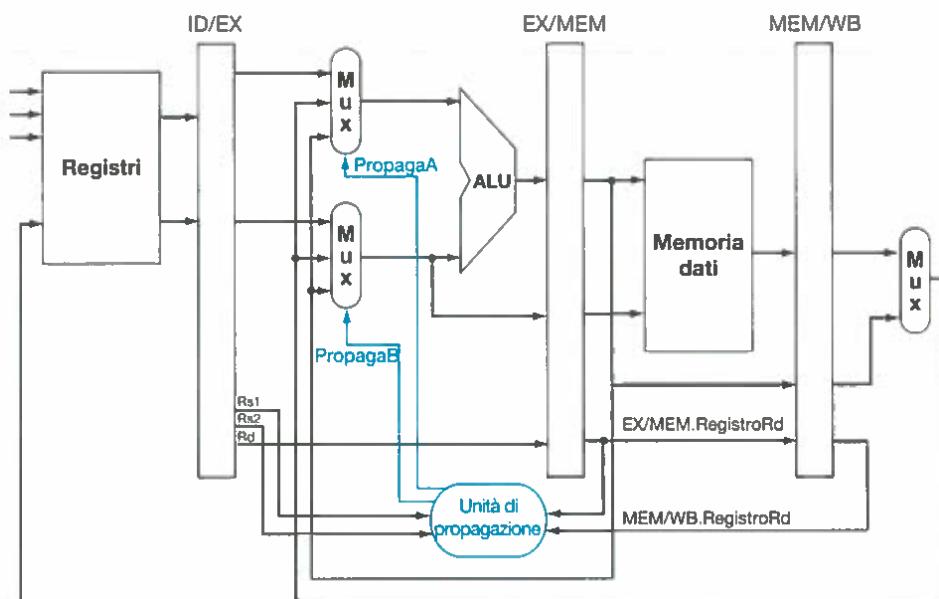
Se possiamo prendere gli input della ALU non solo dal registro ID/EX ma anche dagli *altri* registri di pipeline, allora possiamo propagare i dati corretti. Aggiungendo dei multiplexer sugli ingressi della ALU e utilizzando dei segnali di controllo adeguati, è possibile far funzionare la pipeline alla velocità massima anche in presenza delle dipendenze tra i dati che abbiamo evidenziato.

Faremo inizialmente l'assunzione che la propagazione sia necessaria solamente per le quattro istruzioni del formato R: add, sub, and e or. La Figura 4.52 mostra un dettaglio della ALU e dei registri di pipeline prima e dopo l'aggiunta della propagazione, mentre la Figura 4.53 mostra come le linee di controllo dei multiplexer della ALU selezionino il dato sorgente, prendendolo dal register file oppure da uno dei dati propagati.

Dato che i multiplexer di propagazione all'ingresso della ALU si trovano nello stadio EX, il controllo della propagazione avviene in questo stadio. Di conseguenza, per determinare se propagare o meno i dati è necessario trasportare il numero dei due registri utilizzati come operandi dallo stadio ID al registro di pipeline ID/EX. Prima di considerare la propagazione, non c'era alcun motivo di trasportare anche il numero dei registri Rs1 e Rs2 scrivendoli nel registro ID/EX, mentre ora è diventato necessario per poter verificare se c'è un hazard sui dati: occorre quindi aggiungere questi campi al registro ID/EX.



a. Nessuna propagazione



b. Con propagazione

**Figura 4.52** Nella parte superiore della figura sono mostrate la ALU e i registri di pipeline prima di aggiungere l'hardware che esegue la propagazione. Nella parte inferiore della figura sono stati inseriti dei multiplexer per la selezione dei cammini di propagazione e viene mostrata l'unità di propagazione. L'hardware aggiunto è evidenziato in blu. La figura è semplificata, in quanto non contiene alcuni dettagli dell'unità di elaborazione completa, come il circuito per l'estensione del segnale.

Controllo multiplexer	Sorgente	Spiegazione
PropagaA = 00	ID/EX	Il primo operando della ALU proviene dal register file
PropagaA = 10	EX/MEM	Il primo operando della ALU viene propagato dal risultato della ALU nel ciclo di clock precedente
PropagaA = 01	MEM/WB	Il primo operando della ALU viene propagato dalla memoria dati o da un precedente risultato della ALU
PropagaB = 00	ID/EX	Il secondo operando della ALU proviene dal register file
PropagaB = 10	EX/MEM	Il secondo operando della ALU viene propagato dal risultato della ALU nel ciclo di clock precedente
PropagaB = 01	MEM/WB	Il secondo operando della ALU è propagato dalla memoria dati o da un precedente risultato della ALU

**Figura 4.53** Il valore dei segnali di controllo dei multiplexer di propagazione di Figura 4.52. Il campo immediato dotato di segno che costituisce un ulteriore ingresso della ALU viene descritto nella sezione *Approfondimento* al termine di questo paragrafo.

Si possono ora scrivere sia le condizioni per rilevare gli hazard sia i segnali di controllo per risolvere:

1. *Hazard nello stadio EX:*

```
if (EX/MEM.RegWrite
and (EX/MEM.RegistroRd ≠ 0)
and (EX/MEM.RegistroRd = ID/EX.RegistroRs1)) PropagaA = 10

if (EX/MEM.RegWrite
and (EX/MEM.RegistroRd ≠ 0)
and (EX/MEM.RegistroRd = ID/EX.RegistroRs2)) PropagaB = 10
```

In questo schema di risoluzione degli hazard il risultato che arriva dalla precedente istruzione viene propagato a uno dei due ingressi della ALU; se l'istruzione precedente deve scrivere sul register file e il numero del registro di scrittura coincide con quello del registro da cui viene letto il dato A o B in ingresso alla ALU (purché questo non sia il registro 0), allora il multiplexer viene pilotato in maniera da prelevare il dato dal registro di pipeline EX/MEM.

2. *Hazard nello stadio MEM:*

```
if (MEM/WB.RegWrite
and (MEM/WB.RegistroRd ≠ 0)
and (MEM/WB.RegistroRd = ID/EX.RegistroRs1)) PropagaA = 01

if (MEM/WB.RegWrite
and (MEM/WB.RegistroRd ≠ 0)
and (MEM/WB.RegistroRd = ID/EX.RegistroRs2)) PropagaB = 01
```

Come già accennato in precedenza, non si può verificare un hazard nello stadio WB, poiché si suppone che il register file fornisca in uscita il dato corretto anche quando l'istruzione nello stadio ID deve leggere il contenuto dello stesso registro che viene scritto in quel ciclo di clock dall'istruzione che si trova nello stadio WB. Un register file di questo tipo implementa una particolare forma di propagazione che si verifica al suo interno.

Una possibile fonte di problemi è costituita dagli hazard sui dati che si possono verificare tra il risultato di un'istruzione nello stadio WB, il risultato di un'istruzione nello stadio MEM e l'operando sorgente di un'istruzione nello stadio ALU. Per esempio, se si devono sommare tra loro gli elementi di un vettore e scrivere il risultato in un particolare registro, si ottiene una sequenza di istruzioni che leggono e scrivono nello stesso registro:

```
add x1, x1, x2
add x1, x1, x3
add x1, x1, x4
. . .
```

In questo caso il risultato viene propagato dallo stadio MEM, dal momento che quel risultato è il più recente. Quindi il test da eseguire per rilevare un hazard nello stadio MEM sarà il seguente (le aggiunte sono evidenziate in blu):

```
if (MEM/WB.RegWrite
and (MEM/WB.RegistroRd ≠ 0)
and not (EX/MEM.RegWrite and (EX/MEM.RegistroRd ≠ 0)
        and (EX/MEM.RegistroRd = ID/EX.RegistroRs1)))
and (MEM/WB.RegistroRd = ID/EX.RegistroRs1)) PropagaA = 01

if (MEM/WB.RegWrite
and (MEM/WB.RegistroRd ≠ 0)
```

```

and not(EX/MEM.RegWrite and (EX/MEM.RegistroRd ≠ 0)
       and (EX/MEM.RegistroRd = ID/EX.RegistroRs2))
and (MEM/WB.RegistroRd = ID/EX.RegistroRs2)) Propagab = 01

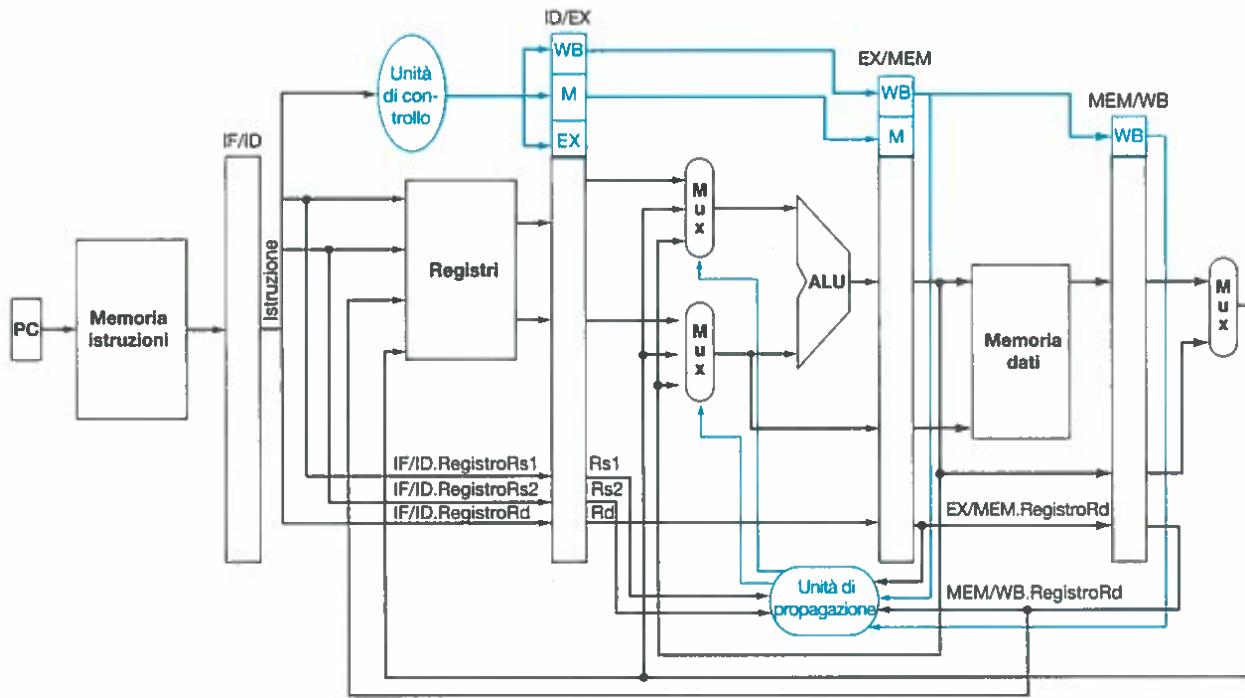
```

La Figura 4.54 mostra l'hardware necessario a implementare la propagazione per le istruzioni che utilizzano nello stadio EX il risultato di istruzioni precedenti. Si noti che il campo EX/MEM.RegistroRd contiene il numero del registro di scrittura sia per le istruzioni che operano sulla ALU sia per le istruzioni di load.

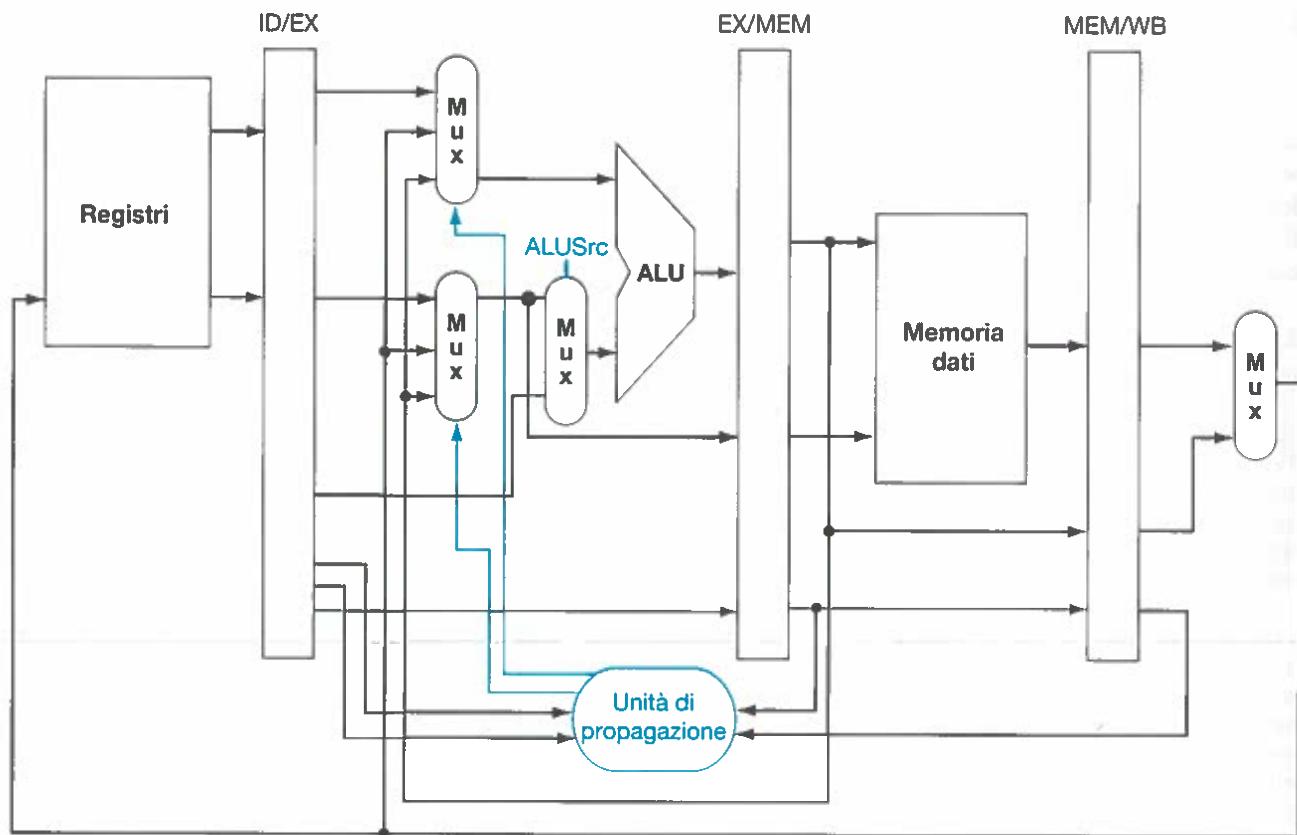
Per i lettori interessati ad altri esempi illustrati mediante diagrammi a singolo ciclo, il paragrafo 4.13 mostra due frammenti di codice RISC-V contenenti hazard che richiedono la propagazione.

**Approfondimento.** La propagazione può essere di aiuto anche per gli hazard che si verificano quando un'istruzione di store è dipendente da altre istruzioni. Poiché le istruzioni di store utilizzano soltanto un dato durante lo stadio MEM, la propagazione è facile. Si consideri, tuttavia, il caso in cui una load sia immediatamente seguita da alcune istruzioni di store (è il caso in cui si vogliono copiare i dati da memoria a memoria nelle architetture RISC-V). Poiché questo tipo di operazione è frequente, è necessario aggiungere altro hardware per la propagazione per far sì che la copia dei dati da memoria a memoria sia eseguita più velocemente. Se in Figura 4.51 sostituissimo le istruzioni sub e and, rispettivamente, con 1d e sd, vedremmo che è possibile evitare lo stallo, poiché il dato che la sd deve scrivere in memoria viene scritto nel registro MEM/WB dall'istruzione 1d precedente, in tempo per essere utilizzato nello stadio MEM dell'istruzione di store. In questo caso, per rendere disponibile il dato alla sd, sarebbe necessario aggiungere la propagazione anche all'interno dello stadio di accesso alla memoria. Lasciamo al lettore questa modifica come esercizio.

Inoltre, nell'unità di elaborazione di Figura 4.54 manca il campo immediato dotato di segno che viene richiesto come secondo ingresso della ALU dalle istruzioni load e store. Dal momento che è l'unità di controllo centrale a decidere se selezionare



**Figura 4.54** L'unità di elaborazione dopo le modifiche apportate per risolvere gli hazard tramite la propagazione. Se si confronta questo schema con l'unità di elaborazione di Figura 4.49, si nota che le aggiunte sono costituite dai multiplexer sugli ingressi della ALU. La figura rappresenta uno schema stilizzato, in cui mancano alcuni dettagli relativi all'implementazione completa che contiene anche la logica dei salti e dell'estensione del segno.



**Figura 4.55** Rappresentazione più dettagliata dell'unità di elaborazione di Figura 4.52; essa contiene un multiplexer 2:1, che è stato aggiunto per potere eventualmente selezionare come secondo operando della ALU il campo immediato dotato di segno.

come secondo operando il contenuto di un registro o il campo immediato (mentre l'unità di propagazione sceglie da quale registro di pipeline prendere il secondo input della ALU quando questo deve essere preso da un registro), la soluzione più semplice consiste nell'aggiungere un multiplexer 2:1 che selezioni il secondo ingresso della ALU tra l'uscita PropagaB del multiplexer e il campo immediato dotato di segno. La Figura 4.55 mostra questa aggiunta.

### Hazard sui dati e stallo

*Se al primo tentativo non hai successo, riformula il concetto di successo.*

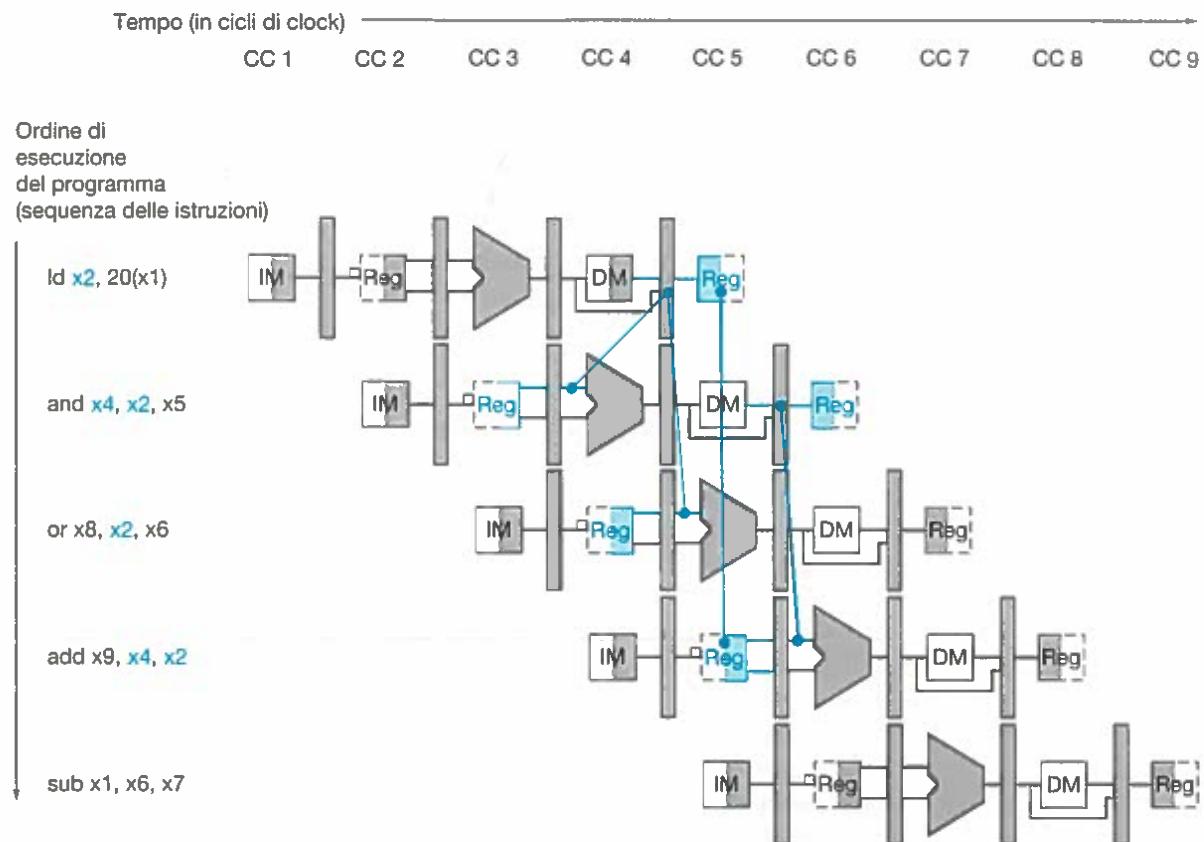
Anonimo

Come si è detto nel paragrafo 4.5, un caso nel quale la propagazione non può risolvere il problema è quello in cui un'istruzione tenta di leggere il registro che verrà scritto da un'istruzione di load che la precede; questa situazione è illustrata in Figura 4.56. Durante il quarto ciclo di clock il dato deve ancora essere letto dalla memoria, ma la ALU si appresta già a eseguire l'operazione prevista per l'istruzione successiva. Occorre quindi mettere in stallo la pipeline quando un'istruzione di load è seguita da un'istruzione che utilizza il dato letto.

È quindi necessaria, oltre all'unità di propagazione, anche un'*unità di rilevamento degli hazard* che, durante lo stadio ID di un'istruzione, possa inserire uno stallo tra la lettura del dato e il suo utilizzo. Per le istruzioni di load, l'unità di rilevamento degli hazard implementa questa semplice condizione logica:

```

if (ID/EX.MemRead and
    ((ID/EX.RegistroRd = IF/ID.RegistroRs1) or
     (ID/EX.RegistroRd = IF/ID.RegistroRs2)))
    metti in stallo la pipeline
  
```



**Figura 4.56 Sequenza di istruzioni nella pipeline.** Poiché la dipendenza tra l'istruzione di load e l'istruzione successiva (and) è a ritroso nel tempo, questo hazard non può essere risolto con la propagazione, ma occorre che l'unità di rilevamento degli hazard generi uno stallo della pipeline.

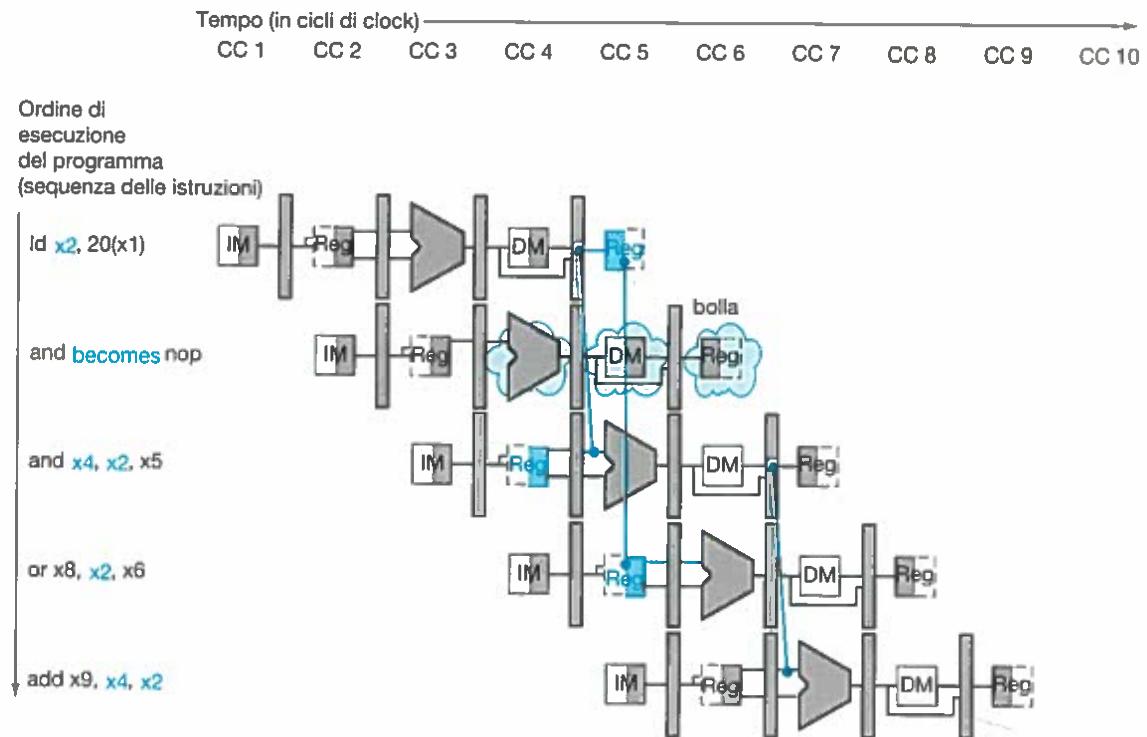
Ricordate che stiamo utilizzando il RegistroRd per indicare il registro specificato nei bit 11:7 sia per le istruzioni di load che di tipo R. La prima condizione stabilisce se l'istruzione nello stadio EX è una load: questa, infatti, è l'unica istruzione che legge dati dalla memoria. Le successive due condizioni controllano se il numero del registro di scrittura della load coincide con uno dei registri sorgente dell'istruzione nello stadio ID. Se la condizione è verificata, l'istruzione nello stadio ID viene messa in stallo per un ciclo di clock. Dopo questo stallo, la dipendenza sui dati può essere gestita dalla logica di propagazione e l'esecuzione può quindi procedere. Se non ci fosse la propagazione le istruzioni di Figura 4.56 richiederebbero lo stallo per un secondo ciclo di clock. Se l'istruzione nello stadio ID viene messa in stallo, anche l'istruzione nello stadio IF deve essere messa in stallo, altrimenti l'istruzione prelevata nella fase di fetch precedente andrebbe persa. Impedire a queste due istruzioni di procedere può essere ottenuto semplicemente impedendo la commutazione del registro PC e del registro di pipeline IF/ID. Se il valore di questi registri viene conservato attraverso il fronte del clock, nel ciclo di clock successivo l'istruzione che viene letta nello stadio IF verrà prelevata dallo stesso indirizzo contenuto nel PC nel ciclo di clock precedente, e il numero dei registri da leggere nello stadio ID, contenuti nei campi del registro IF/ID, sarà lo stesso del ciclo di clock precedente. Riprendendo l'esempio del bucato, è come se si facesse ripartire la lavatrice con gli stessi vestiti e si facesse lavorare l'asciugatrice a vuoto. Naturalmente gli stadi a valle dello stadio ID, a partire dallo stadio EX, continuano a lavorare, ma come nel caso dell'asciugatrice occorre che eseguano istruzioni che non producono effetti: queste istruzioni sono chiamate **nop** (*not operation*).

**Nop:** letteralmente "non un'operazione", è un'istruzione che non esegue alcuna operazione che modifichi lo stato della macchina.

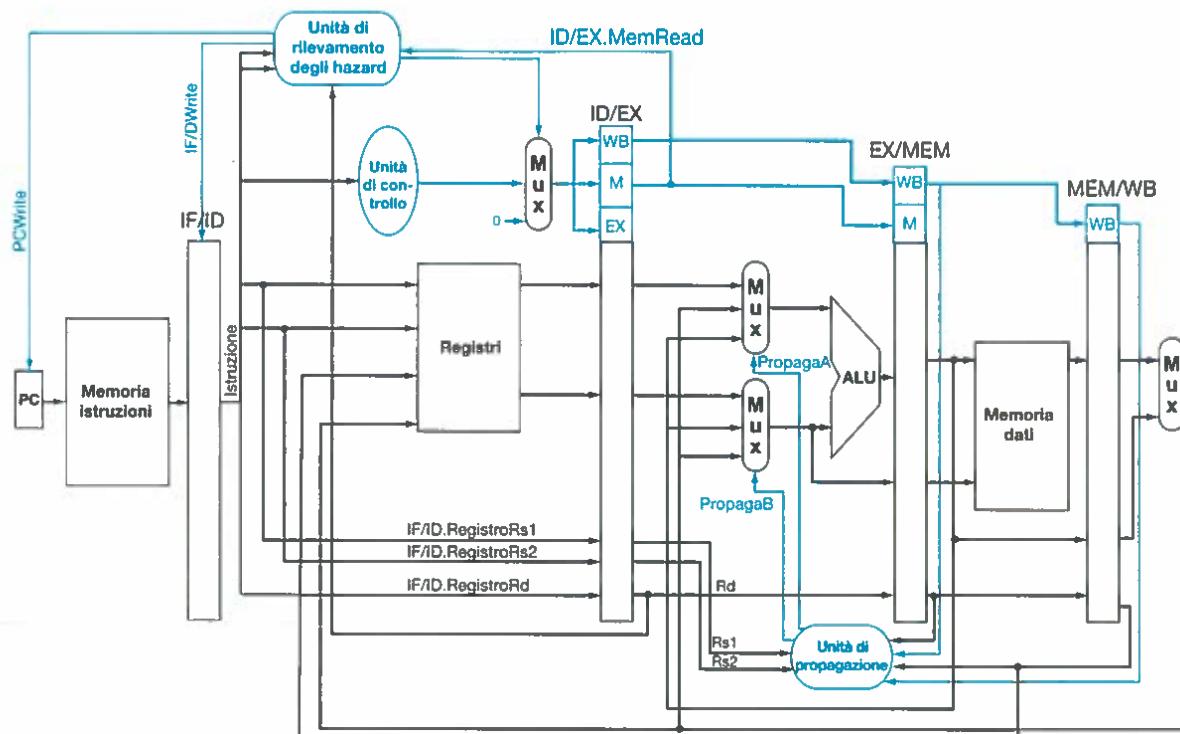
Come possiamo inserire una nop, che si comporta come una "bolla d'aria", all'interno della pipeline? In Figura 4.47 abbiamo visto che impostando a 0 tutti e sette i segnali di controllo degli stadi EX, MEM e WB è possibile creare un'istruzione che non fa nulla, cioè una nop. Se identifichiamo l'hazard nello stadio ID, possiamo inserire nella pipeline una bolla mettendo a 0 i campi EX, MEM e WB del registro di pipeline ID/EX, i quali contengono i segnali di controllo; questi segnali di controllo per così dire sterilizzati vengono poi trasportati in avanti nella pipeline a ogni colpo di clock producendo l'effetto desiderato: se i segnali di controllo sono tutti a 0 non vengono scritti né i registri né la memoria dati.

La Figura 4.57 mostra ciò che avviene effettivamente nell'hardware: l'esecuzione dell'istruzione and nello stadio EX della pipeline è trasformata nell'esecuzione di una nop e l'esecuzione di tutte le istruzioni successive a quella di load, and compresa, viene ritardata di un ciclo di clock. Come una bolla d'aria in un tubo pieno d'acqua, la bolla associata allo stallo ritarda tutto ciò che la segue e procede stadio dopo stadio fino a giungere alla fine della pipeline. Nell'esempio riportato, l'hazard costringe le istruzioni and e or a ripetere nel quarto ciclo di clock ciò che avevano fatto nel terzo ciclo: l'istruzione and legge il contenuto dei registri e viene decodificata, mentre l'istruzione or viene prelevata una seconda volta dalla memoria istruzioni. Uno stallo implica quindi la ripetizione di una certa quantità di lavoro, anche se il suo effetto è di allungare il tempo di esecuzione delle istruzioni and e or e di ritardare il caricamento dell'istruzione add.

La Figura 4.58 mette in evidenza le connessioni dell'unità di rilevamento degli hazard e dell'unità di propagazione all'interno della pipeline. Come in precedenza, l'unità di propagazione controlla i multiplexer all'ingresso della



**Figura 4.57** Modalità con cui gli stalli vengono realmente inseriti in una pipeline. Una "bolla" viene inserita nel quarto ciclo di clock trasformando l'istruzione and in una nop. Si noti che l'istruzione and viene prelevata e poi decodificata rispettivamente nel secondo e terzo ciclo di clock, ma il suo stadio EX viene ritardato fino al quinto ciclo (in assenza di stallo, lo stadio EX verrebbe invece eseguito nel quarto ciclo di clock). Allo stesso modo, l'istruzione or viene prelevata nel terzo ciclo di clock, ma il suo stadio ID viene ritardato fino al quinto ciclo di clock (in assenza di stallo lo stadio ID verrebbe eseguito nel quarto ciclo di clock). Dopo l'inserimento della bolla, tutti i dati che generano le dipendenze si spostano in avanti nel tempo e non si verifica più alcun hazard.



**Figura 4.58 Panoramica sul controllo della pipeline: nella figura compaiono i due multiplexer per la propagazione, l'unità di rilevamento degli hazard e l'unità di propagazione.** Sebbene gli stadi ID ed EX siano stati semplificati (mancano la logica per l'estensione del segno del campo immediato e la logica dei salti), lo schema mostra gli elementi essenziali dell'hardware per la propagazione.

ALU per sostituire eventualmente il valore proveniente da uno dei registri del register file con il contenuto dell'opportuno registro di pipeline. L'unità di rilevamento degli hazard controlla la scrittura del PC e del registro IF/ID oltre al multiplexer che sceglie se scrivere nel registro ID/EX 0 oppure i segnali di controllo dell'istruzione in fase di decodifica. L'unità di rilevamento degli hazard mette in stallo la pipeline e imposta a 0 i segnali di controllo se si verifica la condizione di hazard per l'utilizzo del dato di una load riportata in precedenza. Per i lettori interessati ad approfondire l'argomento, nel paragrafo 4.13 viene riportato un esempio di codice RISC-V contenente hazard che provocano lo stallo della pipeline.

## QUADRO D'INSIEME

Sebbene il compilatore generalmente si basi sull'hardware per risolvere gli hazard e per assicurare l'esecuzione corretta del codice, esso deve comprendere il funzionamento della pipeline per ottenere prestazioni ottimali; altrimenti potrebbero verificarsi stalli inaspettati che riducono le prestazioni del codice compilato. ■

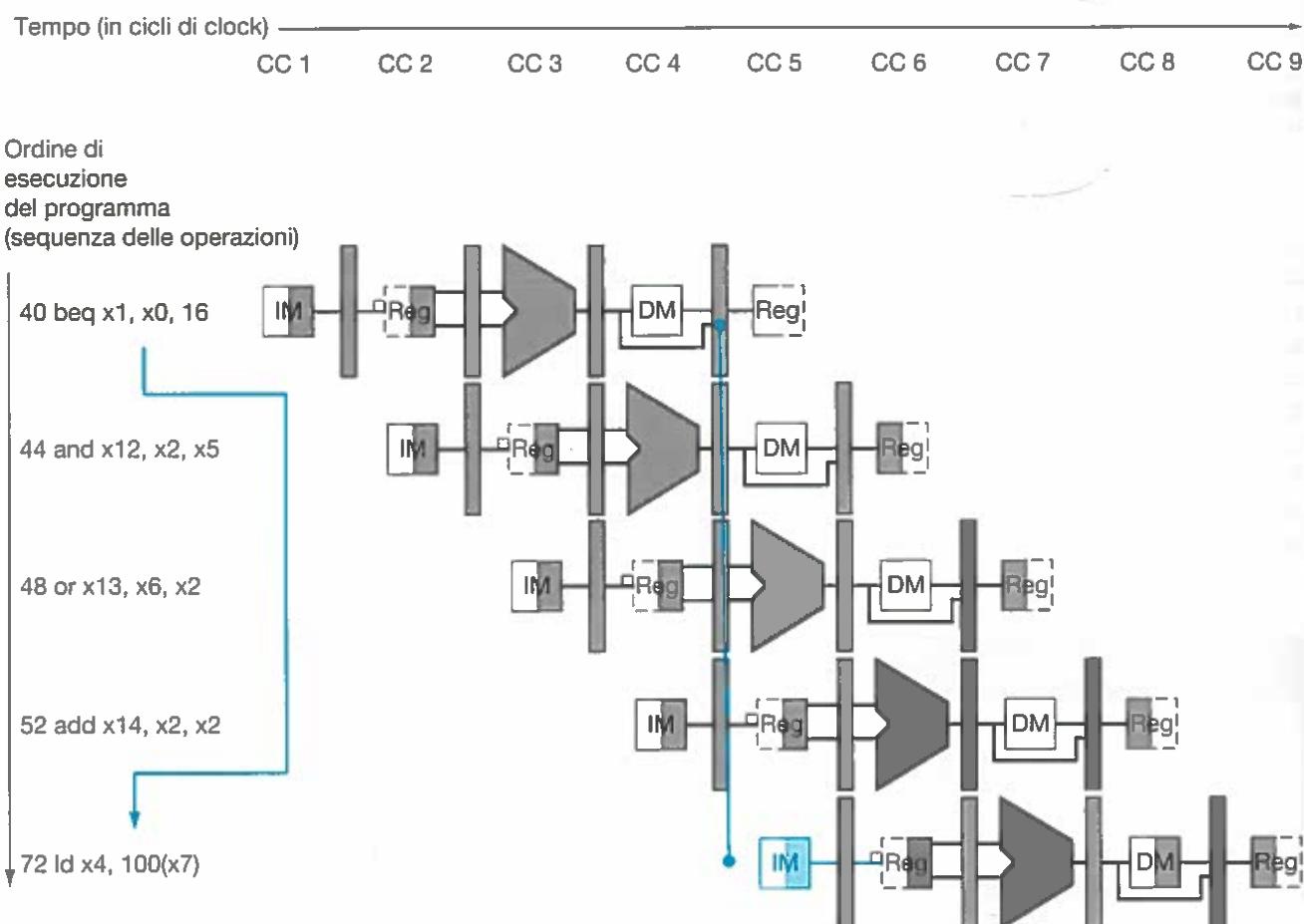
**Approfondimento.** A proposito della precedente osservazione sull'azzeramento dei segnali di controllo, si noti che per evitare la scrittura nei registri o in memoria i segnali RegWrite e MemWrite devono essere impostati a 0, mentre gli altri segnali di controllo possono assumere un valore indifferente.

## 4.8 Hazard sul controllo

*Sono in migliaia ad accanirsi contro le ramificazioni del male, ma è solo uno a colpirne la radice.*  
Henry David Thoreau, *Walden*, 1854

Finora ci siamo preoccupati soltanto degli hazard che coinvolgono le operazioni aritmetiche e di trasferimento dati. Tuttavia, come abbiamo visto nel paragrafo 4.5, ci sono anche hazard della pipeline che riguardano i salti condizionati. La Figura 4.59 mostra una sequenza di istruzioni e indica in quale punto della pipeline verrebbe eseguito il salto. Per alimentare la pipeline è necessario prelevare un'istruzione a ogni ciclo di clock, ma nell'implementazione hardware vista finora le decisioni relative ai salti condizionati non vengono prese prima dello stadio MEM della pipeline. Come abbiamo detto nel paragrafo 4.5, tale ritardo nella determinazione dell'istruzione successiva da prelevare è chiamato *hazard sul controllo* o *hazard sui salti condizionati*, per distinguerlo dagli *hazard sui dati* appena esaminati.

Questo paragrafo è più breve del paragrafo precedente (relativo agli hazard sui dati), poiché gli hazard sul controllo sono relativamente facili da comprendere, si verificano molto meno frequentemente e non c'è nulla di veramente efficace per risolverli. Le modifiche circuitali sono quindi più semplici. In particolare esamineremo due schemi per la risoluzione degli hazard sul controllo e un'ottimizzazione che permette di migliorare tali schemi.



**Figura 4.59 Impatto della pipeline sulle istruzioni di salto condizionato.** I numeri alla sinistra di ciascuna istruzione (40, 44, ...) rappresentano l'indirizzo delle istruzioni. Dato che l'istruzione di salto condizionato decide se occorre saltare nello stadio MEM, corrispondente al quarto ciclo di clock dell'istruzione beq, le tre istruzioni che la seguono vengono comunque prelevate dalla memoria istruzioni e viene avviata la loro esecuzione. Se non si intervenisse, l'esecuzione di queste tre istruzioni inizierebbe prima che beq salti eventualmente alla 1d contenuta all'indirizzo 72. In Figura 4.29 si supponeva che esistesse dell'hardware aggiuntivo per ridurre l'hazard sul controllo a un solo ciclo di clock, mentre in questa figura viene rappresentata l'unità di elaborazione non ottimizzata.

## Ipotizzare che il salto condizionato non sia eseguito

Come abbiamo visto nel paragrafo 4.5, mettere in stallo la pipeline fino a quando l'esecuzione dell'istruzione di salto condizionato non è stata completata fa perdere troppo tempo. Una tecnica frequentemente utilizzata per risolvere il problema consiste nel **predire** che il salto non debba essere eseguito, continuando quindi a eseguire le istruzioni secondo il loro normale flusso. Se poi risultasse che il salto doveva essere eseguito, le istruzioni che si trovano nella fase di fetch, di decodifica e di calcolo devono essere scartate e l'esecuzione deve riprendere a partire dall'istruzione contenuta all'indirizzo di destinazione del salto. Se i salti condizionati non devono essere eseguiti anche solo nella metà dei casi e se costa poco scartare le istruzioni, questa ottimizzazione permette di dimezzare il costo degli hazard sul controllo.

Per scartare le istruzioni bisogna semplicemente cambiare il valore dei segnali di controllo, forzandoli a 0, analogamente a quanto si è fatto per ottenere lo stallo nel caso degli hazard per l'utilizzo del dato di una load. La differenza consiste nel fatto che, quando l'istruzione di salto condizionato raggiunge lo stadio MEM, è necessario anche eliminare le tre istruzioni presenti negli stadi IF, ID e EX.

Se per le situazioni di stallo legate all'utilizzo del dato di una load era sufficiente mettere a 0 il valore dei segnali di controllo nello stadio ID, lasciando che questi si propagassero lungo la pipeline, ora bisogna eliminare le istruzioni presenti negli stadi IF, ID e EX della pipeline; in termini tecnici occorre effettuare un **flush** (svuotamento) di queste istruzioni.



**Flush (di istruzioni):** eliminazione di una o più istruzioni da una pipeline, di solito a causa di un evento inaspettato.

## Ridurre i ritardi associati ai salti condizionati

Un modo per migliorare le prestazioni sui salti condizionati consiste nel ridurre il costo nel caso in cui un salto condizionato venga erroneamente eseguito. Finora abbiamo supposto (per i salti condizionati) che l'indirizzo successivo da caricare nel PC venga preso dallo stadio MEM; ma se anticipassimo la decisione sul salto all'interno della pipeline si potrebbe scartare un numero inferiore di istruzioni. Anticipare la decisione su un salto condizionato richiede che due azioni vengano eseguite anticipatamente: calcolare l'indirizzo di destinazione del salto e valutare la condizione logica sulla base della quale il salto condizionato viene effettuato o meno.

La parte facilmente realizzabile di questa modifica è l'anticipazione del calcolo dell'indirizzo di destinazione del salto; il contenuto del PC e il campo immediato sono già disponibili nel registro IF/ID della pipeline ed è quindi sufficiente spostare il sommatore che calcola l'indirizzo di salto dallo stadio EX allo stadio ID; naturalmente, il calcolo dell'indirizzo di destinazione del salto sarà effettuato per tutte le istruzioni, ma verrà utilizzato solamente quando richiesto.

La parte più difficile consiste invece nel prendere la decisione sul salto. Se il salto è del tipo branch if equal, si deve confrontare il contenuto dei due registri letti nello stadio ID e verificare se il loro contenuto è uguale. Il test di uguaglianza può essere eseguito effettuando dapprima uno XOR bit a bit del contenuto dei due registri e inserendo poi in un NOR l'uscita degli XOR. Spostare il confronto nello stadio ID implica l'aggiunta di altra logica di propagazione e di rilevazione degli hazard; dato che la decisione sul salto dipende dai dati contenuti nella pancia della pipeline, dovremo assicurare che l'istruzione di salto condizionato funzioni correttamente anche con questa ottimizzazione. Per esempio, per implementare l'istruzione di branch if equal (e la sua inversa), abbiamo bisogno di propagare il risultato del test di uguaglianza dello stadio ID. Ci sono due elementi che complicano le cose: vediamoli.

1. Durante lo stadio ID dobbiamo decodificare l'istruzione, decidere se occorra propagare uno dei due operandi all'ingresso del comparatore e valutare l'uguaglianza, in modo che se l'istruzione richiede di saltare si possa scrivere nel PC l'indirizzo di destinazione del salto. Finora la propagazione degli operandi veniva gestita dall'unità di propagazione anche per i salti condizionati, ma l'introduzione del test di uguaglianza nello stadio ID richiede della logica di propagazione aggiuntiva. Si noti che gli operandi di un'istruzione di salto condizionato possono eventualmente essere prelevati dai registri EX/MEM o MEM/WB della pipeline.
2. Poiché i dati su cui fare il confronto sono richiesti già nello stadio ID, ma possono essere prodotti più tardi nella pipeline, è possibile che si verifichi un hazard sui dati e che la pipeline debba essere messa in stallo. Per esempio, se un'istruzione che opera sulla ALU si trova subito prima del salto e produce uno degli operandi per il confronto, ci sarà bisogno di uno stallo, poiché lo stadio EX dell'istruzione che opera sulla ALU è successivo nel tempo allo stadio ID dell'istruzione di salto condizionato. Di conseguenza, se una load è seguita da un'istruzione di salto condizionato e il dato letto dalla memoria è uno dei due operandi, viene richiesto uno stallo di due cicli di clock, dato che il risultato dell'istruzione di load compare nella pipeline solamente alla fine dello stadio MEM, ma viene richiesto all'inizio dello stadio ID dell'istruzione di salto.

Malgrado queste difficoltà, spostare l'esecuzione di un salto condizionato nello stadio ID rappresenta un miglioramento, perché riduce la penalità di un'istruzione di salto condizionato a una sola istruzione da scartare (cioè quella che si trova nella fase di fetch) nel caso in cui il salto debba essere eseguito. L'implementazione dei percorsi di propagazione e di rilevamento degli hazard verrà analizzata dettagliatamente negli esercizi.

Per eliminare l'istruzione in esecuzione dallo stadio IF viene aggiunto un segnale di controllo, denominato IF.Scarta (IF.Flush), che azzerza la parte del registro di pipeline IF/ID che contiene l'istruzione. Azzerando questa parte del registro si trasforma l'istruzione appena caricata in una `nop`, ossia in un'istruzione che non esegue alcuna operazione e non modifica quindi nessun elemento di stato del sistema.

### Salti condizionati nella pipeline

#### ESEMPIO

Mostrare che cosa succede quando all'interno della seguente sequenza di istruzioni il salto condizionato viene eseguito, supponendo che i salti vengano solitamente non eseguiti e che l'esecuzione del salto sia stata spostata nello stadio ID:

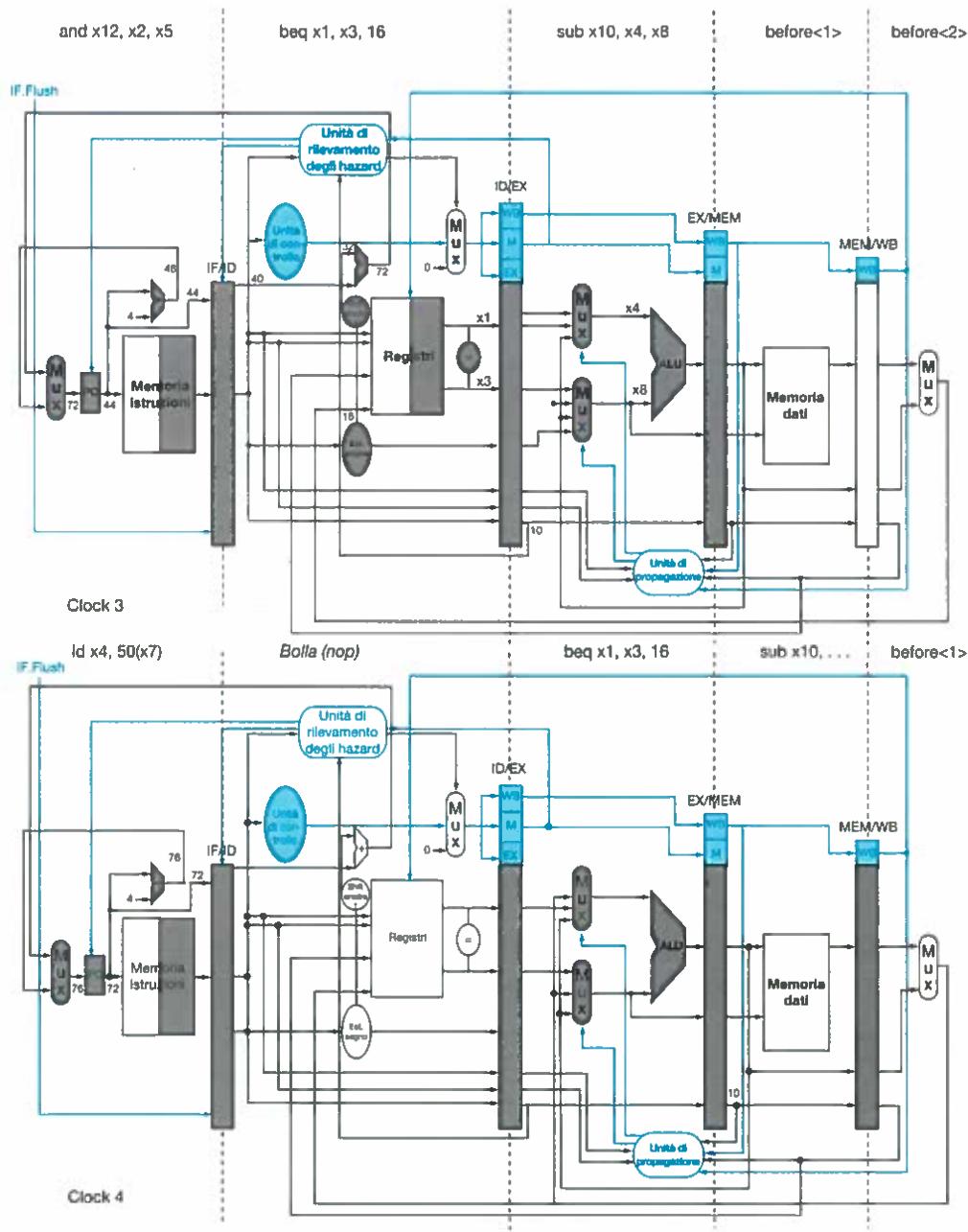
```

36 sub x10, x4, x8
40 beq x1, x3, 16 // Salto relativo al PC: 40 + 16 * 2 = 72
44 and x12, x2, x5
48 or x13, x2, x6
52 add x14, x4, x2
56 sub x15, x6, x7
...
72 ld x4, 50(x7)

```

#### SOLUZIONE

La Figura 4.60 mostra ciò che succede quando viene eseguito il salto. A differenza di Figura 4.59, occorre inserire una sola bolla nella pipeline nel caso in cui il salto venga eseguito.



**Figura 4.60** Nel terzo ciclo di clock lo stadio ID decide che il salto deve essere eseguito, seleziona quindi 72 come successivo indirizzo da caricare nel PC ed elimina l'istruzione caricata nel registro IF/ID. Nel quarto ciclo di clock viene caricata l'istruzione corrispondente all'indirizzo 72 e si inserisce la bolla corrispondente a un'istruzione nop all'interno della pipeline al posto dell'istruzione successiva all'istruzione di salto.

## Predizione dinamica dei salti

Quando si suppone che un salto condizionato non debba essere eseguito si realizza una forma rossa di *predizione del salto*; nel caso in cui la suposizione si riveli errata, occorre eliminare le istruzioni contenute nella pipeline. Per una semplice pipeline a cinque stadi questo approccio, possibilmente accoppiato a qualche tecnica di predizione utilizzata dal compilatore, risulta probabilmente adeguato. Con pipeline più profonde il costo dei salti condizionati, misurato in numero di cicli di clock, aumenta. In maniera simile nei processori a esecuzione parallela (par. 4.10) il costo dei salti condizionati cresce, in termini di numero di istruzioni perse. Ciò significa che in una pipeline “aggressiva” questo



### Predizione dinamica dei salti:

predizione dei salti durante l'esecuzione effettuata utilizzando informazioni raccolte durante l'esecuzione stessa.

**Buffer di predizione dei salti:** detto anche **tavella della storia dei salti**, è una piccola memoria, indicizzata dalla porzione meno significativa dell'indirizzo dell'istruzione di salto, contenente uno o più bit che indicano se in precedenza il salto è stato effettuato.

semplice schema di predizione statica causerebbe probabilmente una perdita di prestazioni eccessiva. Come accennato nel paragrafo 4.5, con più hardware a disposizione è possibile cercare di **predire** il comportamento dei salti durante l'esecuzione stessa del programma.

Un possibile approccio consiste nel verificare se l'ultima volta che un'istruzione di salto è stata eseguita il salto sia stato effettivamente eseguito: in caso positivo vengono caricate le istruzioni a partire da quella presente all'indirizzo di destinazione del salto. Questa tecnica viene chiamata **predizione dinamica dei salti**.

Un modo per implementare questa strategia consiste nell'utilizzare un **buffer di predizione dei salti** (*branch prediction buffer*), detto anche **tavella della storia dei salti**. Un buffer di predizione dei salti è una piccola memoria indicizzata attraverso la parte inferiore dell'indirizzo dell'istruzione di branch: questa memoria contiene un bit che indica se il salto era stato eseguito o meno nell'ultima esecuzione.

Questa predizione utilizza il buffer più semplice che si possa immaginare: infatti non possiamo sapere se la predizione sia corretta o meno; per esempio, la predizione potrebbe riguardare un'altra istruzione di salto condizionato la cui parte inferiore dell'indirizzo coincide con quella dell'istruzione di branch che stiamo analizzando. Tuttavia, un errore sulla predizione non influenza la correttezza dell'esecuzione: la predizione è solamente un suggerimento che si spera risulti corretto, su cui è basato il prelevamento dell'istruzione successiva. Se il suggerimento si rivela sbagliato, le istruzioni caricate erroneamente saranno eliminate, verrà invertito il bit di predizione e l'esecuzione ripartirà dal prelevamento dell'istruzione corretta.

Questo semplice schema di predizione a 1 bit ha un inconveniente: anche se un salto venisse quasi sempre effettuato, la predizione risulterebbe sbagliata almeno due volte invece che una volta sola, quando il salto non viene effettuato. L'esempio che segue illustra questo problema.

### Cicli e predizione

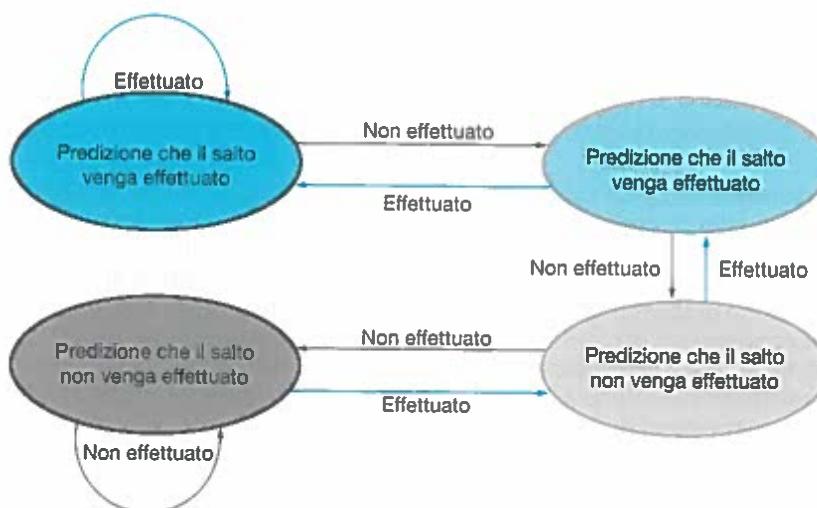
#### ESEMPIO

Si consideri un salto condizionato all'interno di un ciclo e si supponga che il salto venga eseguito nove volte di seguito e poi alla decima volta non venga eseguito. Qual è l'accuratezza fornita dalla previsione per questo salto, assumendo che il bit di previsione nel buffer di predizione non venga alterato dalle altre istruzioni?

#### SOLUZIONE

A regime la previsione non risulterà corretta nella prima e nell'ultima iterazione del ciclo. L'errore sull'ultima iterazione è inevitabile, in quanto la previsione ci dirà che il salto deve essere eseguito, essendo stato eseguito nove volte di seguito fino ad allora. L'errore sulla prima iterazione si verifica perché il bit di predizione commuta in occasione della precedente esecuzione dell'ultima iterazione del ciclo, quando il salto non era stato eseguito. Di conseguenza, l'accuratezza della previsione per questo salto (che viene eseguito nel 90% dei casi) è pari all'80% (due previsioni sbagliate e otto corrette).

In linea di principio si vorrebbe che per questi salti altamente regolari l'accuratezza della predizione corrispondesse alla frequenza di esecuzione del salto. Per rimediare a questa inadeguatezza, si utilizzano spesso degli schemi di previsione a 2 bit. In uno schema a 2 bit la predizione deve risultare sbagliata due volte di seguito prima di essere modificata. La Figura 4.61 mostra la macchina a stati finiti che implementa questo schema.



**Figura 4.61** Stati possibili in uno schema di predizione su 2 bit. Utilizzando 2 bit anziché 1 bit, una branch che abbia una forte tendenza a eseguire il salto o a non eseguirlo, cosa che accade molto frequentemente, avrà una predizione errata soltanto una volta. I 2 bit sono richiesti per codificare i 4 stati della macchina a stati finiti. Lo schema di predizione su 2 bit è un esempio di predizione basata sui contatori. Questi vengono incrementati quando la predizione si rivela accurata e decrementati in caso contrario; ai contatori vengono associati dei predittori, che utilizzano il valore medio di conteggio per decidere se il salto dovrà essere effettuato o meno.

Un buffer di predizione dei salti può essere realizzato come un piccolo buffer speciale, accessibile tramite l'indirizzo dell'istruzione durante lo stadio IF della pipeline. Se si prevede che il salto debba essere eseguito, il caricamento delle istruzioni successive sarà effettuato a partire dall'indirizzo di destinazione del salto non appena questo viene scritto nel PC; come detto in precedenza, questo può avvenire già nello stadio ID. In caso contrario, il prelevamento e l'esecuzione delle istruzioni continueranno in modo sequenziale. Se la predizione si rivela sbagliata, i bit di predizione saranno modificati, come mostrato in Figura 4.61.

**Approfondimento.** Un predittore dei salti ci dice se un salto debba essere effettuato o meno, ma il calcolo dell'indirizzo di destinazione del salto rimane un problema. In una pipeline a cinque stadi questo calcolo richiede un ciclo di clock. Ciò significa che i salti che vengono effettuati avrebbero il costo di un ciclo. Un approccio consiste nell'utilizzare una memoria cache per salvare l'indirizzo di destinazione del salto o l'istruzione di destinazione, detta anche **buffer degli indirizzi di salto** (*branch target buffer*).

Lo schema di predizione dinamica a 2 bit tiene conto solo delle informazioni sul particolare salto condizionato ad essa associato. I ricercatori hanno notato che utilizzare le informazioni sia sul salto considerato sia sul comportamento globale dei salti eseguiti di recente permette di migliorare notevolmente l'accuratezza della predizione, a parità di bit utilizzati. Questi predittori sono chiamati **predittori correlati** (*correlating predictors*). Un tipico predittore correlato contiene due predittori a 2 bit per ogni salto condizionato e la scelta tra i due predittori viene compiuta in base al risultato (effettuato o meno) dell'ultimo salto. Quindi, la storia globale del salto può essere pensata come un insieme ulteriore di indici per l'accesso alla tabella di predizione.

Un altro approccio alla predizione dei salti è rappresentato dai predittori a torneo. Un **predittore di salto a torneo** (*tournament branch predictor*) utilizza più predittori, tenendo traccia per ciascun salto del predittore migliore. Un tipico predittore a torneo contiene due predizioni per un certo salto: una è basata

**Buffer degli indirizzi di salto:** una struttura che memorizza l'indirizzo di destinazione del salto o l'istruzione di destinazione. Di solito viene organizzata come una cache e contiene anche il campo "tag", il che la rende più costosa di un semplice buffer di predizione.

**Predittore correlato:** un predittore di salto che combina il comportamento del salto considerato con informazioni globali associate al comportamento di alcune delle istruzioni di salto condizionato recentemente eseguite.

**Predittore di salto a torneo:** un predittore di salto con predizioni multiple per ciascun salto, dotato di un meccanismo di selezione che permette di scegliere il predittore più accurato per un certo salto.

su informazioni locali, l'altra su informazioni globali. Un circuito di selezione sceglierà quale predittore utilizzare; questo circuito funziona in modo simile a un predittore a 1 o 2 bit e favorisce il predittore che risulta più accurato. Alcuni dei microprocessori più recenti utilizzano questi predittori sofisticati.

**Approfondimento.** Un modo per ridurre il numero di istruzioni di salto condizionato è quello di aggiungere istruzioni di spostamento *condizionato*. Invece di modificare il contenuto del PC con un'istruzione di salto condizionato queste istruzioni modificano in modo condizionato il contenuto del registro di destinazione dell'istruzione di spostamento. Per esempio, l'architettura dell'insieme di istruzioni ARMv8 contiene un'istruzione di selezione condizionata chiamata **CSEL**. Questa istruzione definisce un registro destinazione, due registri sorgente e una condizione. Il registro destinazione assume il valore del primo operando se la condizione è vera e del secondo operando se la condizione è falsa. Quindi **CSEL X8, X11, X4, NE** copia il contenuto del registro 11 nel registro 8 se, in accordo con il codice della condizione specificato, il risultato dell'operazione non è zero, e copia il contenuto del registro 11 nel registro 4 se il risultato è zero. Quindi, i programmi che utilizzano l'insieme delle istruzioni ARMv8 possono contenere meno salti condizionati di quelli scritti in RISC-V.

## Riepilogo sulla pipeline

Siamo partiti dal bucato, mostrando i principi della pipeline applicati a un'attività comune della vita quotidiana. Utilizzando questo esempio abbiamo spiegato l'esecuzione in pipeline, passo dopo passo, a partire dall'unità di elaborazione a singolo ciclo a cui abbiamo aggiunto via via i registri di pipeline, i cammini di propagazione, le unità di rilevamento delle criticità, la predizione dei salti condizionati e l'eliminazione delle istruzioni richiesta quando si verificano eccezioni. La Figura 4.62 mostra l'unità di elaborazione finale con la relativa unità di controllo. Siamo ora pronti ad affrontare il difficile argomento delle eccezioni.

### Autovalutazione

Si considerino tre schemi di predizione dei salti: salti non effettuati, salti effettuati e predizione dinamica. Si supponga che tutti e tre non introducano un costo aggiuntivo (in termini di cicli di clock) quando predicono accuratamente il salto e comportino invece due cicli di clock in più quando sbagliano la predizione. Si supponga, inoltre, che l'accuratezza media della predizione sia del 90%. Quale schema di predizione è migliore nelle seguenti situazioni?

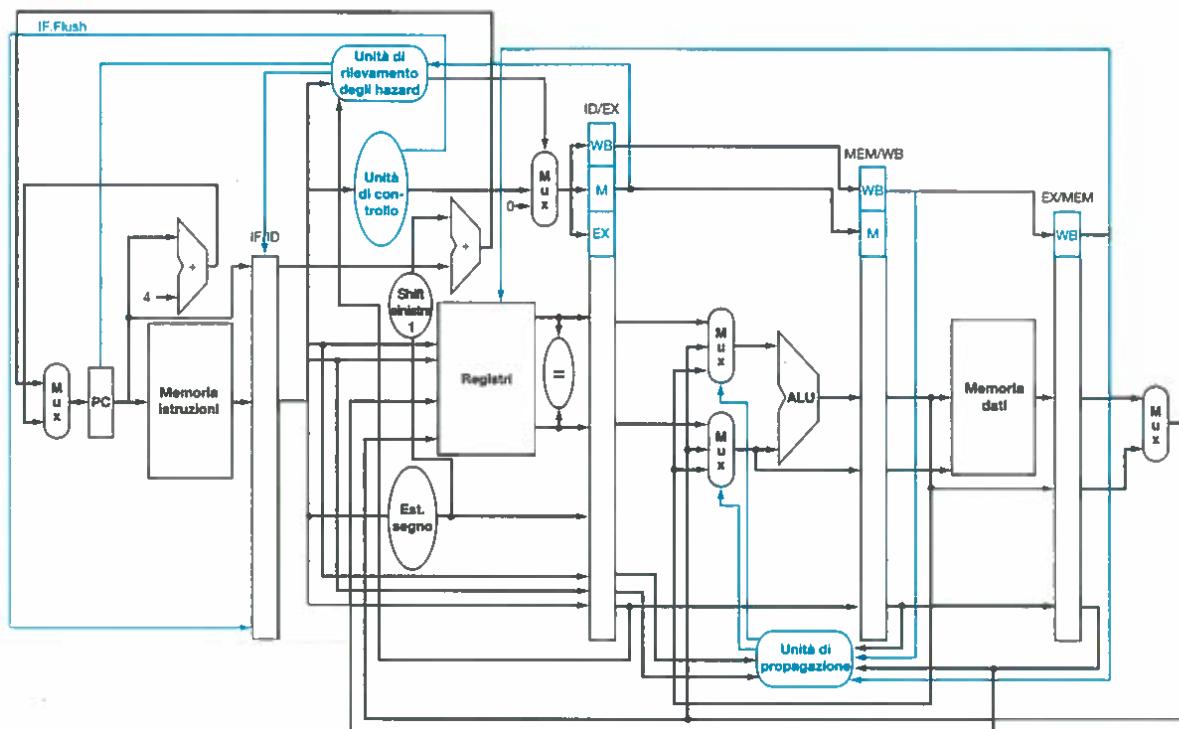
1. Un salto condizionato viene eseguito con una frequenza del 5%.
2. Un salto condizionato viene eseguito con una frequenza del 95%.
3. Un salto condizionato viene eseguito con una frequenza del 70%.

## 4.9 Le eccezioni

L'unità di controllo di un processore è la parte più difficile da far funzionare correttamente ed è anche la più difficile da rendere veloce. Una delle sfide maggiori è rappresentata dalla gestione delle **eccezioni** e degli **interrupt**, eventi distinti dai salti condizionati che alterano il normale flusso sequenziale di esecuzione delle istruzioni. Inizialmente erano stati concepiti per gestire eventi inaspettati che si verificavano all'interno del processore, come l'overflow aritmetico; lo stesso meccanismo di base fu poi esteso alla comunicazione tra i dispositivi di I/O e il processore (vedi Cap. 5).

*Far sì che un calcolatore dotato di una gestione automatica degli interrupt si comportasse [in modo sequenziale] non fu facile, dato che, quando si verifica un interrupt, il numero di istruzioni nei diversi stadi di elaborazione può essere molto grande.*

Fred Brooks Jr., *Planning a Computer System: Project Stretch*, 1962



**Figura 4.62** L'unità di elaborazione finale e la relativa unità di controllo della pipeline vista in questo capitolo. Si noti che questo è uno schema stilizzato e l'unità di elaborazione non è dettagliata; per esempio, manca il multiplexer pilotato dal segnale di controllo ALUSrc presente in Figura 4.55 e i segnali di controllo di Figura 4.49.

Molte architetture e molti ricercatori non fanno distinzione tra interrupt ed eccezioni, utilizzando spesso entrambi i termini per indicare entrambi i tipi di eventi. Per esempio, l'architettura x86 di Intel utilizza il termine interrupt. In questo testo utilizzeremo il termine *eccezione* per riferirci a un *qualsiasi* cambiamento non previsto del flusso di controllo, senza distinguere se esso abbia cause esterne o interne al processore; utilizzeremo invece il termine *interrupt* per riferirci solo agli eventi che hanno cause esterne. La tabella seguente riporta cinque tipi di eventi, distinguendo se vengono generati all'interno o all'esterno del processore e il nome utilizzato dal RISC-V.

**Eccezione:** detta anche **interrupt** (interruzione), è un evento non previsto che interrompe l'esecuzione di un programma. Può essere utilizzata, per esempio, per segnalare le condizioni di overflow.

**Interrupt:** eccezione che ha origine all'esterno del processore. Alcune architetture utilizzano il termine *interrupt* per indicare tutti i tipi di eccezioni.

Tipo di evento	Provenienza	Terminologia RISC-V
Reset di sistema	Esterno	Eccezione
Richiesta di un dispositivo di I/O	Esterno	Interrupt
Chiamata al sistema operativo di un programma utente	Interno	Eccezione
Utilizzo di un'istruzione non definita	Interno	Eccezione
Malfunzionamenti hardware	Entrambi	Eccezione o interrupt

Molti requisiti per il supporto delle eccezioni provengono dal contesto specifico che le causa. Torneremo perciò su questo argomento nel Capitolo 5, dove comprenderemo meglio la necessità di avere funzionalità aggiuntive per la gestione delle eccezioni. In questo paragrafo verrà descritta l'implementazione di un'unità di controllo in grado di rilevare due tipi di eccezioni che si possono verificare nell'insieme delle istruzioni e nell'architettura che le implementa presentata nei paragrafi precedenti.

Il rilevamento di condizioni eccezionali e la scelta dell'azione opportuna si trovano spesso sul cammino critico di un processore, che determina il periodo del clock e quindi le prestazioni. Occorre porre particolare attenzione alle eccezioni durante la progettazione dell'unità di controllo, perché aggiungere a posteriori la gestione delle eccezioni a un'implementazione hardware complessa può comportare una significativa riduzione delle prestazioni; inoltre, diventerebbe più difficile realizzare un'implementazione hardware che funzioni correttamente.

### Gestione delle eccezioni nelle architetture RISC-V

I soli tipi di eccezione che possono essere generati nell'implementazione hardware vista finora sono l'esecuzione di un'istruzione non valida e i malfunzionamenti hardware. In questo paragrafo supporremo che un malfunzionamento hardware si verifichi durante l'esecuzione dell'istruzione `add x11, x12, x11`. L'operazione fondamentale che il processore deve compiere al verificarsi di un'eccezione consiste nel salvare l'indirizzo dell'istruzione che l'ha generata nel *registro causa del supervisore delle eccezioni* (SEPC, *supervision exception cause register*) e trasferire il controllo a un indirizzo specifico del sistema operativo.

Il sistema operativo potrà quindi intraprendere le azioni più opportune, che potrebbero consistere nel fornire alcuni servizi al programma utente, eseguire un'azione predeterminata in risposta al malfunzionamento o terminare l'esecuzione del programma segnalando un errore. Dopo aver compiuto le azioni necessarie per rispondere all'eccezione, il sistema operativo può terminare il programma o riprenderne l'esecuzione utilizzando il SEPC per determinare il punto da cui ripartire. Nel Capitolo 5 esamineremo in maggior dettaglio come viene fatta ripartire l'esecuzione di un programma.

Per poter gestire correttamente un'eccezione, il sistema operativo deve conoscere la causa, oltre a sapere quale istruzione l'ha generata. Ci sono due metodi per comunicare la causa di un'eccezione: il metodo utilizzato nell'architettura RISC-V è quello di prevedere un registro dedicato (detto *registro causa di supervisione dell'eccezione* o SCAUSE, *Supervisor Exception Cause Register*) contenente un campo che indica la causa dell'eccezione.

Un metodo alternativo consiste nell'adottare **interrupt vettorizzati**, nei quali l'indirizzo a cui si deve trasferire il controllo viene determinato dalla causa dell'eccezione stessa, che fornisce un numero che viene sommato eventualmente a un registro base che punta a un'area di memoria contenente le istruzioni di risposta agli interrupt vettorizzati. Per esempio, per gestire i due tipi di eccezioni sopra riportati si potrebbero definire i seguenti indirizzi:

Tipo di eccezione	Indirizzo letto dell'eccezione vettorizzata da aggiungere al Registro Base della Tabella dei Vettori
Istruzione non definita	00 0100 0000 <sub>due</sub>
Errore di sistema (malfunzionamento hardware)	01 1000 0000 <sub>due</sub>

Il sistema operativo capisce il motivo dell'eccezione dall'indirizzo da cui inizia la risposta all'eccezione. Se le eccezioni non sono vettorizzate, come nel caso del RISC-V, il sistema operativo inizia a rispondere a tutte le eccezioni dallo stesso indirizzo e deve decodificare il registro causa per capire che cosa abbia generato l'eccezione. Per le architetture con una vettorizzazione delle eccezioni, gli indirizzi potrebbero essere spaziati, per esempio, di 32 byte ovverosia otto istruzioni, e il sistema operativo deve registrare il motivo dell'eccezione e può eseguire un'elaborazione di portata limitata in questo spazio.

**Interrupt vettorizzato:** un interrupt gestito in modo tale che l'indirizzo del sistema operativo al quale viene trasferito il controllo per la risposta a un'eccezione viene determinato dalla causa stessa dell'eccezione.

Si possono implementare le operazioni richieste per la gestione delle eccezioni aggiungendo all'implementazione base già vista alcuni registri e alcuni segnali di controllo ed estendendo leggermente l'unità di controllo. Supponiamo di implementare il metodo di gestione delle eccezioni con un singolo indirizzo di risposta, l'indirizzo  $0000\ 0000\ 1C09\ 0000_{esa}$  (l'implementazione delle eccezioni vettoriali non porrebbe maggiori difficoltà). Occorrerà introdurre due registri aggiuntivi alla nostra unità di elaborazione RISC-V:

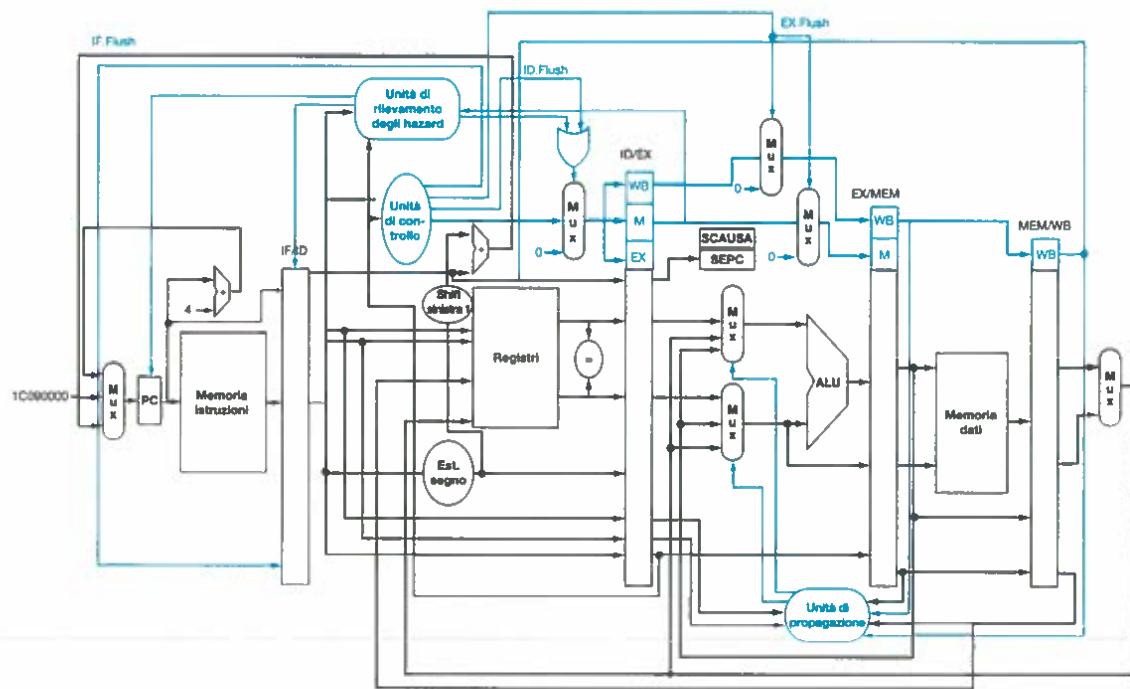
- **SEPC:** è un registro a 64 bit utilizzato per memorizzare l'indirizzo dell'istruzione che ha generato l'eccezione; questo registro serve anche nel caso di eccezioni vettorializzate.
- **SCAUSA:** un registro utilizzato per memorizzare la causa dell'eccezione. Nell'architettura RISC-V tale registro è a 64 bit, sebbene alcuni bit siano, ad oggi, inutilizzati. Si supponga che all'interno del registro esista un campo di cinque bit che codifica le due possibili cause di eccezione sopra citate: il codice 2 sarà associato all'eccezione di istruzione non definita e il codice 12 al malfunzionamento hardware.

## Eccezioni e loro gestione nella pipeline

In un'unità di elaborazione dotata di pipeline le eccezioni vengono trattate come se fossero una forma di hazard sul controllo. Per esempio, supponiamo che si verifichi un overflow aritmetico in un'istruzione di somma; esattamente come facevamo nel paragrafo precedente per i salti condizionati quando venivano presi, dobbiamo eliminare dalla pipeline le istruzioni che seguono la add e iniziare a prelevare le istruzioni dal nuovo indirizzo, che è quello di risposta all'eccezione. Utilizzeremo quindi lo stesso meccanismo dei salti condizionati, con la differenza che in questo caso i segnali di controllo devono essere deasserriti.

Quando dovevamo gestire gli errori sulle predizioni dei salti, abbiamo visto come eliminare un'istruzione nello stadio IF convertendola in una nop. Per eliminare un'istruzione che si trovava nello stadio ID, invece, utilizziamo il multiplexer già presente che mette a 0 i segnali di controllo, realizzando così di fatto lo stall della pipeline. Un nuovo segnale di controllo, chiamato ID.Scarta (ID.Flush), viene messo in OR con il segnale di stall che proviene dall'unità di rilevamento degli hazard per eliminare l'istruzione nello stadio ID. Per eliminare un'istruzione nello stadio EX, viene utilizzato un nuovo segnale di controllo, chiamato EX.Scarta (EX.Flush), che pilota un nuovo multiplexer utilizzato per mettere a 0 i segnali di controllo di questo stadio. Per prelevare l'istruzione successiva dall'indirizzo  $0000\ 0000\ 1C09\ 0000_{esa}$ , che nel RISC-V contiene l'indirizzo della prima istruzione della procedura del sistema operativo di risposta alle eccezioni, occorre semplicemente aggiungere un'altra linea in ingresso, che porta il valore  $0000\ 0000\ 1C09\ 0000_{esa}$ , al multiplexer che seleziona il nuovo valore del PC. Queste modifiche sono mostrate in Figura 4.63.

Questo esempio mostra un problema nella gestione delle eccezioni: se non si interrompe l'istruzione prima della fine della sua esecuzione, il programmatore non sarà più in grado di vedere il valore originale del registro  $x1$ , dato che tale registro verrà sporco, poiché è anche il registro destinazione della stessa istruzione add che ha generato l'eccezione. Tuttavia, se supponiamo che l'eccezione venga riconosciuta nello stadio EX si può utilizzare il segnale EX.Flush per prevenire che l'istruzione ora nello stadio di EX scriva poi il registro destinazione nello stadio WB. Molte eccezioni richiedono di completare normalmente l'esecuzione dell'istruzione che ha causato l'eccezione come se fossero eseguite normalmente. Il modo più semplice per ottenere ciò è eliminare l'istruzione e farla eventualmente ripartire dall'inizio dopo che l'eccezione è stata gestita.



**Figura 4.63** L'unità di elaborazione in grado di gestire le eccezioni, con i relativi segnali di controllo. I cambiamenti più importanti sono l'aggiunta di: un nuovo ingresso con valore 0000 0000 1C09 0000<sub>esa</sub> al multiplexer che fornisce il nuovo indirizzo al PC, di un registro SCAUSA che memorizza la causa dell'eccezione e di un registro SEPC per salvare l'indirizzo dell'istruzione che ha provocato l'eccezione. L'ingresso 0000 0000 1C09 0000<sub>esa</sub> del multiplexer rappresenta l'indirizzo della prima istruzione della procedura del sistema operativo di risposta alle eccezioni. Il segnale di overflow della ALU è uno degli ingressi dell'unità di controllo, anche se in questo schema non viene mostrato.

L'ultimo passo consiste nel salvare l'indirizzo dell'istruzione che ha causato l'eccezione nel *program counter del supervisore delle eccezioni* (SEPC). La Figura 4.63 mostra una versione schematica dell'unità di elaborazione che comprende i circuiti e i componenti necessari per la gestione dei salti condizionati e delle eccezioni.

### Eccezioni in un calcolatore dotato di pipeline

#### ESEMPIO

Supponiamo di avere la seguente sequenza di istruzioni:

```

40esa sub    x11, x2, x4
44esa and    x12, x2, x5
48esa or     x13, x2, x6
4Cesa add    x1, x2, x1
50esa sub    x15, x6, x7
54esa ld     x16, 100(x7)
...

```

e che la procedura a cui viene trasferito il controllo in caso di eccezione inizi nel modo seguente:

```

1C090000esa sd    x26, 1000(x10)
1C090004esa sd    x27, 1008(x10)
...

```

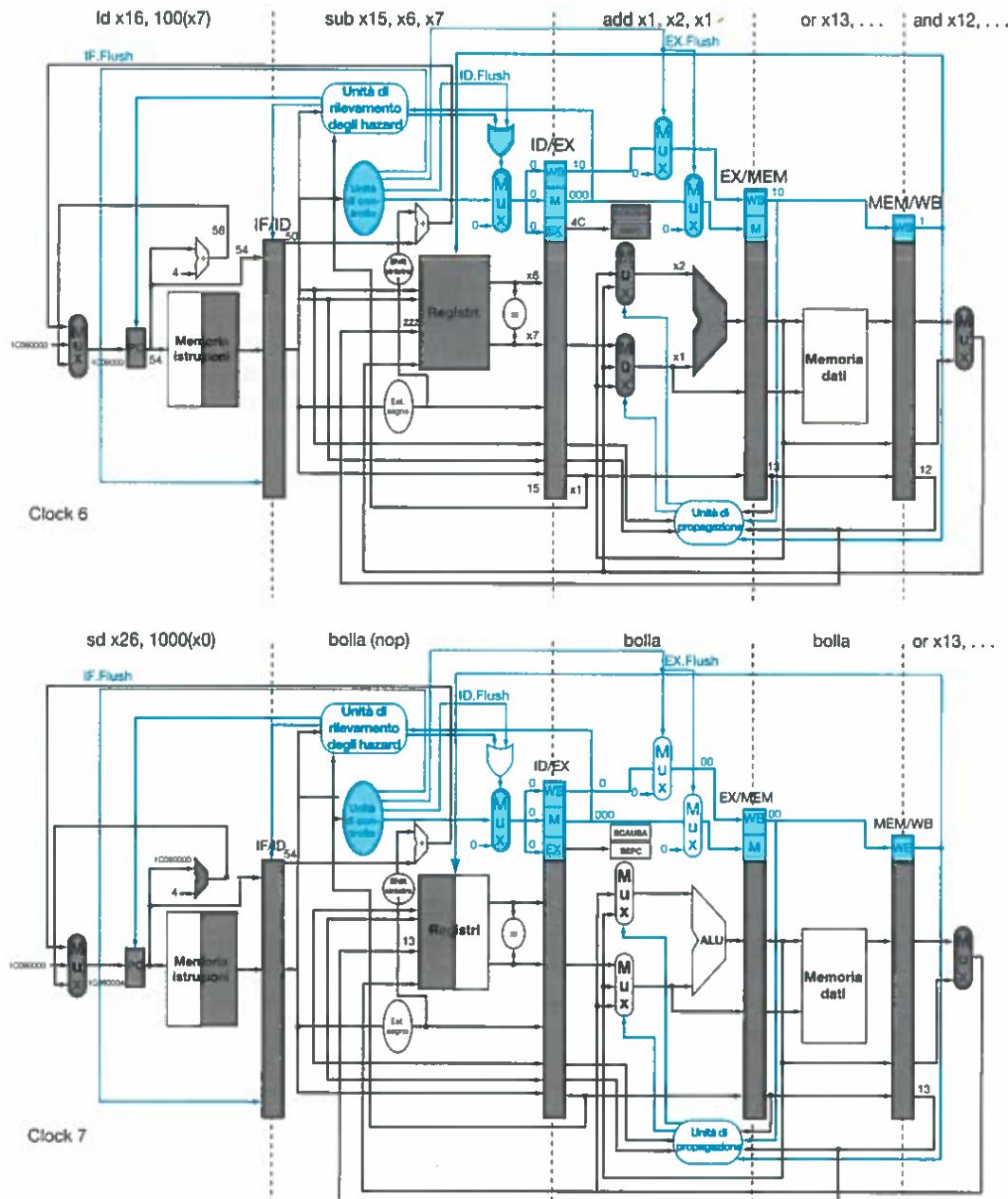
Mostrare che cosa avviene all'interno della pipeline se si verifica un'eccezione di malfunzionamento dell'hardware nell'istruzione add.

(continua)

(continua)

**SOLUZIONE**

La Figura 4.64 mostra gli eventi che si susseguono a partire da quando l'istruzione add si trova nello stadio EX. Si suppone che il malfunzionamento dell'hardware venga rilevato in questo stadio e che il valore 0000 0000 1C09 0000<sub>esa</sub> venga forzato nel PC. Il settimo ciclo di clock mostra che la add e le istruzioni successive vengono scartate e che viene poi prelevata la prima istruzione della procedura di risposta alle eccezioni. Si noti che si salva l'indirizzo dell'istruzione add: 4C<sub>esa</sub>.



**Figura 4.64** Risultato di un'eccezione dovuta a un malfunzionamento dell'hardware in un'istruzione add. L'eccezione viene rilevata nello stadio EX del ciclo di clock numero 6 e viene salvato l'indirizzo dell'istruzione add nel registro SEPC (4C<sub>esa</sub>). Ciò provoca l'asserramento di tutti i segnali di Flush in prossimità della fine di questo ciclo di clock e il deasserramento (impostazione a 0) dei segnali di controllo associati alla add. Nel ciclo di clock numero 7 le istruzioni già presenti nella pipeline risultano convertite in bolle mentre viene caricata in fase di fetch la prima istruzione della funzione di risposta alle eccezioni – sd x26, 1000 (x0) – dall'indirizzo della memoria istruzioni 0000 0000 1C09 0000<sub>esa</sub>. Si noti che l'esecuzione delle istruzioni and e or, caricate prima della add, viene comunque completata.

All'inizio del paragrafo 4.9 abbiamo introdotto alcune possibili eccezioni e vedremo altri esempi nel Capitolo 5. Con cinque istruzioni attive a ogni ciclo di clock, la sfida è associare ciascuna eccezione all'istruzione corretta; inoltre, esiste la possibilità che nello stesso ciclo di clock più eccezioni si verifichino simultaneamente. La soluzione utilizzata consiste nell'associare una priorità alle eccezioni, in modo tale che sia semplice determinare quale eccezione debba essere gestita per prima. Nella maggior parte delle implementazioni RISC-V, l'hardware ordina le eccezioni in modo da interrompere la prima istruzione che ha generato un'eccezione.

Le richieste da parte dei dispositivi di I/O e i malfunzionamenti hardware non sono associati a istruzioni specifiche e quindi, in questi casi, è consentita una maggiore flessibilità su quando interrompere la pipeline. Ne consegue che il meccanismo utilizzato per le eccezioni funziona bene anche in questi casi.

Il registro SEPC mantiene l'indirizzo delle istruzioni interrotte e il registro SCAUSA, quando si verificano più eccezioni in uno stesso ciclo di clock, memorizza l'eccezione a priorità più elevata.

## Interfaccia hardware/software

L'hardware e i sistemi operativi devono lavorare in modo congiunto affinché le eccezioni si comportino nel modo previsto. I compiti dell'hardware sono normalmente i seguenti: bloccare l'istruzione responsabile dell'eccezione prima che venga terminata, far sì che le istruzioni precedenti vengano completate, eliminare tutte le istruzioni successive già caricate nella pipeline, scrivere in un registro la causa dell'eccezione, salvare l'indirizzo dell'istruzione responsabile dell'eccezione e trasferire il controllo a un indirizzo preassegnato. Il compito del sistema operativo consiste invece nell'esaminare la causa dell'eccezione e nell'agire in modo appropriato. Per le eccezioni causate da istruzioni non definite o malfunzionamenti hardware, il sistema operativo normalmente termina (in inglese *kills*) il programma e segnala il motivo dell'interruzione. Nel caso di richieste dei dispositivi di I/O o di chiamate al sistema operativo, quest'ultimo salva lo stato del programma, svolge il compito richiesto e riprende quindi l'esecuzione del programma dal punto in cui era stato lasciato. Nel caso di richieste dei dispositivi di I/O, si sceglie spesso di eseguire un altro programma prima di riprendere l'esecuzione di quello che aveva inviato l'eccezione di I/O, poiché spesso occorre attendere che l'I/O sia stato completato. Questa è la ragione per cui è un passaggio critico salvare e ripristinare lo stato di un qualsiasi programma. Uno degli utilizzi più importanti e frequenti delle eccezioni è la gestione dei page fault e delle eccezioni del TLB che verranno descritte in dettaglio nel Capitolo 5.

**Interrupt imprecisi:** detti anche **eccezioni imprecise**, sono interrupt oppure eccezioni in calcolatori dotati di pipeline che non sono associati all'istruzione che ha realmente causato l'eccezione o l'interrupt.

**Interrupt precisi:** detti anche **eccezioni precise**, sono interrupt oppure eccezioni che vengono sempre associati all'istruzione che ha realmente causato l'eccezione o l'interrupt.

**Approfondimento.** La difficoltà nei calcolatori dotati di pipeline nell'associare sempre correttamente un'eccezione all'istruzione che l'ha provocata ha fatto sì che alcuni progettisti abbiano "ammorbidito" questo vincolo nei casi non critici: si dice che questi processori hanno **interrupt imprecisi** o **eccezioni imprecise**. Nell'esempio precedente, in condizioni normali il PC conterebbe il valore  $58_{esa}$  all'inizio del ciclo di clock successivo al riconoscimento dell'eccezione, anche se l'istruzione che aveva provocato l'eccezione si trovava all'indirizzo  $4C_{esa}$ . Un calcolatore con eccezioni imprecise metterebbe  $58_{esa}$  nel registro SEPC, lasciando al sistema operativo il compito di determinare quale istruzione abbia effettivamente provocato l'eccezione. Il RISC-V e la stragrande maggioranza dei calcolatori di oggi supportano **interrupt** ed **eccezioni precisi**. Uno dei motivi di questa scelta è che i progettisti di pipeline più profonde potrebbero avere la tentazione di memorizzare in SPEC un valore diverso, cosa che creerebbe un bel mal di testa al sistema operativo. Per evitare ciò, è meglio che le pipeline più profonde memorizzino in SPEC lo stesso valore del PC che verrebbe memorizzato nelle pipeline a cinque stadi che abbiamo esaminato. È più semplice per tutti memorizzare il PC dell'istruzione che ha causato l'eccezione. (Un altro motivo è il supporto alla memoria virtuale, come vedremo nel Capitolo 5.)

**Approfondimento.** Il RISC-V utilizza come indirizzo della procedura di risposta alle eccezioni l'indirizzo 0000 0000 1C09 0000<sub>esa</sub>, che è stato scelto in modo arbitrario. Molti calcolatori RISC-V memorizzano l'indirizzo di risposta alle eccezioni in un registro speciale chiamato *vettore delle "trappole" del supervisore* (STVEC, *Supervisor Trap Vector*), nel quale il sistema operativo può memorizzare un indirizzo scelto a suo piacimento.

### Autovalutazione

Quale eccezione verrà riconosciuta per prima in questa sequenza di istruzioni?

1. xxx x11, x12, x11 // istruzione non definita
2. sub x11, x12, x11 // errore hardware

## 4.10 | Parallelismo a livello di istruzioni

Avvisiamo il lettore che questo paragrafo contiene un'introduzione ad argomenti affascinanti, ma molto avanzati. Per un approfondimento consigliamo di consultare il testo, scritto dagli stessi autori, *Computer Architecture: A Quantitative Approach* (quinta edizione), dove gli argomenti qui appena accennati vengono trattati in oltre 200 pagine (comprese le appendici).

La pipeline sfrutta il **parallelismo** potenziale esistente tra le istruzioni. Questa forma di parallelismo viene chiamata, in modo naturale, **parallelismo a livello di istruzioni (ILP, Instruction-Level Parallelism)**. Ci sono due metodi principali per incrementare potenzialmente la quantità di parallelismo a livello di istruzioni. Il primo consiste nell'aumentare la profondità della pipeline per sovrapporre più istruzioni. Utilizzando l'analogia della lavanderia, si può pensare di suddividere il lavaggio tradizionale, supponendo che sia l'operazione più costosa in termini di tempo, in più fasi, utilizzando diverse macchine separate che si occupano di lavare, risciacquare e centrifugare; la pipeline passerebbe in questo caso da quattro a sei stadi e, per ottenere il massimo aumento di velocità, si renderebbe necessario bilanciare nuovamente i restanti stadi in modo tale che abbiano la stessa durata; questo vale sia per il bucato sia per i processori. La quantità di parallelismo è più alta, poiché più istruzioni vengono sovrapposte nel tempo e le prestazioni sono potenzialmente migliori (dato che il ciclo di clock può essere ridotto).

Un altro approccio consiste nel replicare i componenti interni del calcolatore in modo tale che sia possibile lanciare in esecuzione più istruzioni in ogni stadio della pipeline. Il nome generico che si dà a questa tecnica è **esecuzione parallela (multiple issue)**. In una lavanderia a esecuzione parallela, per esempio, la lavatrice e l'asciugatrice di casa verrebbero sostituite con tre lavatrici e tre asciugatrici; inoltre sarebbe necessario trovare altra manodopera per stirare e per riporre negli armadi il triplo della biancheria nello stesso lasso di tempo. Lo svantaggio di questo approccio consiste nel lavoro addizionale necessario a mantenere impegnate tutte le macchine e a trasferire i carichi di bucato da uno stadio della pipeline al successivo.

Lanciare in esecuzione più istruzioni nello stesso stadio fa sì che il tasso di esecuzione delle istruzioni sia maggiore della frequenza di clock o, in altri termini, che il CPI sia minore di 1. In alcuni casi è utile invertire l'unità di misura e utilizzare il *numero di istruzioni per ciclo di clock*, detto anche *IPC*. Quindi, un microprocessore a 3 GHz a esecuzione parallela con quattro vie ha una velocità



**Parallelismo a livello di istruzioni (ILP):** parallelismo tra istruzioni diverse, le quali vengono eseguite in parallelo.



**Esecuzione parallela (multiple issue):** uno schema in base al quale più istruzioni vengono lanciate in esecuzione in uno stesso ciclo di clock.

di esecuzione di picco di 12 miliardi di istruzioni al secondo e quindi un CPI massimo di 0,33, cioè un IPC di 3. Supponendo una pipeline a cinque stadi, un processore di questo genere potrebbe eseguire fino a 20 istruzioni contemporaneamente. Gli attuali calcolatori di fascia alta possono arrivare a eseguire da tre a sei istruzioni in ogni stadio della pipeline e i calcolatori di fascia media puntano a ottenere un IPC di picco pari a 2. Tuttavia, ci sono molti vincoli sulle istruzioni che possono essere eseguite in parallelo e su che cosa fare quando si verificano delle dipendenze.

Esistono due modalità principali per realizzare processori a esecuzione parallela, la cui principale differenza consiste nel modo in cui viene suddiviso il lavoro tra il compilatore e l'hardware: le decisioni possono essere prese staticamente (cioè durante la compilazione) o dinamicamente (cioè durante l'esecuzione). Questi due approcci sono detti **parallelizzazione statica dell'esecuzione** (*static multiple issue*) e **parallelizzazione dinamica dell'esecuzione** (*dynamic multiple issue*). Come vedremo, entrambi questi approcci hanno altri nomi più comuni che possono risultare, però, meno precisi o più restrittivi.

Ci sono due compiti principali, ben distinti, che una pipeline a esecuzione parallela deve essere in grado di eseguire: vediamoli.

1. Impacchettare le istruzioni in **slot di esecuzione** (*issue slot*). Come fa a sapere il processore quali e quante istruzioni possono essere lanciate in esecuzione in un dato ciclo di clock? Nella maggior parte dei processori dotati di parallelizzazione statica, questo processo viene parzialmente gestito dal compilatore. Nelle architetture dotate di parallelizzazione dinamica, in genere, è il processore a decidere durante l'esecuzione stessa, anche se il compilatore spesso cerca comunque di ottimizzare l'ordine delle istruzioni per aumentare la velocità di esecuzione.
2. Gestire gli hazard sui dati e sul controllo. Nei processori dotati di parallelizzazione statica alcune o tutte le conseguenze degli hazard sui dati e sul controllo vengono gestite staticamente dal compilatore, mentre la maggior parte dei processori dotati di parallelizzazione dinamica tenta di attenuarne gli effetti (almeno di alcuni tipi di hazard) utilizzando tecniche hardware che operano durante l'esecuzione stessa.

Benché questi due approcci vengano qui descritti come distinti, in realtà le tecniche di parallelizzazione sono spesso ibride, perché ciascun approccio prende in prestito alcune delle tecniche tipiche dell'altra modalità; in altri termini, nessun approccio alla parallelizzazione è solamente statico o solamente dinamico.

## Concetto di speculazione



**Speculazione:** un approccio all'esecuzione in cui il compilatore o il processore tentano di indovinare il risultato di un'istruzione per rimuovere la dipendenza di altre istruzioni che dovrebbero essere eseguite successivamente.

Uno dei metodi più importanti per scoprire e sfruttare al massimo l'ILP è la speculazione. La **speculazione** è basata sulla grande idea della **predizione** ed è un approccio che permette al compilatore o al processore di "indovinare" le caratteristiche di un'istruzione, in modo da permettere l'inizio dell'esecuzione di altre istruzioni che dipendono da quell'istruzione su cui è stata fatta la speculazione. Per esempio, si può speculare sul risultato del confronto associato a un salto condizionato in modo tale che l'esecuzione delle istruzioni successive possa essere anticipata. Oppure, si può speculare se una store, posta subito prima di una load, operi su un indirizzo di memoria differente, il che permetterebbe alla load di essere eseguita prima della store. Il problema è che la speculazione può rivelarsi sbagliata. Perciò, ogni meccanismo di speculazione deve comprendere anche un meccanismo per verificare se la speculazione è corretta e un meccanismo per poter tornare indietro o eliminare gli effetti delle istruzioni eseguite in base alla speculazione. L'eliminazione degli effetti delle istruzioni eseguite in modo speculativo rende la progettazione della pipeline più complicata.

La speculazione può essere fatta dal compilatore oppure dall'hardware. Per esempio, il compilatore può utilizzare la speculazione per riordinare le istruzioni, spostando un'istruzione dopo un salto condizionato oppure invertendo tra loro una load e una store. Il processore può operare le stesse modifiche durante l'esecuzione mediante tecniche che saranno discusse in seguito.

I meccanismi di correzione utilizzati in caso di speculazione errata sono abbastanza differenti tra loro. Nel caso di speculazione software, il compilatore di solito inserisce istruzioni aggiuntive che controllano l'accuratezza dalla speculazione e fornisce una procedura di riparazione da utilizzare quando si verifica un errore. Nel caso di speculazione hardware, il processore di solito memorizza e mantiene in buffer dedicati i risultati dell'esecuzione speculativa, finché non "capisce" se la speculazione è risultata corretta o meno. Se la speculazione risulta corretta, l'esecuzione delle istruzioni viene completata scrivendo il risultato nei registri o nella memoria. Se invece la speculazione si rivela errata, l'hardware elimina il contenuto dei buffer associati all'istruzione e riprende l'esecuzione con la sequenza corretta di istruzioni. Una speculazione errata tipicamente richiede di svuotare la pipeline o almeno di metterla in stallo, e quindi riduce ulteriormente le prestazioni.

La speculazione introduce un altro genere di problema: essa può introdurre eccezioni che precedentemente non erano presenti. Per esempio, si supponga che sulla base di una speculazione una load venga spostata e che l'indirizzo a cui fa riferimento non sia più valido se la speculazione risultasse sbagliata; in questo caso verrebbe generata un'eccezione di indirizzo non valido che non dovrebbe essere generata. Il problema è reso complicato dal fatto che, quando la load non è oggetto di speculazione, se l'indirizzo risulta non valido, l'eccezione deve essere sicuramente generata! Nella speculazione basata su compilazione, questo problema viene risolto introducendo uno speciale supporto alla speculazione che permette di ignorare queste eccezioni finché non è chiaro se debbano essere sollevate realmente. Nella speculazione hardware, invece, le eccezioni vengono semplicemente memorizzate in un buffer finché la correttezza della speculazione, in base alla quale l'istruzione era stata lanciata in esecuzione, non è stata verificata; solo allora l'esecuzione dell'istruzione può essere terminata. A questo punto l'eccezione viene generata e viene iniziata la procedura di gestione dell'eccezione.

Dato che la speculazione può migliorare le prestazioni quando viene fatta in maniera appropriata, ma può diminuirle quando è implementata in modo poco accurato, un grosso sforzo è dedicato a decidere quando sia opportuno speculare. Verso la fine di questo paragrafo esamineremo alcune tecniche, sia hardware sia software, per la speculazione.

## Parallelizzazione statica dell'esecuzione

Tutti i processori dotati di parallelizzazione statica dell'esecuzione utilizzano il compilatore per formare i pacchetti di istruzioni e gestire gli hazard. In questo tipo di processore le istruzioni che vengono lanciate in esecuzione in parallelo nello stesso ciclo di clock costituiscono un **pacchetto di istruzioni** (*issue packet*), che si può intendere come un'istruzione ampia contenente più istruzioni al suo interno. Questa immagine è più di una semplice analogia. Dal momento che, in un processore dotato di parallelizzazione statica, la combinazione di istruzioni che possono essere lanciate in esecuzione nello stesso ciclo di clock è in generale limitata, è utile pensare al pacchetto di istruzioni come a un'istruzione singola nella quale diverse operazioni vengono codificate in campi predefiniti; questa visione ha prodotto il nome originale di questo approccio: **Very Long Instruction Word (VLIW)**.

La maggior parte dei processori con parallelizzazione statica si appoggia al compilatore per la gestione degli hazard sui dati e sul controllo. Le responsa-

**Pacchetto di istruzioni:** un insieme di istruzioni che vengono lanciate in esecuzione nello stesso ciclo di clock; la composizione del pacchetto può essere determinata sia staticamente dal compilatore sia dinamicamente dal processore.

**Very Long Instruction Word (VLIW):** un tipo di architettura dell'insieme di istruzioni in grado di lanciare in esecuzione più istruzioni, indipendenti tra loro, attraverso un'istruzione codificata su molti bit, la quale tipicamente contiene diversi codici operativi.

bilità del compilatore comprendono la predizione statica dei salti e il riordino del codice per ridurre o prevenire gli hazard. Esaminiamo ora una semplice versione del processore RISC-V dotata di parallelizzazione statica prima di descrivere l'utilizzo di queste tecniche in processori più aggressivi.

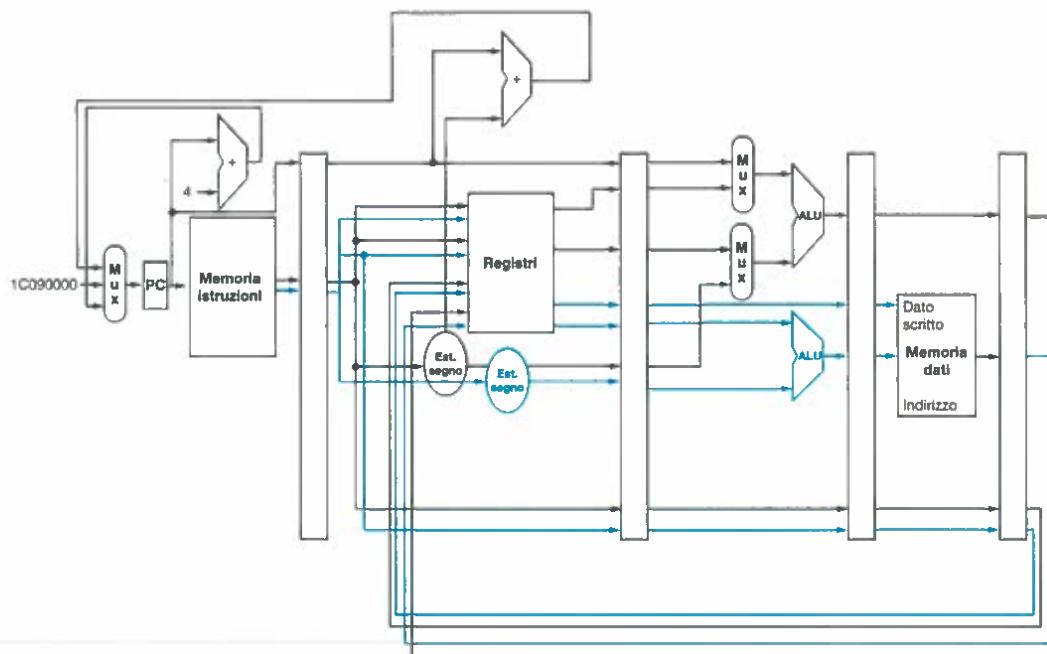
### Un esempio: parallelizzazione statica dell'ISA del RISC-V

Per dare un'idea di come funziona un processore dotato di parallelizzazione statica dell'esecuzione, consideriamo un semplice processore RISC-V dotato di due cammini di esecuzione, detti anche *vie*: il primo cammino può eseguire istruzioni che operano sulla ALU intera o istruzioni di salto condizionato, il secondo può eseguire istruzioni di load e store. Una pipeline di questo genere viene utilizzata in alcuni processori di tipo embedded. Lanciare in esecuzione due istruzioni per ciclo di clock richiede il fetch e la decodifica di istruzioni su 64 bit. In molti processori a parallelizzazione statica (e praticamente in tutti i processori VLIW) esistono dei vincoli su quali istruzioni si possano inserire all'interno della coppia di istruzioni eseguite contemporaneamente, in modo da semplificare la fase di decodifica e di invio all'esecuzione. Quindi, richiederemo che le istruzioni vengano accoppiate e allineate a parole di 64 bit, con l'istruzione che utilizza la ALU intera o l'istruzione di salto per prima. Inoltre, se una delle due istruzioni non può essere utilizzata, essa deve essere rimpiazzata da una nop. Le istruzioni vengono quindi lanciate in esecuzione sempre in coppia, eventualmente con una nop al posto di una delle due. La Figura 4.65 illustra come le istruzioni appaiono quando procedono accoppiate all'interno della pipeline.

I processori dotati di parallelizzazione statica si differenziano in base a come trattano i potenziali hazard sui dati e sul controllo. In alcune architetture il compilatore ha la piena responsabilità di rimuovere *tutti* gli hazard, riordinare il codice e inserire delle istruzioni nop, in modo tale che il codice venga eseguito senza alcun bisogno di rilevare hazard o generare stalli. In altre architetture, è l'hardware a rilevare gli hazard e a generare gli stalli tra due pacchetti di istruzioni consecutivi, mentre viene richiesto al compilatore di eliminare tutte le possibili dipendenze all'interno dello stesso pacchetto di istruzioni. Anche così, un hazard in genere obbliga a mettere in stallo l'intero pacchetto di istruzioni.

Tipo di istruzione	Stadi della pipeline						
Istruzione ALU o di salto	IF	ID	EX	MEM	WB		
Istruzione di load o di store	IF	ID	EX	MEM	WB		
Istruzione ALU o di salto	IF	ID	EX	MEM	WB		
Istruzione di load o di store	IF	ID	EX	MEM	WB		
Istruzione ALU o di salto		IF	ID	EX	MEM	WB	
Istruzione di load o di store		IF	ID	EX	MEM	WB	
Istruzione ALU o di salto			IF	ID	EX	MEM	WB
Istruzione di load o di store			IF	ID	EX	MEM	WB

**Figura 4.65** Unità di elaborazione dotata di parallelizzazione statica a due vie. Le istruzioni che utilizzano la ALU e quelle di trasferimento dati vengono lanciate in esecuzione contemporaneamente. In questa tabella abbiamo supposto una pipeline a cinque stadi simile a quella senza parallelizzazione. Anche se ciò non è strettamente necessario, ha alcuni vantaggi. In particolare, mantenere la scrittura dei registri al termine dell'esecuzione semplifica la gestione delle eccezioni e favorisce l'utilizzo di un preciso modello di gestione delle stesse. Queste funzionalità diventano più complicate nei processori a esecuzione parallela.



**Figura 4.66** Unità di elaborazione dotata di parallelizzazione statica a due vie. Sono evidenziate le aggiunte necessarie per la parallelizzazione dell'esecuzione: 32 bit addizionali dalla memoria istruzioni, due porte di lettura e una di scrittura in più per il register file, e un'altra ALU. Si suppone che la ALU in basso si occupi del calcolo degli indirizzi per il trasferimento dati, mentre quella in alto si occupi di tutto il resto.

che contiene l'istruzione con la dipendenza. Sia che il software debba gestire tutti gli hazard sia che debba solo cercare di ridurre il numero di hazard tra i diversi pacchetti di istruzioni, l'idea di un'istruzione ampia che codifica più operazioni al suo interno ne esce rafforzata. Faremo riferimento a questo secondo approccio nell'esempio seguente.

Per lanciare in esecuzione una coppia di istruzioni, una che esegue un'operazione aritmetico-logica e una di trasferimento dati, il primo componente hardware aggiuntivo, oltre alla logica di rilevamento degli hazard e di messa in stall simile a quella già vista, è una coppia addizionale di porte di lettura del register file (Figura 4.66). In un singolo ciclo di clock, infatti, occorre leggere due registri per l'istruzione che opera sulla ALU e altri due per l'istruzione di trasferimento dati. Servono, inoltre, due porte di scrittura, una per l'operazione della ALU e l'altra per l'operazione di trasferimento dati dalla memoria. Dato che la ALU è impegnata nell'esecuzione dell'operazione aritmetico-logica, occorre anche un sommatore separato per il calcolo dell'indirizzo della memoria per il trasferimento dei dati. Senza queste risorse aggiuntive, il funzionamento di questa pipeline a due vie sarebbe ostacolato da hazard di tipo strutturale.

Chiaramente questo semplice processore a due vie può migliorare le prestazioni al massimo di un fattore 2. Tuttavia, con questo semplice parallelismo, si raddoppia il numero di istruzioni in esecuzione che si sovrappongono nel tempo e questo fa crescere le perdite di prestazioni dovute agli hazard sui dati e sul controllo. Per esempio, nella nostra pipeline a cinque stadi le istruzioni di load hanno una **latenza di utilizzo** di un ciclo di clock, per cui un'altra istruzione non può utilizzare il risultato della load senza uno stall. Questo significa che in una pipeline a cinque stadi e due vie il risultato di un'istruzione di load non può essere utilizzato nel ciclo di clock *successivo* e le *due* istruzioni seguenti non possono utilizzare il risultato della load senza uno stall. Inoltre, le istruzioni che utilizzano la ALU e che non hanno latenza di utilizzo nella pipeline semplice ora presentano una latenza di utilizzo di un'istruzione, poiché il loro risultato non può essere utilizzato da una eventuale load o store con cui sono accoppiate.

**Latenza di utilizzo:** numero di cicli di clock che devono intercorrere tra un'istruzione di load e un'istruzione a cui serve il dato caricato dalla memoria perché si possa evitare che la pipeline debba essere messa in stall.

In un processore a esecuzione parallela, per sfruttare il parallelismo in modo efficace sono necessari compilatori più sofisticati o tecniche di riordinamento hardware del codice. I microprocessori a parallelizzazione statica richiedono che il compilatore svolga questi compiti.

## ESEMPIO

### Riordinamento di un semplice frammento di codice per la parallelizzazione statica dell'esecuzione

Come verrebbe riorganizzato il codice del seguente ciclo in una pipeline RISC-V a due vie a esecuzione parallela?

```
Ciclo: ld x31, 0(x20)      // x31 = elemento di un vettore
       add x31, x31, x21    // somma lo scalare contenuto
                           // in x21
       sd x31, 0(x20)      // memorizza il risultato
       addi x20, x20, -8    // decrementa il puntatore
       blt x22, x20, Ciclo // confronta con il valore
                           // di fine ciclo, esci dal ciclo
                           // se x20 > x22
```

Riordinare le istruzioni in modo da evitare il numero maggiore di stalli della pipeline. Si supponga che il risultato dei salti condizionati possa essere predetto in modo che gli hazard sul controllo vengano gestiti dall'hardware.

## SOLUZIONE

Le prime tre istruzioni, come le due successive, hanno delle dipendenze sui dati. La Figura 4.67 mostra il migliore riordinamento possibile di queste istruzioni. Si noti che solo una coppia di istruzioni viene eseguita in parallelo. Ogni iterazione del ciclo richiede 4 cicli di clock; con 4 cicli di clock per 5 istruzioni si ottiene un CPI pari a 0,8, contro un CPI massimo di 0,5, ovverosia un IPC di 1,25 contro un IPC massimo di 2,0. Si tratta di valori decisamente deludenti. Si noti che nel calcolare il CPI e l'IPC le nop non vengono contate in quanto non sono istruzioni utili. Se contassimo anche le nop aumenteremmo l'IPC ma non le prestazioni!

**Espansione dei cicli:** una tecnica per ottenere prestazioni migliori sui cicli, specialmente se accedono a vettori, in cui il corpo del ciclo viene replicato più volte e le istruzioni provenienti dalle varie repliche possono essere inserite negli stessi pacchetti di istruzioni.

Una tecnica importante adottata dai compilatori per ottenere prestazioni migliori nell'esecuzione dei cicli è l'**espansione dei cicli** (*loop unrolling*). Come suggerisce il nome, con questa tecnica vengono generate più copie del corpo del ciclo e istruzioni appartenenti a iterazioni diverse possono essere avviate all'esecuzione contemporaneamente, inserendole nella stessa VLIW.

Istruzioni ALU o di salto condizionato	Istruzioni trasferimento dati	Ciclo di clock
Ciclo:	ld x31, 0(x20)	1
addi x20, x20, -8		2
add x31, x31, x21		3
blt x22, x20, Ciclo	sd x31, 8(x20)	4

**Figura 4.67** Codice riordinato come apparirebbe in un processore RISC-V con pipeline a due vie. Si noti che avendo spostato la addi prima della sd, abbiamo dovuto aggiungere 8 all'offset della sd.

## Espansione dei cicli per pipeline a esecuzione parallela

Mostrare l'efficacia della tecnica dell'espansione dei cicli nel caso dell'esempio precedente. Per semplicità, si supponga che l'indice del ciclo sia multiplo di 4.

Per riordinare le istruzioni del ciclo in modo che non si verifichino ritardi occorre fare quattro copie del corpo del ciclo. Dopo avere espanso il ciclo e avere eliminato le istruzioni di controllo, che risultano inutili, il ciclo conterrà quattro copie delle istruzioni `ld`, `add` e `sd`, e una copia delle istruzioni `addi` e `blt`. Il codice del ciclo espanso, dopo il riordinamento, è mostrato in Figura 4.68.

Durante il processo di espansione, il compilatore ha introdotto alcuni registri addizionali (`x28`, `x29`, `x30`). L'obiettivo di questa operazione, detta **ridenominazione dei registri** (*register renaming*), è di eliminare le false dipendenze tra i dati che potrebbero generare hazard o limitare il compilatore nel riordino del codice. Si pensi a come apparirebbe il ciclo espanso se si utilizzasse solamente `x31`: si produrrebbero diverse ripetizioni delle istruzioni `ld` `x31, 0(x20)`, `add x31, x31, x21`, seguite da `sd x31, 8(x20)`. Le diverse ripetizioni di queste tre istruzioni, malgrado utilizzino tutte `x31`, sono in realtà completamente indipendenti: non esiste alcun flusso di dati tra una coppia di istruzioni in un'iterazione e la stessa coppia di istruzioni in un'altra iterazione. Si parla in questo caso di **falsa dipendenza**, o **dipendenza nominale**, che implicherebbe un riordinamento delle istruzioni forzato dall'utilizzo di uno stesso nome e non da una reale dipendenza, chiamata anche *dipendenza effettiva*.

Il processo di ridenominazione dei registri durante l'espansione dei cicli permette al compilatore di inserire una dopo l'altra queste istruzioni tra loro indipendenti, in modo da ottenere una migliore riorganizzazione del codice. La ridenominazione dei registri elimina le false dipendenze ma non quelle vere.

Si noti che 12 delle 14 istruzioni del ciclo vengono eseguite in coppia. Occorrono 8 cicli di clock per eseguire 4 iterazioni del ciclo, cioè 2 cicli di clock per iterazione con un ICP pari a  $14/8 = 1,75$ . L'espansione del ciclo e la riorganizzazione del codice in coppie di istruzioni hanno prodotto un miglioramento delle prestazioni di più di un fattore 2, dovuto in parte alla riduzione del numero di istruzioni di controllo del ciclo e in parte all'esecuzione parallela di coppie di istruzioni. Il costo di questo miglioramento consiste nell'uso di quattro registri temporanei, invece che di uno solo, e in un significativo incremento nella dimensione del codice.

### ESEMPIO

### SOLUZIONE

**Ridenominazione dei registri:** è un processo compiuto dal compilatore o dall'hardware per evitare le false dipendenze.

**Falsa dipendenza:** detta anche **dipendenza nominale**, è un riordinamento delle istruzioni causato dall'utilizzo di uno stesso nome, tipicamente di un registro, e non da una vera e propria dipendenza tra i dati delle istruzioni.

## Processori dotati di parallelizzazione dinamica dell'esecuzione

I processori dotati di parallelizzazione dinamica dell'esecuzione sono anche conosciuti come processori superscalari, o semplicemente **superscalari**. Nei processori superscalari più semplici, le istruzioni vengono eseguite in ordine e il processore decide se nessuna, una o più istruzioni possono essere avviate all'esecuzione nello stesso ciclo. Ovviamente, per ottenere buone prestazioni con questo tipo di processore occorre che il compilatore riordini le istruzioni in modo da rimuovere le dipendenze per aumentare il numero di istruzioni che possono essere lanciate in parallelo. Anche con questo supporto del compilatore c'è una differenza importante tra questo semplice processore superscalare e un processore VLIW: è l'hardware che garantisce che il codice, riordinato o meno,

**Superscalare:** una tecnica di pipeline avanzata che permette al processore di eseguire più istruzioni nello stesso ciclo di clock, selezionando le istruzioni da eseguire in parallelo durante l'esecuzione stessa.

	Istruzioni ALU o di salto condizionato	Istruzioni trasferimento dati	Ciclo di clock
Ciclo:	addi x20, x20, -32	ld x28, 0(x20)	1
		ld x29, 24(x20)	2
	add x28, x28, x21	ld x30, 16(x20)	3
	add x29, x29, x21	ld x31, 8(x20)	4
	add x30, x30, x21	sd x28, 32(x20)	5
	add x31, x31, x21	sd x29, 24(x20)	6
		sd x30, 16(x20)	7
	blt x22, x20, Ciclo	sd x31, 8(x20)	8

**Figura 4.68** Il codice di Figura 4.67 dopo l'espansione del ciclo e il riordinamento del codice, come apparirebbe nell'esecuzione su un RISC-V dotato di parallelizzazione statica a due vie. Gli spazi vuoti sono occupati da istruzioni nop. Dato che x20 viene decrementato di 32 dalla prima istruzione del ciclo, gli indirizzi contenuti nel PC saranno rispettivamente il valore originale di x20, x20 meno 8, meno 16 e meno 24.

venga eseguito correttamente. Inoltre, il codice in linguaggio macchina verrà eseguito sempre in modo corretto, indipendentemente da quante istruzioni vengono eseguite in parallelo o dalla struttura della pipeline del processore. In alcuni processori VLIW ciò può non essere vero, e diventa quindi necessario ricompilare i programmi quando ci si sposta su modelli diversi dello stesso processore; in altri processori VLIW, il codice può venire eseguito correttamente anche su modelli diversi, ma con prestazioni così basse da rendere di fatto necessaria una ricompilazione del codice.

Molti processori superscalari estendono la versione base dell'unità che gestisce la parallelizzazione dell'esecuzione in modo da includere un **riordinamento dinamico delle istruzioni** (*dynamic scheduling*). Questo schema permette di decidere quali istruzioni eseguire in un dato ciclo di clock e, allo stesso tempo, di evitare hazard e stalli. Cominciamo con un semplice esempio di hazard sui dati e consideriamo il seguente frammento di codice:

```
ld  x31, 0(x21)
add x1, x31, x2
sub x23, x23, x3
andi x5, x23, 20
```

Anche se l'istruzione sub è pronta per essere eseguita, deve aspettare che l'esecuzione di ld e add sia stata completata, il che può richiedere parecchi cicli di clock se la memoria è lenta (nel Capitolo 5 parleremo delle miss delle cache, che sono il motivo per cui l'accesso alla memoria a volte può essere estremamente lento). La riorganizzazione dinamica del codice nelle pipeline permette di evitare questi hazard completamente o parzialmente.

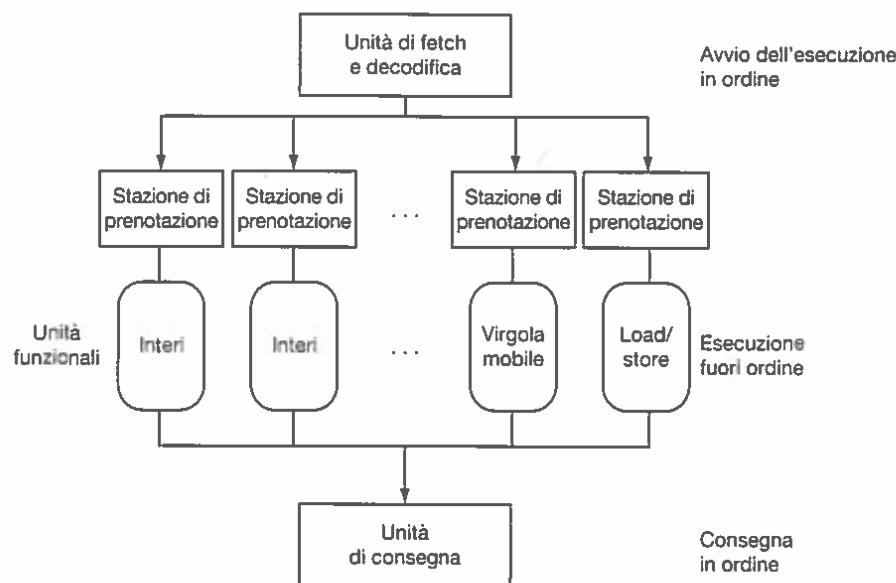
### Riorganizzazione dinamica del codice in una pipeline

Nella riorganizzazione dinamica del codice si sceglie quali istruzioni eseguire nei cicli di clock successivi, eventualmente riordinandole in modo da evitare stalli. In un processore di questo tipo la pipeline viene suddivisa in tre unità principali: un'unità che preleva l'istruzione dalla memoria e la avvia all'esecuzione, un gruppo di più unità funzionali (una decina o più nei processori di fascia alta nel 2015), e un'**unità di consegna** (*commit unit*). La Figura 4.69 mostra il modello corrispondente. La prima unità preleva le istruzioni, le decodifica e invia ciascuna di esse alla corrispondente unità funzionale per

**Riordinamento dinamico delle istruzioni:** supporto hardware al riordino delle istruzioni; viene utilizzato per evitare gli stalli.



**Unità di consegna:** in una pipeline che esegue il codice non nell'ordine scritto, è l'unità che decide quando si può consegnare il risultato di un'istruzione scrivendolo nei registri o nella memoria visibili al programmatore.



**Figura 4.69** Le tre unità fondamentali di una pipeline per il riordinamento dinamico del codice. Il passo finale di aggiornamento dello stato viene chiamato *retirement* (pensionamento) oppure *graduation* (laurea).

l'esecuzione. Ciascuna unità funzionale è dotata di buffer, chiamati **stazioni di prenotazione** (*reservation stations*), che conservano gli operandi e l'operazione. (Nella sezione *Approfondimento* esamineremo una soluzione alternativa alle stazioni di prenotazione che viene impiegata in molti processori recenti.) Non appena la stazione di prenotazione contiene tutti gli operandi richiesti e l'unità funzionale è pronta per l'esecuzione, viene eseguita l'operazione e inviato il risultato alle stazioni di prenotazione che lo stanno attendendo e contemporaneamente all'unità di consegna. Questa conserva il risultato dell'operazione fino a quando non avrà il via libera per salvarlo nel register file o, nel caso di un'istruzione di store, nella memoria. Il buffer presente nell'unità di consegna è chiamato **buffer di riordino** (*reorder buffer*) e viene anche utilizzato per fornire gli operandi, un po' come faceva la logica di propagazione nelle pipeline con smistamento statico delle istruzioni. Una volta scritto il risultato nel register file, esso può essere prelevato da lì come in una normale pipeline.

La combinazione del salvataggio degli operandi nei buffer delle stazioni di prenotazione e dei risultati nel buffer di riordino consente una forma di ride-nominazione dei registri del tutto analoga a quella utilizzata dal compilatore nel caso dell'esempio precedente. Per comprenderne meglio il funzionamento da un punto di vista concettuale, consideriamo i seguenti passi.

1. Quando inizia l'esecuzione di un'istruzione, essa viene copiata nella stazione di prenotazione associata all'unità funzionale appropriata e i suoi operandi disponibili copiati dal register file o dal buffer di riordino. L'istruzione rimane all'interno della stazione di prenotazione finché tutti gli operandi non saranno disponibili e l'unità funzionale non sarà pronta. Non è più necessario copiare dai registri gli operandi per l'istruzione che viene avviata all'esecuzione; inoltre, se occorre scrivere i registri che contengono gli operandi, il loro contenuto può essere tranquillamente sovrascritto.
2. Se un operando non si trova nel register file o nel buffer di riordino, l'istruzione deve aspettare che esso sia prodotto da una delle unità funzionali. In questo caso, il lavoro di quell'unità funzionale viene seguito e il risultato prodotto viene copiato direttamente dall'unità funzionale nella stazione di prenotazione, senza aspettare che venga scritto nel registro di destinazione.

**Stazione di prenotazione:** buffer di un'unità funzionale che conserva gli operandi e l'operazione.

**Buffer di riordino:** il buffer che conserva i risultati prodotti da un processore con riordinamento dinamico del codice finché non è certo che essi possano essere scritti in memoria o nei registri.

**Esecuzione fuori ordine:** la situazione che si verifica in una pipeline quando l'esecuzione di un'istruzione risulta bloccata e le altre istruzioni non dipendenti possono continuare a essere eseguite.

**Completamento in ordine:** una modalità di consegna dei risultati in cui i risultati dell'esecuzione in una pipeline vengono scritti nei registri visibili al programmatore nello stesso ordine con cui le istruzioni associate erano state prelevate.

In questi due passaggi si utilizzano le stazioni di prenotazione e il buffer di riordino per implementare la ridenominazione dei registri.

Concettualmente, possiamo pensare a una pipeline a riordinamento dinamico del codice come a un'entità che analizza la struttura del flusso dei dati di un programma e poi esegue le istruzioni in un ordine che può anche essere diverso, preservando però il flusso dei dati prodotti. Questo stile di esecuzione è chiamato **esecuzione fuori ordine** (*out-of-order execution*), poiché le istruzioni possono essere eseguite in un ordine diverso da quello con cui erano state prelevate dalla memoria istruzioni.

Per fare in modo che i programmi si comportino come se le istruzioni fossero eseguite in ordine, l'unità di fetch e decodifica invia in esecuzione le istruzioni nella sequenza originale, cosa che consente di identificare le dipendenze, e quella di consegna dovrà scrivere i risultati nei registri o in memoria seguendo l'ordine con cui il programma ha prelevato le istruzioni dalla memoria. Questo modello conservativo di esecuzione è detto **completamento in ordine** (*in order commit*). Quindi, nel caso in cui si verifichi un'eccezione, il calcolatore può fare riferimento all'ultima istruzione eseguita, e gli unici registri aggiornati saranno solo quelli scritti dalle istruzioni che precedono l'istruzione che ha generato l'eccezione. Sebbene la parte iniziale della pipeline (fetch e avvio all'esecuzione) e la parte terminale (consegna) eseguano il codice in ordine, le unità funzionali sono invece libere di iniziare e terminare l'esecuzione non appena i dati di cui hanno bisogno si rendono disponibili. Oggi tutte le pipeline dotate di riordinamento dinamico del codice utilizzano il completamento in ordine.

Il riordinamento dinamico del codice viene spesso esteso includendo la speculazione hardware, specialmente per i salti condizionati. Prevedendo la direzione di un salto, un processore con riorganizzazione dinamica del codice può prelevare ed eseguire le istruzioni lungo il ramo di esecuzione scelto. Poiché le istruzioni vengono consegnate in ordine, possiamo capire se un salto è stato previsto correttamente prima che venga consegnato il risultato di una qualsiasi istruzione del ramo di esecuzione scelto. Una pipeline speculativa dotata di riordinamento dinamico del codice può anche supportare la speculazione sugli indirizzi delle load, permettendo il riordino delle istruzioni di load-store; l'unità di consegna verrà utilizzata per evitare di scrivere i risultati prodotti dalle istruzioni eseguite quando la speculazione si rivela errata. Nel prossimo paragrafo esamineremo in che modo si utilizza la riorganizzazione dinamica del codice unita alla speculazione nell'architettura Core i7 di Intel.

## Capire le prestazioni dei programmi



GERARCHIA



PREDIZIONE

Dato che i compilatori possono riordinare il codice anche quando ci sono dipendenze tra i dati, potreste chiedervi perché un processore superscalare debba utilizzare il riordinamento dinamico del codice. I motivi sono principalmente tre. In primo luogo, non tutti gli stalli possono essere predetti. In particolare, le miss delle cache (vedi Cap. 5) nella gerarchia delle memorie possono causare degli stalli non prevedibili. La riorganizzazione dinamica del codice consente al processore di nascondere alcuni di questi stalli continuando a eseguire altre istruzioni in attesa che lo stallo dovuto alla lettura abbia termine.

Il secondo motivo è che se il processore speculasse sulla direzione dei salti condizionati utilizzando la **predizione dinamica** dei salti, non potrebbe comunque sapere qual è l'ordine corretto di esecuzione delle istruzioni all'atto della compilazione, poiché questo dipende dal comportamento predetto e da quello effettivo del salto, durante l'esecuzione. Incorporare la speculazione dinamica per sfruttare un maggiore *parallelismo a livello di istruzioni* (ILP) senza incorporare il riordino dinamico del codice ridurrebbe grandemente i benefici della speculazione.

Infine, poiché la latenza della pipeline e il numero di istruzioni che vengono lanciate in esecuzione in parallelo possono variare molto da un processore all'al-

tro, il modo migliore per compilare una sequenza di istruzioni può cambiare. Per esempio, il modo in cui viene ordinata una sequenza di istruzioni dipendenti è influenzato sia dal numero di istruzioni lanciate in esecuzione in parallelo sia dalla loro latenza. La struttura della pipeline influenza sia il numero di volte che un ciclo deve essere espanso per evitare stalli, sia il processo di ridenominazione dei registri da parte del compilatore. Il riordinamento dinamico del codice consente all'hardware di nascondere la maggior parte di questi dettagli, per cui gli utenti e i distributori di software non devono preoccuparsi di avere versioni diverse dello stesso programma per le differenti implementazioni hardware dello stesso insieme di istruzioni. Analogamente, i vecchi programmi godranno dei miglioramenti delle nuove architetture senza che debba essere ricompilato il codice.

## QUADRO D'INSIEME

Sia la **pipeline** sia l'esecuzione parallela incrementano il throughput di picco delle istruzioni e tentano di sfruttare il **parallelismo** a livello di istruzioni (ILP). Le dipendenze tra i dati e sul controllo nei programmi pongono, tuttavia, un limite superiore a prestazioni elevate, poiché il processore a volte deve aspettare che le dipendenze vengano risolte. Gli approcci basati sul software si affidano all'abilità del compilatore nel trovare e ridurre gli effetti di queste dipendenze, mentre gli approcci basati sull'hardware si affidano a estensioni della pipeline e a meccanismi hardware di parallelizzazione. La speculazione, attuata dal compilatore o dall'hardware, può incrementare, attraverso la **predizione**, la parte di ILP che può essere sfruttata, benché si debba fare attenzione alle speculazioni sbagliate che hanno conseguenze assai dannose sulle prestazioni. ■



I microprocessori moderni ad alte prestazioni sono capaci di lanciare in esecuzione molte istruzioni a ogni ciclo di clock; sfortunatamente, mantenere costante questo tasso di riempimento dei pacchetti di istruzioni è molto difficile. Per esempio, malgrado esistano processori con un numero di istruzioni per pacchetto di esecuzione che va da quattro a sei, ben poche applicazioni riescono a sostenere più di due istruzioni per ciclo di clock. I motivi sono principalmente due.

In primo luogo, nelle pipeline il collo di bottiglia principale è rappresentato dalle dipendenze che non possono essere eliminate; queste riducono il parallelismo tra le istruzioni e quindi il tasso di riempimento dei pacchetti di istruzioni. Anche se si può fare poco contro le dipendenze vere, spesso il compilatore o l'hardware non sono in grado di stabilire con precisione se una dipendenza sia reale o meno, e quindi devono ipotizzare, in modo conservativo, che la dipendenza ci sia. Per esempio, il codice basato sui puntatori, soprattutto quando essi puntano ad aree di memoria parzialmente sovrapposte, porta a dipendenze più complesse da analizzare. Al contrario, la maggiore regolarità dell'accesso agli elementi dei vettori consente spesso al compilatore di capire quando non esistono dipendenze reali. In maniera simile, i salti condizionati che non possono essere predetti accuratamente durante l'esecuzione, o durante la compilazione, limitano fortemente la possibilità di sfruttare l'ILP. Spesso sarebbe possibile utilizzare un ILP maggiore, ma l'abilità del compilatore o dell'hardware nel trovare istruzioni che possono essere molto distanti, a volte anche migliaia di istruzioni, è limitata.

In secondo luogo, le perdite che scaturiscono dalla **gerarchia delle memorie** (argomento del Capitolo 5) limitano anch'esse la capacità di mantenere la pipeline sempre piena. In alcuni sistemi di memoria gli stalli possono essere nascosti, ma un ILP minore limita anche la capacità di nascondere gli stalli.

## Interfaccia hardware/software



## Efficienza energetica e pipeline avanzate

L'incremento del parallelismo a livello di istruzioni mediante parallelizzazione dinamica del codice ha come rovescio della medaglia il potenziale abbassamento dell'efficienza energetica. Ciascuna innovazione discussa in questo paragrafo ha sfruttato il crescente numero di transistor per incrementare le prestazioni, ma spesso ciò è avvenuto in modo scarsamente efficiente. Ora che la barriera dell'energia è stata raggiunta, si cominciano a produrre architetture con più processori sullo stesso chip, nelle quali ciascun processore ha una pipeline meno profonda ed è dotato di una speculazione meno aggressiva rispetto ai precedenti.

Nonostante i processori più semplici non siano tanto veloci quanto i loro sofisticati fratelli, infatti, essi hanno prestazioni migliori per joule e possono quindi garantire prestazioni più elevate per chip quando i vincoli sul progetto sono dettati più dall'energia assorbita che dal numero di transistor.

La Figura 4.70 mostra il numero di stadi di pipeline, l'ampiezza del pacchetto di istruzioni, il livello di speculazione, la frequenza di clock, il numero di core per singolo chip e la potenza assorbita di alcuni processori vecchi e recenti. Si noti la diminuzione del numero di stadi di pipeline e della potenza dissipata con l'avvento dei microprocessori multicore.

**Approfondimento.** Un'unità di consegna controlla l'aggiornamento del register file e della memoria. Alcuni processori dotati di riordinamento dinamico del codice aggiornano il register file immediatamente durante l'esecuzione, utilizzando dei registri aggiuntivi per implementare la ridenominazione e preservare il vecchio contenuto di questi registri fintanto che l'istruzione che ha scritto il register file non viene considerata corretta. Altri processori, invece, memorizzano il risultato in buffer dedicati, che, come abbiamo visto, hanno tipicamente la struttura di un buffer di riordino, e l'aggiornamento effettivo del register file avviene solamente più tardi, durante la fase di consegna del risultato. In questo caso i dati scritti in memoria dalle store devono essere contenuti nel buffer fino al momento della loro consegna; questo buffer può essere un *buffer di scrittura* (vedi Cap. 5) oppure il buffer di riordino. L'unità di consegna consente all'istruzione di store di scrivere in memoria dal buffer non appena al buffer viene associato un indirizzo valido e contiene dati validi, e le dipendenze dai salti condizionati sono state risolte.

**Approfondimento.** Gli accessi a memoria sono avvantaggiati dall'utilizzo di cache *non bloccanti*, che continuano a gestire gli accessi alla cache anche durante la gestione di una miss della cache (vedi Cap. 5). I processori con esecuzione fuori ordine hanno bisogno di circuiti aggiuntivi della cache per consentire di eseguire altre istruzioni durante una miss.

Micropocessore	Anno	Frequenza del clock	Stadi di pipeline	Aampiezza del pacchetto	Esecuzione fuori ordine/Speculazione	Core per chip	Potenza assorbita
Intel 486	1989	25 MHz	5	1	No	1	5 W
Intel Pentium	1993	66 MHz	5	2	No	1	10 W
Intel Pentium Pro	1997	200 MHz	10	3	Sì	1	29 W
Intel Pentium 4 Willamette	2001	2000 MHz	22	3	Sì	1	75 W
Intel Pentium 4 Prescott	2004	3600 MHz	31	3	Sì	1	103 W
Intel Core	2006	2930 MHz	14	4	Sì	2	75 W
Intel Core i5 Nehalem	2010	3300 MHz	14	4	Sì	2-4	87 W
Intel Core i7 Ivy Bridge	2012	3400 MHz	14	4	Sì	8	77 W

**Figura 4.70** Complessità della pipeline, numero di core e potenza dissipata da alcuni microprocessori Intel. Gli stadi della pipeline del Pentium 4 non comprendono gli stadi dell'unità di consegna. Se includessimo anche questi, la profondità della pipeline del Pentium 4 sarebbe ancora maggiore.

## Autovalutazione

Stabilire se le seguenti tecniche e i seguenti componenti sono associati principalmente all'approccio software o hardware di sfruttamento dell'ILP. In alcuni casi la risposta corretta può essere entrambi gli approcci.

1. Predizione dei salti
2. Parallelizzazione dell'esecuzione
3. VLIW
4. Superscalare
5. Parallelizzazione dinamica
6. Esecuzione fuori ordine
7. Speculazione
8. Buffer di riordino
9. Ridenominazione dei registri

## 4.11 Un caso reale: la pipeline del Cortex-A53 ARM e del Core i7 Intel

La Figura 4.71 descrive i due processori analizzati in questo paragrafo, che sono rivolti ai due mercati di riferimento dell'era post-PC.

### Cortex-A53 ARM

Il Cortex-A53 ARM lavora a una frequenza di 1,5 GHz, ha una pipeline di 8 stadi ed esegue l'insieme di istruzioni ARM v8. La sua pipeline utilizza la parallelizzazione dinamica del codice e può lanciare in esecuzione due istruzioni per ciclo di clock. È una pipeline statica e prevede un'esecuzione "in ordine", per cui

Processore	ARM A53	Core i7 920 Intel
Mercato	Dispositivi Mobili Personali	Server, Cloud
Potenza termica di progetto	100 milliWatt (1 core @ 1 GHz)	130 watt
Frequenza di clock	1,5 GHz	2,66 GHz
Core/chip	4 (configurabile)	4
Virgola mobile?	Sì	Sì
Parallelizzazione dell'esecuzione?	Dinamica	Dinamica
Istruzioni/ciclo di clock (picco)	2	4
Stadi della pipeline	8	14
Scheduling della pipeline	Statico in ordine	Dinamico fuori ordine, con speculazione
Predizione dei salti	Ibrida	A 2 livelli
Cache 1° livello/core	16-64 KiB I, 16-64 KiB D	32 KiB I, 32 KiB D
Cache 2° livello/core	128-2048 KiB (condivisa)	256 KiB (per core)
Cache 3° livello (condivisa)	Dipendente dalla piattaforma	2-8 MiB

Figura 4.71 Specifiche del Cortex-A53 ARM e dell'Intel Core i7 920.

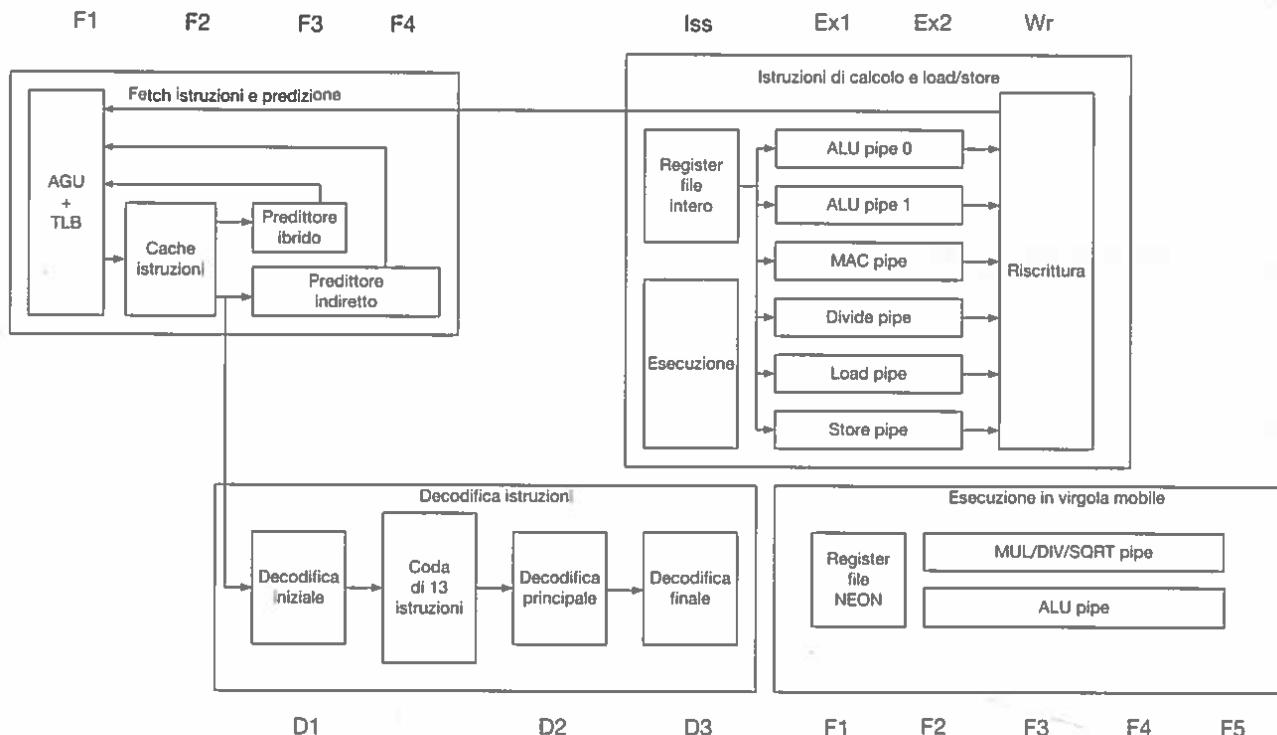
le istruzioni vengono lanciate in esecuzione, eseguite e il risultato consegnato in ordine. La pipeline può essere suddivisa in tre sezioni: fetch, decodifica delle istruzioni ed esecuzione, come mostrato in Figura 4.72.

I primi tre stadi carico dalla memoria due istruzioni alla volta e cercano di mantenere piena la coda di 13 istruzioni. Utilizzano un predittore ibrido di 6k-bit dei salti condizionati, un predittore dei salti indiretti di 256 elementi e uno stack degli indirizzi di ritorno di 8 elementi per predire il ritorno dalle funzioni. La predizione dei salti indiretti richiede un ulteriore stadio della pipeline. Questa scelta progettuale implica l'introduzione di una latenza aggiuntiva se la coda delle istruzioni non può disaccoppiare la fase di fetch da quella di decodifica ed esecuzione, principalmente nel caso del fallimento della predizione di un salto condizionato o di una miss della cache istruzioni. Quando il predittore dei salti condizionati sbaglia, la pipeline viene svuotata e questo risulta in una penalità associata a un errore di predizione di otto cicli di clock.

Gli stadi di codifica della pipeline determinano se ci sono dipendenza tra le coppie di istruzioni, che forzerebbero un'esecuzione sequenziale, e a quale cammino di esecuzione inviare l'istruzione.

La fase di esecuzione occupa principalmente tre stadi di pipeline e fornisce un cammino di esecuzione per le istruzioni di load, una pipeline per le istruzioni di store, due pipeline per le operazioni aritmetiche su interi e pipeline separate per le operazioni di moltiplicazione e divisione. Ciascuna istruzione di una coppia può essere inviata sulle pipeline di load o store. Gli stadi di esecuzione hanno una propagazione completa tra i diversi stadi di pipeline.

Le operazioni in virgola mobile e SIMD aggiungono altri due stadi di pipeline alla fase di esecuzione delle istruzioni e sfruttano una pipeline per le operazioni di moltiplicazione/divisione/radice quadrata e una pipeline per le altre operazioni aritmetiche.



**Figura 4.72 Pipeline del Cortex-A53.** I primi tre stadi leggono le istruzioni della memoria e le inseriscono in una coda di 13 elementi. L'unità di generazione degli indirizzi (AGU, Address Generation Unit) utilizza un predittore ibrido, un predittore indiretto e uno stack degli indirizzi di ritorno per predire i salti condizionati e cercare di mantenere piena la coda delle istruzioni. La decodifica delle istruzioni richiede tre stadi, così come l'esecuzione. Altri due stadi sono riservati alle operazioni in virgola mobile e alle istruzioni SIMD.

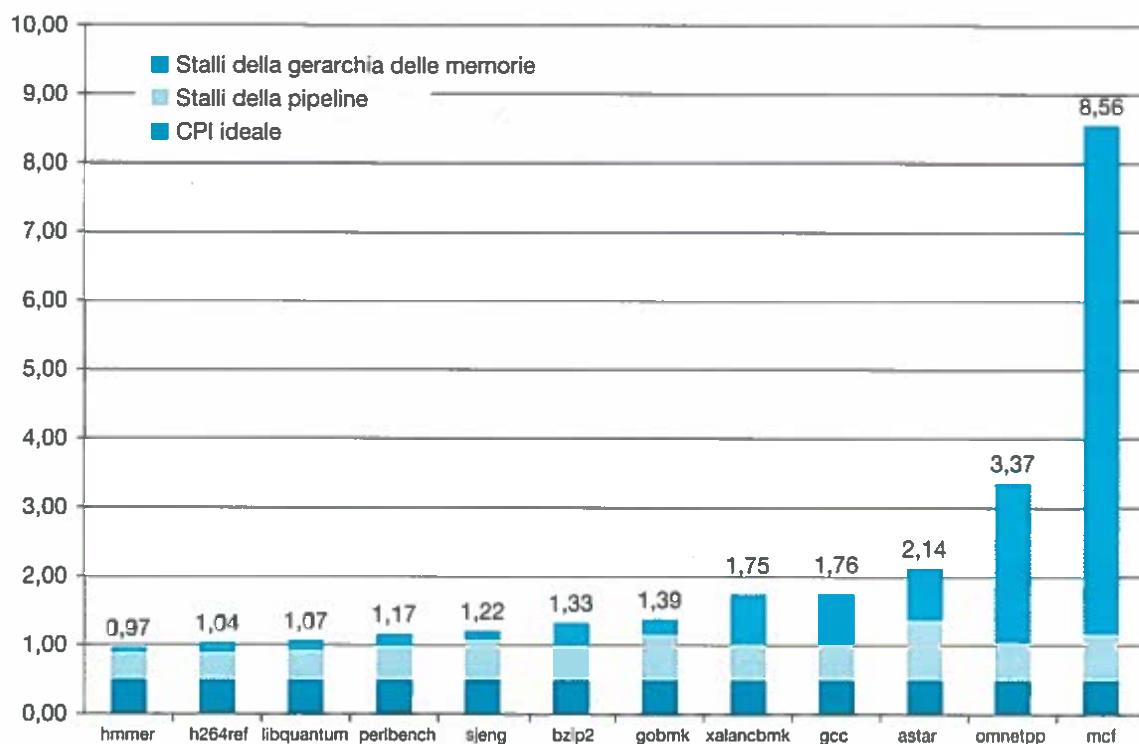


Figura 4.73 CPI del Cortex-A53 ARM per i benchmark Interi SPEC2006.

La Figura 4.73 mostra il CPI del Cortex-A53 su programmi dei benchmark SPEC2006. Anche se il CPI ideale è 0,5, il valore migliore di CPI ottenuto è 1,0, il valore mediano è 1,3 e il valore peggiore 8,6. Per il caso mediano il 60% degli stalli è dovuto a hazard della pipeline e il 40% alla gerarchia delle memorie. Gli stalli della pipeline sono causati da predizioni errate sui salti condizionati, da hazard strutturali e da dipendenze tra i dati di coppie di istruzioni. Data la natura statica della pipeline del Cortex-A53, è compito del compilatore cercare di evitare gli hazard strutturali e le dipendenze tra i dati.

**Approfondimento.** Il Cortex-A53 è un core configurabile che supporta l'architettura dell'insieme di istruzioni ARMv8. Viene distribuito come un core IP (*Intellectual Property* – proprietà intellettuale). I core IP sono la forma principale con cui viene distribuita la tecnologia nel mercato dei dispositivi embedded e mobili e nei mercati associati: miliardi di processori ARM e MIPS sono stati creati a partire da questi core IP.

Osserviamo che i core IP sono diversi dai core alla base dei microprocessori multicore i7 di Intel. Un core IP (che potrebbe anche essere un multicore) è progettato per essere integrato con altri circuiti per produrre un dispositivo (di cui risulta il nucleo) che possa essere utilizzato in processori dedicati ad applicazioni specifiche quali gli encoder e decoder video, le interfacce di I/O e le interfacce della memoria. Questi dispositivi vengono ottimizzati per le funzioni per cui sono sviluppati. Sebbene i processori core siano identici, i dispositivi risultanti sono molto diversi tra loro. Un esempio è la dimensione dalla cache L2, che può variare di un fattore 16.

## Core i7 920 di Intel

I microprocessori x86 utilizzano approcci sofisticati che fanno uso per la loro pipeline a 14 stadi della parallelizzazione dinamica dell'esecuzione, del riordinamento dinamico del codice con esecuzione fuori ordine e della speculazione.

Tuttavia, questi processori devono comunque affrontare la sfida posta dalla gestione del complesso insieme di istruzioni dell'x86 descritto nel Capitolo 2. Gli Intel prelevano dalla memoria le istruzioni x86 e le traducono in istruzioni interne simili alle istruzioni RISC-V; queste istruzioni interne vengono chiamate da Intel *micro-operazioni*. Le micro-operazioni vengono quindi eseguite da una pipeline sofisticata, dotata di riordinamento dinamico del codice e di speculazione, capace di sostenere l'esecuzione fino a sei micro-operazioni per ciclo di clock. Questo paragrafo approfondisce il funzionamento di questa pipeline.

**Microarchitettura:** l'organizzazione interna del processore, comprendente le unità funzionali principali, le interconnessioni e l'unità di controllo.

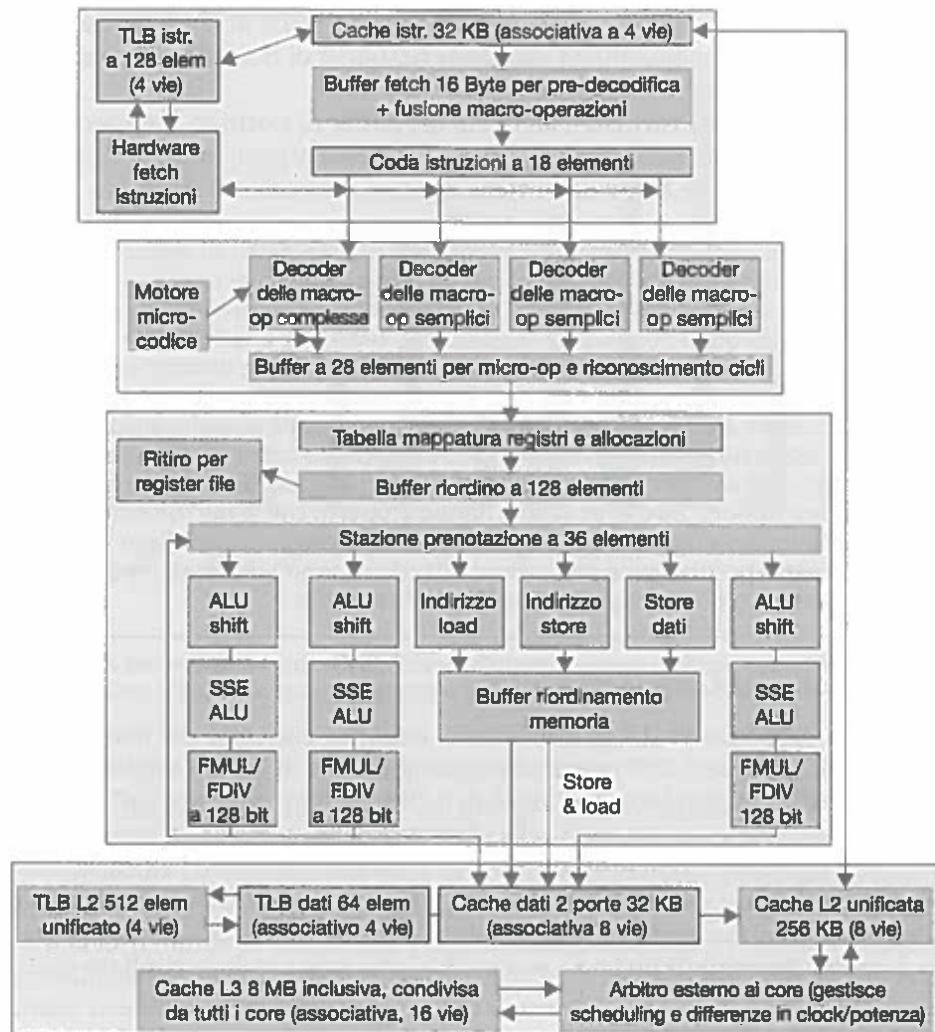
**Registri dell'architettura:** l'insieme dei registri del processore visibili dall'insieme delle istruzioni. Per esempio, nel RISC-V i registri dell'architettura sono i 32 registri interi del register file e i 32 registri in virgola mobile.

Quando si analizza un processore complesso, dotato di riordinamento dinamico del codice, la progettazione delle unità funzionali, della cache e del register file, dell'avvio a esecuzione delle istruzioni e del controllo di tutta la pipeline non può essere fatta separatamente; ciò rende difficile isolare l'elaborazione dei dati dal resto della pipeline. Per questo motivo, molti progettisti e ricercatori hanno adottato il termine **microarchitettura** per indicare l'architettura interna di un processore.

Il Core i7 Intel utilizza uno schema basato sul buffer di riordino e la ridenominazione dei registri per risolvere le false dipendenze e le speculazioni errate. Per operare la ridenominazione dei registri, il processore rinomina esplicitamente i **registri dell'architettura** (16 nel caso della versione a 64 bit dell'architettura x86) nei registri presenti fisicamente nell'architettura, che sono in numero maggiore, in modo da eliminare le false dipendenze. Il processore mantiene una mappa tra i registri dell'architettura e i registri fisici, che indica quale registro fisico contiene la copia più recente di un registro dell'architettura. Tenendo traccia delle ridenominazioni, il processore ha a disposizione un altro approccio per correggere gli errori nelle speculazioni: è sufficiente cancellare le mappature generate a partire dalla prima istruzione sulla quale la speculazione ha fallito. Questo consente al processore di ripristinare lo stato in cui si trovava quando era in esecuzione l'ultima istruzione eseguita correttamente, mantenendo la mappatura corretta tra registri fisici e registri dell'architettura.

La Figura 4.74 mostra l'organizzazione complessiva della pipeline del Core i7. Riportiamo ora gli otto passi richiesti dall'esecuzione di un'istruzione x86.

1. Fetch istruzioni. Il processore utilizza un buffer degli indirizzi di salto a più livelli per ottenere un buon bilanciamento tra velocità e accuratezza delle predizioni. È presente anche uno stack degli indirizzi di ritorno per rendere più veloce il ritorno dalle funzioni. Le predizioni errate provocano una penalità di circa 15 cicli di clock. L'unità di fetch legge dalla cache delle istruzioni 16 byte alla volta, a partire dall'indirizzo fornito dal circuito di predizione.
2. Questi 16 byte vengono inseriti nel buffer di predecodifica delle istruzioni. Lo stadio di predecodifica trasforma i 16 byte in singole istruzioni x86. Questa operazione non è banale poiché la lunghezza delle istruzioni x86 è variabile da 1 a 15 byte e l'unità di predecodifica deve analizzare il contenuto di diversi campi prima di riuscire a determinare la lunghezza di ciascuna istruzione. Le singole istruzioni x86 vengono inserite nella coda delle istruzioni a 18 elementi.
3. Decodifica micro-op. Le singole istruzioni x86 vengono convertite in micro-operazioni (micro-op). Tre decodificatori gestiscono la traduzione delle istruzioni che vengono convertite in un'unica micro-operazione. Per le istruzioni x86 che hanno una semantica più complessa, c'è un motore che converte queste istruzioni in una sequenza di micro-operazioni: produce fino a quattro micro-operazioni per ogni ciclo di clock e procede fino a quando non ha completato la generazione della sequenza di microistruzioni. Le micro-operazioni vengono inserite, secondo l'ordine delle istruzioni x86, nel buffer delle micro-operazioni a 28 elementi.



**Figura 4.74** Pipeline del Core i7 con i componenti della memoria. La profondità totale della pipeline è di 14 stadi, con un costo degli errori di predizione di 17 cicli di clock. Questo processore può inserire nei buffer 48 load e 32 store. I sei cammini di esecuzione indipendenti possono iniziare l'esecuzione di un'operazione RISC pronta, a ogni ciclo di clock.

4. Il buffer delle micro-operazioni effettua il *riconoscimento dei cicli*. Se c'è una piccola sequenza di istruzioni (meno di 28 istruzioni e lunga meno di 256 byte) che contiene un ciclo, il circuito di riconoscimento dei cicli identifica il ciclo e lancia in esecuzione le micro-operazioni direttamente dal buffer, eliminando così la necessità di attivare gli stadi di fetch e decodifica per le istruzioni del ciclo nelle iterazioni successive.
5. Avvio a esecuzione standard. Viene identificata la posizione dei registri nella tabella dei registri, vengono rinominati i registri e allocata una posizione del buffer di riordino, e vengono eventualmente trasferiti dati dai registri o dal buffer di riordino, prima di inviare la micro-operazione alla stazione di prenotazione.
6. L'i7 utilizza una stazione di prenotazione centralizzata a 36 elementi condivisa da sei unità funzionali. Fino a sei micro-operazioni possono essere inviate a esecuzione alla sei unità funzionali a ogni ciclo di clock.
7. Le singole unità funzionali eseguono le micro-operazioni e inviano il risultato alla stazione di prenotazione, se in attesa, e all'unità di consegna dei

dati per il register file, non appena l'istruzione diventa non più speculativa. L'elemento corrispondente all'istruzione nel buffer di riordino viene marcato come completato.

8. Quando una o più istruzioni all'inizio del buffer di riordino vengono marcate come completate, le scritture associate pendenti nell'unità di consegna vengono eseguite e l'istruzione viene rimossa dal buffer di riordino.

**Approfondimento.** L'hardware del secondo e quarto passo di esecuzione può combinare o fondere delle operazioni per ridurre il numero di operazioni da eseguire. La *fusione delle macro-operazioni* nel secondo passo prende combinazioni di istruzioni x86, quali "compare" seguita da "branch", e le fonde in un'unica operazione. La *microfusion* nel quarto passo combina coppie di micro-operazioni quali una load e un'operazione sulla ALU e un'operazione sulla ALU e una store e le invia a esecuzione alla stessa coda di esecuzione (nella quale possono ancora essere lanciate in esecuzione separatamente), incrementando così l'utilizzo del buffer. In uno studio sull'architettura Core di Intel che comprendeva anche l'analisi delle micro- e macrofusioni, Bird et al. [2007] hanno scoperto che le microfusioni hanno un piccolo impatto sulle prestazioni, mentre le macrofusioni sembrano avere un discreto impatto positivo sulle prestazioni quando vengono utilizzati numeri interi ma un impatto piccolo sulle prestazioni in virgola mobile.

### Prestazioni del Core i7 920 Intel

La Figura 4.75 mostra il CPI del Core i7 Intel per ciascuno dei benchmark SPEC2006. Sebbene il CPI massimo raggiungibile sia di 0,25, il miglior valore ottenuto è di 0,44, il valore mediano è di 0,79 e il valore peggiore 2,67.

È difficile distinguere tra stalli della pipeline e della memoria, in una pipeline dinamica con esecuzione fuori ordine, ma possiamo mostrare l'efficacia della predizione dei salti e della speculazione. La Figura 4.76 mostra la percentuale di errori di predizione sui salti e la percentuale di lavoro (misurata come numero di micro-operazioni inviate a esecuzione) che non produce risultato (cioè,

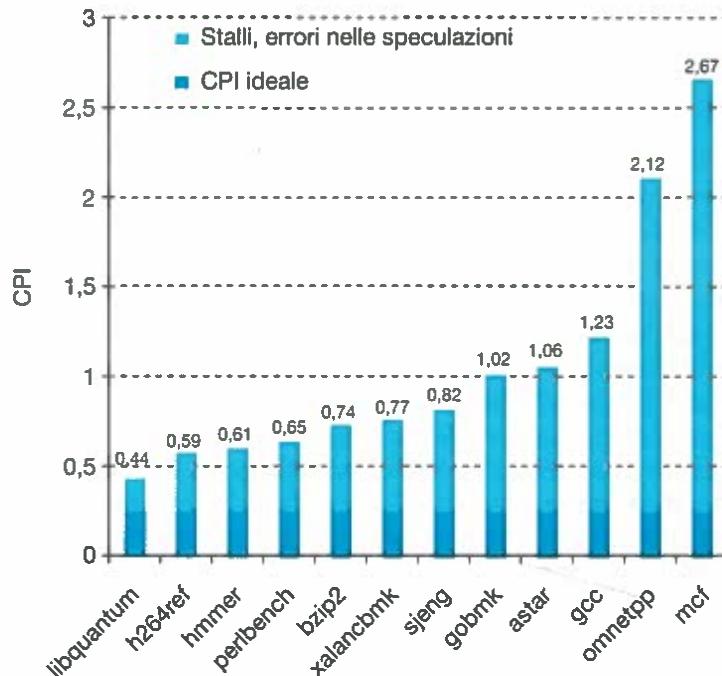
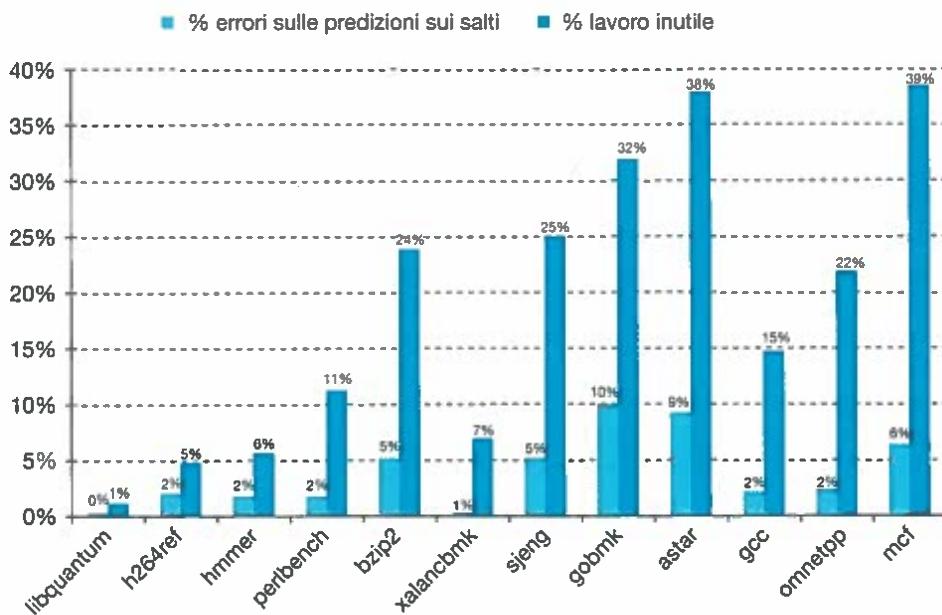


Figura 4.75 CPI del Core i7 920 Intel sui benchmark interi dello SPEC2006.



**Figura 4.76** Percentuale di errori nella predizione delle branch e lavoro inutile causato dalle speculazioni errate sul Core i7 920 Intel durante l'esecuzione dei benchmark interi dello SPEC2006.

il cui risultato è stato annullato) rispetto al numero totale di microistruzioni inviate a esecuzione. La percentuale minima, media e massima degli errori di predizione sui salti è di 0%, 2% e 10%, con una percentuale di lavoro inutile pari a 1%, 18% e 39% rispettivamente.

Il lavoro inutile in alcuni casi è legato strettamente agli errori di predizione sui salti, come per esempio nei benchmark "gobmk" e "astar", per altri benchmark, come "mcf" è molto maggiore. Questa differenza è spiegabile con il comportamento della memoria: con l'alta percentuale di miss della cache generata da "mcf", durante una speculazione sbagliata su una load, il processore invia a esecuzione tante istruzioni, fino a quando le stazioni di prenotazione riescono a contenere i dati che devono essere caricati dalle load. Se tra le istruzioni eseguite in modo speculativo c'è una branch e viene sbagliata la predizione sul salto, le micro-operazioni associate a tutte queste istruzioni devono essere eliminate.

Il Core i7 Intel combina una pipeline a 14 stadi con una parallelizzazione aggressiva dell'esecuzione per ottenere prestazioni elevate. Mantenendo bassa la latenza per le operazioni che vengono eseguite una di seguito all'altra, l'impatto della dipendenza tra i dati viene ridotto. Quali sono i maggiori colli di bottiglia per le prestazioni dei programmi eseguiti su questo processore? L'elenco che segue comprende alcune delle problematiche che hanno un potenziale impatto sulle prestazioni, di cui le ultime tre riguardano, in forme diverse, un qualsiasi processore ad alte prestazioni dotato di pipeline.

- L'utilizzo delle istruzioni x86 che non si possono trasformare in poche semplici micro-operazioni.
- I salti condizionati che sono difficili da predire; questi causano stalli quando la predizione si rivela errata, con conseguente svuotamento e riavvio della pipeline.
- Lunghe dipendenze. Queste sono tipicamente causate da istruzioni con elevato tempo di esecuzione o dalla gerarchia delle memorie e portano a stalli.
- Ritardi nell'esecuzione dovuti all'accesso alla memoria (vedi Cap. 5) che causano lo stallo del processore.

## Capire le prestazioni dei programmi



## 4.12 Come andare più veloci: parallelismo a livello di istruzioni e moltiplicazione di matrici

Ritorniamo all'esempio DGEMM del Capitolo 3 ed esaminiamo l'impatto del parallelismo a livello di istruzioni, espandendo il ciclo così che il processore a parallelizzazione dinamica con esecuzione fuori ordine abbia più istruzioni sulle quali lavorare. La Figura 4.77 mostra l'espansione del ciclo di Figura 3.22, che contiene le funzioni intrinseche del C per generare istruzioni AVX.

Espandiamo il ciclo quattro volte come abbiamo fatto nell'esempio di Figura 4.68. Invece di espandere manualmente il ciclo in C copiando quattro volte ciascuna delle istruzioni intrinseche di Figura 3.22, possiamo sfruttare il compilatore gcc per espandere il ciclo utilizzando l'ottimizzazione -O3. Utilizziamo la costante `ESPANDI` del codice C per controllare il numero di volte per cui il ciclo deve essere eseguito. Circondiamo ciascuna istruzione intrinseca con un semplice ciclo `for` che fa 4 iterazioni (linee 9, 15 e 20) e sostituiamo lo scalare `C0` di Figura 3.22 con un vettore `c[ ]` di 4 elementi (linee 8, 10, 16 e 21).

Il codice assembler associato al codice C nel quale i cicli sono stati espansi è mostrato in Figura 4.78. Come ci si poteva immaginare, in Figura 4.78 ci sono 4 copie di ciascuna delle istruzioni AVX di Figura 3.23, con un'eccezione: abbiamo bisogno di una sola copia dell'istruzione `vbroadcastsd`, dato che possiamo utilizzare le 4 copie dell'elemento di `B`, contenute nel registro `$ymm0`, più volte all'interno del ciclo. Perciò, le 5 istruzioni AVX di Figura 3.22 diventano 17 in Figura 4.78, e le sette istruzioni intere appaiono in entrambi i programmi, sebbene cambino le costanti e gli indirizzi per tenere conto dell'espansione del ciclo. Perciò, nonostante il ciclo sia stato espanso quattro volte, il numero di istruzioni è solamente raddoppiato passando da 12 a 24.

```

1 //include <x86intrin.h>
2 #define UNROLL (4)
3
4 void dgemm (int n, double* A, double* B, double* C)
5 {
6     for ( int i = 0; i < n; i+=UNROLL*4 )
7         for ( int j = 0; j < n; j++ ) {
8             __m256d c[4];
9             for ( int x = 0; x < UNROLL; x++ )
10                 c[x] = _mm256_load_pd(C+i+x*4+j*n);
11
12            for( int k = 0; k < n; k++ )
13            {
14                __m256d b = _mm256_broadcast_sd(B+k+j*n);
15                for (int x = 0; x < UNROLL; x++)
16                    c[x] = _mm256_add_pd(c[x],
17                        _mm256_mul_pd(_mm256_load_pd(A+n*k+x*4+i), b));
18            }
19
20            for ( int x = 0; x < UNROLL; x++ )
21                _mm256_store_pd(C+i+x*4+j*n, c[x]);
22        }
23    }

```

**Figura 4.77** Versione C ottimizzata di DGEMM che utilizza le funzioni intrinseche del C per generare istruzioni AVX. Queste sfruttano il parallelismo a livello di parola dell'x86 (vedi Figura 3.22) e l'espansione dei cicli per creare più opportunità per sfruttare il parallelismo a livello di parola. La Figura 4.78 mostra il codice assembler prodotto dal compilatore per il ciclo interno, che espande il corpo dei tre cicli `for` per evidenziare il parallelismo a livello di istruzioni.

```

1  vmovapd (%r11),%ymm4          // Carica 4 elementi di C in %ymm4
2  mov    %rbx,%rax             // Registro %rax = %rbx
3  xor    %ecx,%ecx             // Registro %ecx = 0
4  vmovapd 0x20(%r11),%ymm3      // Carica 4 elementi di C in %ymm3
5  vmovapd 0x40(%r11),%ymm2      // Carica 4 elementi di C in %ymm2
6  vmovapd 0x60(%r11),%ymm1      // Carica 4 elementi di C in %ymm1
7  vbroadcastsd (%rcx,%r9,1),%ymm0 // Fai 4 copie dell'elemento di B
8  add    $0x8,%rcx              // Registro %rcx = %rcx + 8
9  vmulpd (%rax),%ymm0,%ymm5    // Moltiplicazione parallela %ymm1,4 elementi di A
10 vaddpd %ymm5,%ymm4,%ymm4     // Somma parallela %ymm5, %ymm4
11 vmulpd 0x20(%rax),%ymm0,%ymm5 // Moltiplicazione parallela %ymm1,4 elementi di A
12 vaddpd %ymm5,%ymm3,%ymm3     // Somma parallela %ymm5, %ymm3
13 vmulpd 0x40(%rax),%ymm0,%ymm5 // Moltiplicazione parallela %ymm1,4 elementi di A
14 vmulpd 0x60(%rax),%ymm0,%ymm0 // Moltiplicazione parallela %ymm1,4 elementi di A
15 add    %r8,%rax              // Registro %rax = %rax + %r8
16 cmp    %r10,%rcx              // Compara %r8 e %rax
17 vaddpd %ymm5,%ymm2,%ymm2      // Somma parallela %ymm5, %ymm2
18 vaddpd %ymm0,%ymm1,%ymm1      // Somma parallela %ymm0, %ymm1
19 jne    68 <dgemm+0x68>       // Salta se %r8 != %rax
20 add    $0x1,%esi              // Registro %esi = %esi + 1
21 vmovapd %ymm4,(%r11)         // Memorizza %ymm4 in 4 elementi di C
22 vmovapd %ymm3,0x20(%r11)      // Memorizza %ymm3 in 4 elementi di C
23 vmovapd %ymm2,0x40(%r11)      // Memorizza %ymm2 in 4 elementi di C
24 vmovapd %ymm1,0x60(%r11)      // Memorizza %ymm1 in 4 elementi di C

```

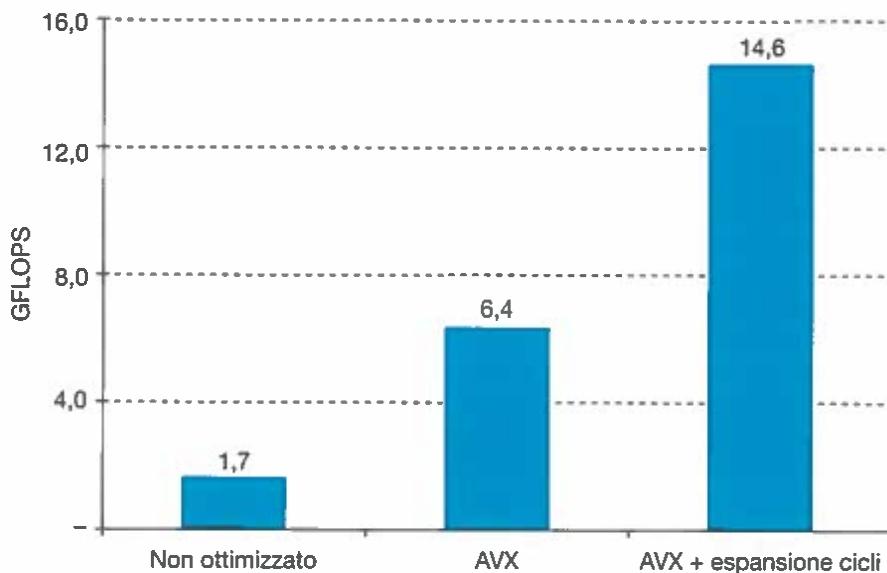
**Figura 4.78** Codice assembler x86 del corpo dei cicli interni generati dal compilatore a partire dal codice C di Figura 4.77, in cui i cicli sono stati espansi.

La Figura 4.79 mostra l'aumento di prestazioni di DGEMM su matrici  $32 \times 32$  passando da una versione non ottimizzata a una versione AVX e poi a una versione AVX con l'espansione dei cicli. L'espansione dei cicli da sola più che raddoppia le prestazioni, che passano da 6,4 GFLOPS a 14,6 GFLOPS. L'applicazione del parallelismo a livello di parola e del parallelismo a livello di istruzioni risulta in un aumento di velocità di 8,8 volte rispetto alla versione DGEMM non ottimizzata di Figura 3.21.



**Approfondimento.** Come abbiamo mostrato nella sezione *Approfondimento* del paragrafo 3.8, questi risultati sono stati ottenuti con la modalità Turbo spenta. Se la inseriamo, come abbiamo fatto nel Capitolo 3, possiamo migliorare tutti i risultati grazie all'aumento temporaneo della frequenza di clock di un fattore  $3,3/2,6 = 1,27$  e ottenere 2,1 GFLOPS per la versione DGEMM non ottimizzata, 8,1 GFLOPS per la versione AVX e 18,6 GFLOPS per la versione AVX con l'espansione dei cicli. Come mostrato nel paragrafo 3.8, la modalità Turbo funziona particolarmente bene in questo caso perché stiamo utilizzando uno solo dei core di un microprocessore a otto core.

**Approfondimento.** Non si verificano stalli della pipeline nonostante il registro %ymm5 venga riutilizzato nelle linee 9-17 di Figura 4.78, perché la pipeline del Core i7 Intel rinomina i registri.



**Figura 4.79** Prestazioni di tre versioni di DGEMM su matrici  $32 \times 32$ . Il parallelismo a livello di parola e di istruzioni ha consentito un aumento di velocità di circa un fattore 9 rispetto al codice non ottimizzato di Figura 3.21.

### Autovalutazione

Stabilire se le seguenti affermazioni sono vere o false.

- Il Core i7 Intel utilizza una pipeline dotata di parallelizzazione dell'esecuzione che esegue direttamente le istruzioni dell'x86.
- Il Cortex-A53 e il Core i7 utilizzano la riorganizzazione dinamica del codice.
- La microarchitettura del Core i7 ha molti più registri di quanti ne richieda l'x86.
- Il Core i7 Intel utilizza meno della metà degli stadi di pipeline del Pentium 4 Prescott (Figura 4.70).

## 4.13 Argomenti avanzati: un'introduzione alla progettazione digitale con un linguaggio di progettazione dell'hardware e un modello di pipeline, e approfondimenti sulla pipeline

La progettazione digitale oggi viene effettuata utilizzando linguaggi di descrizione dell'hardware e strumenti moderni di CAD con cui si possono creare circuiti dettagliati a partire dalla descrizione del loro funzionamento, utilizzando sia librerie sia algoritmi di sintesi dei circuiti logici. Su questi linguaggi e sul loro impiego nella progettazione digitale sono stati scritti libri interi. Questo paragrafo, disponibile online , fornisce una breve introduzione e mostra come un linguaggio di progettazione dell'hardware, il Verilog, possa essere utilizzato per descrivere l'unità di controllo dal punto di vista del suo funzionamento e in una forma adatta alla sintesi hardware. Fornisce anche una serie di modelli in

Verilog del comportamento della pipeline a cinque stadi. Nel primo modello gli hazard vengono ignorati, mentre nei modelli successivi si mettono in evidenza le aggiunte necessarie alla propagazione e alla gestione degli hazard sui dati e sui salti condizionati.

Infine, per i lettori che vogliono approfondire il funzionamento della pipeline, vengono riportati una decina di schemi che mostrano l'esecuzione di alcune istruzioni RISC-V utilizzando la rappresentazione grafica della pipeline a singolo ciclo.

## 4.14 Errori e trabocchetti

*Errore: realizzare una pipeline è facile.*

I diversi libri scritti dagli autori testimoniano la difficoltà nel costruire una pipeline che funzioni correttamente: la prima edizione del testo di architettura degli elaboratori rivolto agli specialisti conteneva un errore nella pipeline; nonostante il testo fosse stato rivisto da più di 100 persone e utilizzato nei corsi di architettura degli elaboratori di 18 università diverse, l'errore saltò fuori soltanto quando qualcuno cercò di costruire effettivamente il calcolatore descritto nel libro. Il fatto che il codice Verilog che descrive una pipeline come quella del Core i7 Intel sia costituito da migliaia di linee è un'indicazione della complessità.

*Errore: le idee su come migliorare le pipeline possono essere implementate in modo indipendente dalla tecnologia.*

Quando il numero di transistor per chip e la loro velocità rendevano la pipeline a cinque stadi la soluzione più conveniente, i salti ritardati (vedi la sezione *Approfondimento* a pagina 242) rappresentavano una soluzione semplice per gli hazard sul controllo. Con pipeline più profonde, dotate di esecuzione super-scalare e predizione dinamica dei salti, tale soluzione è diventata ridondante. Nei primi anni '90 il riordinamento dinamico del codice richiedeva troppe risorse e non veniva quindi impiegato per ottenere prestazioni elevate; con il numero di transistor che continua ad aumentare secondo la Legge di Moore e grazie al fatto che i circuiti logici sono diventati molto più veloci della memoria, unità funzionali multiple e pipeline dinamiche hanno molto più senso. Oggi, l'attenzione verso l'assorbimento di potenza porta a considerare progetti meno aggressivi.



LEGGE DI MOORE

*Trabocchetto: non tenere in giusta considerazione la struttura dell'insieme delle istruzioni può influire in modo negativo sulle pipeline.*

Molte difficoltà nel progettare le pipeline nascono dalle complicazioni causate dagli insiemi di istruzioni. Ecco alcuni esempi.

- Lunghezza e tempo di esecuzione fortemente variabili da istruzione a istruzione possono portare a sbilanciamenti tra i diversi stadi della pipeline; inoltre, complicano di molto il riconoscimento degli hazard quando la progettazione della pipeline viene effettuata sull'insieme delle istruzioni. Questo problema è stato affrontato e superato per la prima volta nei calcolatori DEC VAX 8500 degli anni '80, utilizzando uno schema basato su micro-operazioni e micropipeline simile a quello impiegato oggi dal Core i7 Intel. Ovviamente, rimane il costo aggiuntivo richiesto per tradurre e mantenere la corrispondenza tra le micro-operazioni e le istruzioni.

- Modalità di indirizzamento sofisticate possono condurre a diversi tipi di problemi. Per esempio, modalità di indirizzamento che aggiornano i registri complicano il riconoscimento degli hazard. Altre modalità di indirizzamento, che richiedono accessi multipli alla memoria, rendono significativamente più complessa l'unità di controllo della pipeline e rendono difficile ottenere un flusso costante delle istruzioni nella pipeline.
- L'esempio migliore è forse quello dei DEC Alpha e dei DEC NVAX. Nonostante i due processori fossero basati su tecnologie confrontabili, il nuovo insieme di istruzioni Alpha consentì un'implementazione hardware le cui prestazioni erano doppie in velocità rispetto a quelle del NVAX. Un altro esempio è fornito da Bhandarkar e Clark (1991) che confrontarono le prestazioni del MIPS M/2000 con quelle del VAX 8700 contando i cicli di clock relativi ai benchmark SPEC; essi trovarono che, benché il MIPS M/2000 eseguisse un maggior numero di istruzioni, il VAX utilizzava in media un numero di cicli di clock pari a 2,7 volte quello del MIPS; il MIPS risultò quindi più veloce.

## 4.15 Note conclusive

*La saggezza consiste per nove decimi nell'essere saggi al momento giusto.*

Proverbo popolare americano



**Latenza dell'istruzione:** il tempo di esecuzione intrinseco di un'istruzione.

Come abbiamo visto in questo capitolo, sia l'unità di elaborazione sia l'unità di controllo di un processore possono essere progettate a partire dall'architettura dell'insieme delle istruzioni e da una conoscenza delle caratteristiche di base della tecnologia. In particolare, nel paragrafo 4.3 abbiamo visto come l'unità di elaborazione dei processori RISC-V possa essere costruita a partire dall'architettura delle istruzioni e dalla decisione di costruire un processore che esegua un'istruzione per ogni ciclo di clock. Naturalmente la tecnologia utilizzata influenza molte delle scelte di progetto, suggerendo quali componenti possano essere impiegati nell'unità di elaborazione e se l'implementazione a singolo ciclo abbia senso o meno.

L'utilizzo delle **pipeline** aumenta il throughput ma non il tempo di esecuzione intrinseco delle singole istruzioni, o **latenza dell'istruzione**; per alcune istruzioni, tale latenza ha un valore simile a quello dell'approccio a singolo ciclo di clock. L'esecuzione di più istruzioni in parallelo richiede l'aggiunta di altro hardware in modo tale da poter avviare a esecuzione più istruzioni per ogni ciclo di clock al prezzo di un aumento della latenza effettiva della singola istruzione. La pipeline è stata introdotta per ridurre il tempo di esecuzione rispetto all'unità di elaborazione a singolo ciclo di clock, mentre il parallelismo dell'esecuzione punta chiaramente a ridurre il *numero di cicli di clock per istruzione* (CPI).

La pipeline e il parallelismo dell'esecuzione tentano entrambi di sfruttare il parallelismo a livello di istruzioni. Le dipendenze sui dati e sul controllo (che possono diventare hazard) rappresentano la limitazione principale allo sfruttamento del parallelismo. Il riordinamento del codice e la speculazione, mediante **predizione** sia hardware sia software, sono le tecniche principali utilizzate per ridurre l'impatto delle dipendenze sulle prestazioni.

Abbiamo mostrato che l'espansione del ciclo della funzione DGEMM di 4 volte consente di esporre più istruzioni che possono essere sfruttate dal motore di esecuzione fuori ordine del Core i7 per ottenere più del raddoppio delle prestazioni.

Il passaggio a pipeline più profonde, alla parallelizzazione dinamica dell'esecuzione e al riordinamento dinamico del codice, avvenuto a metà degli anni '90, ha contribuito a mantenere il miglioramento annuale delle prestazioni del 60%, iniziato a partire dagli anni '80. Come si è detto nel Capitolo 1,



PREDIZIONE

i microprocessori hanno mantenuto il modello sequenziale di scrittura del codice, ma si sono scontrati con la barriera dell'energia. Perciò, l'industria è stata costretta a sperimentare i multiprocessori, che sfruttano il parallelismo a un livello più elevato (*vedi Cap. 6*). Questa tendenza ha anche richiesto che i progettisti rivalutassero le implicazioni sul rapporto potenza-prestazioni di alcune delle invenzioni introdotte a partire dalla metà degli anni '90; questo ha portato a una semplificazione delle pipeline nelle versioni più recenti delle microarchitetture.

Considerato il tasso di miglioramento delle prestazioni nei processori paralleli, la Legge di Amdahl suggerisce che un'altra parte del sistema diventerà il collo di bottiglia. Il collo di bottiglia in questione è l'argomento del prossimo capitolo: la gerarchia delle memorie.



## 4.16 | Inquadramento storico e approfondimenti

Questo paragrafo, disponibile online , tratta la storia dei primi processori dotati di pipeline e dei primi processori superscalari. Inoltre, illustra lo sviluppo delle tecniche di esecuzione fuori ordine e di speculazione e gli importanti sviluppi nella tecnologia dei compilatori.

## 4.17 | Esercizi

**4.1** Si consideri la seguente istruzione:

Istruzione: and rd, rs1, rs2

Interpretazione:  $\text{Reg}[\text{rd}] = \text{Reg}[\text{rs1}] \text{ AND } \text{reg}[\text{rs2}]$

**4.1.1** [5] <4.3> Determinare il valore dei segnali di controllo generati per questa istruzione dall'unità di controllo di Figura 4.10.

**4.1.2** [5] <4.3> Quali risorse (blocchi funzionali) eseguono il lavoro utile per l'esecuzione di queste istruzioni?

**4.1.3** [10] <4.3> Quali risorse (blocchi funzionali) producono comunque un output, anche se il loro output non viene poi utilizzato dall'istruzione? Quali risorse non producono alcun output per queste istruzioni?

**4.2** [10] <4.4> Spiegare il significato di tutti i termini "indifferente" riportati in Figura 4.18.

**4.3** Si consideri la seguente combinazione di istruzioni:

Tipo R	Tipo I (non load)	Load	Store	Branch	Jump
24%	28%	25%	10%	11%	2%

**4.3.1** [5] <4.4> Quale percentuale di tutte le istruzioni fa uso della memoria dati?

**4.3.2** [5] <4.4> Quale percentuale di tutte le istruzioni fa uso della memoria istruzioni?

**4.3.3** [5] <4.4> Quale percentuale di tutte le istruzioni fa uso dell'estensione del segno?

**4.3.4** [5] <4.4> Cosa fa l'estensione del segno nei cicli nei quali il suo output non serve?

**4.4** Quando i chip di silicio vengono fabbricati, difetti nei materiali (per es. nel silicio) ed errori di fabbricazione possono fare sì che i circuiti risultino difettati. Un difetto molto comune è la rottura di una pista: il segnale da esso trasportato rimane sempre bloccato a 0, indipendentemente da quale sia il valore logico che dovrebbe assumere. Questo viene spesso chiamato collegamento forzato a 0 ("fault stuck-at-0")

**4.4.1** [5] <4.4> Quali istruzioni non funzionerebbero più correttamente se il bus del segnale MemToReg fosse bloccato a zero?

**4.4.2** [5] <4.4> Quali istruzioni non funzionerebbero più correttamente se il bus del segnale ALUSrc fosse bloccato a zero?

**4.5** In questo esercizio verrà esaminato in dettaglio come viene eseguita un'istruzione da un'unità di elaborazione a singolo ciclo. I problemi di questo esercizio si riferiscono al ciclo di clock nel quale il processore carica dalla memoria istruzioni la seguente parola: 0x00c6ba23.

**4.5.1 [10] <4.4>** Qual è il valore degli ingressi dell'unità di controllo della ALU per questa istruzione?

**4.5.2 [5] <4.4>** Qual è il nuovo indirizzo del PC dopo che l'istruzione è stata eseguita? Evidenziare il cammino che calcola il nuovo indirizzo.

**4.5.3 [10] <4.4>** Mostrare il valore fornito in ingresso e in uscita da ciascun multiplexer durante l'esecuzione di questa istruzione. Riportare il valore in uscita a Reg [xn].

**4.5.4 [10] <4.4>** Determinare quali siano i dati in ingresso alla ALU e ai due addizionatori.

**4.5.5 [10] <4.4>** Determinare tutti i valori in ingresso al register file.

**4.6** Il paragrafo 4.4 non tratta le istruzioni di tipo I quali addi e andi.

**4.6.1 [5] <4.4>** Quali blocchi logici aggiuntivi sono eventualmente necessari per aggiungere le istruzioni di tipo I alla CPU mostrata in Figura 4.21? Aggiungere i blocchi eventualmente necessari alla Figura 4.21 e spiegare la loro funzione.

**4.6.2 [10] <4.4>** Riportare i valori dei segnali generati dall'unità di controllo centrale per le istruzioni di addi. Spiegare il ragionamento che sta dietro alla specifica di "indifferente" per alcuni segnali di controllo.

**4.7** Nei problemi di questo esercizio si suppone che i blocchi logici richiesti per implementare l'unità di elaborazione di un processore abbiano le latenze riportate nella tabella riportata in fondo alla pagina. "Lettura Registro" si riferisce al tempo che intercorre tra il fronte di salita del clock e l'istante in cui il valore contenuto nel registro compare all'uscita del registro stesso. Questo tempo si applica solo al program counter. "Impostazione Registro" si riferisce al tempo per cui il dato in input a un registro deve rimanere stabile prima del fronte di salita del clock. Questo tempo si applica sia al PC sia al Register File.

**4.7.1 [5] <4.4>** Qual è la latenza di un'istruzione di tipo R (cioè quanto deve durare il periodo del clock per essere sicuri che l'istruzione venga eseguita correttamente)?

**4.7.2 [10] <4.4>** Qual è la latenza di un'istruzione di Id? (Fate attenzione: molti studenti inseriscono dei multiplexer in più sul cammino critico).

**4.7.3 [10] <4.4>** Qual è la latenza di un'istruzione di sd? (Fate attenzione: molti studenti inseriscono dei multiplexer in più sul cammino critico).

**4.7.4 [5] <4.4>** Qual è la latenza di un'istruzione di beq?

**4.7.5 [5] <4.4>** Qual è la latenza di un'istruzione di tipo I?

**4.7.6 [5] <4.4>** Qual è il periodo di clock minore possibile per questo PC?

**4.8 [10]** Si supponga che si possa costruire una CPU in cui il periodo di clock sia diverso per le diverse istruzioni. Quale sarà lo l'incremento di velocità (speed-up) di questa CPU rispetto alla CPU presentata in Figura 4.21 per la combinazione di istruzioni qui riportata?

Tipo R/Tipo I (non Id)	Id	sd	beq
52%	25%	11%	12%

**4.9** Si supponga di aggiungere un moltiplicatore alla CPU mostrata in Figura 4.21. Questo moltiplicatore aggiungerebbe 300 ps alla latenza della ALU, ma ridurrebbe il numero di istruzioni del 5% (perché non ci sarebbe più la necessità di emulare le istruzioni di moltiplicazione).

**4.9.1 [5] <4.4>** Quale sarebbe il periodo di clock con e senza questo miglioramento?

**4.9.2 [10] <4.4>** Quale sarebbe l'incremento di velocità ottenuto con questo miglioramento?

**4.9.3 [10] <4.4>** Di quanto potrebbe essere più lenta al massimo la ALU, pur mantenendo un miglioramento delle prestazioni?

**4.10** Quando i progettisti dei processori valutano un possibile miglioramento dell'unità di elaborazione, la

Mem-I/ Mem-D	Register File	Mux	ALU	Addiz	Porta logica singola	Lettura Registro	Impostazione Registro	Estensione Segno	Controllo
250 ps	150 ps	25 ps	200 ps	150 ps	5 ps	30 ps	20 ps	50 ps	50 ps

decisione finale dipende di solito dal rapporto costo/prestazioni. Nei prossimi tre problemi si supponga di partire dall'unità di elaborazione di Figura 4.21, dalle latenze dell'Esercizio 4.7 e dai costi presenti nella tabella riportata in fondo alla pagina. Si consideri che raddoppiare il numero di registri di uso generico da 32 a 64 possa ridurre il numero di istruzioni di `ld` e `sd` del 12% ma aumenti la latenza del register file da 150 ps a 160 ps e raddoppi il costo da 200 a 400. (Si utilizzi la combinazione di istruzioni dell'Esercizio 4.8 e trascurate gli altri effetti sull'ISA discussi nell'Esercizio 2.18).

**4.10.1 [5] <4.4>** Qual è l'incremento di velocità ottenuto con questo miglioramento?

**4.10.2 [10] <4.4>** Confrontare le variazioni in termini di velocità con l'incremento in termini di costo.

**4.10.3 [10] <4.4>** Dati i rapporti costo/velocità appena calcolati, descrivere una situazione in cui abbia senso aggiungere altri registri e descriverne un'altra in cui non abbia senso aggiungere altri registri.

**4.11** Esaminare la difficoltà ad aggiungere un'istruzione di `lwi.d rd, rs1, rs2` ("Load con incremento") al RISC-V.

Interpretazione:

$$\text{Reg}[rd] = \text{Mem}[\text{Reg}[rs1] + \text{Reg}[rs2]]$$

**4.11.1 [5] <4.4>** Quali nuovi blocchi funzionali servirebbero eventualmente per questa nuova istruzione?

**4.11.2 [5] <4.4>** Quali blocchi funzionali esistenti dovrebbero eventualmente essere modificati per questa nuova istruzione?

**4.11.3 [5] <4.4>** Quali nuovi cammini interni all'unità di elaborazione dovrebbero eventualmente essere introdotti per l'esecuzione di questa nuova istruzione?

**4.11.4 [5] <4.4>** Quali nuovi segnali di controllo dovrebbero eventualmente essere introdotti per supportare l'esecuzione di questa nuova istruzione?

**4.12** Esaminare la difficoltà ad aggiungere un'istruzione di `swap rs1, rs2` al RISC-V.

Interpretazione:

$$\text{Reg}[rs2] = \text{Reg}[rs1]; \quad \text{Reg}[rs1] = \text{Reg}[rs2]$$

**4.12.1 [5] <4.4>** Quali nuovi blocchi funzionali servirebbero eventualmente per questa nuova istruzione?

**4.12.2 [10] <4.4>** Quali blocchi funzionali esistenti dovrebbero eventualmente essere modificati per questa nuova istruzione?

**4.12.3 [5] <4.4>** Quali nuovi cammini interni all'unità di elaborazione dovrebbero eventualmente essere introdotti per l'esecuzione di questa nuova istruzione?

**4.12.4 [5] <4.4>** Quali nuovi segnali di controllo dovrebbero eventualmente essere introdotti per supportare l'esecuzione di questa nuova istruzione?

**4.12.5 [5] <4.4>** Modificare la Figura 4.21 per dimostrare una possibile implementazione di questa nuova istruzione.

**4.13** Esaminare la difficoltà ad aggiungere un'istruzione di `ss rs1, rs2` al RISC-V.

**4.13.1 [10] <4.4>** Quali nuovi blocchi funzionali servirebbero eventualmente per questa nuova istruzione?

**4.13.2 [10] <4.4>** Quali blocchi funzionali esistenti dovrebbero eventualmente essere modificati per questa nuova istruzione?

**4.13.3 [5] <4.4>** Quali nuovi cammini interni all'unità di elaborazione dovrebbero eventualmente essere introdotti per l'esecuzione di questa nuova istruzione?

**4.13.4 [5] <4.4>** Quali nuovi segnali di controllo dovrebbero eventualmente essere introdotti per supportare l'esecuzione di questa nuova istruzione?

**4.13.5 [5] <4.4>** Modificare la Figura 4.21 per dimostrare una possibile implementazione di questa nuova istruzione.

**4.14 [5] <4.4>** Per quali istruzioni (se ce ne sono) il blocco funzionale Gen Imm si trova sul cammino critico?

**4.15** L'istruzione di `ld` ha la latenza maggiore nella CPU del paragrafo 4.4. Se modificassimo le istruzioni di `ld` e `sd` in modo che non debba calcolare l'offset (cioè che l'indirizzo da/a cui debba essere trasferito il dato sia calcolato e scritto nel registro `rs1` prima

Mem-I	Register File	Mux	ALU	Addiz	Mem-D	Registro Singolo	Estensione Segno	Porta logica singola	Controllo
1000	200	10	100	30	2000	5	100	1	500

di chiamare l'istruzione di 1d/sd), per cui non ci sono istruzioni che debbano utilizzare sia la ALU sia la memoria dati. Questo ci consentirebbe di ridurre il periodo del clock, ma aumenterebbe il numero di istruzioni dato che molte 1d e sd dovrebbero essere sostituite dalla coppia di istruzioni 1d/add o sd/add.

**4.15.1 [5] <4.4>** Quale sarebbe il nuovo periodo del clock?

**4.15.2 [10] <4.4>** Il programma costituito dalla combinazione di istruzioni presentata nell'Esercizio 4.7 verrebbe eseguito più velocemente o più lentamente da questa nuova CPU? Di quanto? (Per semplicità, si supponga di sostituire ciascuna istruzione di 1d e sd con una coppia di istruzioni 1d/add e sd/add).

**4.15.3 [5] <4.4>** Qual è il fattore principale che determina se un programma viene eseguito più velocemente o più lentamente sulla nuova CPU?

**4.15.4 [5] <4.4>** Pensate che la CPU originaria (quella mostrata in Figura 4.21) sia un progetto complessivamente migliore oppure pensate che sia migliore la nuova CPU? Perché?

**4.16** In questo esercizio verrà esaminato come la pipeline influenzi il periodo del clock di processore. I problemi di questo esercizio suppongono che i singoli stadi dell'unità di elaborazione abbiano la seguente latenza:

IF	ID	EX	MEM	WB
250 ps	350 ps	150 ps	300 ps	200 ps

Si assume inoltre che le istruzioni eseguite dal processore siano suddivise come segue:

ALU/Logiche	Salti incondiz/condiz	Load	Store
45%	20%	20%	15%

**4.16.1 [5] <4.5>** Qual è il periodo del clock di un processore dotato di pipeline e di un processore senza pipeline?

**4.16.2 [10] <4.5>** Qual è la latenza totale di un'istruzione 1d in un processore dotato di pipeline e in uno senza pipeline?

**4.16.3 [10] <4.5>** Se si potesse suddividere uno stadio dell'unità di elaborazione dotata di pipeline in due nuovi stadi, ciascuno avente metà latenza dello stadio originale, quale stadio sarebbe meglio suddividere e quale sarebbe la nuova durata del ciclo di clock del processore?

**4.16.4 [10] <4.5>** Supponendo che non ci siano stalli o hazard, qual è la percentuale di utilizzo della memoria dati?

**4.16.5 [10] <4.5>** Supponendo che non ci siano stalli o hazard, qual è la percentuale di utilizzo della porta di scrittura del register file?

**4.17 [10] <4.5>** Qual è il minimo numero di cicli di clock richiesti per portare a termine l'esecuzione di n istruzioni su una CPU con una pipeline a k stadi? Giustificare la formula ricavata.

**4.18 [5] <4.5>** Si supponga che x11 venga inizializzato a 11 e x12 venga inizializzato a 22. Si supponga di eseguire il codice riportato sotto su una realizzazione della pipeline mostrata nel paragrafo 4.5 che non gestisce gli hazard sui dati (cioè il programmatore è responsabile di gestire gli hazard sui dati inserendo istruzioni NOP quando serve). Quale sarà il contenuto finale dei registri x13 e x14?

```
addi x11, x12, 5
add x13, x11, x12
addi x14, x11, 15
```

**4.19 [10] <4.5>** Si supponga che x11 venga inizializzato a 11 e x12 venga inizializzato a 22. Si supponga di eseguire il codice riportato sotto su una realizzazione della pipeline mostrata nel paragrafo 4.5 che non gestisce gli hazard sui dati (cioè il programmatore è responsabile di gestire gli hazard sui dati inserendo istruzioni NOP quando serve). Quale sarà il contenuto finale del registro x15? Si supponga che il register file venga scritto all'inizio del ciclo di clock e letto alla fine di una ciclo. Quindi, lo stadio ID fornisce in uscita il contenuto dello stadio WB durante lo stesso ciclo di clock. Si esaminino il paragrafo 4.7 e la Figura 4.51 per ulteriori dettagli.

```
addi x11, x12, 5
add x13, x11, x12
addi x14, x11, 15
add x15, x11, x11
```

**4.20 [5] <4.5>** Aggiungere un'istruzione NOP al codice riportato sotto in modo tale che funzioni correttamente su una pipeline che non gestisce gli hazard.

```
addi x11, x12, 5
add x13, x11, x12
addi x14, x11, 15
add x15, x13, x12
```

**4.21** Si consideri una realizzazione della pipeline del paragrafo 4.5 che non gestisce gli hazard sui dati (cioè il programmatore è responsabile di gestire gli hazard sui dati inserendo istruzioni NOP quando serve). Si supponga che (dopo l'ottimizzazione) un programma tipico contenente n istruzioni richieda

altri  $4^*n$  istruzioni NOP per gestire correttamente gli hazard sui dati.

**4.21.1 [5] <4.5>** Si supponga che il periodo del ciclo di clock di questa pipeline senza la propagazione sia di 250 ps. Si supponga anche che l'aggiunta di hardware per la propagazione riduca il numero di NOP da  $0,4^*n$  a  $0,05^*n$ , ma aumenti il periodo del clock a 300 ps. Qual è l'aumento di velocità di questa nuova pipeline rispetto a quella senza la propagazione?

**4.21.2 [10] <4.5>** Programmi diversi richiedono un numero diverso di NOP. Quante NOP (in percentuale delle istruzioni del codice) possono rimanere prima che il tipico programma riportato sopra venga eseguito più lentamente sulla pipeline dotata di propagazione?

**4.21.3 [10] <4.5>** Ripetere l'Esercizio 4.21.2; tuttavia, questa volta considerare  $x$  il numero di NOP rispetto a  $n$  (nell'Esercizio 4.21.2,  $x$  era uguale a 0,4.). La risposta sarà funzione di  $x$ .

**4.21.4 [10] <4.5>** Può un programma contenente solo  $0,075^*n$  NOP essere eseguito più velocemente da una pipeline dotata di propagazione? Sia in caso affermativo che negativo, spiegare il motivo.

**4.21.5 [10] <4.5>** Qual è la percentuale minima, rispetto al numero di istruzioni del codice, di NOP che un programma deve avere prima che possa venire eseguito più velocemente sulla pipeline dotata di propagazione?

**4.22 [10] <4.5>** Si consideri il seguente frammento di codice RISC-V:

```
sd      x29,12(x16)
ld      x29,8(x16)
sub    x17,x15,x14
beqz  x17,Etichetta
add    x15,x11,x14
sub    x15,x30,x14
```

Si supponga di modificare la pipeline in modo che abbia solamente una memoria (in grado di gestire sia i dati sia le istruzioni). In questo caso, si verificherà un hazard strutturale ogni volta che un programma ha bisogno di caricare un'istruzione durante lo stesso ciclo di clock in cui un'altra istruzione ha bisogno di accedere ai dati in memoria.

**4.22.1 [5] <4.5>** Disegnare lo schema di una pipeline dove verrebbe messo in stallo il frammento di codice riportato sopra.

**4.22.2 [5] <4.5>** In generale, è possibile ridurre il numero di stalli/NOP che risultano da questo hazard strutturale riordinando il codice?

**4.22.3 [5] <4.5>** Questo hazard strutturale deve essere gestito dall'hardware? Abbiamo visto che gli hazard sui dati possono venire eliminati aggiungendo un'istruzione NOP al codice. Si può fare la stessa cosa con questo hazard strutturale? Se sì, spiegare come. Se no, spiegare perché.

**4.22.4 [5] <4.5>** Approssimativamente, quanti stalli ci si può aspettare che vengano generati da un programma tipico? (Si utilizzi la combinazione di istruzioni dell'Esercizio 4.8.)

**4.23** Se modifichiamo le istruzioni di load/store per utilizzare un registro (senza l'offset) per l'indirizzo, queste istruzioni non hanno più bisogno di utilizzare la ALU (vedi Esercizio 4.15). Di conseguenza, gli stadi MEM ed EX si possono sovrapporre e la pipeline sarebbe costituita da solo quattro stadi.

**4.23.1 [10] <4.5>** Come viene influenzato il periodo di clock dalla riduzione della profondità della pipeline?

**4.23.2 [5] <4.5>** Come potrebbe questo cambiamento migliorare le prestazioni della pipeline?

**4.23.3 [5] <4.5>** Come potrebbe questo cambiamento degradare le prestazioni della pipeline?

**4.24 [10] <4.7>** Quale dei due diagrammi di esecuzione descrive meglio il funzionamento dell'unità di rilevamento degli hazard?

Scelta 1:

ld x11, 0(x12):	IF ID EX ME WB
add x13, x11, x14	IF ID EX..ME WB
or x15, x16, x17:	IF ID..EX ME WB

Scelta 2:

ld x11, 0(x12):	IF ID EX ME WB
add x13, x11, x14	IF ID..EX ME WB
or x15, x16, x17:	IF..ID EX ME WB

**4.25** Si consideri il ciclo seguente:

CICLO:	ld x10, 0(x13)
	ld x11, 8(x13)
	add x12, x10, x11
	subi x13, x13, 16
	bnez x12, CICLO

Si supponga una predizione esatta delle branch e quindi che non si verifichino stalli dovuti a hazard sul controllo. Si supponga anche che non ci siano delay slot e che la pipeline abbia un supporto completo alla propagazione. Infine, si supponga che i salti condizionati vengano risolti nello stadio EX invece che nello stadio ID.

EX a 1° Solo	MEM a 1° Solo	EX a 2° Solo	MEM a 2° Solo	EX a 1° e EX a 2°
5%	20%	5%	10%	10%

**4.25.1 [10] <4.7>** Mostrare lo schema di esecuzione della pipeline per le prime due iterazioni di questo ciclo.

**4.25.2 [10] <4.7>** Identificare gli stadi della pipeline che non fanno lavoro utile. Con quale frequenza, espressa come percentuale del numero totale di cicli di clock, si verificano cicli di clock in cui tutti e cinque gli stadi della pipeline compiono lavoro utile? (Iniziare dalla situazione in cui la subi si trova nello stadio IF e terminare con il ciclo nel quale la bnez si trova nello stadio IF.)

**4.26** Questo esercizio è stato pensato per aiutarvi a capire il compromesso tra costo, complessità e prestazioni della propagazione in un processore dotato di pipeline. I problemi si riferiscono all'unità di elaborazione dotata di pipeline descritta in Figura 4.53 supponendo che, tra tutte le istruzioni eseguite dal processore, una parte esibisca una particolare dipendenza tra i dati RAW (dati grezzi). Questa dipendenza viene identificata dallo stadio che produce il risultato richiesto (EX o MEM) e dall'istruzione che deve utilizzare questo risultato: la prima istruzione immediatamente successiva all'istruzione che produce il risultato, la seconda, o entrambe. Si suppone che la scrittura dei registri avvenga nella prima metà del periodo di clock e che la lettura avvenga nella seconda metà, per cui la dipendenza tra l'uscita dallo stadio EX e gli operandi della terza istruzione e tra l'uscita dallo stadio MEM e gli operandi della seconda istruzione non viene considerata un hazard sui dati. Inoltre, si supponga che i salti condizionati siano risolti nello stadio EX e che il CPI del processore sia uguale a 1 e che non si verifichino hazard sui dati (tabella in alto alla pagina). Si suppongano le latenze riportate nella tabella in fondo alla pagina per i singoli stadi della pipeline. Per lo stadio EX, la latenza viene fornita separatamente per un processore senza propagazione e per un processore dotato di diversi tipi di propagazione.

**4.26.1 [5] <4.7>** Per ciascuna delle dipendenze GREZZE riportate nella tabella, fornire una sequenza di almeno tre istruzioni assembler che esibiscono quella dipendenza.

**4.26.2 [5] <4.7>** Per ciascuna delle dipendenze GREZZE riportate sopra, quante NOP sarebbe necessario inserire per consentire che il codice dell'Esercizio 4.26.1 venga eseguito correttamente su una pipeline che non sia dotata né di propagazione né di rilevamento degli hazard? Mostrare dove si possano inserire le NOP.

**4.26.3 [10] <4.7>** Analizzando ciascuna istruzione indipendentemente, verrebbe sovrastimato il numero di NOP necessarie per il corretto funzionamento del programma su una pipeline non dotata né di unità di propagazione né di rilevamento degli hazard. Scrivere una sequenza di tre istruzioni assembler tale per cui la somma del numero di stalli quando si analizza ciascuna istruzione separatamente è maggiore del numero di stalli richiesti effettivamente dalla sequenza per evitare hazard sui dati.

**4.26.4 [5] <4.7>** Se si supponesse che non ci siano altri hazard, quale sarebbe il CPI del programma descritto nella tabella riportata sopra quando verrebbe eseguito su una pipeline senza propagazione? Quale sarebbe la percentuale di cicli in cui la pipeline è in stallo? (Per semplicità, assumere che tutti i casi possibili siano riportati nella tabella sopra e che possano essere trattati indipendentemente.)

**4.26.5 [5] <4.7>** Quale sarebbe il CPI se utilizzassimo la propagazione completa (propagazione di tutti i risultati che possono essere propagati)? Qual è la percentuale dei cicli di clock in cui la pipeline è in stallo?

**4.26.6 [10] <4.7>** Si supponga di non potersi permettere un multiplexer a tre input come quelli richiesti dalla propagazione completa. In questo caso occorre decidere se sia meglio propagare il dato dal registro di pipeline EX/MEM (propagazione dal ciclo di clock precedente), oppure dal registro di pipeline MEM/WB (propagazione da due cicli di clock precedenti). Quale soluzione porta a un minore numero di stalli sui dati?

**4.26.7 [5] <4.7>** Per la probabilità degli hazard e la latenza dei diversi stadi di pipeline riportati, quale

IF	ID	EX (no Propag)	EX (Propag completa)	EX (Propag solo da EX/MEM)	EX (Propag solo da MEM/WB)	MEM	WB
120 ps	100 ps	110 ps	130 ps	120 ps	120 ps	120 ps	100 ps

incremento di velocità si otterrebbe con i diversi tipi di propagazione (EX/MEM, MEM/WB, completo) rispetto a una pipeline che sia sprovvista di propagazione?

**4.26.8 [5] <4.7>** Quale sarebbe l'incremento di velocità ulteriore (relativo al processore più veloce dell'Esercizio 4.26.7) se aggiungessimo una propagazione "attraverso il tempo" che eliminasse tutti gli hazard sui dati? Si supponga che il circuito di "propagazione attraverso il tempo", ancora da inventare, aggiunga 100 ps alla latenza dello stadio EX dotato di propagazione completa.

**4.26.9 [5] <4.7>** La tabella dei tipi di hazard riporta separatamente i tempi per le condizioni da EX a 1° e da EX a 1° e 2°. Perché non ci sono valori separati per da MEM a 1° e da MEM a 2°?

**4.27** I problemi di questo esercizio si riferiscono alla seguente sequenza di istruzioni e suppongono che il codice venga eseguito su un'unità di elaborazione dotata di pipeline a cinque stadi:

```
add x15, x12, x11
ld x13, 4(x15)
ld x12, 0(x2)
or x13, x15, x13
sd x13, 0(x15)
```

**4.27.1 [5] <4.7>** Inserire opportunamente delle NOP per assicurare la corretta esecuzione del codice, supponendo che non ci siano né propagazione né rilevamento degli hazard.

**4.27.2 [10] <4.7>** Ora, modificare e/o riordinare il codice per minimizzare il numero di NOP richieste. Si utilizzi il registro x17 per memorizzare il contenuto di variabili temporanee aggiunte nel codice modificato.

**4.27.3 [10] <4.7>** Supponendo che il processore sia dotato di un'unità di propagazione ma non di un'unità di rilevamento degli hazard, che cosa succederebbe durante l'esecuzione della sequenza di istruzioni originaria sopra riportata?

**4.27.4 [20] <4.7>** Supponendo che il processore sia dotato di propagazione, specificare quali segnali vengano asseriti dall'unità di rilevamento delle criticità e dall'unità di propagazione disegnate in Figura 4.59 in ciascuno dei primi sette cicli di clock dell'esecuzione della sequenza di istruzioni riportata sopra.

**4.27.5 [10] <4.7>** Se non ci fosse la propagazione, quali segnali di ingresso e di uscita aggiuntivi sarebbero richiesti per l'unità di riconoscimento degli hazard di Figura 4.59? Utilizzando la sequenza di istruzioni riportata come esempio, spiegare per ciascun segnale il motivo per cui occorre aggiungerlo.

**4.27.6 [20] <4.7>** Per la nuova unità di rilevamento degli hazard dell'Esercizio 4.26.5, specificare quali segnali di uscita vengono asseriti in ciascuno dei primi cinque cicli di clock dell'esecuzione di questa sequenza di istruzioni.

**4.28** L'importanza di avere un buon predittore dei salti condizionati dipende da quanto spesso i salti condizionati vengono eseguiti. Questo esercizio consentirà di valutare sia l'accuratezza dei predittori dei salti sia il tempo che viene perso negli stalli dovuti agli errori nella predizione dei salti condizionati. Si assume che la suddivisione delle istruzioni nelle diverse categorie sia la seguente:

Tipo R	beqz/bnez	jal	ld	sd
40%	25%	5%	25%	5%

Inoltre, si supponga che la predizione dei salti condizionati abbia la seguente accuratezza:

Sempre effettuato	Sempre non effettuato	Su 2 bit
45%	55%	85%

**4.28.1 [10] <4.8>** I cicli di clock associati agli stalli dovuti a errori nella predizione delle branch aumentano il CPI. Qual è l'aumento del CPI per errori nella predizione delle branch quando viene utilizzato un predittore del tipo "salto sempre effettuato"? Si supponga che il risultato del confronto delle branch venga determinato nello stadio di ID e applicato nello stadio EX, che non ci siano hazard sui dati e che non vengano utilizzati delay slot.

**4.28.2 [10] <4.8>** Risolvere l'Esercizio 4.28.1 considerando un predittore che suppone che i salti siano sempre non effettuati.

**4.28.3 [10] <4.8>** Risolvere l'Esercizio 4.28.1 considerando un predittore su 2 bit.

**4.28.4 [10] <4.8>** Quando si utilizza un predittore su 2 bit, quale incremento di velocità si otterrebbe se si convertissero metà delle istruzioni di branch in istruzioni che utilizzano la ALU? Si supponga che le istruzioni di salto condizionato che vengono predette correttamente o non correttamente abbiano la stessa probabilità di essere sostituite.

**4.28.5 [10] <4.8>** Quando si utilizza un predittore su 2 bit, quale incremento di velocità si otterrebbe se si convertissero metà delle istruzioni di branch in una coppia di istruzioni che utilizzano la ALU? Si suppon-

ga che le istruzioni di salto condizionato che vengono predette correttamente o non correttamente abbiano la stessa probabilità di essere sostituite.

**4.28.6 [10] <4.8>** Alcune istruzioni di salto condizionato sono più predicibili di altre. Sapendo che l'80% dei salti condizionati che vengono effettuati sono facili da predire perché saltano indietro all'inizio di un ciclo e vengono predetti quasi sempre correttamente, quale sarebbe l'accuratezza del predittore a 2 bit sul rimanente 20% delle istruzioni di salto condizionato?

**4.29** Questo esercizio esamina l'accuratezza di vari predittori di salti condizionati per la seguente sequenza di risultati ripetuti (per es. cicli) della predizione: T, NT, T, T, NT.

**4.29.1 [5] <4.8>** Qual è l'accuratezza della predizione "sempre preso" e "sempre non preso" per la sequenza di risultati riportata?

**4.29.2 [5] <4.8>** Qual è l'accuratezza del predittore a 2 bit quando il risultato delle prime quattro branch è quello riportato e il predittore, all'inizio, si trova nello stato riportato in basso a sinistra di Figura 4.61 ("predetto non preso")?

**4.29.3 [10] <4.8>** Qual è l'accuratezza del predittore a 2 bit se la sequenza di risultati riportata si ripetesse all'infinito?

**4.29.4 [30] <4.8>** Progettare un predittore che raggiunga un'accuratezza del 100% se la sequenza dei risultati del confronto è ripetuta sempre allo stesso modo. Il predittore dovrà essere costituito da un circuito sequenziale che fornisce in uscita la predizione (1 se il salto deve essere effettuato, 0 altrimenti) e riceve in ingresso solamente il clock e un segnale di controllo che indica quando l'istruzione in esecuzione è un salto condizionato.

**4.29.5 [10] <4.8>** Qual è l'accuratezza del predittore costruito per risolvere l'Esercizio 4.29.4 quando la sequenza dei risultati del confronto, associata alle istruzioni di salto condizionato, è esattamente l'opposto di quella riportata?

**4.29.6 [20] <4.8>** Risolvere l'Esercizio 4.29.4 realizzando un predittore che sia capace, dopo un periodo di prova durante il quale alcune predizioni possono rivelarsi sbagliate, di predire perfettamente il risultato dei confronti sia per la sequenza riportata sia per la sua opposta. Il predittore dovrà avere un input che gli dice quale sia la predizione corretta. Suggerimento: questo segnale deve permettere al predittore di capire quale delle due sequenze sia quella in ingresso.

**4.30** Questo esercizio esplora come la gestione delle eccezioni influenzi il progetto di una pipeline. I primi tre problemi si riferiscono alla seguente coppia di istruzioni:

Istruzione 1	Istruzione 2
beqz x11, ETICHETTA	ld x11, 0(x12)

**4.30.1 [5] <4.9>** Quali eccezioni possono essere generate da queste istruzioni? Specificare in quale studio della pipeline può essere identificata ciascuna di queste eccezioni.

**4.30.2 [10] <4.9>** Se la procedura di risposta alle eccezioni avesse un indirizzo diverso per ogni eccezione, mostrare come si dovrebbe modificare la pipeline per poter gestire ciascuna eccezione. Si può supporre che l'indirizzo di ciascuna procedura sia stato specificato quando il processore è stato progettato.

**4.30.3 [10] <4.9>** Descrivere che cosa succederebbe nella pipeline se la seconda istruzione della coppia venisse prelevata dalla memoria subito dopo la prima istruzione e la prima istruzione generasse la prima eccezione identificata nell'Esercizio 4.30.1. Mostrare il diagramma di esecuzione della pipeline dal momento in cui la prima istruzione viene prelevata dalla memoria istruzioni fino a quando la prima istruzione della procedura di risposta alle eccezioni termina.

**4.30.4 [20] <4.9>** Nella risposta alle eccezioni vettorizzata, la tabella che contiene gli indirizzi delle procedure di risposta alle diverse eccezioni si trova in una locazione nota (prefissata) della memoria. Modificare la pipeline per implementare questo meccanismo di risposta alle eccezioni. Risolvere l'Esercizio 4.30.3 utilizzando questa pipeline modificata e il meccanismo vettorizzato di risposta alle eccezioni.

**4.30.5 [15] <4.9>** Si vuole emulare la risposta vettorizzata alle eccezioni (descritta nell'Esercizio 4.30.4) su una macchina che prevede un unico indirizzo, prefissato, per la procedura di risposta alle eccezioni. Scrivere la sequenza di istruzioni che deve essere poi caricata in memoria a partire da questo indirizzo. Suggerimento: il codice dovrà identificare l'eccezione, leggere l'indirizzo corretto nella tabella del vettore delle eccezioni e trasferire l'esecuzione alla procedura corretta per la gestione di quell'eccezione.

**4.31** Questo esercizio confronterà le prestazioni dei processori senza parallelizzazione dell'esecuzione a una via con quelli dotati di parallelizzazione a due vie tenendo conto delle modifiche richieste ai programmi per ottimizzare l'esecuzione su un processore

che esegua in parallelo due istruzioni. I problemi si riferiscono a questo ciclo scritto in C:

```
for(i=0;i!=j;i+=2)
    b[i]=a[i]-a[i+1];
```

Un compilatore facendo poco o nulla per ottimizzare il codice potrebbe produrre la sequenza di istruzioni RISC-V seguente:

```
        li      x12, 0
        jal    ENT
INIZIO: slli   x5, x12, 3
        add    x6, x10, x5
        ld     x7, 0(x6)
        ld     x29, 8(x6)
        sub    x30, x7, x29
        add    x31, x11, x5
        sd     x30, 0(x31)
        addi   x12, x12, 2
ENT:    bne   x12, x13, INIZIO
```

Questo codice utilizza i registri seguenti:

i	j	a	b	Valori temporanei
x12	x13	x10	x11	x5-x7, x29-x31

Si supponga che il processore a 2 vie, con smistamento statico delle istruzioni, abbia per questo esercizio le seguenti proprietà:

1. Un'istruzione deve essere un'operazione sulla memoria; l'altra può essere un'istruzione aritmetico-logica o un salto condizionato.
2. Il processore possiede tutti i possibili cammini di propagazione tra i diversi stadi (compresi i cammini dallo stadio ID per la risoluzione degli hazard associati ai salti condizionati).
3. Il processore possiede una predizione perfetta dei salti condizionati.
4. Due istruzioni non possono risiedere nello stesso pacchetto di esecuzione se una dipende dall'altra (par. 4.10.)
5. Se uno stall diventa necessario, entrambe le istruzioni del pacchetto di esecuzione incriminato devono andare in stall (par. 4.10.).

Svolgendo questi esercizi, si osservi quale sforzo venga richiesto per generare del codice che produca un incremento delle prestazioni vicino a quello ottimale.

**4.31.1 [30] <4.10>** Disegnare una pipeline che mostri come il codice RISC-V riportato sopra venga eseguito sul processore a due vie. Si supponga che il programma esca dal ciclo dopo due iterazioni.

**4.31.2 [10] <4.10>** Qual è il guadagno di prestazioni passando da un processore a una via a un processore a due vie? (Si supponga che il ciclo venga eseguito migliaia di volte.)

**4.31.3 [10] <4.10>** Riordinare/riscrivere il codice RISC-V riportato sopra per ottenere prestazioni migliori su un processore a una via. Suggerimento: utilizzare l'istruzione `beqz x13, TERMINATO` per saltare completamente il ciclo quando  $j = 0$ .

**4.31.4 [20] <4.10>** Riordinare/riscrivere il codice RISC-V riportato sopra per ottenere prestazioni migliori su un processore a due vie (senza srotolare i cicli.)

**4.31.5 [30] <4.10>** Ripetere l'Esercizio 4.31.1, ma questa volta utilizzare il codice ottimizzato dell'Esercizio 4.31.4.

**4.31.6 [10] <4.10>** Qual è l'incremento di velocità nel passaggio dal processore a una via al processore a due vie quando si utilizza il codice ottimizzato degli Esercizi 4.31.3 e 4.31.4?

**4.31.7 [10] <4.10>** Srotolare il codice RISC-V dei cicli dell'Esercizio 4.31.3 in modo tale che ciascuna iterazione del ciclo srotolato contenga due iterazioni del ciclo originario. Poi riordinare/riscrivere il codice srotolato per ottenere prestazioni migliori sul processore a una via. Si può supporre che  $j$  sia un multiplo di 4.

**4.31.8 [20] <4.10>** Srotolare il codice RISC-V dei cicli dell'Esercizio 4.31.4 in modo tale che ciascuna iterazione del ciclo srotolato contenga due iterazioni del ciclo originario. Poi riordinare/riscrivere il codice srotolato per ottenere prestazioni migliori sul processore a due vie. Si può supporre che  $j$  sia un multiplo di 4. (Suggerimento: riorganizzare il ciclo in modo tale che i calcoli compaiano al di fuori del ciclo e alla fine del ciclo. Si può supporre che i valori contenuti nei registri temporanei non siano richiesti dopo il ciclo.)

**4.31.9 [10] <4.10>** Qual è l'incremento di velocità quando si passa da un processore a una via a un processore a due vie quando viene eseguito il codice non srotolato e non ottimizzato dell'Esercizio 4.31.7 e 4.31.8?

**4.31.10 [30] <4.10>** Ripetere gli Esercizi 4.31.8 e 4.31.9, ma questa volta supponendo che il processore a due vie possa eseguire contemporaneamente due istruzioni aritmetico-logiche. (In altre parole, la prima istruzione di un pacchetto può essere di tipo qualsiasi, ma la seconda deve essere un'istruzione aritmetica o un'istruzione logica. Due operazioni sulla memoria non possono essere avviate a esecuzione nello stesso istante.)

**4.32** Questo esercizio riguarda l'efficienza energetica e la sua relazione con le prestazioni. I problemi di questo esercizio ipotizzano il seguente consumo di energia da parte di queste unità funzionali: memoria

istruzioni, registri e memoria dati. Si può supporre che gli altri componenti funzionali dell'unità di elaborazione consumino una quantità trascurabile di energia ("Lettura registro" e "Scrittura registro" si riferiscono solo al register file).

Mem-I	Lettura 1 registro	Scrittura registro	Lettura Mem-D	Scrittura Mem-D
140 pJ	70 pJ	60 pJ	140 pJ	120 pJ

Si supponga che i componenti dell'unità di elaborazione abbiano le seguenti latenze (si può supporre che gli altri componenti dell'unità di elaborazione abbiano latenze trascurabili):

Mem-J	Unità controllo	Lettura/Scrittura registri	ALU	Mem-D L/S
200 ps	150 ps	90 ps	90 ps	250 ps

**4.32.1** [10] <4.3, 4.6, 4.14> Quanta energia viene richiesta per eseguire un'istruzione add in un sistema a singolo ciclo e in una pipeline a cinque stadi?

**4.32.2** [10] <4.6, 4.14> Qual è l'istruzione RISC-V peggiore in termini di energia consumata, e a quanto ammonta l'energia richiesta per la sua esecuzione?

**4.32.3** [10] <4.6, 4.14> Se la riduzione di energia fosse l'elemento più importante, come dovrebbe essere modificata la pipeline? Quale sarebbe la quantità di energia consumata da un'istruzione 1d dopo la modifica?

**4.32.4** [10] <4.6, 4.14> Quali altre istruzioni potrebbero potenzialmente beneficiare dei cambiamenti discussi nell'Esercizio 4.32.3?

**4.32.5** [10] <4.6, 4.14> Qual è l'impatto sulle prestazioni di una CPU dotata di pipeline delle modifiche introdotte nell'Esercizio 4.32.3?

**4.32.6** [10] <4.6, 4.14> Si può eliminare il segnale di controllo di lettura della memoria dati (MemRead) e quindi leggere la memoria dati a ogni ciclo di clock; avremo cioè MemRead = 1 permanentemente. Spiegare perché il processore funzionerebbe ancora correttamente dopo questa modifica. Qual è l'impatto di questa modifica sul periodo di clock e sull'energia consumata?

**4.33** Quando i chip di silicio vengono fabbricati, difetti nei materiali (per es. nel silicio) ed errori di fabbricazione possono fare sì che i circuiti risultino difettati. Un difetto molto comune è chiamato "cross-talk" (interferenza) e si verifica quando un segnale che transita su una pista influenza un altro segnale che transita su un'altra pista. Un tipo particolare di malfunzionamento dovuto al cross-talk si manifesta

quando uno dei due segnali è trasportato da una pista erroneamente collegata a un valore logico costante, per esempio all'alimentazione elettrica. In questo caso si ottiene un segnale che rimane sempre bloccato a 0 o a 1, indipendentemente da quale sia il valore logico che dovrebbe assumere. I problemi successivi si riferiscono al bit 0 del Registro di Scrittura del Register File di Figura 4.21.

**4.33.1** [10] <4.3, 4.4> Si supponga che la verifica del funzionamento del processore venga effettuata scrivendo dei dati, a piacimento, nel PC, nei registri, nella memoria dati e nella memoria istruzioni, per poi eseguire una sola istruzione. Al termine dell'esecuzione verrà letto il contenuto del PC, delle memorie e dei registri, e questi valori verranno esaminati per determinare se si sia verificato un qualche malfunzionamento. Progettare un test, cioè determinare un valore opportuno da scrivere nel PC, nelle memorie e nei registri, che consenta di determinare se il bit analizzato è bloccato a 0.

**4.33.2** [10] <4.3, 4.4> Risolvere l'Esercizio 4.33.1 nel caso in cui si voglia determinare se il bit analizzato è bloccato a 1. È possibile utilizzare un unico test per verificare se il bit è bloccato a 0 o a 1? Se sì, spiegare come; se non è possibile, spiegare perché.

**4.33.3** [10] <4.3, 4.4> Se sapessimo che in un processore il bit analizzato è bloccato a 1, potremmo comunque utilizzare il processore? Per potere utilizzare comunque questo processore occorrerebbe convertire i programmi che vengono eseguiti su un RISC-V normale in programmi che possano essere eseguiti su questo processore, tenendo conto del malfunzionamento. Si supponga che ci sia abbastanza memoria istruzioni per aggiungere delle istruzioni e abbastanza memoria dati libera per memorizzare altri dati.

**4.33.4** [10] <4.3, 4.4> Risolvere l'Esercizio 4.33.1 supponendo che il malfunzionamento da identificare sia sul segnale di controllo MemRead, il quale diventa 0 se il segnale branch è anch'esso 0; non ci sono malfunzionamenti in tutti gli altri casi.

**4.33.5** [10] <4.3, 4.4> Risolvere l'Esercizio 4.33.1 supponendo che il malfunzionamento da identificare sia sul segnale di controllo MemRead, il quale diventa 1 se il segnale di controllo RegRd diventa 1; non ci sono malfunzionamenti in tutti gli altri casi. Suggerimento: questo problema richiede delle conoscenze sui sistemi operativi. Identificate cosa possa causare un errore di segmentazione della memoria (*segmentation fault*).

## Risposte alle domande di autovalutazione

**Paragrafo 4.1, pagina 209** – 3 o 5: unità di controllo, unità di elaborazione, memoria. Mancano input e output.

**Paragrafo 4.2, pagina 213** – Falso. Gli elementi di stato sincronizzati sui fronti rendono possibile e non ambigua la lettura e scrittura simultanea.

**Paragrafo 4.3, pagina 219** – 1. a. 2. c.

**Paragrafo 4.4, pagina 230** – Si, i segnali di Branch e ALUOP0 sono identici. Inoltre si può utilizzare la flessibilità offerta dai bit il cui valore è indifferente per raggruppare assieme altri segnali. ALUSrc e MemReg possono essere raggruppati in un unico segnale se si impostano a 1 e 0 rispettivamente i due bit associati al segnale MemtoReg. ALUOp1 e MemtoReg possono essere raggruppati dopo avere negato uno dei due, se si impostano i bit indifferenti di MemtoReg a 1. Non è necessario un inverter: è sufficiente cambiare l'ordine degli input nel multiplexer associato a MemtoReg!

**Paragrafo 4.5, pagina 243** – 1. Stallo sul risultato della ld. 2. Evitare lo stallo nella terza istruzione per l'hazard sul dato per la lettura dopo la scrittura su x11 propagando il risultato della add. 3. Non occorrono stalli, anche senza propagazione.

**Paragrafo 4.6, pagina 257** – Le affermazioni 2 e 4 sono corrette; le altre due sono errate.

**Paragrafo 4.8, pagina 278** – 1. Predetto non preso. 2. Predetto preso. 3. Predizione dinamica.

**Paragrafo 4.9, pagina 285** – La prima istruzione, dato che viene eseguita prima delle altre.

**Paragrafo 4.10, pagina 297** – 1. Entrambi. 2. Entrambi. 3. Software. 4. Hardware. 5. Hardware. 6. Hardware. 7. Entrambi. 8. Hardware. 9. Entrambi.

**Paragrafo 4.12, pagina 306** – Le prime due affermazioni sono false, mentre le ultime due affermazioni sono vere.

# 5

## Grande e veloce: la gerarchia delle memorie

*Idealmente si desidererebbe possedere una memoria indefinitamente grande tale che ogni particolare ... parola risulti immediatamente disponibile ... Siamo ... costretti a riconoscere la possibilità di costruire una gerarchia delle memorie, ognuna delle quali abbia una capacità maggiore rispetto alla precedente, ma sia accessibile meno velocemente.*

**A.W. Burks, H.H. Goldstine e J. von Neumann,**  
*Discussione preliminare sul progetto logico di uno strumento elettronico di computazione, 1946*

### 5.1 | Introduzione

Fin dal principio dell'era dell'informatica i programmati hanno sperato di avere a disposizione una quantità illimitata di memoria che fosse al contempo veloce. Tutti gli argomenti che saranno illustrati in questo capitolo sono incentrati sulle tecniche che permettono di dare ai programmati quest'illusione, ossia di poter usufruire di una memoria che sia contemporaneamente veloce e di dimensioni illimitate. Prima di esaminare come venga effettivamente realizzata questa illusione, consideriamo una semplice analogia che illustra i principi chiave e i meccanismi che vengono utilizzati.

Immaginiamo uno studente che debba scrivere una relazione sulle tappe principali della storia dell'architettura dei calcolatori. Lo studente è seduto a una scrivania della biblioteca di ingegneria o di informatica e sta esaminando una serie di testi che ha prelevato dagli scaffali. A un certo punto si accorge che molti dei calcolatori che deve descrivere sono illustrati nei libri che ha davanti, ma non trova alcuna informazione sull'EDSAC; così, si alza e torna a guardare sugli scaffali, terminando la ricerca solo dopo aver trovato un testo dedicato ai primi calcolatori inglesi che descrive l'EDSAC. Dopo che lo studente ha selezionato un certo numero di libri e li ha riposti sulla sua scrivania, ci sarà un'alta probabilità che la maggior parte degli argomenti a cui è interessato sia trattata in questi libri; egli passerà perciò la maggior parte del tempo a consultare i libri raccolti, senza dover tornare ogni volta agli scaffali. Avere a disposizione molti libri sulla scrivania consente allo studente di risparmiare molto tempo rispetto

ad avere un solo libro e doversi continuamente alzare per riporlo sullo scaffale e prenderne un altro.

Lo stesso principio ci permette di creare l'illusione di avere a disposizione in un calcolatore una memoria di grandi dimensioni a cui poter accedere con la stessa velocità con cui si accede a una memoria di piccole dimensioni. Proprio come lo studente non ha bisogno di dover consultare, con la stessa probabilità, tutti i libri della biblioteca, un programma non deve accedere con la stessa probabilità a tutte le sue istruzioni e a tutti i suoi dati contemporaneamente. Se così non fosse, sarebbe impossibile pensare di rendere veloce l'accesso alla memoria e, allo stesso tempo, avere una memoria di grandi dimensioni, così come è impossibile per lo studente mettere tutti i libri della biblioteca sulla sua scrivania e pensare ancora di riuscire a trovare velocemente le informazioni che gli interessano.

Il **principio di località** sta alla base del comportamento dei programmi in un calcolatore ed è del tutto simile al modo di cercare informazioni in una biblioteca. Questo principio afferma che un programma, in un certo istante di tempo, accede soltanto a una porzione relativamente piccola del suo spazio di indirizzamento, proprio come lo studente accede solo a una piccola porzione dei libri nella biblioteca. Esistono due diversi tipi di località:

- **località temporale** (località nel tempo): quando si fa riferimento a un elemento, c'è la tendenza a fare riferimento allo stesso elemento dopo poco tempo; se lo studente ha appena portato un libro sulla scrivania per consultarlo, è probabile che abbia bisogno di consultarla ancora poco dopo;
- **località spaziale** (località nello spazio): quando si fa riferimento a un elemento, c'è la tendenza a fare riferimento poco dopo ad altri elementi che hanno l'indirizzo vicino a esso. Per esempio, quando lo studente ha prelevato il libro sui primi calcolatori costruiti in Inghilterra per trovare informazioni sull'EDSAC, ha anche notato che il testo a fianco riguardava le prime calcolatrici meccaniche; così, ha deciso di portare sulla sua scrivania anche questo libro, in cui più tardi potrà trovare delle informazioni utili. I bibliotecari dispongono i libri su uno stesso argomento sullo stesso scaffale per aumentare la località spaziale. Più avanti in questo capitolo vedremo come la località spaziale venga sfruttata dalle gerarchie delle memorie.

La località emerge in modo naturale dalle strutture di controllo semplici e tipiche dei programmi, proprio come l'accesso ai libri della biblioteca mostra un comportamento naturalmente locale. Per esempio, la maggior parte dei programmi contiene dei cicli; le istruzioni e i dati che si trovano all'interno di queste strutture di controllo vengono letti ripetutamente dalla memoria, dimostrando così un elevato livello di località temporale. Dato che le istruzioni di un programma normalmente vengono caricate in sequenza dalla memoria, i programmi presentano anche un'elevata località spaziale. Anche l'accesso ai dati è intrinsecamente caratterizzato da una località spaziale; per esempio, è naturale che l'accesso agli elementi consecutivi di un vettore o di una struttura dati abbia un elevato grado di località spaziale.

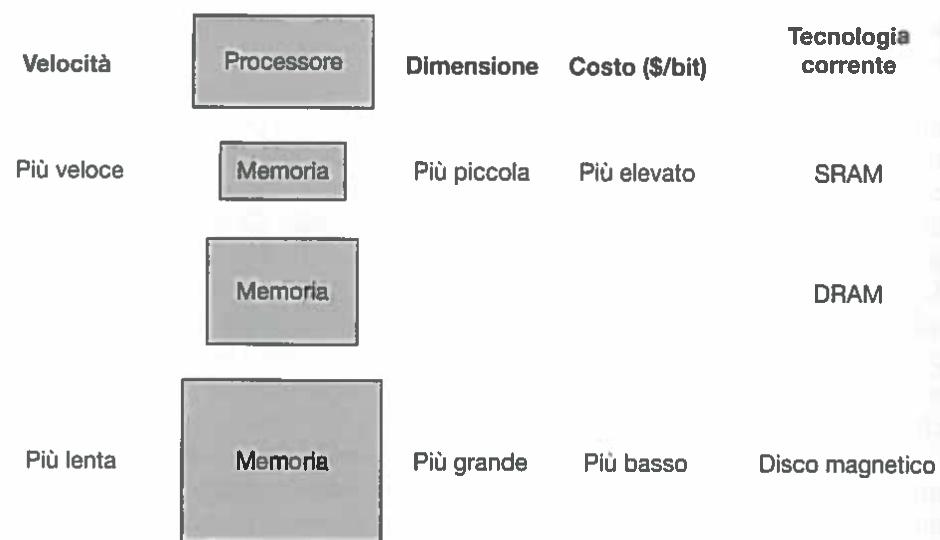
Il principio di località viene sfruttato strutturando la memoria di un calcolatore in forma gerarchica. La **gerarchia delle memorie** consiste in un insieme di livelli di memoria, ciascuno caratterizzato da una diversa velocità e dimensione: a parità di capacità, le memorie più veloci hanno un costo più elevato per singolo bit di quelle più lente, perciò, di solito, sono più piccole.

La Figura 5.1 mostra che la memoria più veloce è posta vicino al processore e quella più lenta, meno costosa, è posizionata più lontano. L'obiettivo è di fornire all'utente una quantità di memoria pari a quella disponibile nella tecnologia più economica, consentendo allo stesso tempo una velocità di accesso pari a quella garantita dalla memoria più veloce.

**Località temporale:** principio secondo cui se si accede a una determinata locazione di memoria, è molto probabile che vi si acceda di nuovo dopo poco tempo.

**Località spaziale:** principio secondo cui se si accede a una determinata locazione di memoria, è molto probabile che si acceda alle locazioni vicine ad essa dopo poco tempo.

**Gerarchia delle memorie:** una struttura che utilizza più livelli di memoria; all'aumentare della distanza dal processore, crescono sia la dimensione sia il tempo di accesso alla memoria.



**Figura 5.1 Struttura base di una gerarchia delle memorie.** Realizzando il sistema di memoria in maniera gerarchica, l'utente ha l'impressione di utilizzare una memoria grande come quella del livello contenente la memoria di dimensioni maggiori, ma con una velocità di accesso pari a quella della memoria realizzata con la tecnologia più veloce. Le memorie flash hanno sostituito i dischi in molti dispositivi embedded e potrebbero portare in breve tempo all'introduzione di un nuovo livello di memoria nei calcolatori desktop e nei server (par. 5.2).

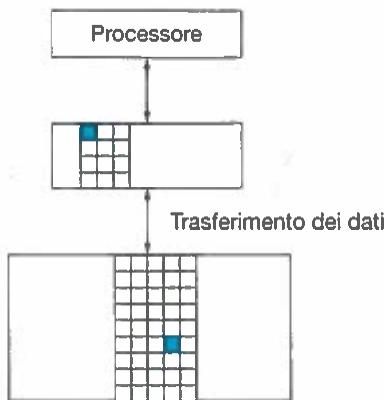
Anche i dati sono organizzati in modo gerarchico: un livello più vicino al processore contiene in generale un sottoinsieme dei dati memorizzati in ognuno dei livelli sottostanti, e tutti i dati si trovano memorizzati nel livello più basso della memoria. Si consideri ancora l'analogia con la biblioteca: i libri sulla scrivania rappresentano il sottoinsieme dei libri della biblioteca su cui lo studente sta lavorando; a sua volta la biblioteca è un sottoinsieme di tutte le biblioteche dell'università. Quanto più ci si allontana dal processore, tanto più aumenta il tempo necessario per accedere ai dati, proprio come accade per la consultazione dei libri in una gerarchia di biblioteche universitarie.

**Blocco:** detto anche **linea**, è l'unità minima di informazione che può essere presente o assente in una cache.

**Hit rate:** detto anche **frequenza delle hit**, è la frazione degli accessi alla memoria nei quali la parola cercata viene trovata in un certo livello della gerarchia.

Una gerarchia delle memorie può essere composta da più livelli, ma i dati vengono di volta in volta trasferiti solo tra due livelli vicini; possiamo perciò concentrare la nostra attenzione su due soli livelli di memoria. Quello superiore, più vicino al processore, è più piccolo e veloce del livello inferiore, poiché utilizza una tecnologia più costosa. La più piccola quantità di informazione che può essere presente o assente in questa gerarchia su due livelli è denominata **blocco**, o **linea**, come illustrato in Figura 5.2; nell'analogia con la biblioteca, il blocco di informazione corrisponde a un libro.

Se il dato richiesto dal processore è contenuto in uno dei blocchi presenti nel livello superiore, si dice che la richiesta ha avuto successo e si denota questo evento con il termine inglese *hit* (letteralmente, "colpito"); una hit corrisponde allo studente che trova le informazioni che cerca in uno dei libri che ha a disposizione sulla scrivania. Se il dato non viene trovato nel livello superiore della gerarchia, si dice che la richiesta fallisce e si indica questo evento con il termine *miss* (letteralmente, "mancato"). In questo secondo caso, per trovare il blocco che contiene il dato richiesto occorre accedere al livello inferiore della gerarchia. Proseguendo con l'analogia della biblioteca, lo studente si alza dalla scrivania e va a cercare il libro che gli serve sugli scaffali. La frequenza di hit, chiamata **hit rate** (letteralmente, "frequenza delle hit"), è la frazione degli accessi alla memoria nei quali il dato desiderato è stato trovato nel livello superiore; spesso l'hit rate viene utilizzato come indice delle prestazioni della gerarchia delle memorie. La frequenza delle miss, indicata con il termine



**Figura 5.2** Ogni coppia di livelli in una gerarchia delle memorie può essere vista come formata da un livello superiore e da un livello inferiore. All'interno di ciascun livello, l'unità di informazione che è presente o assente è chiamata *blocco*, o *linea*. Di solito trasferiamo un intero blocco quando copiamo qualcosa fra i due livelli.

**miss rate**, è pari a  $(1 - \text{hit rate})$  ed è la frazione degli accessi alla memoria nei quali il dato desiderato non è stato trovato nel livello superiore.

Poiché il motivo principale che ha condotto all'organizzazione gerarchica delle memorie è l'aumento delle prestazioni, la velocità degli accessi è importante sia in caso di successo sia in caso di fallimento. Il **tempo di hit** è il tempo di accesso al livello superiore della gerarchia delle memorie e comprende anche il tempo necessario a stabilire se il tentativo di accesso si risolve in un successo o in un fallimento, cioè se produce una hit o una miss. Il tempo di hit è analogo al tempo necessario allo studente per passare in rassegna i libri che ha sulla scrivania. La **penalità di miss** è il tempo necessario a sostituire un blocco del livello superiore con un nuovo blocco, caricato dal livello inferiore della gerarchia, e a trasferire i dati contenuti in questo blocco al processore. La penalità di miss è il tempo che serve per prendere un nuovo libro dagli scaffali e metterlo sulla scrivania. Siccome il livello superiore è più piccolo ed è costruito utilizzando componenti più veloci, il tempo di hit sarà molto inferiore al tempo necessario ad accedere al secondo livello della gerarchia, che rappresenta la componente principale della penalità di miss. Proseguendo con l'analogia della biblioteca, il tempo richiesto per esaminare i libri sulla scrivania è molto inferiore a quello impiegato per alzarsi e andare a prelevare un altro libro dagli scaffali della biblioteca.

Come vedremo più avanti, i concetti utilizzati per realizzare i sistemi di memoria influenzano molti altri aspetti di un calcolatore, come le modalità di gestione della memoria e delle periferiche di I/O che possono essere adottate dal sistema operativo, la generazione del codice da parte dei compilatori e perfino la modalità di utilizzo del calcolatore da parte delle diverse applicazioni. Ovviamente, poiché tutti i programmi spendono molto del loro tempo di esecuzione in accessi alla memoria, il sistema di memoria è uno degli elementi cruciali che determinano le prestazioni del calcolatore nel suo complesso. La dipendenza delle prestazioni dall'utilizzo delle gerarchie delle memorie ha fatto sì che i programmati – che inizialmente consideravano la memoria come un contenitore dal quale prendere i dati con lo stesso tempo di accesso – debbano ora comprendere i meccanismi di funzionamento della memoria per ottenere prestazioni più elevate. Illustreremo più avanti quanto sia importante comprendere il funzionamento della gerarchia delle memorie con degli esempi: in Figura 5.18 e nel paragrafo 5.14 mostreremo come raddoppiare le prestazioni della moltiplicazione di matrici in doppia precisione.

**Miss rate:** detto anche **frequenza delle miss**, è la frazione degli accessi alla memoria nei quali la parola cercata non viene trovata in un certo livello della gerarchia.

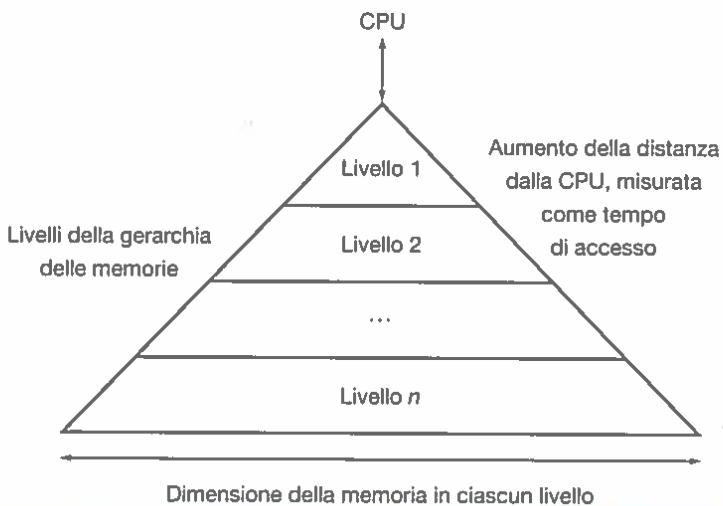
**Tempo di hit:** il tempo richiesto per accedere a un livello della gerarchia delle memorie, compreso il tempo necessario a determinare se l'accesso ha avuto successo o meno, cioè se ha prodotto una hit o una miss.

**Penalità di miss:** il tempo richiesto per caricare un blocco dal livello inferiore della gerarchia delle memorie, comprendente il tempo per accedere al blocco, portarlo al livello superiore (che ha generato la miss) e scriverlo in quel livello di memoria.

Proprio perché i sistemi di memoria sono così critici per le prestazioni, sono stati oggetto di una grande attenzione da parte dei progettisti, i quali hanno sviluppato meccanismi molto sofisticati per migliorarli. In questo capitolo esamineremo le principali problematiche che stanno alla base dei sistemi di memoria; abbiamo scelto di adottare diverse semplificazioni e astrazioni per rendere i concetti il più possibile concisi e comprensibili.

### QUADRO D'INSIEME

I programmi mostrano un comportamento locale sia dal punto di vista temporale, ossia hanno la tendenza a riutilizzare i dati che sono stati caricati di recente dalla memoria, sia da quello spaziale, cioè tendono a richiedere dati che si trovano in prossimità di quelli che sono stati appena letti o scritti. Le gerarchie delle memorie sfruttano la località temporale mantenendo i dati utilizzati più di recente nel livello più vicino al processore e sfruttano la località spaziale spostando verso i livelli superiori della gerarchia blocchi composti da più parole aventi indirizzi di memoria adiacenti. La Figura 5.3 mostra che una gerarchia delle memorie impiega la tecnologia più veloce con una capacità ridotta per realizzare i livelli più vicini al processore; in questo modo, gli accessi ai dati che si trovano nei livelli più alti, quando hanno successo, possono essere effettuati velocemente. Gli accessi che falliscono richiedono di caricare i dati dai livelli inferiori della gerarchia, che hanno dimensioni maggiori ma sono caratterizzati da un maggior tempo d'accesso. Se la frequenza di hit è sufficientemente alta, la gerarchia delle memorie ha un tempo medio di accesso vicino a quello del livello più alto (e più veloce) e una dimensione complessiva pari a quella del livello più basso (che ha una capacità maggiore). In molti sistemi, la memoria è una gerarchia vera e propria, nel senso che un dato non può essere presente nel livello  $i$  a meno che non si trovi anche nel livello  $i + 1$ . ■



**Figura 5.3** Questo schema mostra la struttura di una gerarchia delle memorie: al crescere della distanza dal processore, cresce la dimensione della memoria. Questa struttura, con gli appropriati meccanismi di gestione, permette al processore di avere un tempo di accesso che è determinato principalmente dal livello 1 della gerarchia, e allo stesso tempo di avere una memoria la cui dimensione è pari a quella del livello  $n$ . Il modo in cui si crea questa illusione è l'oggetto di questo capitolo. Sebbene il disco locale costituisca solitamente il livello inferiore della gerarchia, alcuni sistemi utilizzano nastri o file server su reti LAN (*Local Area Network*) come livelli addizionali della gerarchia.

## Autovalutazione

Quali delle seguenti affermazioni sono generalmente vere?

1. Le gerarchie delle memorie sfruttano la località temporale.
2. Il valore restituito da una lettura dipende da quali blocchi si trovano nella cache.
3. La maggior parte del costo di una gerarchia delle memorie appartiene al livello più alto.
4. La maggior parte della capacità di una gerarchia delle memorie si trova nel livello più basso.

## 5.2 | Tecnologie delle memorie

Al giorno d'oggi vengono utilizzate quattro principali tecnologie per costruire le gerarchie delle memorie. La memoria principale viene realizzata utilizzando DRAM (*Dynamic Random Access Memory*, memoria dinamica ad accesso casuale), mentre i livelli più vicini al processore (cache) sono composti da SRAM (*Static Random Access Memory*, memoria statica ad accesso casuale). La DRAM ha un costo per bit più basso della SRAM, ma il suo accesso è decisamente più lento. La differenza di costo dipende dal fatto che le memorie DRAM utilizzano meno transistor per ogni bit da memorizzare e riescono così a raggiungere una capacità superiore a parità di porzione di silicio utilizzata; la differenza di velocità è dovuta invece ad altri fattori, descritti nel paragrafo A.9 dell'Appendice A . La terza tecnologia è quella delle memorie flash. Questa memoria non volatile viene utilizzata nei dispositivi mobili. La quarta tecnologia viene impiegata per implementare il livello di memoria più capiente e lento della gerarchia, ed è solitamente rappresentata dai dischi magnetici. Il tempo di accesso e il costo per bit variano ampiamente per ciascuna di queste tecnologie, come si vede dalla seguente tabella che riporta i costi e i tempi di accesso tipici, riferiti all'anno 2012:

Tecnologia di memoria	Tempo di accesso tipico	\$ per GiB nel 2012
Memoria a semiconduttore SRAM	0,5-2,5 ns	\$500-\$1000
Memoria a semiconduttore DRAM	50-70 ns	\$10-\$20
Memoria flash a semiconduttore	5000-50 000 ns	\$0,75-\$1
Dischi magnetici	5 000 000-20 000 000 ns	\$0,05-\$0,1

Descriviamo ora le diverse tecnologie.

### Tecnologia SRAM

Le SRAM sono semplicemente dei circuiti integrati organizzati come vettori di memoria che (di solito) hanno una sola porta di accesso che può fornire sia la lettura sia la scrittura. Le SRAM hanno uno stesso tempo di accesso per tutti i dati, anche se i tempi di accesso in lettura e scrittura possono essere diversi.

Le SRAM non hanno bisogno di essere rinfrescate (*refresh* è il termine inglese di uso comune), per cui il loro tempo di accesso è molto vicino al periodo del clock; utilizzano da sei a otto transistor per bit per evitare che l'informazione



possia essere disturbata quando viene letta. Le SRAM hanno bisogno di una potenza minima per mantenere la carica quando si trovano in modalità stand-by.

In passato, la maggior parte dei PC e dei server utilizzava chip di SRAM diversi per i diversi livelli di cache (primo, secondo e, in alcuni casi, terzo livello). Oggi, grazie alla **Legge di Moore**, le cache di tutti i livelli sono integrate nel chip del processore, per cui il mercato delle SRAM singole è praticamente svanito.

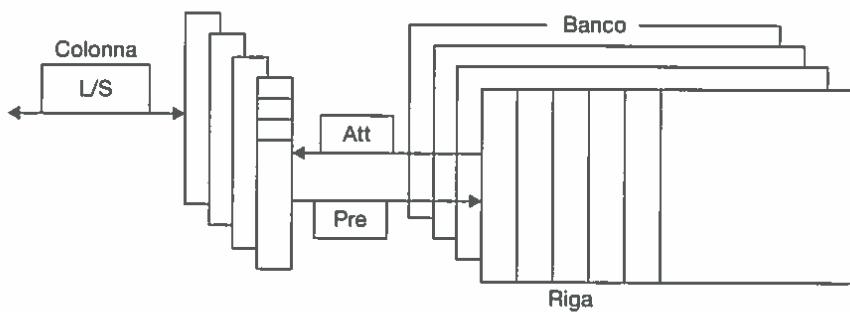
## Tecnologia DRAM

In una memoria SRAM il dato rimane memorizzato per tutto il tempo in cui l'alimentazione è attiva. In una memoria RAM *dinamica* (DRAM), invece, il dato viene memorizzato come carica in un condensatore e un solo transistor è sufficiente per leggere il dato o per sovrascriverlo. Dato che le DRAM utilizzano un solo transistor per bit memorizzato, sono molto più dense e il costo per bit è inferiore a quello delle SRAM. Tuttavia, dato che nelle DRAM l'informazione viene memorizzata in un condensatore, non rimane indefinitamente e occorre rinfrescarla periodicamente. Questo è il motivo per cui questo tipo di memoria è chiamato dinamico, per distinguere dalle memorie statiche delle celle SRAM.

Per rinfrescare una cella, occorre semplicemente leggerne il contenuto e riscriverlo. La carica può essere mantenuta per pochi millisecondi, per cui, se dovessimo leggere e riscrivere i diversi bit di una DRAM uno per uno, passeremmo tutto il tempo a fare il refresh della DRAM, senza avere il tempo per accedervi. Fortunatamente, le DRAM utilizzano un'architettura a due livelli che consente di rinfrescare un'intera *riga*, mediante un ciclo di lettura seguito immediatamente da un ciclo di riscrittura.

La Figura 5.4 mostra l'organizzazione interna di una DRAM, mentre la Figura 5.5 mostra come siano cambiati negli anni la densità, il costo e il tempo di accesso delle DRAM.

L'organizzazione per righe aiuta nel refresh della memoria e aiuta anche nelle prestazioni: le DRAM utilizzano un buffer per memorizzare una riga per consentire accessi ripetuti. Il buffer si comporta come una SRAM: cambiando opportunamente l'indirizzo si può accedere ai diversi bit del buffer fino a quando non si deve accedere ai dati di un'altra riga. Questa funzionalità migliora significativamente il tempo di accesso, dato che il tempo di accesso a bit diversi della stessa riga è molto più basso. Anche costruire il chip più ampio migliora l'ampiezza di banda del chip. Quando la riga si trova nel buffer, il suo conte-



**Figura 5.4 Organizzazione interna di una DRAM.** Le DRAM moderne sono organizzate in banchi, tipicamente quattro per le DDR3. Ciascun banco consiste in una serie di righe. Inviano un comando di precaricamento (Pre) si apre o chiude un banco. L'indirizzo di una riga viene inviato con un comando di attivazione (Att), che attiva il trasferimento del contenuto della riga in un buffer. Quando la riga si trova nel buffer, il suo contenuto può essere trasferito specificando gli indirizzi successivi di gruppi di bit qualsiasi sia l'ampiezza della DRAM (tipicamente 4, 8 o 16 bit per le DDR3) o specificando l'indirizzo di partenza e che si intende trasferire l'intero blocco. Ciascun comando e il comando di trasferimento di un blocco sono sincronizzati con il clock.

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	\$4 000	60 ns	10 ns
2000	256 Mebibit	\$1 000	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

**Figura 5.5** La dimensione delle DRAM è quadruplicata approssimativamente ogni tre anni fino al 1996, mentre è cresciuta molto meno negli anni successivi. Il miglioramento del tempo di accesso è stato meno marcato ma continuo. L'andamento del costo segue da vicino l'incremento della densità, sebbene il costo sia spesso influenzato da altri elementi quali la disponibilità e la richiesta. Il costo per Gibibyte non tiene conto dell'inflazione.

nuto può essere trasferito specificando gli indirizzi successivi di gruppi di bit qualsiasi sia l'ampiezza della DRAM (tipicamente 4, 8 o 16 bit) o specificando l'indirizzo di partenza e che si intende trasferire l'intero blocco.

Per migliorare ulteriormente l'interfaccia con i processori, alle DRAM è stato aggiunto il clock, per cui queste memorie sono chiamate più propriamente DRAM sincrone o SDRAM (*Synchronous DRAM*). Il vantaggio delle SDRAM sta nel fatto che il clock elimina il tempo necessario a sincronizzare la memoria con il processore. L'aumento di velocità delle DRAM sincrone proviene dalla possibilità di trasferire gruppi di bit adiacenti a raffica (*in burst*) senza dovere specificare altri indirizzi. La versione attuale più veloce è chiamata DDR (*Double Data Rate* – frequenza doppia dei dati) SDRAM. Questo termine indica che il trasferimento dei dati avviene sia sul fronte di salita sia di discesa del clock, ottenendo così il doppio della larghezza di banda rispetto alla banda calcolata a partire dall'ampiezza dei dati e dalla frequenza del clock. La versione più recente di questa tecnologia viene chiamata DDR4. Una DRAM DDR4-3200 può effettuare 3200 milioni di trasferimenti al secondo, che corrisponde a una frequenza di clock di 1600 MHz.

Sostenere una tale larghezza di banda richiede un'organizzazione particolarmente intelligente all'interno della DRAM. Invece di un buffer di riga più veloce, si possono organizzare internamente le DRAM per leggere o scrivere da diversi *banchi*, ciascuno con il suo buffer di riga. Inviare un indirizzo a più banchi consente di leggere o scrivere simultaneamente. Per esempio, con quattro banchi occorre un unico tempo di accesso alla memoria, dopodiché si possono indirizzare sequenzialmente a rotazione i quattro banchi, ottenendo così il quadruplo della larghezza di banda. Questo schema di accesso sequenziale a rotazione viene chiamato *indirizzamento con interleaving* (letteralmente, "intreccio").

Sebbene i dispositivi mobili come gli iPad (vedi Cap. 1) utilizzino DRAM singole, per i server vengono comunemente vendute delle schedine chiamate DIMM (*Dual Inline Memory Modules* – moduli di memoria in linea doppi, cioè fronte e retro). Una DIMM tipicamente contiene 4-16 DRAM ed è organizzata in modo da offrire un'ampiezza di 8 byte per i server. Una DIMM che utilizza

le SDRAM DDR4-3200 può arrivare a trasferire  $8 \times 3200 = 25\,600$  megabyte al secondo e prende il nome dalla sua ampiezza di banda; PC25600. Dato che una DIMM può avere così tanti chip che solo alcuni sono interessati da ogni trasferimento, occorre un termine per indicare il sottoinsieme di chip di una DRAM che condividono le linee di indirizzamento. Per evitare confusione con i termini che indicano la struttura interna, ovvero righe e banchi, utilizzeremo il termine *rango della memoria* (*memory rank*) per indicare questo sottoinsieme di chip di una DIMM.

**Approfondimento.** Un modo per misurare le prestazioni del sistema di memoria che sta dietro le cache è il *benchmark Stream* (cioè con flusso di dati continuo) [McCalpin, 1995], che misura le prestazioni di operazioni effettuate su vettori molto lunghi. Queste operazioni non presentano località temporale e accedono a vettori di dimensioni maggiori di quelle delle cache del calcolatore che viene esaminato.

## Memorie flash

Le memorie flash sono un tipo di memoria a sola lettura, cancellabile elettricamente e programmabile (EEPROM, *Electrically Erasable Programmable Read-Only Memory*).

A differenza dei dischi e delle DRAM, ma come le altre tecnologie EEPROM, i bit delle memorie flash si deteriorano dopo un certo numero di scritture. Per far fronte a questa limitazione, la maggior parte dei dispositivi che utilizzano memorie flash contiene un controllore che distribuisce le scritture consentite rimappando i blocchi di memoria che sono stati scritti più spesso sui blocchi che sono stati scritti meno di frequente. Questa tecnica è chiamata *livellamento dell'usura* (*wear leveling*). Adottando questa tecnica è difficile che nei dispositivi mobili si possa superare il limite sul numero massimo di scritture per singola cella di memoria flash. I controllori che implementano il livellamento dell'usura fanno diminuire le prestazioni teoriche di una memoria flash, ma sono necessari, a meno che non sia il software di controllo ad alto livello a verificare l'usura dei blocchi di memoria. I controllori delle memorie flash possono anch'essi migliorare le prestazioni, identificando le celle di memoria costruite male e quindi mal funzionanti.

## Memorie a disco

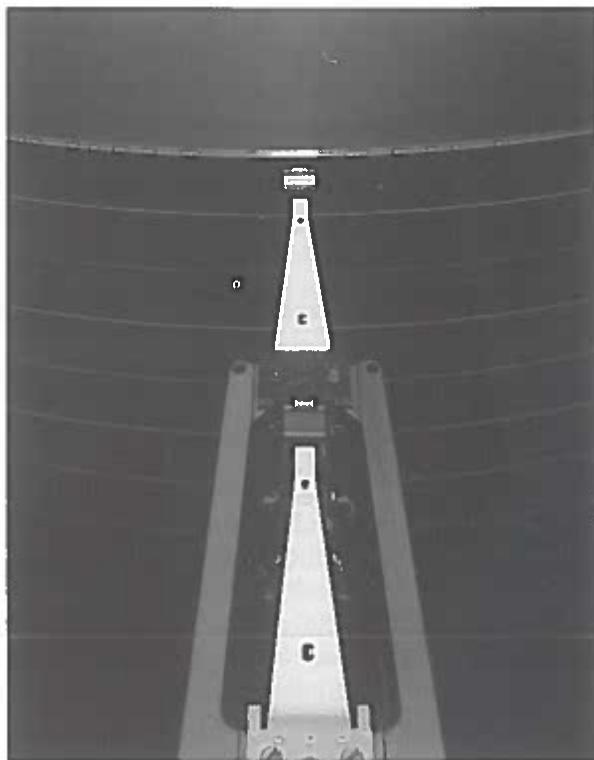
Come mostrato in Figura 5.6, un disco magnetico è formato da un gruppo di dischi, detti *piatti*, che ruotano solidali a una velocità compresa tra i 5400 e i 15 000 giri al minuto (GPM). Ciascun piatto metallico è ricoperto su entrambe le facce da materiale magnetico registrabile simile a quello delle cassette audio o delle videocassette. Per leggere e scrivere dati su un disco rigido, poco al di sopra di ognuna delle due superfici di ogni disco è posizionato un *braccio mobile* contenente una piccola bobina elettromagnetica chiamata *testina di lettura/scrittura*. L'intera unità è sigillata per controllare i parametri fisici dell'ambiente interno; ciò consente anche di avvicinare molto di più le testine alla superficie del disco.

La superficie di ogni disco è divisa in cerchi concentrici chiamati **tracce** (*tracks*): ci sono tipicamente decine di migliaia di tracce per ogni lato di un piatto<sup>1</sup>. Ogni traccia è a sua volta suddivisa in **settori** (*sectors*), che sono l'unità di memorizzazione delle informazioni: ciascuna traccia può contenere migliaia

**Traccia:** uno delle migliaia di cerchi concentrici che compongono la superficie di un disco magnetico.

**Settore:** uno dei segmenti che compongono una traccia di un disco magnetico; un settore è la più piccola quantità di informazione letta o scritta sul disco.

<sup>1</sup> La struttura del singolo piatto e la testina di lettura assomigliano ai giradischi degli anni Settanta [N.d.T.].



**Figura 5.6** Disco rigido costituito da 10 piatti paralleli e dalle testine di lettura/scrittura. Il diametro dei dischi è oggi di 2,5 o 3,5 pollici e ci sono tipicamente uno o due piatti per ogni disco.

di settori. Ogni settore ha tipicamente una dimensione variabile da 512 a 4096 byte. La sequenza di informazioni registrata sul materiale magnetico è costituita dal numero del settore, seguito da uno spazio vuoto di separazione, dalle informazioni memorizzate nel settore, comprendenti il codice di correzione degli errori (par. 5.5), quindi da un altro spazio vuoto di separazione, dal numero del settore successivo ecc.

Le testine di lettura/scrittura associate a ognuna delle due superfici dei piatti sono collegate assieme e si muovono congiuntamente, così che ogni testina si trova in corrispondenza della stessa traccia su tutte le superfici. Il termine *cilindro* viene utilizzato per riferirsi a tutte le tracce dei diversi piatti che si trovano sotto le testine in un certo istante di tempo.

Per accedere ai dati, il sistema operativo deve guidare il disco attraverso un processo in tre passi. Il primo passo consiste nel posizionare la testina sulla traccia giusta. Questa operazione è chiamata **seek** (ricerca), e il tempo necessario per spostare la testina sopra la traccia desiderata viene chiamato *tempo di ricerca (seek time)*.

I costruttori di dischi riportano nei loro manuali il tempo minimo di ricerca, quello massimo e quello medio. I primi due sono facili da misurare, ma quello medio è aperto a diverse interpretazioni, perché dipende da quanto la testina si deve spostare. I costruttori hanno deciso di calcolare il tempo medio di ricerca come la somma dei tempi di ricerca di tutti i possibili settori diviso per il numero dei possibili settori. Il tempo medio di ricerca che viene di solito dichiarato dai costruttori è compreso tra i 3 e i 13 ms, ma, a seconda dell'applicazione e dell'ordine con cui arrivano le richieste di accesso a disco, il tempo medio effettivo di una ricerca si può ridurre fino a un valore compreso tra il 25 e il 33% del tempo di ricerca medio dichiarato, grazie alla località degli accessi.

**Seek:** il processo di posizionamento della testina di lettura e scrittura su una data traccia del disco.

**Latenza di rotazione:** detta anche **ritardo di rotazione**, è il tempo necessario perché, ruotando il disco, il settore desiderato passi sotto la testina di lettura/scrittura; in genere si suppone che sia pari a metà del tempo di rotazione.

al disco. Questa località può derivare sia dal fatto che gli accessi sono per parti adiacenti di uno stesso file, sia dal riordinamento degli accessi imposto dal sistema operativo.

Quando la testina di lettura/scrittura ha raggiunto la traccia corretta, deve aspettare che il settore desiderato passi sotto. Questo tempo di attesa viene chiamato **latenza di rotazione** o **ritardo di rotazione**. La latenza media corrisponde a mezzo giro del disco. I dischi ruotano a una velocità compresa tra i 5400 e i 15 000 GPM (giri/minuto). La latenza di rotazione media a 5400 GPM sarà:

$$\text{Latenza di rotazione media} = \frac{0,5 \text{ rotazione}}{5400 \text{ GPM}} = \frac{0,5 \text{ rotazione}}{5400 \text{ GPM} / \left( 60 \frac{\text{secondi}}{\text{minuto}} \right)}$$

$$= 0,0056 \text{ secondi} = 5,6 \text{ ms}$$

L'ultima componente del tempo di accesso al disco, il *tempo di trasferimento*, è il tempo impiegato per trasferire un blocco di bit; questo tempo è funzione della dimensione del settore, della velocità di rotazione e della densità di memorizzazione delle tracce. Le velocità di trasferimento dei dati nel 2012 erano tipicamente comprese tra i 100 e i 200 MB/s. Una delle complicazioni è che il controllore del disco contiene al suo interno una cache dove vengono memorizzati i dati letti dai diversi settori del disco e la velocità di trasmissione da questa cache è tipicamente più elevata e poteva arrivare a 750 MB/s, ovvero 6 Gbit/s, nel 2012. Il risultato è che la disposizione dei blocchi non è più intuitiva. L'ipotesi sottostante il modello basato su settori-tracce-cilindri è che blocchi contigui siano posizionati sulla stessa traccia, che i blocchi sullo stesso cilindro richiedano un tempo di accesso minore dato che in questo caso il tempo di ricerca è nullo, e che alcune coppie di tracce siano più vicine tra loro di altre. Il motivo per cui questo modello sta perdendo consistenza è l'aumento delle prestazioni delle interfacce. Per incrementare la velocità degli accessi sequenziali, queste interfacce intelligenti organizzano i dischi in una forma più vicina a quella dei dispositivi a nastro che a quella delle memorie a semiconduttore. I blocchi logici sono disposti secondo una serpentina sulla superficie del disco, in modo da percorrere tutti i settori, i quali mantengono la stessa densità di bit. Di conseguenza blocchi in sequenza possono trovarsi su tracce diverse.

Riassumendo, le due differenze principali tra i dischi magnetici e le memorie a semiconduttore sono che i dischi hanno un tempo di accesso superiore perché sono dispositivi meccanici – le memorie flash sono 1000 volte più veloci e le DRAM 100 000 volte più veloci dei dischi – ma hanno un costo per bit più basso perché offrono una grandissima capacità di memoria a un costo modesto – i dischi sono da 10 a 100 volte più economici. I dischi magnetici sono non volatili come le memorie flash, ma a differenza delle memorie flash non ci sono problemi di usura. Tuttavia le memorie flash sono molto più robuste e quindi più adatte a sopportare gli urti inerenti all'utilizzo dei dispositivi mobili.

*Cache: un posto sicuro dove nascondere o riporre le cose.*

Webster's New World Dictionary of the American Language, Third College Edition, 1988

### 5.3 Principi base delle memorie cache

Nel nostro esempio della biblioteca la scrivania si comporta come una cache, ossia un posto sicuro dove riporre le cose (i libri) che dobbiamo esaminare. Cache<sup>2</sup> era il nome scelto per indicare il livello della memoria gerarchica che si trova

<sup>2</sup> Il nome proviene dal termine francese *caché*, che significa "nascosto", e deriva dal fatto che la memoria cache e il suo utilizzo sono generalmente trasparenti agli occhi del programmatore, e quindi nascosti [N.d.T.].

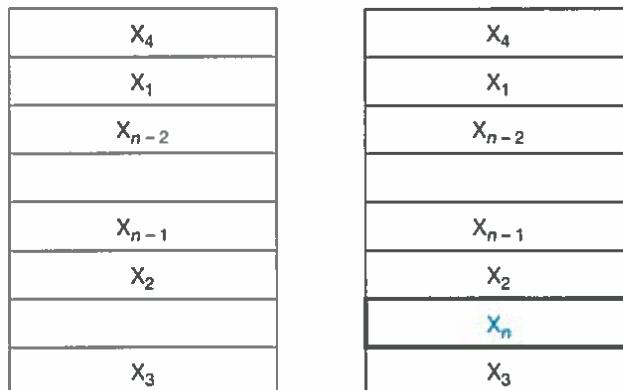
tra il processore e la memoria principale nel primo calcolatore commerciale che utilizzava una gerarchia delle memorie. Le memorie contenute nell'unità di elaborazione descritta nel Capitolo 4 possono essere semplicemente sostituite da memorie cache. Oggi, anche se questo rimane l'utilizzo più comune, il termine cache viene usato anche per indicare i sistemi di memoria gestiti in modo tale da ottenere i massimi benefici dalla località degli accessi. La cache apparve per la prima volta in architetture sperimentali progettate nei primi anni '60 e nei calcolatori commerciali costruiti verso la fine di quel decennio; quasi tutti i calcolatori prodotti oggi, dai server ai processori embedded a basso consumo, contengono memorie cache.

In questo paragrafo esamineremo una cache molto semplice, dalla quale il processore carica una sola parola per volta e in cui i blocchi sono costituiti da una singola parola. I lettori che hanno già familiarità con i principi base delle cache possono passare direttamente alla lettura del paragrafo 5.4. La Figura 5.7 mostra questa semplice cache, prima e dopo la richiesta di un elemento che, all'inizio, non è ivi presente. Prima della richiesta, la cache contiene un insieme di elementi, quelli utilizzati più di recente:  $X_1, X_2, \dots, X_{n-1}$ . Quando il processore richiede la parola  $X_n$ , che non è presente in cache, si produce una miss e la parola  $X_n$  viene prelevata dal livello inferiore della memoria e portata nella cache.

Osservando lo scenario rappresentato in Figura 5.7, è naturale chiedersi: come si fa a sapere se un dato è presente nella cache? E se è contenuto nella cache, come facciamo a trovarlo? Le risposte a queste domande sono collegate. Se ogni parola può essere scritta in una sola posizione della cache, allora sappiamo dove trovarla, ammesso che la parola sia presente nella cache. La maniera più semplice per associare una sola locazione della cache a ogni parola della memoria consiste nel definire una corrispondenza tra l'*indirizzo in memoria* della parola e la locazione nella cache. Quest'organizzazione della cache è detta **a mappatura diretta** (*direct mapped cache*), dato che ogni locazione della memoria principale corrisponde, in modo univoco, a una locazione della cache. Per le cache a mappatura diretta, la corrispondenza tra indirizzi della memoria e locazioni della cache è di solito molto semplice. Per esempio, quasi tutte le cache a mappatura diretta utilizzano la seguente operazione per trovare il blocco che corrisponde a un dato indirizzo della memoria principale:

(Indirizzo del blocco) modulo (numero di blocchi nella cache)

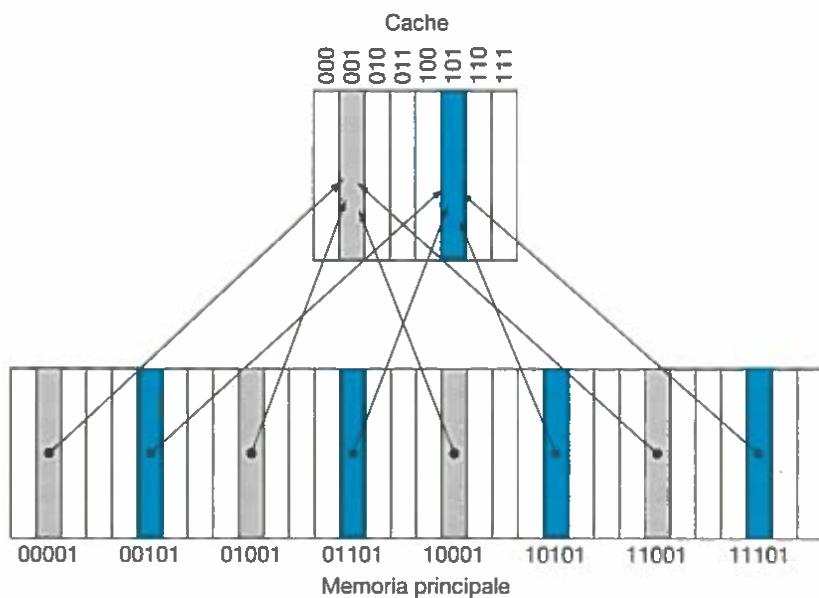
**Cache a mappatura diretta:** una cache in cui a ogni locazione della memoria principale corrisponde una (e una sola) locazione della cache.



a. Prima di richiedere il dato  $X_n$

b. Dopo aver richiesto il dato  $X_n$

**Figura 5.7** Stato della cache immediatamente prima e immediatamente dopo la richiesta di una parola  $X_n$ , che non è presente inizialmente nella cache. Questa richiesta provoca una miss che ha come conseguenza il caricamento della parola  $X_n$  dalla memoria principale e la sua scrittura nella cache.



**Figura 5.8** Esempio di cache a mappatura diretta contenente otto elementi e corrispondenza tra gli indirizzi della memoria principale, riferiti alla parola e compresi tra 0 e 31, e le locazioni della cache. Dato che la cache contiene 8 parole, un indirizzo  $X$  della memoria principale viene mappato su l'elemento  $X$  modulo 8 della cache a mappatura diretta. Cioè, i  $\log_2(8) = 3$  bit meno significativi dell'indirizzo vengono utilizzati come indice della cache. Perciò gli indirizzi 00001<sub>due</sub>, 01001<sub>due</sub>, 10001<sub>due</sub> e 11001<sub>due</sub> vengono mappati tutti sull'elemento 001<sub>due</sub> della cache, mentre gli indirizzi 00101<sub>due</sub>, 01101<sub>due</sub>, 10101<sub>due</sub> e 11101<sub>due</sub> vengono mappati tutti sull'elemento 101<sub>due</sub>.

Se il numero di elementi della cache è una potenza di 2, l'operazione di modulo può essere eseguita semplicemente considerando il logaritmo in base due ( $\log_2$ ) dei bit meno significativi dell'indirizzo, che corrispondono alla dimensione della cache in numero di blocchi. Quindi, una cache di 8 blocchi utilizzerà i 3 bit meno significativi per indirizzare il blocco, dato che  $8 = 2^3$ . Per esempio, la Figura 5.8 mostra come gli indirizzi della memoria principale compresi tra 1<sub>dec</sub> (00001<sub>due</sub>) e 29<sub>dec</sub> (11101<sub>due</sub>) vengano mappati nelle locazioni comprese tra 1 (001<sub>due</sub>) e 5 (101<sub>due</sub>) di una cache a mappatura diretta contenente otto parole.

Poiché ogni elemento della cache può contenere dati provenienti da diverse locazioni della memoria principale, come si riesce a capire se il dato presente nella cache corrisponde effettivamente alla parola desiderata? In altre parole, come facciamo a sapere se la parola richiesta si trova nella cache oppure no? È possibile risolvere questo problema aggiungendo alla cache un insieme di bit che costituiscono il campo **tag** (etichetta). I tag contengono le informazioni necessarie a verificare se una parola della cache corrisponda o meno alla parola cercata. Un tag contiene solamente la parte superiore dell'indirizzo della parola nella memoria principale, in particolare i bit dell'indirizzo che non vengono utilizzati come indice per individuare il blocco all'interno della cache. Per esempio, nella cache di Figura 5.8 occorre memorizzare nel campo tag solamente i 2 bit più significativi dei 5 bit totali dell'indirizzo, dato che i 3 bit inferiori servono per selezionare il blocco all'interno della cache. I bit di indice non vengono scritti nel campo tag perché sarebbero ridondanti, visto che, per definizione, l'indice associato a un qualsiasi blocco della cache coincide con il numero del blocco.

È necessario anche disporre di un metodo per sapere quando un blocco della cache non contiene informazioni valide. Per esempio, quando un processore viene avviato, la cache è vuota e i numeri contenuti nei campi tag non hanno alcun significato. Anche dopo aver eseguito molte istruzioni, alcune delle locazioni della cache possono ancora essere vuote, come mostrato in Figura 5.7.

**Tag:** un campo all'interno di una tabella utilizzata in una gerarchia delle memorie che contiene l'informazione necessaria per capire se il blocco di memoria associato al tag contenga o meno la parola richiesta.

Occorre perciò sapere quando il campo tag associato a queste locazioni deve essere ignorato. La procedura più comune consiste nell'aggiungere un **bit di validità** (*valid bit*) per indicare se il corrispondente elemento della cache contiene dati validi. Se il bit non è impostato a 1, la richiesta di lettura dell'elemento associato non può avere successo, perché il contenuto del blocco della cache sarebbe privo di significato.

Nel resto di questo paragrafo vedremo come viene gestita la lettura di una memoria cache. In generale, la lettura è un'operazione più semplice della scrittura, dato che essa non modifica i dati contenuti nella cache. Dopo aver esaminato i principi base del funzionamento della lettura e di come possano essere gestite le miss di una cache, analizzeremo le diverse organizzazioni di cache che si possono trovare nei calcolatori reali e illustreremo in dettaglio le modalità che queste adottano per gestire la scrittura.

**Bit di validità:** un campo in una tabella di una gerarchia delle memorie che indica se il blocco di memoria associato contiene un dato valido.

## QUADRO D'INSIEME

La memoria cache è forse l'esempio più importante della grande idea della **predizione**. Si basa sul principio di località, cerca di trovare i dati desiderati nei livelli più alti della gerarchia delle memorie e fornisce i meccanismi necessari per garantire che, quando una predizione si rivela errata, trovi e utilizzi i dati corretti presi dai livelli inferiori della gerarchia. La frequenza di hit della predizione nelle cache dei calcolatori moderni è spesso più alta del 95% (Figura 5.46). ■



PREDIZIONE

## Accesso alla cache

Nella tabella seguente viene riportata una sequenza di nove richieste di dati a una memoria cache di otto blocchi, inizialmente vuota, assieme all'azione associata a ogni richiesta. La Figura 5.9 mostra come il contenuto della cache cambi dopo ogni miss. Dato che nella cache ci sono otto blocchi, i tre bit meno significativi dell'indirizzo forniscono il numero del blocco:

Indirizzo decimale del dato nella memoria principale	Indirizzo binario del dato nella memoria principale	Hit o miss nell'accesso alla cache	Blocco della cache corrispondente (dove trovare o scrivere il dato)
22	10110 <sub>due</sub>	Miss (5.9b)	(10110 <sub>due</sub> mod 8) = 110 <sub>due</sub>
26	11010 <sub>due</sub>	Miss (5.9c)	(11010 <sub>due</sub> mod 8) = 010 <sub>due</sub>
22	10110 <sub>due</sub>	Hit	(10110 <sub>due</sub> mod 8) = 110 <sub>due</sub>
26	11010 <sub>due</sub>	Hit	(11010 <sub>due</sub> mod 8) = 010 <sub>due</sub>
16	10000 <sub>due</sub>	Miss (5.9d)	(10000 <sub>due</sub> mod 8) = 000 <sub>due</sub>
3	00011 <sub>due</sub>	Miss (5.9e)	(00011 <sub>due</sub> mod 8) = 011 <sub>due</sub>
16	10000 <sub>due</sub>	Hit	(10000 <sub>due</sub> mod 8) = 000 <sub>due</sub>
18	10010 <sub>due</sub>	Miss (5.9f)	(10010 <sub>due</sub> mod 8) = 010 <sub>due</sub>
16	10000 <sub>due</sub>	Hit	(10000 <sub>due</sub> mod 8) = 000 <sub>due</sub>

Essendo la cache inizialmente vuota, alcuni dei primi accessi in lettura costituiscono delle miss; le azioni che vengono intraprese per ciascun accesso alla memoria sono descritte nella didascalia di Figura 5.9. All'ottavo accesso, si verifica un conflitto sul contenuto di un blocco: la parola contenuta all'indirizzo 18 (10010<sub>due</sub>) della memoria principale dovrebbe essere scritta nel blocco 2

Indice	V	Tag	Dati
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

a. Lo stato iniziale della cache dopo l'accensione del calcolatore

Indice	V	Tag	Dati
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	S	10 <sub>due</sub>	Memoria (10110 <sub>due</sub> )
111	N		

b. Dopo avere gestito una miss all'indirizzo (10110<sub>due</sub>)

Indice	V	Tag	Dati
000	N		
001	N		
010	S	11 <sub>due</sub>	Memoria (11010 <sub>due</sub> )
011	N		
100	N		
101	N		
110	S	10 <sub>due</sub>	Memoria (10110 <sub>due</sub> )
111	N		

c. Dopo avere gestito una miss all'indirizzo 11010<sub>due</sub>

Indice	V	Tag	Dati
000	S	10 <sub>due</sub>	Memoria (10000 <sub>due</sub> )
001	N		
010	S	11 <sub>due</sub>	Memoria (11010 <sub>due</sub> )
011	N		
100	N		
101	N		
110	S	10 <sub>due</sub>	Memoria (10110 <sub>due</sub> )
111	N		

d. Dopo avere gestito una miss all'indirizzo 10000<sub>due</sub>

Indice	V	Tag	Dati
000	S	10 <sub>due</sub>	Memoria (10000 <sub>due</sub> )
001	N		
010	S	11 <sub>due</sub>	Memoria (11010 <sub>due</sub> )
011	S	00 <sub>due</sub>	Memoria (00011 <sub>due</sub> )
100	N		
101	N		
110	S	10 <sub>due</sub>	Memoria (10110 <sub>due</sub> )
111	N		

e. Dopo avere gestito una miss all'indirizzo 00011<sub>due</sub>

Indice	V	Tag	Dati
000	S	10 <sub>due</sub>	Memoria (10000 <sub>due</sub> )
001	N		
010	S	10 <sub>due</sub>	Memoria (10010 <sub>due</sub> )
011	S	00 <sub>due</sub>	Memoria (00011 <sub>due</sub> )
100	N		
101	N		
110	S	10 <sub>due</sub>	Memoria (10110 <sub>due</sub> )
111	N		

f. Dopo avere gestito una miss all'indirizzo 10010<sub>due</sub>

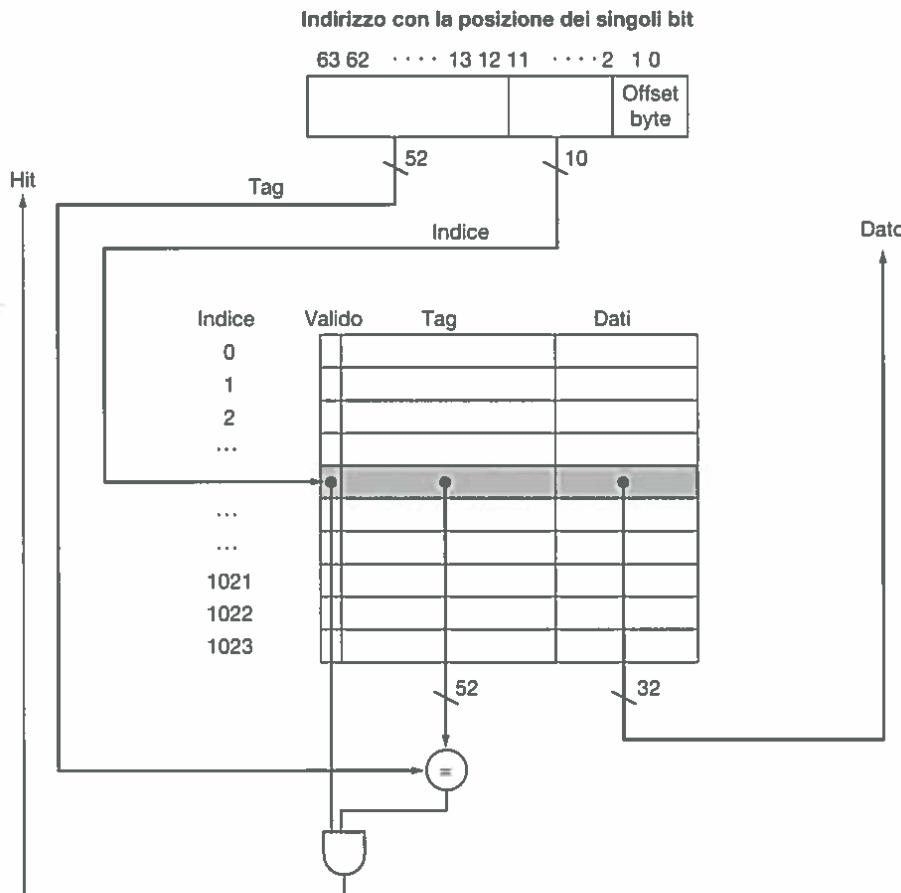
**Figura 5.9** Contenuto della cache dopo ogni accesso che provoca una miss all'interno della sequenza riportata nella tabella della pagina precedente; i campi indice e tag sono in binario. All'inizio la cache è vuota e tutti i bit di validità (campo V) sono impostati a 0 (N). Il processore richiede i dati relativi ai seguenti indirizzi della memoria principale: 10110<sub>due</sub> (miss), 11010<sub>due</sub> (miss), 10110<sub>due</sub> (hit), 11010<sub>due</sub> (hit), 10000<sub>due</sub> (miss), 00011<sub>due</sub> (miss), 10000<sub>due</sub> (hit) e 10010<sub>due</sub> (miss). Le tabelle mostrano il contenuto della cache dopo la gestione di ognuna delle miss che si verificano nella sequenza di richieste. Quando il processore richiede il dato all'indirizzo 10010<sub>due</sub> (18<sub>dec</sub>) della memoria principale, il contenuto dell'elemento della cache corrispondente, che contiene in quel momento il dato con indirizzo 11010<sub>due</sub> (26<sub>dec</sub>), deve essere sostituito dal contenuto dell'indirizzo 10010<sub>due</sub>; una successiva richiesta del dato contenuto all'indirizzo 11010<sub>due</sub> provocherà una nuova miss. Si noti che il campo tag contiene soltanto la parte più significativa dell'indirizzo. L'indirizzo completo di una parola nel blocco /esimo di questa cache a cui è associato un tag j, si ottiene come 8 × j + i, oppure, in maniera equivalente, lo si ricava dalla concatenazione del campo tag, j, e del campo indice, / Per esempio, nella cache raffigurata nel pannello f al dato con indice 010<sub>due</sub> è associata l'etichetta 10<sub>due</sub>, per cui l'indirizzo del dato nella memoria principale sarà 10010<sub>due</sub>.

della cache ( $010_{\text{due}}$ ), sostituendo quindi la parola con indirizzo 26 ( $11010_{\text{due}}$ ) scritta in precedenza. Questo comportamento consente alla cache di sfruttare la località temporale: le parole richieste più di recente sostituiscono quelle richieste in precedenza.

Questa situazione corrisponde esattamente a quella in cui lo studente ha bisogno di prendere un altro libro dagli scaffali della biblioteca, ma non ha più spazio sulla scrivania: alcuni dei testi sulla scrivania devono essere prima riportati sugli scaffali. In una cache a mappatura diretta c'è una sola locazione in cui poter scrivere il nuovo dato richiesto e quindi è uno solo il blocco della cache che può essere sostituito.

Sappiamo quale blocco della cache corrisponde a ogni indirizzo della memoria principale: i bit meno significativi dell'indirizzo vengono utilizzati per individuare l'unico elemento della cache a cui l'indirizzo può corrispondere. La Figura 5.10 illustra come l'indirizzo della memoria principale venga suddiviso in:

- un *campo tag*, che viene confrontato con il contenuto del campo tag della cache;
- un *campo indice*, utilizzato per selezionare il blocco della cache.



**Figura 5.10** Per questa cache, la parte bassa dell'indirizzo viene utilizzata per selezionare un elemento della cache che consiste in un dato largo una parola e in un campo tag. Questa cache contiene 1024 parole o 4 KiB e in questo capitolo supponiamo indirizzi a 64 bit. Il tag della cache viene confrontato con la parte alta dell'indirizzo al fine di verificare se l'elemento della cache corrisponda all'elemento richiesto. Poiché la cache contiene  $2^{10}$  (cioè 1024) parole e la dimensione di un blocco è di una parola, 10 bit dell'indirizzo della memoria principale vengono utilizzati per l'indirizzamento interno della cache, mentre  $64 - 10 - 2 = 52$  bit vengono confrontati con il contenuto del campo tag. Se il campo tag è uguale ai 52 bit più significativi dell'indirizzo e se il bit di validità è impostato a 1, allora la ricerca del dato nella cache ha avuto successo, cioè si verifica una hit e il dato viene fornito al processore. Altrimenti scatta una miss.

L'indice di un blocco della cache, insieme al contenuto del campo tag, specifica in maniera univoca l'indirizzo della memoria principale associato alla parola contenuta in quel blocco della cache. Dato che i bit del campo indice vengono utilizzati come indirizzo interno alla cache e che un campo di  $n$  bit contiene  $2^n$  valori diversi, il numero totale degli elementi contenuti in una memoria cache a mappatura diretta deve essere una potenza di 2. Poiché nell'architettura MIPS le parole sono allineate a multipli di quattro byte e gli indirizzi sono riferiti ai byte, i due bit meno significativi dell'indirizzo specificano i byte all'interno di una parola e quindi non vengono utilizzati per selezionare le parole all'interno di un blocco della cache.

Il numero complessivo di bit necessari per realizzare una cache è funzione della dimensione della cache e della dimensione degli indirizzi, poiché la cache deve contenere spazio di memoria sia per i dati sia per i tag.

La dimensione dei blocchi considerata finora era di una singola parola, ma normalmente un blocco è costituito da più parole. Si consideri la seguente situazione:

- indirizzo su 64 bit;
- cache a mappatura diretta;
- dimensione della cache di  $2^n$  blocchi, per cui  $n$  bit vengono utilizzati per l'indice;
- dimensione del blocco della cache di  $2^m$  parole, ossia  $2^{m+2}$  byte, per cui  $m$  bit vengono utilizzati per individuare una parola all'interno di un blocco, mentre due bit servono a individuare un byte all'interno di una parola.

In questo caso, la dimensione del campo tag è:

$$64 - (n + m + 2)$$

Il numero totale di bit contenuti in una cache a mappatura diretta è quindi:

$$2^n \times (\text{dimensione\_blocco} + \text{dimensione\_tag} + \text{bit\_validità})$$

Tenendo presente che la dimensione del blocco è di  $2^m$  parole, ossia  $2^{m+2}$  byte che corrispondono a  $2^{m+5}$  bit, e considerando il bit di validità, il numero di bit contenuti in questa cache è:

$$2^n \times [2^m \times 32 + (64 - n - m - 2) + 1] = 2^n \times (2m \times 32 + 63 - n - m).$$

Sebbene questa sia la dimensione reale della cache, per convenzione si escludono dal calcolo i bit del campo tag e il bit di validità, per cui si considera solo la dimensione dei dati. Quella di Figura 5.10 è considerata quindi una cache da 4 KiB.

### Bit di una cache

#### ESEMPIO

Da quanti bit è costituita una cache a mappatura diretta contenente 16 KiB di dati e avente blocchi di 4 parole, ipotizzando un indirizzo su 64 bit?

#### SOLUZIONE

(continua)

Sappiamo che 16 KiB equivalgono a 4096 parole (cioè  $2^{12}$  parole) e che la dimensione del blocco è di 4 parole ( $2^2$ ); la cache conterrà quindi 1024

(continua)

blocchi ( $2^{10}$  blocchi). Ogni blocco ha dimensione  $4 \times 32 = 128$  bit di dati più la dimensione del campo tag, che è pari a  $64 - 10 - 2 - 2 = 50$  bit, più il bit di validità. Quindi, la dimensione totale della cache è:

$$2^{10} \times [4 \times 32 + (64 - 10 - 2 - 2) + 1] = 2^{10} \times 179 = 179 \text{ Kibibit}$$

vale a dire 22,4 KiB per una cache che contiene 16 KiB di dati. In questo caso, il numero totale di bit della cache è circa 1,4 volte lo spazio richiesto per memorizzare i dati.

### Mappatura di un indirizzo in una cache con blocchi composti da più parole

Si consideri una cache con 64 blocchi di 16 byte ciascuno. A quale numero di blocco corrisponde l'indirizzo 1200, espresso in byte, della memoria principale?

In base alla relazione introdotta all'inizio del paragrafo 5.3, il blocco contenente il dato viene identificato come:

(Indirizzo del blocco) modulo (numero dei blocchi nella cache)

dove l'indirizzo del blocco è:

$$\frac{\text{Indirizzo del dato in byte}}{\text{Byte del blocco}}$$

Si noti che le parole appartenenti al blocco individuato hanno indirizzo compreso tra:

$$\left[ \frac{\text{Indirizzo del dato in byte}}{\text{Byte del blocco}} \right] \times \text{Byte per blocco}$$

e

$$\left[ \frac{\text{Indirizzo del dato in byte}}{\text{Byte del blocco}} \right] \times \text{Byte per blocco} + (\text{Byte per blocco} - 1)$$

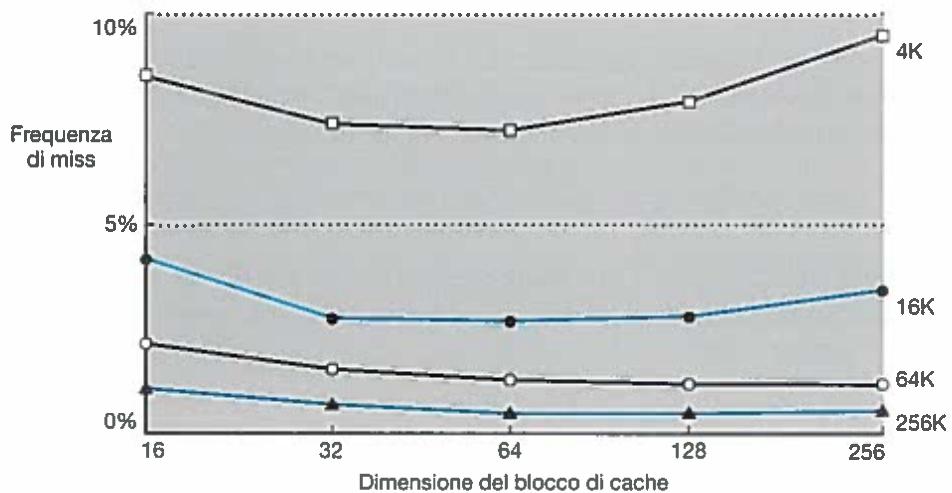
Perciò, dato che ci sono 16 byte in ogni blocco, l'indirizzo 1200 corrisponde all'indirizzo di blocco:

$$\left[ \frac{1200}{16} \right] = 75$$

che corrisponde al blocco della cache avente indice (75 modulo 64) = 11. Infatti, sul blocco 11 vengono mappati tutti i dati contenuti negli indirizzi compresi tra 1200 e 1215 della memoria principale.

ESEMPIO

SOLUZIONE



**Figura 5.11** Frequenza delle miss in funzione della dimensione del blocco della cache. Si noti che la frequenza delle miss, in realtà, cresce se la dimensione del blocco è troppo grande rispetto alla dimensione della cache. Ogni linea rappresenta una cache di dimensioni diverse. Questa figura è indipendente dal grado di associatività di cui discuteremo più avanti. Purtroppo, occorrerebbe troppo tempo per raccogliere i dati relativi ai benchmark SPEC CPU2000 includendo anche la dimensione del blocco della cache tra i parametri, e quindi questi dati sono basati sugli SPEC92.

Blocchi di dimensioni maggiori sfruttano maggiormente la località spaziale e diminuiscono la frequenza delle miss. Come mostrato in Figura 5.11, all'aumentare della dimensione dei blocchi di solito diminuisce la frequenza delle miss. Tuttavia la frequenza delle miss può tornare a crescere se la dimensione dei blocchi diventa troppo grande rispetto alla dimensione della cache, perché in questo caso il numero dei blocchi che possono essere memorizzati nella cache diventa piccolo, mentre cresce la competizione per occuparli. Di conseguenza, i blocchi vengono scaricati dalla cache prima ancora che molti dei dati in essi contenuti siano stati utilizzati. In altre parole, la località spaziale tra le parole di un blocco diminuisce con blocchi di dimensioni troppo grandi e, di conseguenza, il miglioramento legato alla frequenza delle miss si riduce.

Un problema ancora più serio, associato all'aumento della dimensione dei blocchi, è la crescita del costo di una miss. La penalità di una miss, infatti, è determinata dal tempo necessario a prelevare un blocco dal livello inferiore della gerarchia e a scriverlo nella cache. Il tempo impiegato per prelevare il blocco è costituito da due parti: la latenza per ottenere la prima parola del blocco e il tempo di trasferimento del resto del blocco. Chiaramente, a meno che non si modifichi il sistema di memoria, il tempo di trasferimento, e quindi la penalità di miss, cresce con le dimensioni del blocco. Inoltre, il miglioramento nella frequenza delle miss comincia a diminuire man mano che aumenta la dimensione del blocco. Il risultato è che l'incremento della penalità di miss ha un impatto superiore alla diminuzione della frequenza delle miss e quindi le prestazioni della cache si abbassano. Ovviamente, se si progetta la memoria per trasferire in maniera più efficiente blocchi di dimensioni maggiori, è possibile aumentare la dimensione del blocco e ottenere un ulteriore incremento delle prestazioni della cache. Discuteremo questo aspetto nel prossimo paragrafo.

**Approfondimento.** Sebbene sia difficile fare qualcosa contro l'aumento della durata della componente preponderante della penalità di miss per blocchi di grandi dimensioni, è possibile nascondere una parte del tempo di trasferimento in modo da rendere effettivamente più breve la penalità di miss. Il metodo più semplice per

ottenere tale risultato viene chiamato metodo della *ripartenza anticipata (early restart)*, e consiste semplicemente nel riprendere l'esecuzione non appena è stata trasferita la parola richiesta, invece di aspettare che sia stato trasferito l'intero blocco. Molti processori utilizzano questo meccanismo dove funziona meglio, ossia nell'accesso alla cache delle istruzioni. L'accesso alle istruzioni, infatti, è in gran parte sequenziale: ciò significa che se il sistema di memoria è in grado di trasferire una parola a ogni ciclo di clock, il processore può riprendere l'attività non appena la parola richiesta è stata trasferita, mentre il sistema di memoria continua a trasferire le altre istruzioni mettendole a disposizione in tempo. Questa tecnica, di solito, è meno efficace per le cache contenenti i dati, perché è facile che le parole di un blocco vengano richieste in maniera meno prevedibile e perché è alta la probabilità che il processore abbia bisogno di altre parole presenti in blocchi diversi della cache prima che il trasferimento di tutto il blocco sia stato completato. Se il processore non può accedere alla cache dei dati perché c'è un trasferimento in corso, allora deve generare uno stall e fermare l'esecuzione.

Uno schema ancora più sofisticato consiste nell'organizzare la memoria in maniera tale che la parola richiesta venga trasferita per prima dalla memoria alla cache. Il resto del blocco viene poi trasferito a partire dall'indirizzo successivo a quello della parola richiesta, per poi tornare alla prima parola del blocco dopo avere trasferito l'ultima parola del blocco stesso. Questa tecnica, nota come *tecnica della prima parola richiesta (requested word first)*, può risultare un po' più veloce della ripartenza anticipata, ma è limitata dagli stessi problemi che affliggono la ripartenza anticipata.

## Gestione delle miss della cache

Prima di esaminare la struttura della cache di un sistema reale, vediamo in che modo l'unità di controllo gestisce le **miss della cache** (fallimenti della cache). Descriveremo in dettaglio il funzionamento del controllore della cache nel paragrafo 5.9. L'unità di controllo deve riconoscere il fallimento del tentativo di accesso e provvedere a risolverlo andando a prelevare i dati dalla memoria principale, oppure, come vedremo più avanti, da una cache di livello inferiore. Se invece il tentativo di accesso va a buon fine, il calcolatore continua l'elaborazione utilizzando il dato come se nulla fosse accaduto.

Modificare l'unità di controllo del processore per gestire una hit è un'operazione banale. La gestione delle miss, invece, richiede del lavoro aggiuntivo e un'unità di controllo separata che collabora con il processore. Questa si occupa dell'accesso alla memoria principale e del "rifornimento" della cache. La risoluzione di una miss richiede lo stall della pipeline (vedi Cap. 4) ed è differente da un interrupt o da un'eccezione, che implicherebbe il salvataggio dello stato di tutti i registri. A ogni miss della cache possiamo mettere in stall l'intero processore, essenzialmente bloccando il contenuto dei registri temporanei e di quelli visibili al programmatore, per tutto il tempo necessario a caricare i dati dalla memoria principale. I processori più sofisticati, dotati di esecuzione fuori ordine, consentono l'esecuzione di altre istruzioni mentre attendono la risoluzione di una miss della cache, ma in questo paragrafo considereremo solamente i processori che eseguono il codice in sequenza; questi devono essere messi in stall in seguito a una miss.

Vediamo più da vicino come vengono gestite le miss sulle istruzioni; gli stessi meccanismi si possono estendere alla gestione delle miss sui dati. Se l'accesso a un'istruzione si traduce in una miss, il contenuto del registro istruzioni non sarà più valido. Per caricare l'istruzione corretta nella cache dobbiamo poter dire al livello inferiore della gerarchia delle memorie di eseguire un'operazione di lettura. Dal momento che il program counter viene incrementato di 4 nel primo ciclo di clock di esecuzione, l'indirizzo dell'istruzione che ha generato la miss sarà uguale al contenuto del PC meno 4. Una volta costruito l'indirizzo corretto, si può chiedere alla memoria principale di effettuare la lettura dell'istruzione corrispondente; si deve poi attendere che la memoria abbia terminato la lettura,

**Miss della cache:** la richiesta di un dato alla cache che non può essere soddisfatta perché il dato non è presente nella cache.

dato che l'accesso richiede più cicli di clock. Infine, occorre attendere che la parola contenente l'istruzione desiderata venga scritta nella cache.

Riassumiamo schematicamente i passi che devono essere eseguiti quando si verifica una miss nella cache istruzioni:

1. inviare il valore originale del program counter alla memoria;
2. comandare alla memoria principale di eseguire un'operazione di lettura e attendere che la memoria completi la lettura;
3. scrivere la parola che proviene dalla memoria nella posizione opportuna del blocco della cache, aggiornare il campo tag corrispondente scrivendovi i bit più significativi dell'indirizzo presi direttamente dalla ALU e impostare a 1 il bit di validità;
4. far ripartire l'esecuzione dell'istruzione dall'inizio, ripetendo la fase di fetch, che questa volta troverà l'istruzione all'interno della cache.

Le operazioni di controllo svolte da una cache per la lettura dei dati sono essenzialmente identiche: in caso di miss, il processore viene messo in stallo finché la memoria non risponde, restituendo il dato richiesto.

## Gestione della scrittura

**Write-through:** uno schema secondo il quale un elemento viene scritto sia in cache sia nel livello inferiore della gerarchia delle memorie, assicurando così che i dati presenti nelle due memorie siano sempre coerenti tra di loro.

La scrittura funziona in modo un po' diverso. Supponiamo che un'operazione di store scriva il dato solamente nella cache dei dati, senza modificare la memoria principale; al termine della scrittura, la memoria principale avrebbe un contenuto diverso da quello della cache. In questo caso, si dice che la memoria e la cache sono *incoerenti*. La maniera più semplice per conservare la coerenza tra memoria e cache consiste nello scrivere sempre il dato in entrambe le memorie. Questo schema viene chiamato **write-through** ("scrivere attraverso").

L'altro elemento fondamentale della scrittura è rappresentato dalla gestione delle miss in scrittura. Dapprima occorre caricare dalla memoria principale le parole appartenenti al blocco interessato. Dopo avere caricato il blocco e averlo scritto in cache, possiamo sovrascrivere la parola del blocco che aveva causato la miss; questa parola viene scritta anche nella memoria principale utilizzando il suo indirizzo completo.

Sebbene il meccanismo sopra descritto consenta di gestire le scritture con molta semplicità, non offre buone prestazioni. L'utilizzo dello schema write-through comporta che a ogni scrittura la parola venga salvata anche nella memoria principale. Questa operazione richiede molto tempo, almeno 100 cicli di clock del processore, e può quindi rallentare il sistema in maniera considerevole. Per esempio, supponendo che il 10% delle istruzioni siano istruzioni di store, che il CPI senza miss sia pari a 1,0 e che ogni miss richieda 100 cicli di clock in più per ciascuna scrittura, il CPI del processore diventerebbe  $1,0 + 100 \times 10\% = 11$ , riducendo le prestazioni di più di un fattore 10.

Una possibile soluzione a questo problema consiste nell'utilizzare una memoria tampone, chiamata **buffer di scrittura**. Il buffer di scrittura memorizza i dati in attesa che essi vengano scritti in memoria: dopo aver salvato il dato nella cache e nel buffer di scrittura, il processore può proseguire l'esecuzione. Una volta completata la scrittura di un dato nella memoria principale, il corrispondente spazio nel buffer di scrittura viene liberato. Se il buffer di scrittura è pieno e il processore deve eseguire un'operazione di scrittura, il processore viene messo in stallo finché non si libera spazio nel buffer. Ovviamente, se la velocità con cui la memoria principale completa le scritture è inferiore a quella con cui il processore genera i dati da scrivere, per quanto grande possa essere il buffer di scrittura, prima o poi il processore dovrà comunque essere messo in stallo, poiché i dati da scrivere vengono generati più velocemente di quanto il sistema di memoria sia in grado di accettarli.

**Buffer di scrittura:** una coda che contiene i dati che aspettano di essere scritti nella memoria principale.

Anche se la frequenza con cui vengono effettuate le scritture fosse *inferiore* alla frequenza con cui la memoria può gestirle, si possono comunque generare degli stalli. Ciò può succedere quando le scritture devono essere effettuate a raffica (*in burst*). Per ridurre il numero di stalli, i processori sono dotati solitamente di un buffer di scrittura in grado di contenere più di un singolo elemento.

Lo schema alternativo al write-through è chiamato **write-back**. In questo schema, quando si verifica una scrittura, il dato viene scritto solamente nel blocco corrispondente della cache e il blocco modificato viene salvato nel livello inferiore della gerarchia solo quando deve essere rimpiazzato. Lo schema write-back può produrre un miglioramento delle prestazioni, specialmente quando il processore può generare le scritture velocemente o più velocemente di quanto le scritture possano essere gestite dalla memoria principale; il write-back è comunque più complesso da implementare del write-through.

Nel resto di questo paragrafo descriveremo le cache di alcuni processori reali ed esamineremo come esse gestiscano sia le letture sia le scritture. Nel paragrafo 5.8 esamineremo più dettagliatamente la gestione delle scritture.

**Approfondimento.** La scrittura in cache introduce alcune problematiche che non sono presenti nella lettura. In questa sezione di approfondimento ne esamineremo due: la strategia adottata per le miss in scrittura e l'implementazione efficiente della scrittura in una cache con write-back.

Consideriamo una miss di una cache con write-through. La strategia adottata in molte cache di questo tipo, chiamata *write-allocate* (allocazione per la scrittura), è di allocare un blocco intero della cache. Il blocco intero viene prelevato dalla memoria e solamente la parte appropriata del blocco viene sovrascritta. Una strategia alternativa, chiamata *no-write-allocate*, è quella di aggiornare solo la parte appropriata del blocco direttamente nella memoria principale senza aggiornare la cache. Questi due metodi derivano dall'osservazione che a volte i programmi scrivono blocchi interi di dati, per esempio quando il sistema operativo inizializza a zero una pagina della memoria principale. In questi casi, il prelevamento del dato associato alla prima miss può risultare inutile. Alcuni calcolatori consentono una politica di write-allocate che può essere cambiata in base alla pagina di memoria considerata.

Implementare la scrittura in modo efficiente in una cache di tipo write-back è più complicato rispetto a una cache write-through. In una cache write-through possiamo scrivere il dato nella cache e leggere il campo tag; se questo non coincide con il corrispondente campo dell'indirizzo del dato, si verifica una miss. Dato che la cache è di tipo write-through, la sovrascrittura del blocco della cache non ha effetti catastrofici, perché la memoria principale contiene già il valore corretto. In una cache con write-back, quando un dato deve essere scritto nella cache e si ha una miss, occorre prima scrivere il blocco nella memoria principale. Se sovrascrivessimo semplicemente il blocco a seguito di un'istruzione di store prima ancora di sapere se l'accesso produca una hit (come si potrebbe fare nel caso di una cache di tipo write-through), distruggeremmo il contenuto del blocco, perché esso non è stato ancora copiato nel livello inferiore della gerarchia delle memorie.

In una cache di tipo write-back, poiché non si possono sovrascrivere i blocchi, ci sono due possibili soluzioni. Nella prima, le istruzioni di store richiedono due cicli di clock: un ciclo per verificare se si ottiene una hit, seguito da un secondo ciclo per eseguire effettivamente la scrittura. In alternativa occorre avere un buffer di scrittura per conservare temporaneamente il dato, consentendo così alla scrittura di impiegare un solo ciclo e inserendo, di fatto, l'operazione di scrittura in pipeline. Quando si utilizza un buffer di scrittura, il processore accede alla cache e scrive il dato nel buffer durante il normale ciclo di accesso alla cache. Nel caso in cui si verifichi una hit, il nuovo dato viene copiato dal buffer di scrittura in cache al primo ciclo di clock successivo in cui la cache non viene utilizzata.

In una cache di tipo write-through, invece, le scritture possono essere sempre eseguite in un solo ciclo di clock, nel quale viene letto il campo tag dell'indirizzo e viene salvato il dato nella porzione corrispondente del blocco della cache. Se il campo tag coincide con quello del blocco presente in cache, il processore può proseguire normalmente l'esecuzione, essendo stato aggiornato il blocco corretto;

**Write-back:** uno schema in cui un elemento viene scritto solo nel blocco di cache; il blocco corrispondente del livello inferiore della gerarchia viene scritto solamente quando il blocco della cache viene sostituito.

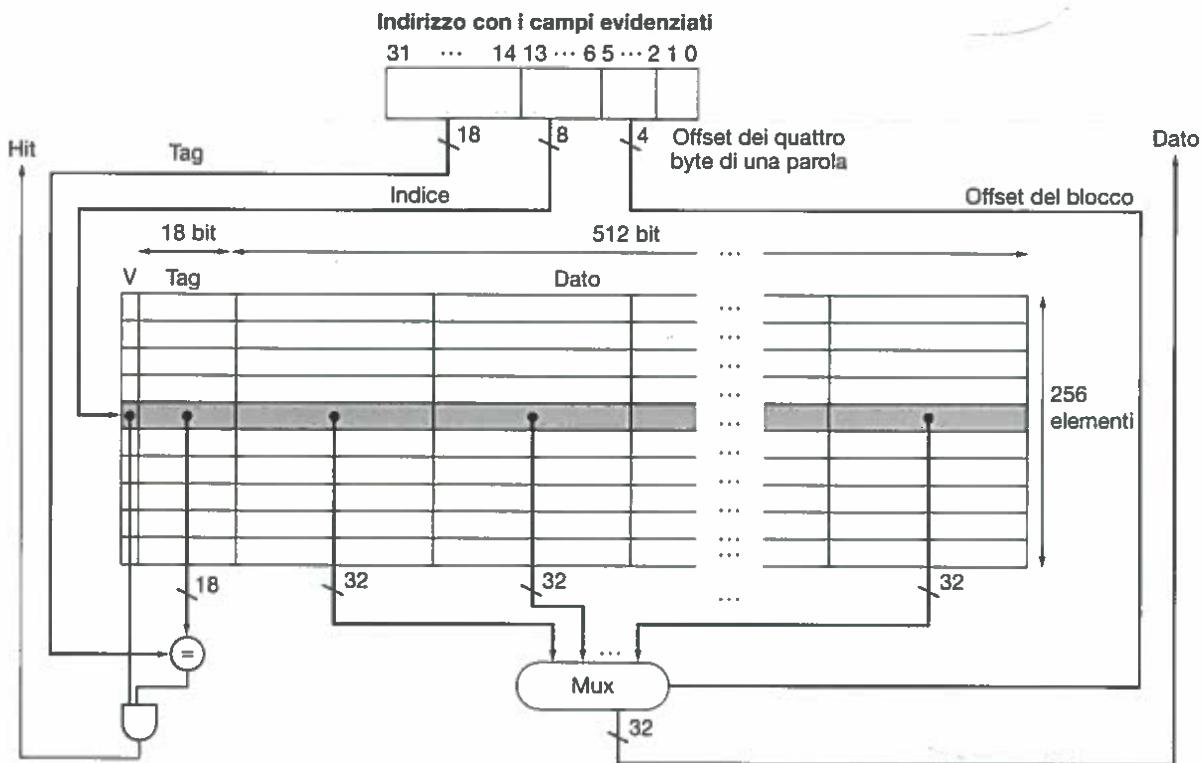
altrimenti, il processore genera una miss in scrittura per caricare le altre parole del blocco corrispondenti a quell'indirizzo.

Molte cache di tipo write-back contengono anche dei buffer di scrittura che vengono utilizzati per ridurre la penalità di miss dovuta alla sostituzione di un blocco che è stato modificato. In questi casi, il blocco modificato viene spostato in un buffer di write-back associato alla cache mentre il blocco richiesto viene letto dalla memoria. Il contenuto del buffer di write-back viene successivamente copiato nella memoria principale. Supponendo che una seconda miss non avvenga immediatamente, questa tecnica dimezza la penalità di miss associata alla sostituzione di un blocco modificato.

### Un esempio di memoria cache: il processore FastMATH Intrinsity

Il processore FastMATH Intrinsity è un microprocessore embedded veloce che utilizza l'architettura MIPS e una semplice implementazione della cache. Iniziamo da questo esempio semplice (ma reale) per ragioni didattiche; più avanti esamineremo le strutture più complesse adottate per le cache dei microprocessori ARM e Intel. La Figura 5.12 mostra l'organizzazione della cache dati del processore FastMATH Intrinsity. Si noti che la dimensione dell'indirizzo in questo calcolatore è di soli 32 bit e non 64 bit come nel resto del libro.

Questo processore ha una pipeline a 12 stadi. Quando opera a pieno regime, il processore può richiedere sia una parola di dati sia una parola di istruzioni a ogni ciclo di clock. Per soddisfare queste richieste senza generare stalli nella pipeline, vengono utilizzate una cache dati e una cache istruzioni separate. Ciascuna cache è di 16 KiB, ossia di 4096 parole, con 16 parole per blocco.



**Figura 5.12** Le due cache del processore FastMATH Intrinsity hanno una dimensione di 16 KiB: ognuna contiene 256 blocchi con 16 parole ciascuno. Il campo tag è di 18 bit, il campo indice è di 8 bit, mentre un campo di 4 bit (i bit dal 5 al 2) viene utilizzato per indicizzare il blocco e selezionare la parola all'interno del blocco mediante un multiplexer 16:1. In pratica, per eliminare il multiplexer, le cache utilizzano una RAM separata di grandi dimensioni per i dati e una RAM più piccola per i tag, mentre l'offset del blocco fornisce i bit aggiuntivi per indirizzare i dati all'interno della RAM più grande. In questo caso, la RAM dati ha parole di 32 bit ciascuna e deve contenere 16 parole per ogni blocco.

Frequenza di miss per le istruzioni	Frequenza di miss per i dati	Frequenza di miss totale
0,4%	11,4%	3,2%

**Figura 5.13** Frequenza approssimativa delle miss sulle istruzioni e sui dati in un processore FastMATH Intrinsity quando vengono eseguiti i benchmark SPEC CPU2000. La frequenza di miss totale è la frequenza di miss effettiva rilevata per una configurazione che prevede una cache dati da 16 KiB e una cache istruzioni da 16 KiB. Questa viene ottenuta pesando la frequenza di miss delle istruzioni e dei dati con la frequenza di accesso alla cache dati e alla cache istruzioni.

Le richieste di lettura dalla cache non destano problemi. Poiché la cache dati e la cache istruzioni sono separate, occorrono segnali diversi per leggere e scrivere in ognuna delle due memorie; si ricordi, inoltre, che occorre aggiornare la cache istruzioni ogni volta che si verifica una miss. Quindi, i passi associati a una richiesta di lettura, validi per entrambe le memorie, sono i seguenti.

1. Inviare l'indirizzo alla cache appropriata. L'indirizzo si ricava dal PC (per le istruzioni) oppure dalla ALU (per i dati).
2. Se la cache segnala una hit, la parola richiesta è disponibile sulle linee dati. Poiché ci sono 16 parole all'interno del blocco selezionato, occorre scegliere la parola giusta. Il campo offset del blocco viene utilizzato per controllare il multiplexer mostrato in basso nella figura, il quale seleziona la parola richiesta tra le 16 parole del blocco selezionato.
3. Se la cache genera una miss, l'indirizzo della parola viene inviato alla memoria principale. Quando la memoria restituisce il dato richiesto, esso viene scritto nella cache e quindi trasferito dalla cache al processore per soddisfare la richiesta.

Per le operazioni di scrittura, il FastMATH Intrinsity offre sia la modalità write-through sia la modalità write-back, lasciando al sistema operativo il compito di decidere quale strategia utilizzare per l'applicazione corrente. La dimensione del buffer di scrittura è di un singolo elemento.

Quale percentuale di miss ci possiamo aspettare da una cache con una struttura come quella utilizzata dal FastMATH Intrinsity? La Figura 5.13 mostra la frequenza di miss per le istruzioni e per i dati. La frequenza di miss totale è la frequenza di miss effettiva degli accessi a cache calcolata per ciascun programma tenendo conto delle differenti frequenze di accesso ai dati e alle istruzioni.

Sebbene la frequenza delle miss sia una caratteristica importante da considerare durante la progettazione di una cache, l'indicatore fondamentale è l'impatto del sistema di memoria sul tempo di esecuzione dei programmi; vedremo tra poco come frequenza delle miss e tempo di esecuzione siano collegati tra loro.

**Approfondimento.** Una cache unica di dimensioni pari alla somma delle due **split cache** (cache separate) di solito ha una frequenza di hit superiore. Questo accade perché, utilizzando una memoria unica, i blocchi non sono separati rigidamente tra istruzioni e dati. Nonostante ciò, oggi quasi tutte le architetture adottano cache separate per dati e istruzioni, al fine di aumentare la *larghezza di banda* della cache, per soddisfare le esigenze delle moderne pipeline. Inoltre, questa struttura permette di ottenere un numero minore di miss conflittuali (par. 5.8).

Di seguito riportiamo alcune misure ottenute su cache tipiche del processore FastMATH Intrinsity e su cache unificate che abbiano dimensioni pari alla somma delle due cache:

- dimensione complessiva della cache: 32 KiB;
- frequenza effettiva delle miss per le cache separate: 3,24%;
- frequenza delle miss per la cache unica: 3,18%.

La frequenza delle miss è di poco inferiore nel caso di cache separate.

**Split cache:** uno schema in cui un livello della gerarchia delle memorie è composto da due cache indipendenti che operano in parallelo: una cache gestisce le istruzioni e l'altra i dati.

antaggio ottenuto raddoppiando la larghezza di banda della cache, conseguenza del supporto all'accesso simultaneo ai dati e alle istruzioni, compensa abbondantemente lo svantaggio di avere una frequenza di miss leggermente più elevata. Questa è un'ulteriore conferma del fatto che non è possibile utilizzare solo la frequenza delle miss come indice delle prestazioni di una cache, come verrà mostrato nel paragrafo 5.4.

## Riepilogo

Abbiamo iniziato il paragrafo esaminando la più semplice delle cache: una cache a mappatura diretta con blocchi di ampiezza pari a quella di una parola. In una cache di questo tipo sia le hit sia le miss sono semplici da gestire, visto che ogni parola può andare esattamente in una sola locazione della cache e per ciascuna di esse è previsto un tag separato. Per mantenere la coerenza tra memoria principale e cache, si può utilizzare uno schema write-through nel quale ogni scrittura in cache provoca l'aggiornamento anche della memoria principale. L'alternativa allo schema write-through è lo schema write-back, nel quale un blocco della cache viene ricopiato nella memoria principale solamente quando deve essere sostituito; quest'ultimo schema verrà discusso in modo più approfondito nei prossimi paragrafi.

Per sfruttare la località spaziale, una cache deve avere blocchi di dimensioni superiori a una parola. L'utilizzo di blocchi di dimensioni maggiori consente di diminuire la frequenza delle miss e di migliorare l'efficienza della cache, riducendo la quantità di memoria dedicata ai campi tag a favore dei dati. Anche se l'utilizzo di blocchi di dimensioni maggiori riduce la frequenza delle miss, questa strategia può causare un aumento della penalità di miss. Se la penalità di miss crescesse linearmente con la dimensione dei blocchi, l'utilizzo di blocchi di dimensioni maggiori porterebbe rapidamente a un peggioramento delle prestazioni.

Per evitare una caduta delle prestazioni, la banda della memoria principale viene aumentata in modo che il trasferimento dei blocchi dalla cache avvenga nel modo più efficiente possibile.

I due metodi più utilizzati per aumentare la larghezza di banda sono: rendere la memoria più larga oppure organizzarla in blocchi interallacciati. Nel tempo i progettisti delle DRAM hanno continuamente migliorato l'interfaccia tra processore e memoria per aumentare la larghezza di banda dei trasferimenti a raffica (in modalità *burst*), in modo da diminuire la latenza del trasferimento di blocchi di cache più ampi.

## Autovalutazione

La velocità del sistema di memoria condiziona il progettista nel decidere quale dimensione adottare per il blocco della cache. Quali delle seguenti linee guida possono essere considerate generalmente valide nella progettazione di una cache?

1. Più piccola è la latenza della memoria, più piccolo sarà il blocco della cache.
2. Più piccola è la latenza della memoria principale, più grande sarà il blocco della cache.
3. Più larga è la banda della memoria, più piccolo sarà il blocco della cache.
4. Più larga è la banda della memoria, più grande sarà il blocco della cache.

## 5.4 Come misurare e migliorare le prestazioni di una cache

Iniziamo questo paragrafo esaminando diversi modi per misurare e analizzare le prestazioni di una cache; esploreremo quindi due diverse tecniche per incrementare le prestazioni. La prima è basata sulla riduzione della frequenza delle miss, ottenuta riducendo la probabilità che due differenti blocchi di memoria principale entrino in conflitto per la stessa locazione della cache. La seconda tecnica riduce la penalità delle miss introducendo nella gerarchia un livello aggiuntivo. Questa secondo metodo, chiamato *cache multilivello (multilevel caching)*, fu introdotto per la prima volta nel 1990 in calcolatori di fascia alta, venduti a una cifra superiore ai 100 000 dollari; da allora è divenuto comune anche nei dispositivi mobili, venduti per alcune centinaia di dollari!

Il tempo di CPU può essere suddiviso nei cicli di clock che la CPU spende per eseguire un programma e in quelli che la CPU trascorre in attesa di una risposta dal sistema di memoria. Di solito si suppone che il costo degli accessi alla cache che producono una hit rientri nei cicli di clock della normale esecuzione della CPU. Si ottiene quindi:

$$\text{Tempo di CPU} = (\text{Cicli di esecuzione CPU} + \text{Cicli di stallo della memoria}) \times \text{Durata del ciclo di clock}$$

I cicli di stallo della memoria sono dovuti principalmente alle miss della cache. In questo capitolo considereremo valida questa ipotesi e restringeremo la discussione a un modello semplificato del sistema di memoria. Nei processori reali, gli stalli generati da operazioni di lettura e scrittura possono risultare anche molto complessi e una previsione accurata delle prestazioni richiede di solito che siano effettuate simulazioni molto dettagliate del processore e del sistema di memoria.

I cicli di stallo dovuti alla memoria possono essere definiti come la somma degli stalli provocati dalle letture e di quelli provocati dalle scritture:

$$\text{Cicli di stallo della memoria} = (\text{Cicli di stallo in lettura} + \text{Cicli di stallo in scrittura})$$

I cicli di stallo dovuti alle letture possono, a loro volta, essere definiti in termini di numero di accessi in lettura fatti dal programma, di penalità di miss per ogni operazione di lettura (in cicli di clock) e di frequenza di miss per le letture:

$$\text{Cicli di stallo in lettura} = \frac{\text{Letture}}{\text{Programma}} \times \text{Frequenza di miss in lettura} \times \times \text{Penalità di miss in lettura}$$

Il discorso diventa più complesso per la scrittura. In uno schema write-through, sono due le cause che generano uno stallo del processore: le miss in scrittura, che di solito richiedono di prelevare dalla memoria principale l'intero blocco prima di continuare la scrittura (per maggiori dettagli sulla gestione delle scritture vedi *Approfondimento* nel sottoparagrafo *Gestione della scrittura* del paragrafo 5.3), e la saturazione del buffer di scrittura, che si verifica quando viene richiesta una scrittura ma il buffer di scrittura risulta pieno. Quindi, il numero di cicli di stallo provocati dalle scritture si ottiene come somma di queste due componenti:

$$\text{Cicli di stallo in scrittura} = \left[ \frac{\text{Scritture}}{\text{Programma}} \times \text{Frequenza di miss in scrittura} \times \times \text{Penalità di miss in scrittura} \right] + \text{Numero di stalli del buffer di scrittura}$$

Poiché gli stalli del buffer di scrittura dipendono da quanto ravvicinate nel tempo sono le scritture, e non solo dalla loro frequenza, non è possibile fornire una semplice equazione per calcolare il numero di stalli. Fortunatamente, nei sistemi provvisti di buffer di scrittura sufficientemente profondi (per es. quattro o più parole) e di una memoria capace di scrivere dati a una frequenza molto superiore (per es. di un fattore 2) alla frequenza media delle richieste di scrittura dei programmi, gli stalli dovuti al buffer saranno molto pochi e si possono quindi trascurare. Se un sistema non dovesse soddisfare questi criteri, sarebbe un sistema progettato male: il progettista avrebbe dovuto utilizzare un buffer di scrittura più profondo, oppure adottare una politica di scrittura di tipo write-back.

Gli schemi di scrittura di tipo write-back hanno un'ulteriore possibile causa di stallo, che nasce dalla necessità di scrivere un blocco della cache nella memoria principale quando il blocco della cache viene rimpiazzato. Discuteremo più in dettaglio questo problema nel paragrafo 5.8.

Nell'organizzazione di molte cache di tipo write-through, la penalità delle miss in lettura e in scrittura è uguale e corrisponde al tempo richiesto per prelevare un blocco dalla memoria principale. Se supponiamo che gli stalli provocati dal buffer di scrittura siano trascurabili, possiamo considerare una sola frequenza e una sola penalità di miss per le letture e le scritture:

$$\text{Cicli di stallo della memoria} = \frac{\text{Accessi alla memoria}}{\text{Programma}} \times \text{Frequenza di miss} \times \\ \times \text{Penalità di miss}$$

che si può scrivere anche come:

$$\text{Cicli di stallo della memoria} = \frac{\text{Istruzioni}}{\text{Programma}} \times \frac{\text{Miss}}{\text{Istruzioni}} \times \text{Penalità di miss}$$

Consideriamo ora un semplice esempio per comprendere l'impatto delle prestazioni delle cache sulle prestazioni complessive di un processore.

### Determinazione delle prestazioni di una cache

#### ESEMPIO

Si supponga una frequenza di miss del 2% per una cache istruzioni e del 4% per una cache dati. Si supponga anche che il processore abbia un CPI di 2 quando non si verificano stalli dovuti alla memoria e che la penalità di miss sia pari a 100 cicli di clock per tutte le miss. Determinare di quanto sarebbe più veloce il processore se fosse dotato di una cache ideale che non provochi mai miss. Si ipotizzi anche che la percentuale delle istruzioni di load e di store sul totale delle istruzioni sia del 36%.

#### SOLUZIONE

Il numero di cicli di clock spesi per gestire le miss per le istruzioni, espresso in numero di istruzioni, I, è:

$$\text{Cicli di miss per le istruzioni} = I \times 2\% \times 100 = 2,00 \times I$$

Dato che la percentuale delle istruzioni di load e di store è del 36%, possiamo calcolare il numero di cicli di clock spesi per le miss per l'accesso ai dati:

$$\text{Cicli di miss per i dati} = I \times 36\% \times 4\% \times 100 = 1,44 \times I$$

(continua)

(continua)

Il numero totale di cicli di clock di stallo dovuti alla memoria è quindi pari a  $2,0 \times I + 1,44 \times I = 3,44 \times I$ , che è più di 3 cicli di clock per istruzione. Di conseguenza, il CPI totale, che tiene conto anche degli stalli dovuti alla memoria, è pari a  $2 + 3,44 = 5,44$ . Dato che non muta né il numero di istruzioni né la frequenza di clock, il rapporto tra i tempi di esecuzione della CPU è:

$$\frac{\text{Tempo di CPU con stalli}}{\text{Tempo di CPU con cache ideale}} = \frac{I \times \text{CPI}_{\text{stallo}} \times \text{Periodo di clock}}{I \times \text{CPI}_{\text{ideale}} \times \text{Periodo di clock}} = \frac{\text{CPI}_{\text{stallo}}}{\text{CPI}_{\text{ideale}}} = \frac{5,44}{2}$$

Perciò, le prestazioni del processore dotato di una cache ideale saranno migliori di un fattore pari a:

$$\frac{5,44}{2} = 2,72$$

Che cosa succede se il processore è reso più veloce ma il sistema di memoria rimane invariato? In questo caso, la quantità di tempo spesa negli stalli dovuti alla memoria aumenterà e diventerà la parte preponderante del tempo totale di esecuzione. La Legge di Amdahl, che abbiamo esaminato nel Capitolo 1, spiega questo meccanismo. Pochi semplici esempi consentono di mostrare quanto questo problema sia importante. Immaginiamo di accelerare il calcolatore considerato nell'esempio precedente riducendo il suo CPI da 2 a 1 senza modificare la frequenza di clock (cosa che può essere ottenuta migliorando la pipeline). Il sistema, considerando le miss della cache, avrà ora un CPI pari a  $1 + 3,44 = 4,44$ , mentre il sistema dotato di una cache ideale sarà:

$$\frac{4,44}{1} = 4,44 \text{ volte più veloce}$$

La quantità di tempo di esecuzione spesa in stalli dovuti alla memoria salirebbe quindi da:

$$\frac{3,44}{5,44} = 63\%$$

a:

$$\frac{3,44}{4,44} = 77\%$$

Analogamente, aumentare la frequenza di clock senza modificare il sistema di memoria aumenterebbe la perdita di prestazioni dovuta alle miss della cache.

Gli esempi e le equazioni viste in precedenza sono basati sull'ipotesi che il tempo di hit non sia un fattore significativo nel calcolo delle prestazioni di una cache. Chiaramente, se il tempo di hit aumentasse, il tempo totale richiesto per accedere a una parola del sistema di memoria aumenterebbe, facendo crescere quindi verosimilmente la durata del periodo di clock del processore. Anche se analizzeremo più avanti altre situazioni in cui aumenta il tempo di hit, introduciamo ora uno degli elementi cruciali: la dimensione della cache. Una cache più grande, infatti, avrà bisogno di un tempo di accesso maggiore; tornando all'analogia con la biblioteca, è come se la scrivania a cui siede lo studente fosse molto grande: occorrerebbe più tempo per trovare un libro rispetto a una scrivania più

piccola, perché ci sarebbero più libri da consultare. Un aumento del tempo di hit è facile che comporti la necessità di aggiungere un altro stadio alla pipeline, poiché occorrerebbero più cicli di clock per accedere e trasferire i dati da una cache a seguito di una hit. Sebbene sia più difficile calcolare l'impatto sulle prestazioni di una pipeline più profonda, ci sarà comunque una dimensione della cache oltre la quale l'aumento del tempo di hit supera il miglioramento ottenuto nella frequenza di hit, portando a un peggioramento delle prestazioni globali del processore.

Per tenere conto del fatto che il tempo di accesso ai dati influenza le prestazioni sia nel caso delle hit sia delle miss, i progettisti utilizzano a volte il *tempo medio di accesso alla memoria* (AMAT, Average Memory Access Time) per esaminare e confrontare tra loro cache diverse. Questo tempo è dato dalla seguente relazione:

$$\text{AMAT} = \text{Durata di una hit} + \text{Frequenza di miss} \times \text{Penalità di miss}$$

### Calcolo del tempo medio di accesso alla memoria

#### ESEMPIO

Determinare l'AMAT di un processore avente un periodo di clock di 1 ns, una penalità di miss di 20 cicli di clock, una frequenza di miss di 0,05 miss per istruzione e un tempo di accesso alla cache, comprendente la generazione del segnale di hit, di 1 ciclo di clock. Si supponga che la penalità di miss sia la stessa per le letture e le scritture e si ignorino gli altri stalli in scrittura.

#### SOLUZIONE

Il tempo medio di accesso alla memoria per singola istruzione è

$$\begin{aligned}\text{AMAT} &= \text{Durata di una hit} + \text{Frequenza di miss} \times \text{Penalità di miss} \\ &= 1 + 0,05 \times 20 \\ &= 2 \text{ cicli di clock}\end{aligned}$$

ossia 2 ns.

Di seguito analizzeremo le possibili diverse organizzazioni di una cache che consentono di diminuire la frequenza di miss; il prezzo da pagare, a volte, è l'aumento del tempo di hit. Discuteremo altri esempi nel paragrafo 5.16.

### Riduzione delle miss di una cache utilizzando un posizionamento più flessibile dei blocchi

**Cache completamente associativa:** una cache in cui un blocco della memoria principale può essere scritto o letto da un qualsiasi blocco della cache.

Finora abbiamo utilizzato uno schema semplice per determinare la posizione di un blocco che deve essere caricato in cache in cui il blocco può trovarsi in una sola posizione. Come già detto in precedenza, questo schema viene chiamato schema a *mappatura diretta*, poiché esiste una mappatura diretta degli indirizzi della memoria principale sulle singole locazioni del livello superiore della gerarchia. In realtà, esiste un ventaglio di schemi per determinare la posizione di un blocco. A un'estremità sta la mappatura diretta, all'estremità opposta si trova uno schema nel quale un blocco della memoria principale può essere scritto in una *qualsiasi* locazione della cache. Una cache che implementi tale schema è chiamata **cache completamente associativa**, perché un blocco della memoria può essere associato a un qualsiasi blocco della cache. Per trovare un blocco in una cache di questo tipo, la ricerca deve essere effettuata su tutti gli elementi della cache, dato che il blocco può trovarsi in una qualsiasi posizione. Per ren-

derla efficace, la ricerca viene svolta in parallelo utilizzando un comparatore per ogni blocco della cache. Tali comparatori incrementano significativamente il costo dell'hardware, rendendo realizzabile lo schema totalmente associativo, di fatto, solo per cache con un numero ridotto di blocchi.

Le cache che adottano schemi intermedi tra la mappatura diretta e quella completamente associativa vengono chiamate **cache set-associative**. In una cache set-associativa ogni blocco della memoria principale può essere caricato in un numero prefissato di posizioni alternative (almeno due), e una cache con  $n$  possibili scelte è chiamata set-associativa a  $n$  vie. Una cache set-associativa a  $n$  vie è costituita da un certo numero di linee (o insiemi), ciascuna costituita da  $n$  blocchi. Ciascun blocco della memoria principale viene mappato su un'unica linea della cache, individuata dal campo indice, e il blocco di dati può essere scritto in *uno qualsiasi* dei blocchi costituenti la linea. Quindi una cache set-associativa combina la mappatura diretta con il posizionamento completamente associativo di un blocco sulla linea: un blocco di memoria principale viene associato a una linea della cache in maniera diretta, e tutti i blocchi della linea vengono esaminati per verificare se contengono l'elemento cercato. La Figura 5.14 mostra dove può essere scritto il blocco 12 della memoria principale in una cache costituita da otto blocchi in totale secondo i tre diversi schemi di posizionamento.

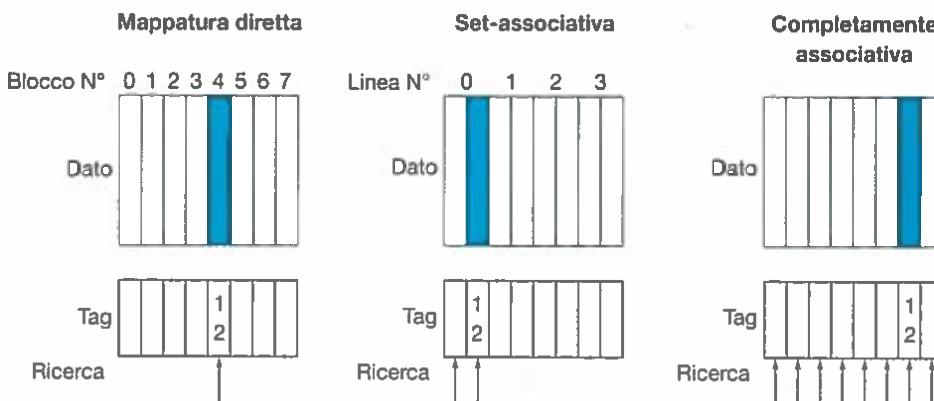
Come abbiamo già visto, in una cache a mappatura diretta la posizione di un blocco della memoria principale è data da:

(Numero del blocco) modulo (Numero dei *blocchi* nella cache)

Nella cache set-associativa, la linea che contiene il blocco viene individuata da:

(Numero del blocco) modulo (Numero delle *linee* della cache)

Dato che il blocco di dati può essere scritto in un qualsiasi blocco di una linea, la ricerca deve essere effettuata *analizzando i campi tag di tutti i blocchi della linea*. In una cache completamente associativa, il blocco può essere messo in una linea qualsiasi della cache, pertanto deve essere analizzato il campo tag di tutti i blocchi della cache durante la ricerca.



**Figura 5.14** La posizione del blocco della memoria principale avente indirizzo 12 in una cache a otto blocchi varia a seconda che la cache sia a mappatura diretta, set-associativa o completamente associativa. Nella cache a mappatura diretta, il blocco 12 può essere scritto in un solo blocco della cache e il numero di questo blocco si ottiene dall'operazione  $(12 \text{ modulo } 8) = 4$ . La cache set-associativa a due vie è costituita da quattro linee e il blocco 12 della memoria dovrà trovarsi nella linea  $(12 \text{ modulo } 4) = 0$ , cioè in uno dei due blocchi che costituiscono questa linea. Nella cache completamente associativa, il blocco di indirizzo 12 può trovarsi in uno qualsiasi degli otto blocchi.

**Cache set-associativa:** una cache che possiede un numero prestabilito di locazioni alternative (almeno due) nelle quali può essere scritto o letto ogni blocco della memoria principale.

## Cache set-associativa a una via

(a mappatura diretta)

Blocco	Tag	Dato
0		
1		
2		
3		
4		
5		
6		
7		

## Cache set-associativa a due vie

Linea	Tag	Dato	Tag	Dato
0				
1				
2				
3				

## Cache set-associativa a quattro vie

Linea	Tag	Dato	Tag	Dato	Tag	Dato	Tag	Dato
0								
1								

## Cache set-associativa a otto vie (completamente associativa)

Tag	Dato												

**Figura 5.15** Cache di otto blocchi configurata come cache a mappatura diretta, set-associativa a due vie, set-associativa a quattro vie e completamente associativa. La dimensione complessiva della cache in termini di blocchi è uguale al numero di linee moltiplicato per il grado di associatività. Perciò, definita la dimensione della cache, al crescere dell'associatività diminuisce il numero delle linee, mentre cresce il numero dei blocchi sulla stessa linea. Con otto blocchi, una cache set-associativa a otto vie è equivalente a una cache completamente associativa.

I diversi schemi di corrispondenza si possono anche vedere come varianti della politica set-associativa. La Figura 5.15 mostra le possibili strutture, con diverso grado di associatività, che possono essere adottate per una cache costituita da otto blocchi. Una cache a mappatura diretta è semplicemente una cache set-associativa a una sola via: ogni elemento della cache contiene un unico blocco, così come ogni linea contiene un solo blocco. Una cache di  $m$  elementi, completamente associativa, è semplicemente una cache set-associativa a  $m$  vie: ha un'unica linea contenente  $m$  blocchi e un elemento può trovarsi in uno qualsiasi dei blocchi.

Il vantaggio dell'aumento del grado di associatività è che, in genere, diminuisce la frequenza delle miss, come mostra l'esempio seguente. Lo svantaggio principale, del quale discuteremo più avanti, è il potenziale aumento del tempo di hit.

**Miss e associatività nelle cache****ESEMPIO**

Consideriamo tre cache di piccole dimensioni, ciascuna costituita da quattro blocchi di una parola: la prima cache è completamente associativa, la seconda è set-associativa a due vie e la terza è a mappatura diretta. Determinare il numero di miss nei tre casi quando vengono richiesti i blocchi con la seguente sequenza di indirizzi: 0, 8, 0, 6, 8.

(continua)

Il caso della mappatura diretta è il più semplice. Innanzitutto occorre determinare a quale blocco della cache corrispondano i diversi indirizzi:

(continua)

**SOLUZIONE**

Indirizzo del blocco	Blocco della cache
0	(0 modulo 4) = 0
6	(6 modulo 4) = 2
8	(8 modulo 4) = 0

Ora è possibile compilare il contenuto della cache dopo l'indirizzamento di ciascun blocco, come è mostrato nella seguente tabella. I blocchi non validi non contengono nulla; il testo in blu mostra il blocco appena inserito nella cache, mentre il testo in nero mostra gli elementi precedentemente inseriti nella cache:

Indirizzo del blocco richiesto della memoria principale	Hit o miss	Contenuto dei blocchi della cache dopo l'indirizzamento			
		0	1	2	3
0	Miss	Memoria[0]			
8	Miss	Memoria[8]			
0	Miss	Memoria[0]			
6	Miss	Memoria[0]	Memoria[6]		
8	Miss	Memoria[8]	Memoria[6]		

La cache a mappatura diretta genera cinque miss per i cinque accessi.

La cache set-associativa è costituita da due linee, aventi indice 0 e 1, ciascuna contenente due blocchi. Come prima cosa occorre individuare a quale linea corrisponda ciascun indirizzo di blocco:

Indirizzo del blocco	Linea della cache
0	(0 modulo 2) = 0
6	(6 modulo 2) = 0
8	(8 modulo 2) = 0

Dato che, in caso di miss, si può scegliere quale blocco della linea sostituire, è necessario stabilire una modalità di sostituzione. Di solito le cache set-associative sostituiscono il blocco della linea utilizzato meno di recente; discuteremo in dettaglio altri metodi di sostituzione più avanti. Utilizzando questa regola di sostituzione, il contenuto della cache set-associativa dopo ogni accesso risulterà essere:

Indirizzo del blocco richiesto della memoria principale	Hit o miss	Contenuto dei blocchi della cache dopo l'indirizzamento			
		Linea 0	Linea 0	Linea 1	Linea 1
0	Miss	Memoria[0]			
8	Miss	Memoria[0]	Memoria[8]		
0	Hit	Memoria[0]	Memoria[8]		
6	Miss	Memoria[0]	Memoria[6]		
8	Miss	Memoria[8]	Memoria[6]		

(continua)

(continua)

Si noti che quando viene indirizzato il blocco 6, questo va a sostituire il blocco 8, poiché il blocco 8 è stato indirizzato meno di recente rispetto al blocco 0. La cache set-associativa a due vie genera in totale 4 miss, una in meno della cache a mappatura diretta.

La cache completamente associativa ha quattro blocchi, distribuiti su una singola linea: ogni blocco della memoria principale può quindi essere memorizzato in un qualsiasi blocco della cache. Le prestazioni di questo schema sono le migliori, dato che si verificano solamente 3 miss, come si può vedere dalla tabella seguente:

Indirizzo del blocco richiesto della memoria principale	Hit o miss	Contenuto dei blocchi della cache dopo l'accesso			
		Blocco 0	Blocco 1	Blocco 2	Blocco 3
0	Miss	Memoria[0]			
8	Miss	Memoria[0]	Memoria[8]		
0	Hit	Memoria[0]	Memoria[8]		
6	Miss	Memoria[0]	Memoria[8]	Memoria[6]	
8	Hit	Memoria[0]	Memoria[8]	Memoria[6]	

Per questa sequenza di accessi, 3 miss è il risultato migliore che si può ottenere, poiché vengono richiesti dati aventi 3 indirizzi di blocco distinti. Si noti che se la cache avesse avuto otto blocchi, nella struttura set-associativa a due vie non sarebbero state necessarie sostituzioni (come potete verificare voi stessi) e si sarebbe ottenuto un numero di miss pari a quello della cache completamente associativa. Analogamente, se i blocchi fossero stati sedici, tutte e tre le cache avrebbero avuto lo stesso numero di miss. Questo esempio, sebbene molto semplice, dimostra che la dimensione della cache e il suo grado di associatività non sono parametri indipendenti nel determinarne le prestazioni.

Di quanto si può ridurre la frequenza di miss sfruttando l'associatività? La Figura 5.16 mostra il miglioramento che si ottiene in un'architettura con una cache dati da 64 KiB e blocchi di 16 parole per diversi gradi di associatività, passando dalla mappatura diretta all'organizzazione set-associativa a otto vie. Il passaggio del grado di associatività da una a due vie migliora già di circa il 15% la frequenza di miss, la quale migliora leggermente aumentando ancora il grado di associatività.

Associatività	Frequenza di miss
1	10,3%
2	8,6%
4	8,3%
8	8,1%

**Figura 5.16** Frequenza di miss in una cache dati con associatività variabile da una via a otto vie, avente un'organizzazione simile a quella del processore FastMATH Intrinsicity, valutata utilizzando come benchmark gli SPEC CPU2000. Questi dati, relativi a 10 programmi del benchmark SPEC CPU2000, provengono da Hennessy e Patterson [2003].

Tag	Indice	Offset di blocco
-----	--------	------------------

**Figura 5.17** Le tre parti di un indirizzo in una cache set-associativa o a mappatura diretta. L'indice viene utilizzato per selezionare la linea, mentre per selezionare il blocco all'interno della linea viene comparato il tag dell'indirizzo con i tag dei blocchi. L'offset di blocco rappresenta la posizione del dato desiderato all'interno del blocco.

## Come trovare un blocco nella cache

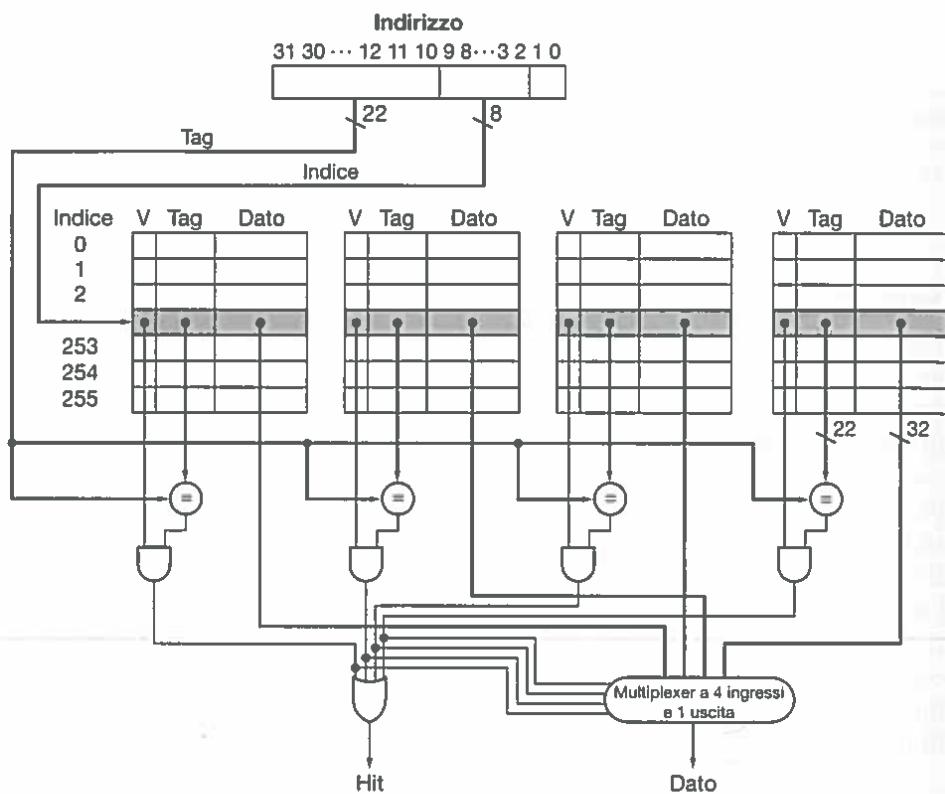
Consideriamo ora il problema di come si possa trovare un blocco all'interno di una cache set-associativa. Proprio come in una cache a mappatura diretta, ogni blocco di una cache set-associativa contiene un campo tag che permette di individuarne l'indirizzo riferito alla memoria principale. Il campo tag di ciascuno dei blocchi contenuti sulla stessa linea della cache viene controllato per verificare se corrisponda a quello dell'indirizzo del dato o dell'istruzione richiesta dal processore. La Figura 5.17 mostra come si scomponga l'indirizzo. Il contenuto del campo indice serve a selezionare la linea che può contenere l'elemento cercato e viene esaminato il campo tag di tutti i blocchi appartenenti alla linea selezionata. Dato che la velocità di esecuzione è fondamentale, il contenuto del campo tag viene esaminato in parallelo per tutti i blocchi della linea selezionata: una ricerca sequenziale renderebbe eccessivo il tempo di hit di una cache set-associativa.

Se si mantiene costante la dimensione totale della cache, al crescere del grado di associatività aumenta il numero dei blocchi contenuti in una singola linea, e quindi anche il numero di confronti che devono essere condotti simultaneamente per effettuare la ricerca in parallelo: ogni volta che raddoppia il grado di associatività, raddoppia anche il numero di blocchi contenuti in una linea e si dimezza il numero di linee. Analogamente, ogni incremento dell'associatività di un fattore 2 fa diminuire di 1 bit la dimensione del campo indice e fa aumentare di 1 bit la dimensione del campo tag. Perciò, in una cache completamente associativa c'è una sola linea e tutti i blocchi devono essere esaminati in parallelo: il campo indice non esiste e tutto l'indirizzo, a parte l'offset del blocco, deve essere confrontato con il tag di ogni blocco. In altre parole, la ricerca viene effettuata su tutta la cache, senza indicizzazione.

Una cache a mappatura diretta richiede un unico comparatore, dato che l'elemento cercato può trovarsi in una sola posizione; in questo caso, si accede alla cache utilizzando solamente il campo indice. In una cache set-associativa a quattro vie, come quella mostrata in Figura 5.18, sono necessari quattro comparatori, oltre a un multiplexer con quattro ingressi e un'uscita, richiesto per scegliere uno dei quattro possibili blocchi della linea selezionata. L'accesso alla cache, quindi, si effettua utilizzando l'indice per individuare la linea appropriata, e poi si esaminano i diversi tag della linea. Il costo di una cache set-associativa consiste nei comparatori aggiuntivi e nel ritardo imposto dalla necessità di confrontare i campi tag con l'indirizzo e di selezionare l'elemento desiderato tra quelli della linea.

In ogni gerarchia delle memorie, la scelta tra lo schema a mappatura diretta, lo schema set-associativo e quello completamente associativo dipende dal confronto tra il costo di una miss e il costo necessario a realizzare l'associatività, dal punto di vista sia del tempo di elaborazione sia dei componenti hardware aggiuntivi.

**Approfondimento.** Una memoria indirizzabile per contenuto (CAM, Content Addressable Memory) è un circuito che combina la comparazione e la memorizzazione in un unico dispositivo. Invece di ricevere un indirizzo e leggere una parola, come nelle RAM, una CAM riceve un dato, cerca al suo interno una copia del dato



**Figura 5.18** La realizzazione di una cache set-associativa a quattro vie richiede quattro comparatori e un multiplexer con 4 ingressi e 1 uscita. I comparatori individuano quale elemento della linea, se esiste, corrisponde al tag dell'indirizzo. L'uscita dei comparatori viene utilizzata per selezionare il dato da uno dei quattro blocchi della linea selezionata, attraverso un multiplexer pilotato dai segnali di hit dei diversi blocchi. In alcune realizzazioni della cache, il segnale di abilitazione dell'uscita, che si trova nella parte del chip dedicata ai dati, può essere utilizzato per selezionare il blocco della linea che deve essere portato in uscita. Il segnale di abilitazione dell'uscita proviene direttamente dai comparatori e consente di portare in uscita il blocco il cui tag corrisponde a quello dell'indirizzo del dato richiesto. Questa struttura elimina la necessità del multiplexer.

in ingresso e restituisce l'indice della riga nella quale ha trovato il dato (se lo trova). Le memorie CAM consentono ai progettisti di realizzare un grado di associatività ancora maggiore di quello che si riesce a ottenere utilizzando SRAM e comparatori. Nel 2013, l'esistenza di CAM più grandi e potenti ha fatto sì che le memorie associative a due e quattro vie fossero costruite ancora con comparatori e blocchi di SRAM standard, mentre le cache a otto vie o con un grado di associatività ancora superiore adottarono memorie CAM.

### Come scegliere il blocco da sostituire

Quando si verifica una miss in una cache a mappatura diretta, il blocco richiesto può essere scritto in un'unica posizione e quello che occupava precedentemente questa posizione deve essere sostituito. In una cache associativa, possiamo scegliere dove scrivere il blocco richiesto e quindi dobbiamo decidere quale sia il blocco da sostituire. In una cache completamente associativa, tutti i blocchi costituiscono un potenziale candidato per la sostituzione. Se la cache è set-associativa, possiamo scegliere tra i blocchi contenuti nella stessa linea della cache.

Lo schema di sostituzione più comunemente utilizzato è detto **utilizzo meno recente (LRU)**, *Least Recently Used* ed è quello che abbiamo adottato nel precedente esempio. Nello schema LRU, il blocco sostituito è quello che è rimasto inutilizzato più a lungo; nel nostro esempio abbiamo quindi sostituito il blocco Memoria(0) invece del blocco Memoria(6).

**Utilizzo meno recente (LRU):** uno schema di sostituzione nel quale il blocco sostituito è quello che è rimasto inutilizzato per più tempo.

La sostituzione tramite schema LRU viene implementata tenendo traccia della sequenza di impiego dei vari blocchi della linea. Per una cache set-associativa a due vie, è sufficiente un unico bit per ogni linea per memorizzare quale dei due blocchi è stato utilizzato più di recente; questo bit viene impostato in maniera opportuna ogni volta che si utilizza uno dei due blocchi della linea. All'aumentare del grado di associatività, l'implementazione dello schema LRU diventa sempre più difficile; nel paragrafo 5.8 analizzeremo uno schema di sostituzione alternativo.

### Dimensione del campo tag e grado di associatività

Incrementare l'associatività di una cache comporta più comparatori e più bit per il campo tag di ogni blocco della cache. Considerando una cache costituita da 4096 blocchi, con blocchi da quattro parole e 64 bit di indirizzo, determinare il numero totale di linee e il numero di bit del campo tag per una cache a mappatura diretta, una cache set-associativa a due vie e a quattro vie e per una cache completamente associativa.

Dato che ci sono  $16 (= 2^4)$  byte per blocco e l'indirizzo è su 64 bit,  $64 - 4 = 60$  bit devono essere utilizzati per l'indice e il tag. La cache a mappatura diretta ha un numero di linee pari al numero di blocchi; occorrono quindi 12 bit per l'indice, essendo  $\log_2(4096) = 12$ . Il numero totale di bit per il campo tag è perciò di  $(60 - 12) \times 4096 = 48 \times 4096 = 197$  K.

Ogni volta che raddoppia il grado di associatività, diminuisce il numero di linee di un fattore 2; decresce quindi di 1 il numero di bit utilizzati per l'indice e aumenta di 1 il numero di bit del campo tag. Quindi, per una cache set-associativa a due vie, ci sono 2048 linee e il numero totale di bit del campo tag è pari a  $(60 - 11) \times 2 \times 2048 = 98 \times 2048 = 401$  Kbit. Per una cache set-associativa a quattro vie, il numero totale di linee è pari a 1024 e il numero totale di bit del campo tag è  $(60 - 10) \times 4 \times 1024 = 100 \times 1024 = 205$  K.

In una cache completamente associativa esiste un'unica linea costituita da 4096 blocchi e il campo tag è costituito da 60 bit; pertanto il numero totale di bit utilizzati per il campo tag è  $60 \times 4096 \times 1 = 246$  K.

### ESEMPIO

### SOLUZIONE

### Ridurre la penalità di miss utilizzando una cache multilivello

Tutti i moderni calcolatori fanno uso delle cache. Per ridurre ulteriormente la differenza tra l'elevata frequenza di clock dei processori moderni e il tempo di accesso relativamente lungo delle DRAM, molti microprocessori supportano un livello addizionale di cache. Questo secondo livello si trova di solito sullo stesso chip del processore e viene utilizzato quando si verifica una miss all'interno della cache primaria. Se il secondo livello di cache contiene il dato desiderato, la penalità di miss per il primo livello di cache sarà essenzialmente il tempo di accesso al secondo livello di cache, molto inferiore al tempo di accesso alla memoria principale. Se né la cache primaria né quella secondaria contendono il dato, viene richiesto un accesso alla memoria principale e, inevitabilmente, si incorre in una penalità di miss elevata.

Quanto può essere significativo l'incremento di prestazioni dovuto all'utilizzo di una cache secondaria? Rispondiamo a questa domanda attraverso il prossimo esempio.

## Prestazioni di una cache multilivello

### ESEMPIO

Supponiamo di disporre di un processore con frequenza di clock pari a 4 GHz e con un CPI di base pari a 1,0, misurato quando tutti gli accessi alla cache primaria producono una hit. La memoria principale ha un tempo di accesso di 100 ns che comprende anche tutto il tempo richiesto per gestire una miss. Supponiamo anche che la frequenza di miss della cache primaria sia pari al 2% delle istruzioni. Di quanto sarà più veloce il processore se introduciamo una cache secondaria con tempo di accesso di 5 ns, sia per le hit sia per le miss, grande a sufficienza da ridurre la frequenza di miss della memoria principale allo 0,5%?

### SOLUZIONE

La penalità di miss della memoria principale è pari a:

$$\frac{100 \text{ ns}}{0,25 \text{ ns}} = \frac{400 \text{ cicli di clock}}{\text{ciclo di clock}}$$

Il CPI effettivo, considerando un livello di cache, è dato da:

$$\text{CPI totale} = \text{CPI base} + \text{Numero di cicli di stall per istruzione}$$

Per il processore in questione, con un livello di cache si ha:

$$\begin{aligned}\text{CPI totale} &= 1,0 + \text{Numero di cicli di stall per istruzione} \\ &= 1,0 + 2\% \times 400 = 9\end{aligned}$$

Con due livelli di cache, una miss della cache primaria (o di primo livello) può essere soddisfatta dalla cache secondaria o dalla memoria principale. La penalità di miss per un accesso al secondo livello di cache è pari a:

$$\frac{5 \text{ ns}}{0,25 \text{ ns}} = \frac{20 \text{ cicli di clock}}{\text{ciclo di clock}}$$

Se la miss viene risolta dalla cache di secondo livello, questo tempo costituisce la penalità di miss totale; se la miss richiede invece di accedere alla memoria principale, il tempo totale di miss sarà dato dalla somma del tempo di accesso alla cache secondaria e di quello alla memoria principale. Quindi, per una cache a due livelli, il CPI totale è la somma del CPI di base e dei cicli di stall introdotti da entrambi i livelli di cache:

$$\begin{aligned}\text{CPI totale} &= 1 + \text{stalli primari per istruzione} + \text{stalli secondari per istruzione} \\ &= 1 + 2\% \times 20 + 0,5\% \times 400 = 1 + 0,4 + 2 = 3,4\end{aligned}$$

Perciò, il processore con il secondo livello di cache è più veloce di un fattore:

$$\frac{9,0}{3,4} = 2,6$$

In alternativa, avremmo potuto calcolare il numero di cicli di stall totali sommando i cicli di stall dovuti agli accessi alla cache secondaria che hanno prodotto una hit  $[(2\% - 0,5\%) \times 20 = 0,3]$  e gli accessi alla memoria principale ai quali bisogna aggiungere anche il costo degli accessi alla cache secondaria  $[0,5\% \times (20 + 400) = 2,1]$ . La somma  $1,0 + 0,3 + 2,1$  dà ancora 3,4.

Le considerazioni da fare in fase di progettazione sono molto diverse per le cache primaria e secondaria perché la presenza di una seconda cache modifica le specifiche. In particolare, una cache a due livelli permette di progettare la cache primaria focalizzandosi sulla minimizzazione del tempo di hit, in modo da ottenere un periodo di clock più breve oppure un numero minore di stadi di pipeline e permette di progettare la cache secondaria focalizzandosi sulla frequenza delle miss per ridurre la penalità dovuta ai lunghi tempi di accesso alla memoria principale.

L'effetto di queste modifiche sulle cache a due livelli si può osservare confrontando ciascuna cache con la cache ottimale a singolo livello. Rispetto a questa, la cache primaria di una **cache multilivello** è spesso di dimensioni ridotte e utilizza di solito blocchi di grandezza inferiore, che consentono una dimensione complessiva inferiore della cache e una penalità di miss ridotta. La cache secondaria, invece, è molto più grande della cache a un solo livello, dato che il tempo di accesso di una cache secondaria è meno critico; avendo una dimensione maggiore, potrà utilizzare, in generale, blocchi più ampi di quelli adatti a una cache a singolo livello. Inoltre, può utilizzare un grado più elevato di associatività della cache primaria, poiché l'obiettivo della cache secondaria è la riduzione della frequenza di miss.

**Cache multilivello:** una gerarchia delle memorie composta da più livelli di cache, invece che da una sola cache e dalla memoria principale.

Il problema dell'ordinamento (*sorting*) di un insieme di elementi è stato ampiamente studiato e ha portato all'ideazione di algoritmi ad ampia diffusione, tra cui: Bubble Sort, Quicksort, Radix Sort e così via. Il numero di istruzioni eseguite per ogni elemento cercato è mostrato in Figura 5.19a per gli algoritmi Quicksort e Radix Sort. Come si poteva immaginare, per vettori di grandi dimensioni Radix Sort mostra un vantaggio nei confronti di Quicksort in termini di numero di operazioni. La Figura 5.19b mostra invece il tempo impiegato per ogni ordinamento. Vediamo che le curve hanno un andamento simile a quello di Figura 5.19a quando gli elementi da ordinare sono pochi, ma la curva relativa a Radix Sort, dopo una discesa iniziale, tende a salire man mano che gli elementi aumentano. Come si spiega questo comportamento? La Figura 5.19c fornisce la risposta, mostrando il numero di miss della cache per ogni elemento che viene ordinato: Quick Sort genera in maniera consistente meno miss per ogni elemento ordinato.

Purtroppo l'analisi standard degli algoritmi ignora spesso l'impatto della gerarchia delle memorie. Come l'incremento della frequenza di clock e la **Legge di Moore** consentono ai progettisti dei calcolatori di ottenere le migliori prestazioni dall'esecuzione di una sequenza di istruzioni, così l'uso efficiente della gerarchia delle memorie è un elemento cruciale per il miglioramento delle prestazioni del sistema. Come abbiamo visto nell'introduzione, comprendere il comportamento della gerarchia delle memorie è quindi fondamentale per valutare le prestazioni dei programmi in un calcolatore moderno.

## Capire le prestazioni dei programmi

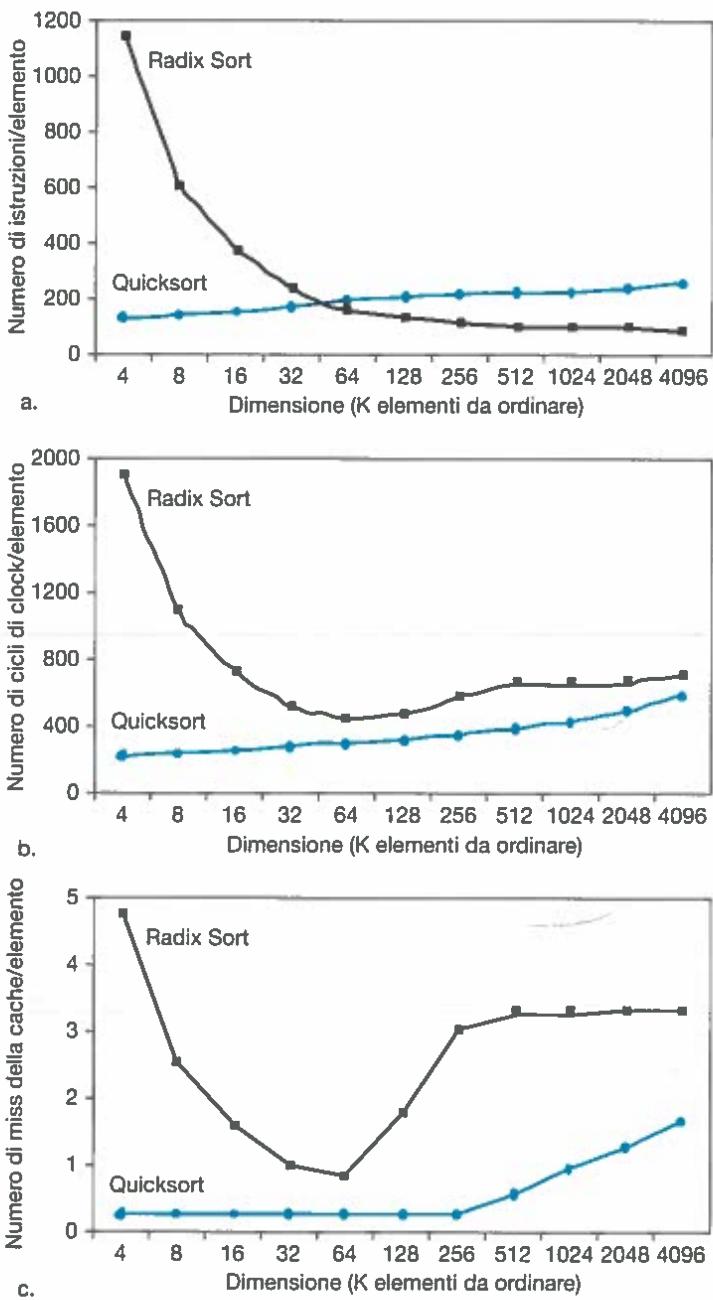


LEGGE DI MOORE

## Ottimizzazione software mediante elaborazione a blocchi

Dato l'impatto delle gerarchie delle memorie sulle prestazioni dei programmi, molte ottimizzazioni software sono state inventate per migliorare notevolmente le prestazioni riutilizzando i dati già contenuti in cache, riducendo la frequenza di miss attraverso l'aumento della località temporale.

Quando si devono utilizzare delle matrici, si possono ottenere buone prestazioni dal sistema di memoria se si memorizzano le matrici in modo tale che l'accesso in memoria ai loro elementi sia sequenziale. Supponiamo di dovere gestire delle matrici multidimensionali e che ad alcune di esse si acceda per righe e ad altre per colonne. La memorizzazione delle matrici per riga o per colonna



**Figura 5.19** Confronto tra Quicksort e Radix Sort in termini di numero di istruzioni eseguite per elemento ordinato (a), tempo per elemento ordinato (b) e miss della cache per elemento ordinato (c). Questi dati provengono da un articolo di LaMarca e Ladner [1996]. Anche se i numeri potranno essere diversi per i nuovi calcolatori, il concetto non cambia. In seguito a questi risultati sono state realizzate nuove versioni dell'algoritmo Radix Sort, che tengono in considerazione anche la gerarchia delle memorie al fine di riguadagnare il suo vantaggio algoritmico (par. 5.15). L'idea base delle ottimizzazioni sull'uso della cache è quella di utilizzare tutti i dati all'interno del blocco ripetutivamente prima della sua sostituzione in seguito a una miss.

non risolve il problema, perché sia le righe sia le colonne vengono utilizzate in ogni iterazione del ciclo.

Invece di operare su intere righe o intere colonne della matrice, gli algoritmi *a blocchi* operano su sottomatrici o *blocchi*. L'obiettivo è massimizzare l'accesso ai dati caricati nella cache prima che debbano essere sostituiti, cioè migliorare la località temporale per ridurre le miss della cache.

Per esempio, i cicli interni della funzione DGEMM (dalla linea 4 alla linea 9 di Figura 3.21 sono:

```
for (int j = 0; j < n; j++)
{
    double cij = C[i+j*n]; /* cij = C[i][j] */
    for (int k = 0; k < n; k++)
        cij += A[i+k*n]*B[k+j*n]; /* cij += A[i][k]*B[k][j] */
    C[i+j*n] = cij; /* C[i][j] = cij */
}
```

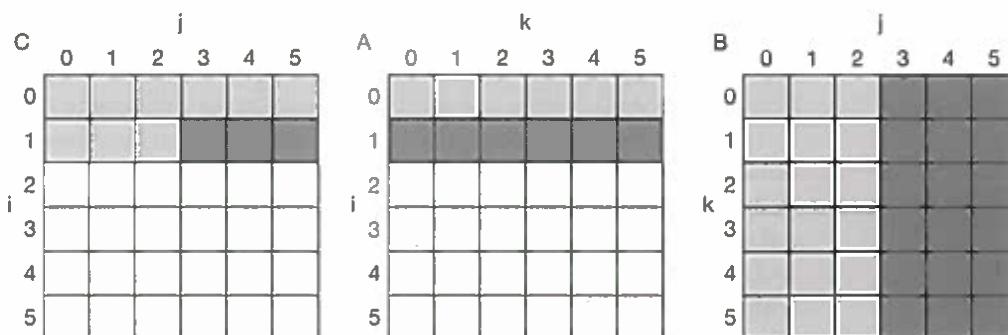
Questa parte di codice legge tutti gli  $N$  per  $N$  elementi di  $B$ , legge tutti gli elementi della  $i$ -esima riga di  $A$  e scrive il risultato nella  $i$ -esima riga di  $N$  elementi di  $C$ . (I commenti rendono facile capire a quali righe e colonne ci stiamo riferendo.) La Figura 5.20 fotografa l'accesso alle tre matrici: le celle grigio scuro indicano un accesso recente, le celle grigio chiaro un accesso precedente e le celle bianche indicano che non ci sono stati ancora accessi.

Il numero delle miss causate dalla dimensione dipende chiaramente da  $N$  e dalla dimensione della cache. Se la cache può contenere tutte e tre le matrici di dimensioni  $N$  per  $N$ :  $A$ ,  $B$  e  $C$ , allora non ci sono conflitti sulla cache. Abbiamo opportunamente scelto  $N$  uguale a 32 negli esempi dei Capitoli 3 e 4 per rientrare in questa situazione. Ciascuna matrice è in questo caso di  $32 \times 32 = 1024$  elementi, ciascuno di 8 byte; le tre matrici occupano quindi 24 KiB, che possono essere contenuti comodamente nella cache dati da 32 KiB del Core i7 Intel (Sandy Bridge).

Se la cache può contenere una matrice  $N$  per  $N$  e una riga di  $N$  elementi, allora almeno l' $i$ -esima riga di  $A$  e la matrice  $B$  possono essere contenute in cache. Se la cache non può contenere neppure questi, si possono verificare miss anche per le matrici  $B$  e  $C$ . Nel caso peggiore, si accede a  $2N^3 + N^2$  parole di memoria per  $N^3$  operazioni.

Per essere sicuri che gli elementi a cui si accede siano contenuti nella cache, modifichiamo il codice originale per effettuare i calcoli su sottomatrici. Chiameremo quindi ripetutamente la versione della funzione DGEMM riportata in Figura 4.78 su matrici di dimensioni  $\text{DIM\_BLOCCO}$  per  $\text{DIM\_BLOCCO}$ , dove  $\text{DIM\_BLOCCO}$  è chiamato *fattore di blocco* (*blocking factor*).

La Figura 5.21 mostra la versione a blocchi della funzione DGEMM. La funzione `esegui_blocco` è uguale alla funzione di Figura 3.22 con tre parametri aggiuntivi: `si`, `sj` e `sk`, che specificano la posizione iniziale delle sottomatrici



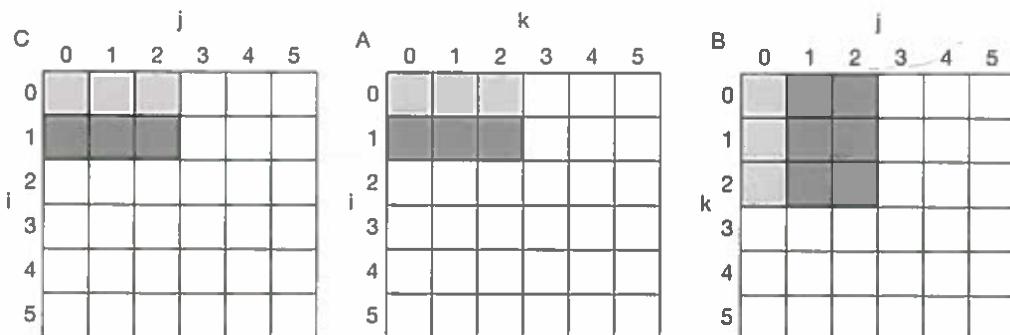
**Figura 5.20** Fotografia delle tre matrici  $C$ ,  $A$  e  $B$  quando  $N = 6$  e  $i = 1$ . Il livello di grigio degli elementi delle matrici rappresenta quanto di recente si è verificato l'ultimo accesso a quell'elemento: grigio chiaro indica accessi più remoti, grigio scuro indica accessi più recenti. Gli elementi bianchi non hanno ancora ricevuto accessi. In confronto alla Figura 5.22, gli elementi di  $A$  e  $B$  vengono letti ripetutamente per calcolare il nuovo elemento di  $C$ . Le variabili  $i$ ,  $j$  e  $k$  vengono mostrate sulle righe e colonne a seconda di come vengono utilizzate per accedere agli elementi delle matrici.

```

1 #define DIM_BLOCCO 32
2 void esegui_blocco (int n, int si, int sj, int sk, double *A, double
3 *B, double *C)
4 {
5     for (int i = si; i < si+DIM_BLOCCO; ++i)
6         for (int j = sj; j < sj+DIM_BLOCCO; ++j)
7             {
8                 double cij = C[i+j*n];/* cij = C[i][j] */
9                 for( int k = sk; k < sk+DIM_BLOCCO; k++ )
10                     cij += A[i+k*n] * B[k+j*n];/* cij+=A[i][k]*B[k][j] */
11                     C[i+j*n] = cij; /* C[i][j] = cij */
12             }
13 }
14 void dgemm (int n, double* A, double* B, double* C)
15 {
16     for ( int sj = 0; sj < n; sj += DIM_BLOCCO )
17         for ( int si = 0; si < n; si += DIM_BLOCCO )
18             for ( int sk = 0; sk < n; sk += DIM_BLOCCO )
19                 esegui_blocco(n, si, sj, sk, A, B, C);
20 }

```

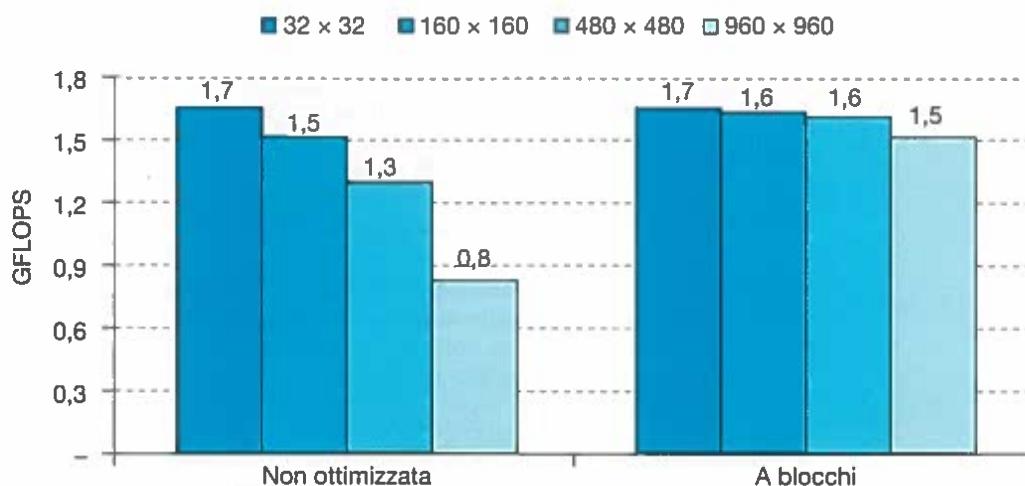
**Figura 5.21** Versione a blocchi della funzione DGEMM di Figura 3.21. Si supponga che la matrice **C** sia stata inizializzata a zero. La funzione **esegui\_blocco** è uguale alla funzione di Figura 3.21 con tre parametri aggiuntivi che specificano la posizione iniziale delle sottomatrici di dimensione **DIM\_BLOCCO**. L'ottimizzatore del compilatore gcc rimuove ogni aggiunta richiesta dalle chiamate a funzione inserendo la funzione direttamente nel codice.



**Figura 5.22** Tempo dall'ultimo accesso per le matrici **C**, **A** e **B** quando **DIM\_BLOCCO** = 3. Si noti che, a differenza di Figura 5.20, pochi elementi hanno ricevuto accessi.

in **A**, **B** e **C**. I due cicli interni di **esegui\_blocco** eseguono ora i calcoli su **DIM\_BLOCCO** elementi invece che su tutta la lunghezza delle righe e colonne di **B** e **C**. L'ottimizzatore del compilatore gcc rimuove ogni aggiunta richiesta dalle chiamate a funzione inserendo la funzione direttamente nel codice (“Inlining” della funzione) per evitare le istruzioni che gestiscono il passaggio dei parametri e il ritorno alla procedura chiamante.

La Figura 5.22 illustra gli accessi alle tre matrici utilizzando l'elaborazione a blocchi. Analizzando solamente le miss causate dalla dimensione, il numero totale di parole a cui si accede in memoria è sceso a  $2 N^3 / \text{DIM\_BLOCCO} + N^2$ , che costituisce un miglioramento di un fattore **DIM\_BLOCCO** approssimativamente. Quindi, l'elaborazione a blocchi sfrutta una combinazione di località spaziale e temporale, dato che **A** beneficia della località spaziale e **B** di quella temporale.



**Figura 5.23** Le prestazioni della funzione DGEMM non ottimizzata (Figura 3.21) confrontata con la versione eseguita su blocchi di dati (Figura 5.21), al variare della dimensione delle matrici da  $32 \times 32$  (dimensione per la quale tutte e tre le matrici possono essere contenute nella memoria cache primaria) a  $960 \times 960$ .

Sebbene il nostro obiettivo fosse ridurre le miss della cache, l'elaborazione a blocchi aiuta anche nell'allocazione dei registri. Scegliendo una dimensione del blocco tale per cui i blocchi possono essere contenuti nei registri, possiamo minimizzare il numero di operazioni del programma di trasferimento da e verso la memoria, cosa che migliora ulteriormente le prestazioni.

La Figura 5.23 mostra l'impatto dell'elaborazione a blocchi sulla funzione DGEMM all'aumentare della dimensione delle matrici oltre la dimensione per la quale tutte e tre le matrici possono essere contenute in cache. Le prestazioni della versione non ottimizzata sono la metà rispetto alla versione a blocchi per matrici  $960 \times 960$ , 900 volte più grandi delle matrici  $32 \times 32$  utilizzate nei Capitoli 3 e 4. La versione a blocchi è invece solamente un 10% più lenta.

**Approfondimento.** Le cache multilivello creano diverse complicazioni. Anzitutto vengono introdotti altri tipi di miss con una loro frequenza di miss. Nell'esempio a pagina 356, abbiamo introdotto la frequenza di miss della cache primaria e la **frequenza di miss globale**, cioè la percentuale di accessi che falliscono in tutti i livelli di cache. C'è anche una frequenza di miss della cache secondaria, che è il rapporto tra tutti i fallimenti nella cache secondaria e il numero totale di accessi ad essa. Questa frequenza di miss è chiamata **frequenza di miss locale** della cache secondaria. Poiché la cache primaria filtra gli accessi, specialmente quelli con una buona località spaziale e temporale, la frequenza di miss locale della cache secondaria è molto più alta della frequenza di miss globale. Per l'esempio a pagina 356, la frequenza di miss locale della cache secondaria risulta pari a  $0,5\%/2\% = 25\%$ ! Fortunatamente è la frequenza di miss globale a indicare quanto spesso dobbiamo accedere alla memoria principale.

**Frequenza di miss globale:**  
la percentuale degli accessi alla memoria che falliscono in tutti i livelli di una cache multilivello.

**Frequenza di miss locale:**  
la percentuale degli accessi a un livello di cache che falliscono; viene utilizzata in una gerarchia multilivello.

**Approfondimento.** Nei processori dotati di esecuzione fuori ordine (vedi Cap. 4) il calcolo delle prestazioni diviene più complesso, dal momento che essi possono eseguire delle istruzioni anche durante la penalità di miss. Invece della frequenza di miss per istruzioni e dati, si utilizza il numero di miss per istruzione calcolato attraverso la seguente formula:

$$\frac{\text{Cicli di stallo dovuti alla memoria}}{\text{Istruzione}} = \\ = \frac{\text{Miss}}{\text{Istruzione}} \times (\text{Latenza totale delle miss} - \text{Latenza delle miss sovrapposte})$$

Non esiste un modo generale per calcolare la latenza delle miss sovrapposte, quindi la valutazione della gerarchia delle memorie per i processori fuori ordine richiede inevitabilmente la simulazione del processore e della gerarchia delle memorie. Solamente osservando l'esecuzione di un processore durante ogni miss possiamo capire se il processore rimane in stallo, in attesa del dato, oppure semplicemente trova altro lavoro da compiere. In generale, la regola è che un processore riesce spesso a nascondere la penalità di miss della cache L1 se il dato viene trovato nella cache L2, ma raramente nasconde una miss della cache L2.

**Approfondimento.** La sfida, per quel che concerne le prestazioni degli algoritmi, è rappresentata dal fatto che le gerarchie delle memorie variano nelle differenti implementazioni hardware della stessa architettura per dimensioni delle cache, grado di associatività, dimensione del blocco e numero di livelli di cache. Per gestire questa variabilità, alcune recenti librerie numeriche rendono parametrici i loro algoritmi ed esplorano lo spazio dei parametri durante l'esecuzione per trovare la miglior combinazione adatta al particolare tipo di calcolatore utilizzato. Questo approccio viene chiamato *autocalibrazione (autotuning)*.

### Autovalutazione

Quale delle seguenti affermazioni è generalmente vera per quel che riguarda la progettazione di cache multilivello?

1. Le cache di primo livello sono focalizzate principalmente sul tempo di accesso, mentre le cache di secondo livello sono focalizzate sulla frequenza di miss.
2. Le cache di primo livello sono focalizzate principalmente sulla frequenza di miss, mentre le cache di secondo livello sono focalizzate sul tempo di accesso.

### Riepilogo

In questa prima parte del capitolo abbiamo discusso tre argomenti principali: le prestazioni delle cache, utilizzando il grado di associatività per ridurre la frequenza delle miss, l'utilizzo di gerarchie delle memorie costituite da cache multilivello per ridurre la penalità di miss e le ottimizzazioni software per migliorare l'efficacia delle cache. Il sistema di memoria ha un impatto significativo sul tempo di esecuzione dei programmi. Il numero totale di cicli di stallo dovuti alla memoria dipende sia dalla frequenza sia dalla penalità di miss. La sfida, come vedremo nel paragrafo 5.8, è ridurre uno di questi due fattori critici della gerarchia delle memorie senza alterare significativamente l'altro.

Per ridurre la frequenza di miss abbiamo esaminato gli schemi associativi di corrispondenza tra memoria principale e cache. Tali schemi possono ridurre la frequenza delle miss di una cache permettendo una più flessibile collocazione dei blocchi all'interno della cache. Gli schemi completamente associativi permettono di scrivere i blocchi della memoria principale in un qualsiasi blocco della cache, ma richiedono anche di controllare tutti i blocchi della cache a ogni accesso; il costo elevato che ne risulta rende praticamente irrealizzabili cache completamente associative di grandi dimensioni. Le cache set-associative sono un'alternativa praticabile, dal momento che qui la ricerca deve essere effettuata solamente tra i blocchi di una stessa linea, selezionata attraverso il campo indice. Le cache set-associative hanno una frequenza di miss più elevata ma con tempi di accesso inferiori. Il grado di associatività che produce le migliori prestazioni dipende sia dalla tecnologia sia dal tipo di implementazione.

Infine, abbiamo visto come le cache multilivello costituiscano una tecnica che può essere utilizzata per ridurre la penalità di miss, permettendo a una cache secondaria di dimensioni maggiori di gestire le miss della cache primaria. Le cache di secondo livello sono divenute comuni da quando i progettisti si sono resi conto che il limitato spazio disponibile sul chip e la necessità di mantenere elevata la frequenza del clock impediscono alle cache primarie di essere troppo grandi. La cache secondaria, che ha dimensioni 10 o più volte superiori di una cache primaria, gestisce molti degli accessi falliti alla cache primaria. In tal caso, la penalità di miss è pari al tempo di accesso alla cache secondaria (tipicamente minore di 10 cicli di clock del processore), che è molto inferiore al tempo di accesso alla memoria principale (tipicamente maggiore di 100 cicli di clock del processore). Come per il grado di associatività, il compromesso tra dimensione e tempo di accesso nella progettazione di una cache secondaria dipende da un gran numero di fattori legati all'implementazione della gerarchia delle memorie.

Infine, data l'importanza delle gerarchie delle memorie per le prestazioni, abbiamo esaminato come modificare un algoritmo per migliorare il comportamento delle cache: l'elaborazione a blocchi è una tecnica molto importante quando si devono elaborare matrici di grandi dimensioni.

## 5.5 | Affidabilità delle gerarchie delle memorie

L'ipotesi alla base delle discussioni dei paragrafi precedenti è che le memorie non dimentichino. Avere una memoria veloce ma non affidabile non è per nulla attraente. Come abbiamo visto nel primo capitolo, una grande idea per ottenere sistemi affidabili è la ridondanza. In questo paragrafo esamineremo innanzitutto i termini associati all'affidabilità per definire delle misure associate ai malfunzionamenti; vedremo anche come la ridondanza possa rendere praticamente infinita la durata delle memorie.



### Definizione di malfunzionamento

Partiamo dall'ipotesi che sappiate quali siano le specifiche di un servizio adeguato. Gli utenti potranno valutare il sistema in base a due possibili condizioni di erogazione del servizio rispetto alle specifiche del servizio stesso:

1. *erogazione corretta del servizio*, quando il servizio viene erogato secondo le specifiche;
2. *interruzione del servizio*, quando il servizio fornito differisce dalle specifiche definite per quel servizio.

Le transizioni dalla condizione 1 alla condizione 2 sono causate da *malfunzionamenti*, mentre le transizioni dalla condizione 2 alla condizione 1 sono chiamate *ripristino*. I malfunzionamenti possono essere permanenti o intermittenti, e i secondi costituiscono il caso più difficile: è più difficile diagnosticare un problema quando un sistema oscilla tra i due stati. I malfunzionamenti permanenti sono invece più facili da diagnosticare. Queste definizioni conducono alla definizione di due termini collegati: affidabilità e disponibilità.

L'*affidabilità* è una misura della continuità con cui viene fornito il servizio a partire da un certo istante; ossia, in maniera equivalente, è il tempo che intercorre prima che si verifichi il primo malfunzionamento. Quindi, il *tempo medio di funzionamento prima di un malfunzionamento* (MTTF, Mean Time To Failure) è una misura di affidabilità. Una misura collegata è la *frequenza annua media dei malfunzionamenti* (AFR, Annual Failure Rate) che rappresenta la percentuale di dispositivi che, in un anno, presentano un malfunzionamento, dato un certo MTTF. Quando l'MTTF diventa grande può essere fuorviante e l'AFR rappresenta una misura migliore.

## MTTF e AFR dei dischi

### ESEMPIO

Alcuni dischi recenti dichiarano un MTTF di 1 000 000 di ore. Dato che 1 000 000 di ore corrisponde a  $1\,000\,000 / (365 \times 24) = 114$  anni, sembrerebbe che questi dischi non si guastino mai. I grandi centri di calcolo che forniscono i servizi Internet come la ricerca su web possono essere dotati anche di 50 000 server. Si supponga che ciascun server abbia 2 dischi. Utilizzare l'AFR per calcolare quanti dischi ci si può aspettare che presentino malfunzionamenti in un anno.

### SOLUZIONE

Un anno è costituito da  $365 \times 24 = 8760$  ore. Un MTTF di 1 000 000 di ore corrisponde a un AFR di  $8760 / 1\,000\,000 = 0,876\%$ . Avendo il centro di calcolo 100 000 dischi, possiamo attenderci che in un anno presentino malfunzionamenti 876 dischi, cioè in media più di due dischi al giorno!

L'interruzione di un servizio viene misurata attraverso il *tempo medio di riparazione* (MTTR, Mean Time To Repair), mentre il *tempo medio tra due guasti* (MTBF, Mean Time Between Failures) è semplicemente la somma di MTTF e MTTR. Sebbene l'MTBF sia la misura più diffusa, l'MTTF spesso è il termine più adeguato. La *disponibilità* è quindi una misura relativa della corretta erogazione di un servizio ed è data rispetto all'alternanza delle condizioni di erogazione corretta e interruzione del servizio. La disponibilità viene misurata, in modo statistico, come:

$$\text{Disponibilità} = \frac{\text{MTTF}}{(\text{MTTF} + \text{MTTR})}$$

Si noti che affidabilità e disponibilità sono quantità misurabili in pratica e non solo sinonimi di sicurezza. Diminuire l'MTTR può fare aumentare l'affidabilità, così come aumentare l'MTTF. Per esempio, strumenti per l'identificazione dei guasti, per la diagnosi e la riparazione possono aiutare a ridurre i tempi di riparazione e quindi migliorare la disponibilità.

Solitamente si vuole che l'affidabilità sia molto alta. Un modo per valutarla è attraverso il numero di "nove" nella percentuale di affidabilità per anno. Per esempio, un servizio Internet molto buono oggi offre da 4 a 5 "nove di affidabilità". Vediamo come questa indicazione si traduce in numeri. Considerando che ci sono 365 giorni in un anno e quindi  $365 \times 24 \times 60 = 526\,000$  minuti, otteniamo:

Un nove:	90%	$\rightarrow 36,5$ giorni di riparazione/anno
Due nove:	99%	$\rightarrow 3,65$ giorni di riparazione/anno
Tre nove:	99,9%	$\rightarrow 526$ minuti di riparazione/anno
Quattro nove:	99,99%	$\rightarrow 52,6$ minuti di riparazione/anno
Cinque nove:	99,999%	$\rightarrow 5,26$ minuti di riparazione/anno

ecc.

Per incrementare l'MTTF occorre migliorare la qualità dei componenti o del progetto del sistema, in modo da poter continuare a svolgere le operazioni anche quando un componente presenta dei malfunzionamenti; i malfunzionamenti

devono perciò essere definiti rispetto al loro contesto: il guasto di un componente potrebbe non implicare il malfunzionamento dell'intero sistema. Perché sia chiara questa distinzione, useremo in seguito il termine *guasto (fault)* per indicare un malfunzionamento permanente o transitorio di un componente. Ci sono tre modi per migliorare l'MTTF:

1. *Evitare il guasto*: costruire il dispositivo in modo da evitare il verificarsi del guasto stesso.
2. *Tolleranza ai guasti*: utilizzare la ridondanza per fare sì che il servizio soddisfi le specifiche anche quando si verifica un guasto.
3. *Previsione dei guasti*: prevedere la presenza di guasti o la possibilità che questi si generino, prima che si verifichino, in modo da sostituire il componente prima che si guasti.



PREDIZIONE

## Codice di Hamming per la correzione di errori singoli e identificazione di errori doppi (SEC/DED)

Richard Hamming ha inventato uno schema molto utilizzato per ridondare la memoria, per il quale ha ricevuto il Premio Turing nel 1968. Per capire i codici ridondanti, occorre definire una misura di "vicinanza" tra due insiemi di bit: definiamo *distanza di Hamming* il minimo numero di bit che devono essere commutati perché i due insiemi risultino uguali. Per esempio, la distanza tra 011011 e 001111 è due. Che cosa succede se la distanza tra due insiemi di bit è due e si ha un errore su un bit? In questo caso l'insieme sul quale si è verificato l'errore non sarà più un numero valido. Perciò, se possiamo riconoscere se un numero è valido o no, possiamo rilevare gli errori su un singolo bit e possiamo dire di avere un **codice di rilevamento degli errori** su un singolo bit.

Hamming ha utilizzato il *codice di parità* per la rilevazione degli errori. In un codice di parità, viene contato il numero di 1 in una parola; si dice che la parola ha parità dispari se il numero di 1 è dispari e che ha parità pari se il numero di 1 è pari. Quando una parola viene scritta in memoria, viene aggiunto il bit di parità (1 per un numero dispari di 1, 0 per un numero pari di 1), di modo che la parità della parola di  $N + 1$  bit sarà sempre pari. Quando la parola viene letta, viene letto e confrontato il bit di parità con il numero di 1 presente nei primi  $N$  bit. Se il numero di 1 è pari e il bit di parità indica un numero dispari di 1, o viceversa, significa che si è verificato un errore.

**Codice di rilevamento degli errori:** un codice che consente la rilevazione di un errore nei dati. Non consente però di identificare la posizione dell'errore all'interno di una parola e quindi di correggerlo.

Calcolare la parità associata a un byte che contiene il numero  $31_{dec}$  e mostrare la stringa di bit che viene memorizzata. Si supponga che il bit di parità sia sulla destra. Si supponga che il bit più significativo commuti mentre il byte risiede in memoria e che leggiate successivamente il byte. Siete in grado di rilevare l'errore? Che cosa succede se i primi due bit più significativi commutano?

ESEMPIO

$31_{dec}$  corrisponde a  $00011111_{due}$ , che contiene cinque 1. Per scrivere un numero pari di 1, dobbiamo scrivere un 1 nel bit di parità e quindi:  $00011111_{due}$ . Se il bit più significativo cambia segno, quando leggiamo il byte otterremo:  $10011111_{due}$ , che contiene sette 1. Dato che ci aspettavamo una parità pari mentre il numero ha una parità dispari, siamo in grado di rilevare l'errore. Se invece commutassero i primi due bit più significativi, leggeremmo  $11011111_{due}$ , che contiene otto 1, cioè un numero pari di 1, e non saremmo in grado di rilevare l'errore.

SOLUZIONE

Se si verifica un errore su 2 bit, quindi, uno schema di parità a 1 bit non sarebbe in grado di rilevare nessun errore, dato che il bit di parità codifica la stessa parità quanto il numero contenente due errori. (In realtà il codice di parità a 1 bit consente di rilevare un numero dispari di errori, ma la probabilità di avere tre errori è molto inferiore della probabilità di averne due, per cui, in pratica, il codice di parità su 1 bit è limitato al riconoscimento di errori sul singolo bit.)

Naturalmente il codice di parità non può correggere gli errori, mentre Hamming voleva arrivare a questo. Se utilizzassimo un codice di errore che abbia una distanza 3 tra due numeri validi, allora ogni errore sul singolo bit produrrebbe un numero più vicino al numero corretto che a ogni altro numero valido. Sulla base di queste osservazioni, Hamming propose una trasformazione dei dati in un codice a distanza 3 chiamato *codice di correzione degli errori* (ECC, *Error Correction Code*) di Hamming, in suo onore. Utilizzeremo i bit di parità aggiunti per consentire di identificare la posizione in cui si è verificato l'errore singolo. Questi sono i passi per calcolare la codifica ECC di Hamming.

1. Iniziare a contare i bit a partire da sinistra, invece di contarli da destra come si fa solitamente.
2. Riservare tutte le posizioni che sono potenze di 2 ai bit di parità (posizioni 1, 2, 4, 8, 16, ...).
3. Tutte le altre posizioni possono essere utilizzate per memorizzare i bit del dato (posizioni 3, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 17...)
4. I bit inseriti nelle posizioni dei bit di parità determinano la sequenza dei bit dei dati da controllare (la Figura 5.24 mostra questo meccanismo graficamente per un numero su otto bit):
  - Bit 1 ( $0001_{\text{due}}$ ) controlla i bit (1, 3, 5, 7, 9, 11, ...) che sono i bit il cui indirizzo contiene un 1 nel bit più a destra ( $0001_{\text{due}}$ ,  $0011_{\text{due}}$ ,  $0101_{\text{due}}$ ,  $0111_{\text{due}}$ ,  $1001_{\text{due}}$ ,  $1011_{\text{due}}$ , ...)
  - Bit 2 ( $0010_{\text{due}}$ ) controlla la parità dei bit (2, 3, 6, 7, 10, 11, 14, 15, ...), che sono i bit per i quali il secondo bit a destra dell'indirizzo è un 1.
  - Bit 4 ( $0100_{\text{due}}$ ) controlla la parità dei bit (4-7, 12-15, 20-23, ...), che sono i bit per i quali il terzo bit a destra dell'indirizzo è un 1.
  - Bit 8 ( $1000_{\text{due}}$ ) controlla la parità dei bit (8-15, 24-31, 40-47, ...), che sono i bit per i quali il quarto bit a destra dell'indirizzo è un 1.

Si noti che ciascun bit del dato è coperto da uno o più bit di parità.

5. Impostare il bit di parità per ottenere una parità pari in ciascun gruppo. Sembra una magia, ma ciò vi consente di determinare se i bit sono commutati controllando i bit di parità. Utilizzando il codice a 12 bit di Figura 5.24, se il valore dei quattro bit di parità ( $p_8, p_4, p_2, p_1$ ) è 0000 vuol dire che non si sono verificati errori. Se, invece, i quattro bit assumono, per esempio, il valore 1010, cioè  $10_{\text{dec}}$ , allora il codice ECC di Hamming ci dice che l'errore è sul bit 10 (d6). Possiamo correggere l'errore semplicemente commutando il valore del bit 10.

Posizione bit	1	2	3	4	5	6	7	8	9	10	11	12
Bit codificati	p1	p2	d1	p4	d2	d3	d4	p8	d5	d6	d7	d8
Copertura del bit di parità	p1	X		X		X		X		X		X
	p2		X	X			X	X			X	X
	p4				X	X	X	X				X
	p8								X	X	X	X

Figura 5.24 I bit di dati ( $d_i$ ) e i bit di parità ( $p_i$ ) con la loro copertura per il codice ECC di Hamming su dati di otto bit.

Si supponga che un dato di ampiezza un byte sia  $10011010_{\text{dec}}$ . Mostrare il codice ECC di Hamming per il dato. Invertire poi il bit  $10_{\text{dec}}$  e mostrare che il codice ECC è in grado di correggere l'errore.

**ESEMPIO**

Lasciando lo spazio per i bit di parità, il codice di parità su 12 bit è:

   1    0 0 1    1 0 1 0

La posizione 1 verifica la parità dei bit 1, 3, 5, 7, 9 e 11, che evidenziamo:    1    0 0 1    1 0 1 0. Per rendere pari il numero di 1 in questo gruppo, occorre che il bit di parità assuma valore 0.

**SOLUZIONE**

La posizione 2 verifica la parità dei bit 2, 3, 6, 7, 10 e 11, che evidenziamo: 0    1    0 0 1    1 0 1 0. Per rendere pari il numero di 1 in questo gruppo, occorre che il bit di parità assuma valore 1.

La posizione 4 verifica la parità dei bit 4, 5, 6, 7 e 12: 011   0 0 1    1 0 1. Per rendere pari il numero di 1 in questo gruppo, occorre che il bit di parità assuma valore 1.

La posizione 8 verifica la parità dei bit 8, 9, 10, 11 e 12: 0111001   1 0 1 0. Per rendere pari il numero di 1 in questo gruppo, occorre che il bit di parità assuma valore 0.

Il codice finale è 011100101010. Invertiamo ora il bit 10; otteniamo: 011100101110.

Il bit di parità 1 è 0 (011100101110 contiene quattro 1 nel gruppo associato, parità pari, e quindi i bit di questo gruppo non contengono errori).

Il bit di parità 2 è 1 (011100101110 contiene cinque 1 nel gruppo associato, parità dispari, e quindi si è verificato un errore in uno di questi bit).

Il bit di parità 4 è 1 (011100101110 contiene due 1 nel gruppo associato, parità pari, e quindi i bit di questo gruppo non contengono errori).

Il bit di parità 8 è 1 (011100101110 contiene tre 1 nel gruppo associato, parità dispari, e quindi si è verificato un errore in uno di questi bit).

I bit di parità 2 e 8 segnalano un errore, che si trova in posizione 10, cioè  $8 + 2$ . Possiamo quindi invertire il contenuto del bit 10 e correggere l'errore: 011100101010.

Hamming non si fermò al codice di errore per la correzione dell'errore su un bit. Con l'aggiunta di un altro bit possiamo aumentare a un minimo di quattro la distanza tra due numeri validi e possiamo quindi *rilevare l'errore su due bit* e correggere l'errore su un bit. L'idea è quella di aggiungere un bit di parità che viene calcolato sull'intera parola. Utilizziamo un dato su quattro bit come esempio. Il codice richiede in questo caso 7 bit per il rilevamento di un errore. I bit di parità dell'ECC, H ( $p_1, p_2$  e  $p_3$ ), vengono calcolati (supponendo una parità pari) e viene anche calcolato il bit di parità globale,  $p_4$ :

1	2	3	4	5	6	7	8
p <sub>1</sub>	p <sub>2</sub>	d <sub>1</sub>	p <sub>3</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	<u>p<sub>4</sub></u>

L'algoritmo per correggere un errore e per identificarne due richiede solamente di calcolare la parità sui gruppi ECC (H) come visto in precedenza e sull'intero gruppo, p4. Si distinguono quattro casi:

1. H è pari e p4 è pari: non si è verificato nessun errore;
2. H è dispari e p4 è dispari: si è verificato un errore singolo che può essere corretto (p4 calcola la parità dispari quando si è verificato 1 errore);
3. H è pari e p4 è dispari: si è verificato un errore in p4 ma non negli altri bit, quindi occorre commutare p4;
4. H è dispari e p4 è pari: si sono verificati due errori (p4 calcola la parità pari quando si sono verificati 2 errori).

Oggi, la correzione di errori singoli e il rilevamento di due errori (SEC/DED, *Single Error Correction/Double Error Detection*) è comune nelle memorie dei server: il SEC/DED per blocchi di dati ampi otto byte richiede un solo byte in più e questo è il motivo per cui molte DIMM hanno un'ampiezza di 72 bit.

**Approfondimento.** Per calcolare quanti bit occorrono per il SEC, supponiamo che  $p$  sia il numero totale dei bit di parità,  $d$  il numero di bit dei dati e  $p + d$  l'ampiezza delle parole di memoria. Se  $p$  bit servono per identificare e correggere gli errori su un bit ( $p + d$  combinazioni) e un bit serve per indicare che non si sono verificati errori, occorrono:

$$2^p \geq p + d + 1 \text{ bit e quindi } p \geq \log(p + d + 1)$$

Per esempio, per dati su 8 bit,  $d = 8$  e  $2^p \geq p + 8 + 1$ , e quindi  $p = 4$ . Analogamente,  $p = 5$  per dati su 16 bit, 6 per dati su 32 bit, 7 per dati su 64 bit ecc.

**Approfondimento.** Nei sistemi di grandissime dimensioni, la probabilità di errori multipli e del malfunzionamento di tutto un chip di memoria diventa significativa. L'IBM ha introdotto *chipkill* per ovviare a questo problema, e molti altri sistemi di grosse dimensioni hanno adottato questo sistema. (Intel ha chiamato la sua versione SDDC.) È un approccio simile ai dischi RAID esaminati nel paragrafo 5.11. Chipkill distribuisce i dati e le informazioni sull'ECC, in modo tale che il malfunzionamento di un intero chip di memoria possa essere gestito ricostruendo i dati mancanti a partire dai dati contenuti negli altri chip di memoria. Supponendo di avere un cluster di 10 000 processori con 4 GiB di memoria per processore, IBM ha calcolato la seguente frequenza di errori di memoria *irrecuperabili* nell'arco di tre anni:

- solo parità: circa 90 000, ovvero un errore irrecuperabile (o non rilevato) ogni 17 minuti;
- SEC/DED solamente: circa 3500, ovvero circa un errore irrecuperabile (o non rilevato) ogni 7,5 ore;
- chipkill: 6, ovvero un errore irrecuperabile (o non rilevato) ogni 2 mesi.

Quindi la funzionalità chipkill è diventata un requisito fondamentale per i calcolatori per i grandi centri di calcolo.

**Approfondimento.** Se gli errori su singolo bit o su due bit sono tipici per i sistemi di memoria, le reti possono avere errori a raffica. Una soluzione è il controllo periodico: *Cyclic Redundancy Check* (controllo periodico della ridondanza). Per un blocco di  $k$  bit, un trasmettitore genera una sequenza di  $n-k$  bit di controllo: trasmette un numero su  $n$  bit che risulta esattamente divisibile per un certo numero. Il ricevitore divide il numero ricevuto per questo numero e controlla il resto: se è uguale a zero non si sono verificati errori. Se il resto è diverso da zero, il ricevitore rigetta il messaggio e chiede la ritrasmissione. Come abbiamo visto nel Capitolo 3, è facile calcolare la divisione per una potenza di due mediante un registro a scorrimento, cosa che ha permesso la diffusione del codice CRC anche quando l'hardware era più prezioso. Facendo un ulteriore passo avanti, i codici di Reed-Solomon utilizzano i campi di Galois per correggere errori di trasmissione su più bit; in questo caso i dati sono considerati i coefficienti di un polinomio e lo spazio dei codici diventa lo spazio dei valori del polinomio. Il calcolo del codice di Reed-Solomon è notevolmente più complicato della divisione binaria!

## 5.6 | Macchine virtuali

Le *macchine virtuali* (VM, *Virtual Machines*) sono state sviluppate a partire dagli anni '60, e negli anni successivi hanno sempre costituito una parte importante dei processi di calcolo nei calcolatori mainframe. Sebbene siano state ignorate negli anni '80 e '90 nei PC a uso singolo, recentemente stanno guadagnando popolarità per i seguenti motivi:

- l'aumento dell'importanza di isolamento e sicurezza nei sistemi moderni;
- i problemi di sicurezza e affidabilità dei sistemi operativi standard;
- la condivisione di uno stesso calcolatore tra diversi utenti non collegati tra loro, in particolare per il cloud computing;
- l'aumento esponenziale della velocità dei processori, che rende il sovraccarico di lavoro delle VM più accettabile.

Nella definizione più ampia di VM sono compresi praticamente tutti i metodi di emulazione che forniscono un'interfaccia software standard, come la VM di Java. In questo paragrafo esamineremo le VM che forniscono un ambiente di sistema completo, in grado di gestire, a livello binario, una certa implementazione architettonica di un insieme di istruzioni (ISA). Sebbene alcune VM possano emulare diverse ISA sul loro hardware nativo, supporremo che una VM corrisponda sempre a un particolare processore. Queste VM sono chiamate *macchine virtuali di sistema* (operativo); alcuni esempi sono il VM/370 di IBM, VMware, il Server ESX e Xen.

Le macchine virtuali di sistema forniscono agli utenti l'illusione di avere a loro disposizione l'intero calcolatore, compresa una copia del sistema operativo. Un singolo calcolatore può eseguire più VM e può supportare diversi *sistemi operativi* (SO). Su una piattaforma convenzionale, un singolo SO "possiede" tutte le risorse hardware, ma con una VM più SO possono condividere le stesse risorse.

Il software che supporta una VM è chiamato *monitor della macchina virtuale* (VMM, *Virtual Machine Monitor*), o *hypervisor*; il VMM è il cuore della tecnologia delle macchine virtuali. La piattaforma su cui gira un VMM è chiamata *host* (letteralmente, "ospitante"), e le sue risorse vengono condivise dalle diverse VM *guest* (gli "ospiti"). Il VMM determina come mappare le risorse virtuali nelle risorse fisiche: una risorsa fisica può essere condivisa in istanti diversi di tempo (*in time sharing*), partizionata o addirittura emulata dal software. Il VMM è molto più piccolo di un sistema operativo tradizionale: il software di un VMM può contenere circa 10 000 linee di codice. Sebbene il nostro interesse in questo capitolo sia focalizzato sull'utilizzo delle VM per aumentare la protezione del calcolatore, una VM fornisce altri due vantaggi significativi dal punto di vista commerciale: vediamoli.

1. *Gestione del software*. Le VM forniscono un'astrazione che può eseguire programmi appartenenti a tutti gli strati del software e può eseguire perfino vecchi sistemi operativi, come il DOS. In una installazione tipica, alcune VM possono emulare una vecchia versione di alcuni SO, altre VM possono emulare la versione stabile attuale del SO e alcune VM possono persino emulare versioni beta del SO.
2. *Gestione dell'hardware*. Uno dei motivi per cui conviene utilizzare server multipli è la possibilità di eseguire ciascuna applicazione con la versione appropriata del sistema operativo su calcolatori diversi, dato che la separazione aumenta l'affidabilità. Le VM consentono di eseguire indipendentemente le diverse applicazioni con la loro versione del SO utilizzando lo stesso hardware, riducendo così il numero di server richiesti. Un'altra applicazione di un VMM è il supporto alla migrazione di una VM in esecuzione su un altro calcolatore, per bilanciare il carico di lavoro o per abbandonare un server che si sta per guastare.

## Interfaccia hardware/software



LEGGE DI MOORE

Amazon fornisce i suoi servizi Web EC2 (AWS, *Amazon Web Services*) utilizzando le macchine virtuali della sua “nuvola” per cinque motivi.

1. Consente al sistema AWS di proteggere gli utenti l’uno dall’altro mentre condividono lo stesso server.
2. Semplifica la distribuzione del software all’interno dei sistemi della dimensione di un grosso centro di calcolo. Il cliente installa l’immagine di una macchina virtuale configurata mediante un software adeguato e AWS lo distribuisce a tutti i calcolatori che un cliente vuole utilizzare.
3. I clienti (e AWS) possono affidabilmente terminare una VM per controllare l’utilizzo delle risorse quando i clienti hanno completato il loro lavoro.
4. Le macchine virtuali nascondono l’identità dell’hardware sul quale il cliente sta eseguendo i propri programmi; questo vuol dire che AWS può continuare a utilizzare i vecchi server e anche introdurre nuovi server più performanti. Il cliente, per esempio, si può aspettare che le prestazioni siano equivalenti alle specifiche riportate per le Unità di Calcolo EC2, che, secondo AWS, “forniscono prestazioni equivalenti a una CPU AMD Opteron del 2007 da 1,0-1,2 GHz o di uno Xeon Intel del 2007”. Grazie alla Legge di Moore, i nuovi server possono chiaramente offrire un numero maggiore di Unità di Calcolo EC2 dei server più vecchi, ma AWS può continuare a noleggiare i server vecchi fino a quando rimangono più economici.
5. I monitor delle macchine virtuali possono controllare la frequenza con cui una VM utilizza un processore, la rete e lo spazio su disco, cosa che consente all’AWS di offrire prezzi diversi per esecuzioni di tipo diverso che pure vengono eseguite dallo stesso server fisico. Per esempio, nel 2012, AWS offriva 14 tipi di esecuzioni, da piccoli programmi standard che costavano \$0,08 l’ora, a programmi molto grandi con I/O intensivo che potevano arrivare a \$3,10 l’ora.

In generale, il costo della virtualizzazione di un processore dipende dal carico di lavoro. I programmi utente che utilizzano principalmente il processore non hanno costi aggiuntivi di virtualizzazione; dato che il SO viene raramente richiamato, tutto continua a essere eseguito alla velocità originale. Invece, i programmi che prevedono un flusso di I/O rilevante, in genere, usano in maniera intensiva anche il sistema operativo, dato che devono eseguire molte chiamate di sistema e istruzioni privilegiate; queste possono produrre un elevato costo aggiuntivo di virtualizzazione. D’altra parte, se il carico di lavoro dell’I/O contiene principalmente *operazioni di I/O*, il costo della virtualizzazione del processore può essere completamente nascosto, poiché il processore si trova spesso nella condizione di attesa dei segnali di I/O. Il costo aggiuntivo dipende sia dal numero di istruzioni che devono essere emulate dal VMM sia dal tempo richiesto per emularle. Perciò, quando la VM ospite esegue la stessa ISA dell’host, come supponiamo qui, l’obiettivo dell’architettura e del VMM è di eseguire praticamente tutte le istruzioni direttamente sull’hardware nativo.

### Requisiti del monitor di una macchina virtuale

Che cosa deve fare il monitor di una VM? Deve fornire un’interfaccia software al software ospitato, deve isolare uno dall’altro lo stato dei diversi software ospitati e deve proteggere se stesso dal software ospitato (compresi i SO ospitati). I requisiti qualitativi sono:

- il software ospitato dovrebbe comportarsi sulla VM esattamente come se venisse eseguito sull’hardware nativo, fatta eccezione per le caratteristiche legate alle prestazioni o alle limitazioni delle risorse prefissate che sono condivise tra le diverse VM;
- il software ospitato non dovrebbe essere in grado di modificare direttamente l’allocazione delle risorse reali del sistema.

Per "virtualizzare" un processore, il VMM deve poter controllare praticamente tutto: accessi allo stato privilegiato, traduzione degli indirizzi, I/O, eccezioni e interrupt, anche quando la VM ospitata e il SO in esecuzione li stanno utilizzando temporaneamente. Per esempio, se si verifica un interrupt del timer, il VMM sosponderà la VM ospitata correntemente in esecuzione, salverà lo stato di questa VM, gestirà l'interrupt, determinerà quale VM tra quelle ospitate deve riprendere l'esecuzione e infine caricherà lo stato della nuova VM. Le VM ospitate, il cui funzionamento è legato agli interrupt del timer, saranno dotate di un timer virtuale e riceveranno dal VMM un interrupt che emula l'interrupt del timer.

Per poter svolgere questa attività, il VMM deve avere un livello di privilegio superiore a quello della VM ospitata, che generalmente viene eseguita in modalità utente; questo assicura anche che l'esecuzione di tutte le istruzioni privilegiate sia gestita dal VMM. I requisiti di base delle macchine virtuali di sistema sono:

- almeno due modalità di esecuzione del processore: sistema e utente;
- un insieme privilegiato di istruzioni, disponibili solamente nella modalità di sistema, che provocano un'eccezione di trap se vengono eseguite in modalità utente; tutte le risorse di sistema devono essere controllabili solamente attraverso queste istruzioni.

### Mancanza di supporto alle macchine virtuali da parte dell'architettura dell'insieme di istruzioni

Se le VM fossero già previste durante la progettazione dell'ISA, sarebbe relativamente semplice sia ridurre il numero di istruzioni che devono essere eseguite da un VMM, sia aumentare la velocità di emulazione. Un'architettura che consente di eseguire una VM direttamente in hardware viene detta *virtualizzabile*; per esempio, le architetture IBM 370 e RISC-V sono di questo tipo.

Tuttavia, dato che le VM sono state introdotte solamente di recente nelle applicazioni per i PC e per i server, la maggior parte degli insiemi di istruzioni, compreso l'x86 e molte architetture RISC (tra le quali gli ARMv7 e i MIPS), è stata creata senza pensare alla virtualizzazione.

Dato che il VMM deve assicurare che il sistema ospitato interagisca solamente con le risorse virtuali, un tipico SO ospitato deve essere eseguito come un programma in modalità utente sopra il VMM. Quindi, se un sistema operativo ospitato cerca di accedere o di modificare informazioni legate alle risorse hardware utilizzando le istruzioni privilegiate, per esempio per leggere o scrivere il bit che abilita gli interrupt, verrà generata un'eccezione di trap che verrà gestita dal VMM. Sarà quindi il VMM a eseguire le modifiche richieste alle risorse fisiche interessate.

Se tutte le istruzioni che cercano di leggere o scrivere queste informazioni delicate provocano un'eccezione di trap quando vengono eseguite in modalità utente, il VMM può intercettare e può supportare una versione virtuale di queste informazioni delicate, gestendole come le gestirebbe il SO ospitato.

Senza questo supporto, occorrerebbe prendere altri provvedimenti: il VMM dovrebbe adottare precauzioni particolari per individuare tutte le istruzioni problematiche e assicurarsi che si comportino correttamente quando vengono eseguite da un SO ospitato, aumentando quindi la complessità del VMM e riducendo le prestazioni della VM in esecuzione.

### Protezione e architettura dell'insieme delle istruzioni

La protezione deriva da uno sforzo congiunto dell'architettura e del sistema operativo, ma i progettisti, quando si sono diffuse le memorie virtuali, hanno dovuto modificare alcuni dettagli problematici delle architetture degli insiemi di istruzioni esistenti.

Per esempio, l'istruzione POPF dell'x86 carica i registri di flag dalla cima dello stack in memoria. Uno dei flag è il flag di abilitazione degli interrupt (IE, *Interrupt Enable*). Quando l'istruzione di POPF viene eseguita in modalità utente, invece di generare un'eccezione di trap modifica semplicemente il valore di tutti i flag ad eccezione di IE, mentre se viene eseguita in modalità di sistema modifica anche IE. I sistemi operativi ospitati in una VM vengono eseguiti in modalità utente e questo rappresenta un problema, poiché ci aspetteremmo di vedere modificato anche il flag IE.

Storicamente, si possono identificare tre tappe nello sviluppo dell'hardware dei mainframe IBM e dei VMM per migliorare le prestazioni delle macchine virtuali:

1. riduzione del costo della virtualizzazione del processore;
2. riduzione dei costi addizionali degli interrupt dovuti alla virtualizzazione;
3. riduzione dei costi degli interrupt pilotando gli interrupt verso la VM appropriata invece che richiedere l'intervento del VMM.

Nel 2006 nuove proposte di AMD e Intel hanno cercato di affrontare il primo punto, riducendo il costo della virtualizzazione del processore. Sarà interessante vedere quante generazioni di architetture e quante modifiche al VMM vedranno la luce prima che questi tre aspetti siano affrontati insieme, e quanto tempo dovrà passare prima che le macchine virtuali del XXI secolo diventino efficienti come i mainframe IBM e i VMM degli anni '70.

**Approfondimento.** I RISC-V intercettano tutte le istruzioni privilegiate quando funzionano in modalità utente, per cui supportano la *virtualizzazione classica*, nella quale il SO ospitato funziona in modalità utente mentre il VMM funziona in modalità supervisore.

## 5.7 Memoria virtuale

*È stato progettato un sistema per fare in modo che la memoria centrale appaia al programmatore come una memoria a un solo livello, dato che la gestione dei trasferimenti richiesti viene effettuata automaticamente.*

Kilburn et al., *One-level storage system*, 1962

**Memoria virtuale:** una tecnica che utilizza la memoria principale come se fosse una "cache" per la memoria di massa.

Nei paragrafi precedenti abbiamo visto come le cache forniscano un accesso veloce alle porzioni dei dati e del codice di un programma utilizzate più di recente. Allo stesso modo, la memoria principale può agire come una cache della memoria di massa, che di solito è costituita da dischi magnetici; questa tecnica è chiamata **memoria virtuale**. Storicamente esistono due motivazioni principali che hanno portato alla realizzazione della memoria virtuale: consentire una condivisione efficiente e sicura della memoria da parte di più programmi e liberare chi scrive i programmi dall'onere di dover gestire una memoria principale piccola e limitata. Quattro decenni dopo la sua invenzione, la prima motivazione è oggi predominante.

Naturalmente, affinché più macchine virtuali possano condividere la stessa memoria, dobbiamo essere in grado di proteggere le macchine virtuali una dall'altra, assicurando che un programma possa leggere e scrivere solamente le porzioni di memoria principale che gli sono state assegnate. La memoria principale conterrà solamente le parti attive delle diverse macchine virtuali, proprio come la cache contiene soltanto le parti attive di un singolo programma. Quindi, il principio di località giustifica la presenza sia della memoria virtuale sia della cache, e la memoria virtuale ci permette una condivisione efficiente del processore e della memoria principale.

Non possiamo sapere quali macchine virtuali condivideranno la memoria con altre macchine virtuali, quando le compiliamo. Infatti, le macchine virtuali che condividono la memoria principale cambiano dinamicamente durante l'esecuzione. A causa di questa interazione dinamica, sarebbe opportuno compilare ciascun programma nel suo *spazio di indirizzamento privato*: un insieme

di locazioni di memoria riservate, accessibili soltanto a quel programma. La memoria virtuale implementa la traduzione degli indirizzi appartenenti allo spazio di indirizzamento privato di un programma in **indirizzi fisici**. Questo processo di traduzione garantisce la **protezione** dello spazio di indirizzamento di un programma dalle altre macchine virtuali.

Il secondo motivo che giustifica l'uso della memoria virtuale è consentire ai programmi utente di superare la dimensione fisica della memoria principale. Precedentemente, se un programma diventava troppo grande da poter essere contenuto nella memoria fisica, era compito del programmatore modificarlo in modo che ci potesse stare. I programmatore suddividevano il programma in porzioni e poi identificavano le parti tra loro indipendenti. Questi moduli, detti *overlay*, ossia «sovraposti» in memoria, venivano caricati o scaricati dalla memoria principale sotto il controllo del programma durante la sua esecuzione; era compito del programmatore garantire che il programma non cercasse mai di accedere a un modulo che non fosse stato già caricato in memoria e che i moduli caricati non superassero mai la dimensione complessiva della memoria. I moduli di overlay erano tradizionalmente organizzati in modo da contenere sia il codice sia i dati, e la chiamata di una procedura appartenente a un modulo diverso da quelli presenti in memoria richiedeva di caricare quel modulo in memoria a spese di uno dei moduli già presenti.

Come si può ben immaginare, questa operazione comportava un ingente carico di lavoro per il programmatore; la memoria virtuale fu inventata per sollevare i programmatori da questo carico, gestendo automaticamente i due livelli della gerarchia costituiti dalla memoria principale, a volte chiamata *memoria fisica* (per distinguerla dalla memoria virtuale), e dalla memoria di massa.

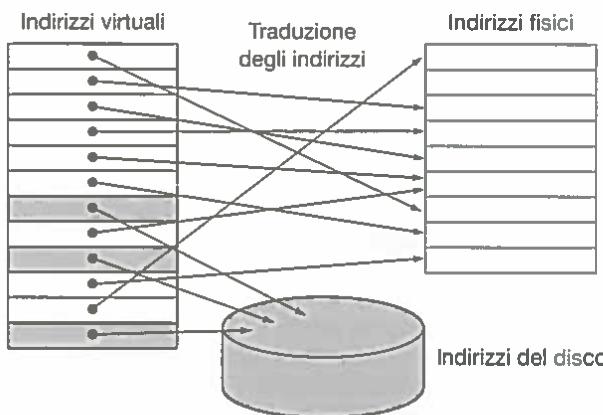
Sebbene i concetti che stanno alla base della memoria virtuale e delle cache siano gli stessi, le diverse origini storiche hanno portato all'utilizzo di una differente terminologia. Un blocco della memoria virtuale è chiamato *pagina* e le miss di una memoria virtuale vengono chiamate **page fault** (letteralmente, "mancanza di pagina"). Quando si utilizza la memoria virtuale, il processore genera un **indirizzo virtuale**; questo viene tradotto da una combinazione di componenti hardware e software in un *indirizzo fisico*, che può essere utilizzato per accedere alla memoria principale. La Figura 5.25 mostra la memoria indi-

**Indirizzo fisico:** l'indirizzo di una cella della memoria principale.

**Protezione:** un insieme di meccanismi che garantiscono che più processi che condividono lo stesso processore, la stessa memoria o gli stessi dispositivi di I/O non interferiscano tra loro, intenzionalmente o involontariamente, evitando così che un programma possa leggere o scrivere i dati di un altro programma. Questi meccanismi consentono anche di isolare il sistema operativo dai processi utente.

**Page fault:** un evento che si verifica quando la pagina a cui si accede non è presente nella memoria principale.

**Indirizzo virtuale:** un indirizzo che corrisponde a una locazione nello spazio virtuale convertito in un indirizzo fisico mediante il processo di mappatura su un indirizzo fisico quando si accede alla memoria.



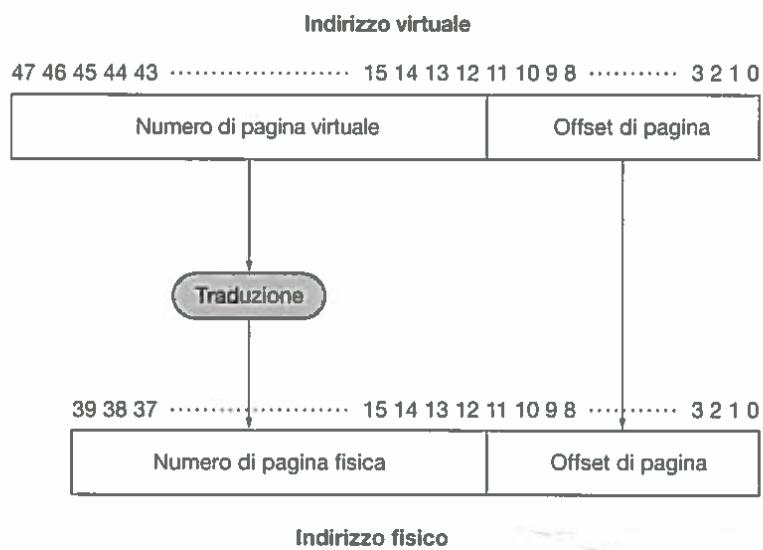
**Figura 5.25** Nella memoria virtuale i blocchi della memoria (chiamati *pagine*) vengono mappati da un insieme di indirizzi (detti *virtuali*) a un altro insieme di indirizzi (detti *fisici*). Il processore produce indirizzi virtuali mentre l'accesso alla memoria avviene utilizzando gli indirizzi fisici. Sia la memoria virtuale sia quella fisica sono suddivise in pagine, in modo tale che una pagina virtuale viene mappata su una pagina fisica. Naturalmente è possibile che una pagina virtuale non sia presente nella memoria principale e che quindi non possa essere mappata su una pagina fisica; in tal caso si troverà sul disco. Le pagine fisiche possono essere condivise facendo in modo che due indirizzi virtuali puntino allo stesso indirizzo fisico. Questa possibilità permette a due diversi programmi di condividere dati e parti di codice.

**Traduzione dell'indirizzo:** detta anche **mappatura dell'indirizzo**, è il processo tramite il quale un indirizzo virtuale viene mappato sull'indirizzo utilizzato per accedere alla memoria.

rizzata in modo virtuale e la mappatura delle pagine virtuali sulle pagine della memoria principale. Questo processo prende il nome di **traduzione dell'indirizzo** o **mappatura dell'indirizzo**. Attualmente, i due livelli della gerarchia delle memorie che vengono controllati utilizzando la memoria virtuale sono di solito le DRAM e le memorie a stato solido per i dispositivi mobili e le DRAM e i dischi magnetici nei server (par. 5.2). Considerando di nuovo l'analogia con la biblioteca, possiamo pensare che l'indirizzo virtuale sia il titolo di un libro e che l'indirizzo fisico sia la posizione di quel libro sullo scaffale.

La memoria virtuale semplifica anche la fase di caricamento di un programma prima della sua esecuzione, fornendo un meccanismo di *rilocazione*. La rilocazione mappa gli indirizzi virtuali utilizzati da un programma nei corrispondenti indirizzi fisici prima ancora che gli indirizzi vengano effettivamente utilizzati per accedere alla memoria. Questa mappatura consente il caricamento dei programmi in una qualsiasi locazione della memoria principale. Inoltre, tutti i meccanismi di gestione della memoria virtuale attualmente in uso rilocano il programma su un insieme di blocchi di dimensione prefissata, detti pagine, eliminando così la necessità di individuare un unico grande blocco di memoria contigua sufficientemente ampio da potervi caricare il programma; il sistema operativo deve solamente trovare un numero sufficiente di pagine della memoria fisica disponibili.

In una memoria virtuale, l'indirizzo viene suddiviso in due parti: *numero della pagina virtuale* e *offset nella pagina*. La Figura 5.26 mostra la traduzione del numero della pagina virtuale nel corrispondente *numero della pagina fisica*. Anche se il RISC-V ha indirizzi su 64 bit, i 16 bit più significativi non vengono utilizzati, per cui l'indirizzo da mappare è su 48 bit. La figura suppone che la memoria fisica sia di 1 TiB o  $2^{40}$  byte, che richiedono un indirizzo su 40 bit. I bit più significativi del numero della pagina fisica definiscono i bit più significativi dell'indirizzo fisico, mentre l'offset della pagina, che non viene modificato, definisce i bit meno significativi dell'indirizzo fisico. Il numero di bit riservati al campo offset determina la dimensione della pagina. Non è necessario che il numero di pagine indirizzabili attraverso l'indirizzo virtuale corrisponda al numero di pagine indirizzabili attraverso l'indirizzo fisico: avere un numero di



**Figura 5.26** Mappatura di un indirizzo virtuale in un indirizzo fisico. La dimensione della pagina è  $2^{12} = 4$  KiB. Il numero di pagine fisiche che possono risiedere in memoria è  $2^{28}$ , dato che il numero di pagina fisica è codificato su 28 bit. Questo significa che la memoria principale può contenere al massimo 1 TiB, mentre lo spazio di indirizzamento virtuale è di 256 TiB. Il RISC-V consente di avere una memoria fisica grande fino a 1 TiB, che è sufficientemente grande per la maggior parte dei calcolatori nel 2016.

pagine virtuali maggiore di quello delle pagine fisiche sta alla base dell'illusione di possedere una memoria virtuale praticamente illimitata.

Nei sistemi di memoria virtuale, molte scelte di progetto sono motivate dall'elevato costo di una miss, che nella memoria virtuale viene tradizionalmente chiamata *page fault*. Un page fault richiede milioni di cicli di clock per essere risolto: la tabella riportata all'inizio del paragrafo 5.2 mostra che la memoria principale è circa 100 000 volte più veloce di un disco rigido. Questa enorme penalità di miss, che per dimensioni normali della pagina è costituita principalmente dal tempo necessario a caricare la prima parola del blocco, condiziona diverse scelte chiave nella progettazione dei sistemi di memoria virtuale.

- Le pagine devono essere sufficientemente grandi da ammortizzare il lungo tempo di accesso. Attualmente le dimensioni tipiche variano tra 4 KiB e 64 KiB. I nuovi sistemi desktop e server vengono progettati per supportare pagine da 32 KiB a 64 KiB, ma i nuovi sistemi embedded stanno andando nella direzione opposta, cioè verso pagine da 1 KiB.
- Le strutture che riducono la frequenza dei page fault sono molto attraenti. La tecnica principale che viene utilizzata consiste nel permettere un posizionamento flessibile delle pagine nella memoria principale mediante tecniche completamente associative.
- I page fault possono essere gestiti via software, perché l'incremento del tempo di elaborazione è piccolo se confrontato con il tempo di accesso al disco. Inoltre, il software può consentire l'utilizzo di algoritmi più intelligenti nella scelta del posizionamento delle pagine, poiché anche una piccola diminuzione della frequenza delle miss ripaga il tempo speso per eseguire questi algoritmi.
- La tecnica di scrittura di tipo write-through non funziona per la memoria virtuale, perché richiederebbe troppo tempo; i sistemi di memoria virtuale utilizzano quindi tecniche di scrittura in write-back.

Esamineremo questi elementi di progetto di una memoria virtuale nei prossimi paragrafi.

**Approfondimento.** Abbiamo introdotto la necessità di avere una memoria virtuale attraverso diverse macchine virtuali che devono condividere la stessa memoria, ma la memoria virtuale è stata inventata originariamente affinché programmi diversi potessero condividere un calcolatore, utilizzandolo per porzioni di tempo, in time sharing. Dato che molti lettori oggi non hanno esperienza dei sistemi in time sharing, abbiamo utilizzato le macchine virtuali per introdurre questo paragrafo.

**Approfondimento.** I RISC-V supportano una varietà di configurazioni di memoria virtuale. Oltre allo schema di indirizzamento virtuale su 48 bit, che è particolarmente adatto ai server di grandi dimensioni nel 2016, quest'architettura può supportare spazi di indirizzamento virtuale su 39 bit o 57 bit. Tutte queste configurazioni utilizzano una dimensione della pagina di 4 kibibyte.

**Approfondimento.** I processori con 32 bit di indirizzamento sono problematici per i server e anche per i PC. Sebbene siamo propensi a pensare che lo spazio degli indirizzi virtuali sia molto più grande dello spazio degli indirizzi fisici, si può verificare la situazione opposta quando l'ampiezza dell'indirizzo, imposta dal processore, è piccola rispetto alla tecnologia corrente della memoria. In questo caso, il singolo programma non può trarre beneficio da una maggiore memoria fisica, ma un insieme di programmi, eseguiti contemporaneamente in parallelo, può beneficiare del fatto che i programmi non debbano essere caricati e scaricati dalla memoria, o che possano essere eseguiti su processori paralleli.

**Approfondimento.** In questo libro la discussione relativa alla memoria virtuale è focalizzata sulla tecnica della paginazione, nella quale vengono utilizzati blocchi di dimensione prefissata. Esiste anche uno schema con blocchi di dimensione va-

**Segmentazione:** uno schema di mappatura a dimensione variabile in cui un indirizzo viene suddiviso in due parti: un numero di segmento, che viene mappato in un indirizzo fisico, e un offset interno al segmento.

riabile, detto **segmentazione**. Nella segmentazione un indirizzo è formato da due parti: un numero di segmento e un offset interno al segmento. Per determinare l'indirizzo fisico effettivo, il contenuto del registro di segmento viene mappato in un indirizzo fisico a cui viene poi sommato l'offset. Dato che i segmenti possono avere dimensioni diverse, è necessario verificare che l'offset non superi i limiti del segmento. L'utilizzo principale della segmentazione è il supporto a tecniche più potenti per la protezione e la condivisione di uno spazio di indirizzamento comune. Molti libri sui sistemi operativi contengono una descrizione estesa della segmentazione (confrontata con la paginazione) e del suo utilizzo per condividere lo spazio di indirizzamento a livello logico. Lo svantaggio principale della segmentazione consiste nel dover suddividere lo spazio di indirizzamento in due parti logicamente separate, che vengono rappresentate e gestite utilizzando due campi diversi dell'indirizzo: il numero di segmento e l'offset. La paginazione invece rende invisibili al programmatore e al compilatore i confini tra il numero di pagina e l'offset della pagina.

I segmenti sono stati utilizzati anche come metodo per estendere lo spazio di indirizzamento, senza dover modificare la dimensione della parola del calcolatore. Questi tentativi non hanno avuto successo a causa della difficoltà di implementazione e della diminuzione delle prestazioni, dovute alla gestione di un indirizzo formato da due parti, di cui i programmatori e i compilatori devono necessariamente tenere conto.

Molte architetture dividono lo spazio di indirizzamento in grossi blocchi di dimensioni prefissate, in modo da rendere più semplice la protezione tra sistema operativo e programmi utente e aumentare l'efficienza nell'implementazione della paginazione. Sebbene queste partizioni vengano spesso chiamate «segmenti», questo meccanismo è molto più semplice della segmentazione in blocchi di dimensioni variabili e non è visibile ai programmi utente; esamineremo in dettaglio questa soluzione più avanti.

## Come individuare la posizione di una pagina e come ritrovarla

A causa dell'altissima penalità dei page fault, i progettisti cercano di ridurne la frequenza ottimizzando la posizione delle pagine. Se si consente di mappare una pagina virtuale su una qualsiasi pagina fisica, il sistema operativo potrà scegliere di sostituire una qualsiasi pagina fisica quando si verifica un errore di page fault. Per esempio, il sistema operativo può utilizzare degli algoritmi sofisticati associati a strutture dati complesse per seguire l'utilizzo delle diverse pagine fisiche e cercare di scegliere di sostituire quella pagina fisica che si presume non verrà richiesta per molto tempo. L'utilizzo di uno schema di sostituzione intelligente e flessibile riduce la frequenza dei page fault e semplifica l'utilizzo di uno schema di posizionamento delle pagine completamente associativo.

Come menzionato nel paragrafo 5.4, la difficoltà che si incontra nell'utilizzare uno schema completamente associativo è nell'individuazione dell'elemento, dato che esso può trovarsi ovunque all'interno del livello superiore della gerarchia e una ricerca esaustiva è impraticabile. In un sistema di memoria virtuale la pagina fisica corrispondente a una pagina virtuale viene individuata attraverso una tabella che indica la memoria; questa tabella viene chiamata **tabella delle pagine** (*page table*) e risiede nella memoria principale. La tabella delle pagine ha come indice il numero di pagina contenuto nell'indirizzo virtuale e specifica il numero della corrispondente pagina fisica. Ogni programma possiede la propria tabella delle pagine, che mappa lo spazio di indirizzamento virtuale del programma nella memoria fisica. Tornando all'analogia con la biblioteca, la tabella delle pagine corrisponde al catalogo, che contiene la corrispondenza tra i titoli dei libri e la loro posizione sugli scaffali. Come il catalogo della biblioteca può contenere informazioni relative anche ai libri che stanno in altre biblioteche, così la tabella delle pagine può contenere informazioni relative a pagine che non sono presenti nella memoria principale. Per individuare la posizione della tabella delle pagine nella memoria, l'architettura prevede un

**Tabella delle pagine:** la tabella che contiene la traduzione degli indirizzi virtuali in indirizzi fisici in un sistema di memoria virtuale. Essa è memorizzata nella memoria principale e viene tipicamente indicizzata tramite il numero della pagina virtuale: ogni elemento della tabella contiene il numero della pagina fisica corrispondente alla pagina virtuale, sempre che la pagina sia presente nella memoria principale.

registro hardware che indica l'indirizzo iniziale della tabella; questo registro è chiamato *registro della tabella delle pagine*. Ipotizzeremo per il momento che la tabella delle pagine si trovi in un'area di memoria prefissata e contigua.

La tabella delle pagine, insieme al program counter e ai registri, specifica lo stato di una macchina virtuale; se si vuole consentire a un'altra macchina virtuale di utilizzare il processore, è necessario salvare lo stato, cioè il contenuto di questi componenti. In seguito, dopo avere ripristinato lo stato, il programma potrà continuare l'esecuzione dal punto in cui era stato interrotto. Spesso si chiama *processo* questo particolare stato; il processo viene considerato *attivo* quando è padrone del processore, altrimenti si dice che è *inattivo*. Il sistema operativo può rendere attivo un processo caricandone lo stato, incluso il program counter, in modo da far riprendere l'esecuzione del programma a partire dall'indirizzo contenuto nel program counter, precedentemente salvato.

Lo spazio di indirizzamento del processo è quindi l'insieme di tutti i dati della memoria a cui il processo può accedere ed è definito dalla sua tabella delle pagine, che risiede anch'essa in memoria. Invece di salvare l'intera tabella delle pagine, il sistema operativo carica semplicemente solo il registro della tabella delle pagine, in modo tale da potere indirizzare la tabella appartenente al processo che vuole rendere attivo. Ogni processo ha la propria tabella delle pagine, dato che processi differenti possono utilizzare gli stessi indirizzi virtuali. Il sistema operativo è responsabile dell'allocazione della memoria fisica e dell'aggiornamento della tabella delle pagine, che deve essere fatta in modo tale che gli spazi virtuali di due processi differenti non entrino in conflitto. Come vedremo tra breve, l'utilizzo di tabelle separate per processi differenti permette di proteggere un processo dagli altri.

In Figura 5.27 è mostrato il modo in cui l'hardware forma l'indirizzo fisico a partire dal registro della tabella delle pagine, la relativa tabella delle pagine e l'indirizzo virtuale. In ogni elemento della tabella è presente anche un bit di validità, proprio come nelle cache. Se il bit è a 0, la pagina non è presente nella memoria fisica e si verifica un errore di page fault; se invece il bit assume il valore 1, la pagina è presente e il corrispondente elemento della tabella delle pagine contiene il numero della pagina fisica.

Poiché la tabella delle pagine contiene la traduzione di ogni possibile pagina virtuale, non c'è bisogno di aggiungere un campo tag. Utilizzando la terminologia delle cache, l'indice utilizzato per accedere alla tabella delle pagine è l'indirizzo completo di un blocco e corrisponde al numero della pagina virtuale.

## Page fault

Se il bit di validità associato a una pagina virtuale è impostato a 0, si verifica un errore di page fault e il controllo deve essere trasferito al sistema operativo. Il trasferimento viene effettuato attraverso il sollevamento di un'eccezione, come verrà descritto più avanti in questo paragrafo (*vedi* anche Cap. 4). Quando il sistema operativo prende il controllo, deve scoprire dove si trova la pagina richiesta nel livello inferiore della gerarchia, solitamente costituita da una memoria flash o da un disco magnetico, e decidere dove copiare la pagina nella memoria principale. L'indirizzo virtuale da solo non consente di individuare immediatamente la posizione della pagina nella memoria secondaria. Tornando all'analogia con la biblioteca, non è possibile individuare la posizione di un libro sugli scaffali conoscendone soltanto il titolo, ma è necessario consultare il catalogo per ricavare la sua posizione. Analogamente, in un sistema di memoria virtuale occorre tenere traccia della posizione nella memoria secondaria di ogni pagina definita nello spazio di indirizzamento virtuale.

## Interfaccia hardware/software

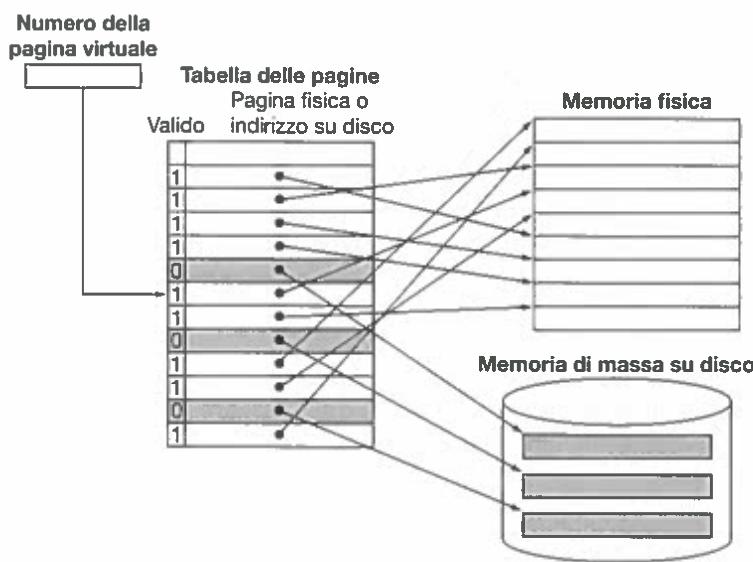


**Figura 5.27** La tabella delle pagine ha come indice il numero della pagina virtuale e contiene la parte corrispondente (più significativa) dell'indirizzo fisico. Si suppone che gli indirizzi siano su 48 bit. L'indirizzo di partenza della tabella delle pagine è contenuto nel registro della tabella delle pagine. Nella figura, la dimensione della pagina è di  $2^{12}$  byte o 4 KiB. Lo spazio di indirizzamento virtuale è di  $2^{48}$  byte, cioè 256 TiB, mentre lo spazio di indirizzamento fisico è di  $2^{40}$  byte, che consente alla memoria fisica di arrivare fino a 1 TiB. Se il RISC-V utilizzasse una singola tabella delle pagine come mostrato in figura, il numero di elementi per pagina sarebbe  $2^{36}$  cioè circa 64 miliardi. (Vedremo a breve la strategia utilizzata nei RISC-V per ridurre il numero di elementi.). Il bit di validità associato a ogni elemento indica se la mappatura è valida. Se il bit di validità è impostato a 0, la pagina richiesta non è presente nella memoria principale. Sebbene ogni elemento della tabella delle pagine qui mostrata richieda solo 29 bit, l'ampiezza degli elementi viene tipicamente arrotondata a una potenza di 2 per ottenere un indirizzamento più facile. Gli elementi della tabella delle pagine del RISC-V sono su 64 bit. I bit aggiuntivi possono essere utilizzati per memorizzare informazioni supplementari relative a ogni pagina, come la protezione..

**Spazio di swap:** lo spazio che viene riservato su disco all'intero spazio di memoria virtuale di un processo.

Dato che non possiamo sapere a priori quando una pagina della memoria principale verrà selezionata per essere sostituita, di solito il sistema operativo riserva sulla memoria flash o sul disco uno spazio sufficiente per tutte le pagine di un processo nel momento in cui crea il processo stesso; questo spazio viene chiamato **spazio di swap** (spazio di scambio). Contemporaneamente, il SO crea anche una struttura dati in cui registra la posizione su disco di ciascuna pagina virtuale. Questa struttura dati può far parte della tabella delle pagine, oppure può essere una struttura dati a se stante, indicizzata come la tabella delle pagine. La Figura 5.28 mostra l'organizzazione della memoria virtuale nel caso in cui si utilizzi una sola tabella che contiene il numero della pagina fisica o l'indirizzo della memoria secondaria, corrispondente a ogni pagina virtuale.

Il sistema operativo crea inoltre una struttura dati che tiene traccia dei processi e degli indirizzi virtuali associati a ciascuna pagina fisica. Quando si verifica un errore di page fault, se tutte le pagine della memoria fisica sono in uso il sistema operativo deve scegliere una pagina da sostituire. Poiché si vogliono minimizzare i page fault, molti sistemi operativi cercano di scegliere una pagi-



**Figura 5.28** La tabella delle pagine mappa ciascuna pagina della memoria virtuale in una pagina della memoria fisica oppure in una pagina memorizzata su disco, che è il livello successivo della gerarchia. Il numero della pagina virtuale viene utilizzato come indice della tabella delle pagine. Se il bit di validità è impostato a 1, la tabella delle pagine fornisce il numero della pagina fisica corrispondente alla pagina virtuale, ossia l'indirizzo della memoria fisica a partire dal quale è memorizzata la pagina virtuale. Se invece il bit di validità è impostato a 0, la pagina in quel momento si trova solamente sul disco, all'indirizzo specificato. In molti sistemi, le tabelle degli indirizzi delle pagine fisiche e delle pagine su disco, sebbene costituiscano una sola tabella dal punto di vista logico, vengono memorizzate in due strutture dati separate. L'utilizzo di tabelle separate è in parte giustificato dal fatto che è necessario conservare gli indirizzi su disco di tutte le pagine, anche di quelle che sono correntemente presenti nella memoria fisica. Ricordiamo che le pagine della memoria principale e le pagine su disco hanno dimensioni identiche.

na che suppongono di non dover utilizzare poco dopo. Analizzando il passato per predire il futuro, i sistemi operativi seguono uno schema di sostituzione di tipo LRU (*Least Recently Used*, già introdotto nel paragrafo 5.4). Il sistema operativo cerca la pagina utilizzata meno di recente, ipotizzando che una pagina inutilizzata da molto tempo abbia una minore probabilità di essere utilizzata rispetto alle altre. Le pagine sostituite vengono salvate nello spazio di swap del disco. Forse vi state domandando come venga trattato il sistema operativo; esso è analogo a ogni altro processo e le tabelle utilizzate dal sistema operativo per gestire la memoria virtuale sono residenti nella memoria principale; una spiegazione dettagliata di questa apparente contraddizione verrà data tra breve.

Implementare uno schema di tipo LRU completamente corretto è troppo dispendioso, poiché richiederebbe l'aggiornamento della struttura dati a ogni accesso alla memoria. Per questo motivo molti sistemi operativi adottano una strategia LRU approssimata, tenendo traccia delle pagine che sono state utilizzate da poco. Per aiutare il sistema operativo a individuare la pagina richiesta meno di recente, alcuni calcolatori forniscono un **bit di utilizzo** (*use bit*), detto anche **bit di indirizzamento** (*reference bit*), che viene impostato a 1 ogni volta che si accede alla pagina corrispondente. Il sistema operativo azzerà periodicamente tutti i bit di utilizzo, per poi impostare a 1 il bit di utilizzo associato a ogni pagina indirizzata da quel momento in poi. Analizzando questo bit, il sistema operativo può determinare quali pagine siano state toccate durante un certo lasso di tempo e può quindi utilizzare questa informazione per selezionare una pagina tra quelle richieste meno di recente, cioè tra quelle il cui bit di utilizzo è impostato a 0. Se questo bit non viene fornito dall'hardware, il sistema operativo deve trovare un altro metodo per individuare quali siano le pagine richieste di recente.

## Interfaccia hardware/software

**Bit di utilizzo (o di indirizzamento):** è un campo che viene impostato a 1 ogni volta che si accede alla pagina associata e che viene utilizzato per implementare la politica LRU o altri schemi di sostituzione.

## Memoria virtuale per un insieme ampio di indirizzi virtuali

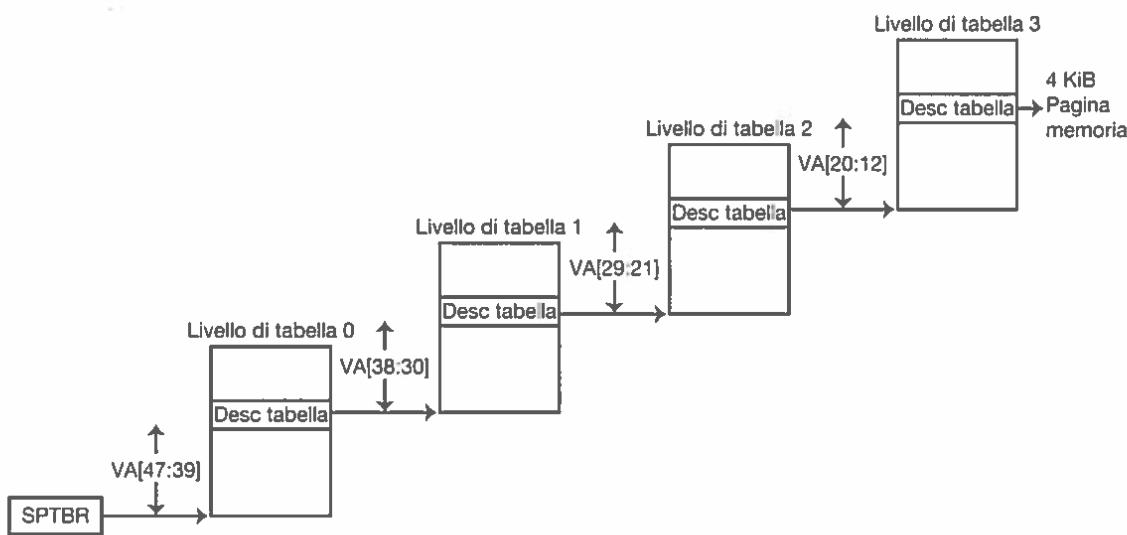
La legenda di Figura 5.27 sottolinea che una tabella delle pagine su un solo livello sarebbe costituita da 64 miliardi di elementi per un indirizzo su 48 bit con pagine di 4 KiB. Dato che ciascun elemento della tabella delle pagine nel RISC-V è di 8 byte occorrerebbero 0,5 TiB solamente per mappare gli indirizzi virtuali in indirizzi fisici! Inoltre, ci potrebbero essere centinaia di processi in esecuzione, ciascuno con la propria tabella delle pagine. Questa enorme quantità di memoria dedicata alla traduzione non sarebbe possibile neppure nei sistemi più grandi.

Esistono molte tecniche per ridurre la quantità di memoria richiesta da una tabella delle pagine. I seguenti cinque metodi hanno lo scopo di ridurre la quantità totale di memoria richiesta e di minimizzare la quantità di memoria principale dedicata alla tabella delle pagine.

1. La tecnica più semplice consiste nell'utilizzare un registro limite che riduca la dimensione della tabella delle pagine di un certo processo. Se il numero di pagine virtuali cresce oltre il valore contenuto nel registro limite, occorre aggiungere nuovi elementi alla tabella delle pagine. Questa tecnica permette di aumentare la dimensione della tabella delle pagine mano a mano che un processo richiede più spazio in memoria. In questo modo, la tabella diventa grande solo se il processo utilizza molte pagine dello spazio di indirizzamento virtuale. Questa tecnica richiede che lo spazio di indirizzamento si espanda in una sola direzione (crescente).
2. Consentire la crescita della memoria in una sola direzione non è sufficiente, poiché la maggior parte dei linguaggi di programmazione richiede almeno due aree di memoria che possano crescere secondo necessità: un'area che contiene lo stack e l'altra che contiene i dati dinamici ed è chiamata *heap*. Per questo motivo, è conveniente suddividere la tabella delle pagine in due parti e lasciare crescere verso il basso la parte superiore e verso l'alto la parte inferiore. Questo significa che ci saranno due tabelle delle pagine separate e due limiti distinti. L'utilizzo di due tabelle delle pagine divide lo spazio di indirizzamento in due segmenti: il bit più significativo di un indirizzo virtuale determina generalmente il segmento e quindi quale delle due tabelle delle pagine debba essere utilizzata per tradurre quell'indirizzo. Dato che il segmento viene specificato dal bit più significativo dell'indirizzo, ognuno dei due segmenti può occupare al massimo metà dello spazio di indirizzamento. Un registro limite associato a ogni segmento contiene la dimensione corrente del segmento, che cresce di una unità per ogni nuova pagina. A differenza della segmentazione precedentemente descritta, questo tipo di segmentazione è invisibile ai programmi applicativi ma non al sistema operativo. Il principale svantaggio di questo schema è che non funziona bene quando lo spazio degli indirizzi viene utilizzato in modo "sparso" e non come un insieme contiguo di indirizzi virtuali.
3. Un altro approccio per ridurre la dimensione della tabella delle pagine consiste nell'applicare all'indirizzo virtuale una funzione di corrispondenza (*hashing function*), in modo tale che la tabella delle pagine abbia dimensione pari soltanto al numero delle pagine fisiche della memoria principale. Una struttura di questo tipo viene chiamata *tabella inversa delle pagine*. Ovviamente il processo di ricerca in una tabella inversa delle pagine è leggermente più complesso, perché non è più sufficiente utilizzare il numero della pagina virtuale come indice della tabella.
4. Per ridurre la quantità di memoria principale consumata dalle tabelle delle pagine, la maggior parte dei sistemi moderni prevede che la tabella delle pagine sia, a sua volta, suddivisa in pagine. Sebbene questo sembri un trucco, funziona utilizzando gli stessi concetti che sono alla base della memoria

virtuale, a patto di consentire alla tabella delle pagine di risiedere nello spazio di indirizzamento virtuale. Inoltre, occorre evitare alcuni problemi piccoli ma critici, come la generazione di una catena infinita di errori di page fault. Il modo in cui questi problemi vengono superati è molto particolare ed è di solito strettamente legato alla macchina. In breve, questi problemi si evitano posizionando tutte le tabelle delle pagine all'interno dello spazio di indirizzamento del sistema operativo e inserendo almeno qualcuna delle tabelle delle pagine utilizzate dal sistema operativo in una porzione della memoria principale che viene indirizzata fisicamente, le cui pagine sono sempre presenti in memoria e quindi non si troveranno mai nella memoria secondaria.

- Per ridurre la quantità totale di memoria richiesta per memorizzare una tabella delle pagine, si possono utilizzare anche tabelle delle pagine su più livelli, e questa è la soluzione adottata dal RISC-V per ridurre l'occupazione di memoria dovuta alla traduzione degli indirizzi. La Figura 5.29 mostra i quattro livelli di traduzione richiesti per passare da un indirizzo virtuale su 48 bit a un indirizzo fisico su 40 bit per pagine di 4 KiB. La traduzione di un indirizzo inizia guardando nella tabella di livello 0, utilizzando i bit più significativi dell'indirizzo. Se questo indirizzo si trova nella tabella, viene utilizzato il gruppo successivo di bit più significativi per indicizzare la tabella delle pagine indicata dall'elemento della tabella indicante il segmento, e così via. In questo modo, la tabella di livello 0 mappa l'indirizzo virtuale in regioni di 512 GB ( $2^{39}$  byte). La tabella di livello 1 a sua volta mappa l'indirizzo virtuale su una regione di 1 GB ( $2^{30}$  byte). La mappa del livello successivo mappa questo indirizzo in una regione di 2 MB ( $2^{21}$  byte). L'ultima tabella mappa l'indirizzo virtuale nella pagina di memoria di 4 KiB ( $2^{12}$  byte). Questo schema consente di utilizzare lo spazio di indirizzamento in maniera sparsa; possono cioè essere attivi segmenti multipli non contigui, senza dover allocare l'intera tabella delle pagine. Tali schemi sono particolarmente utili con spazi di indirizzamento molto ampi e in sistemi software che richiedano un'allocazione non contigua delle pagine. Il principale svantaggio di questa tecnica a più livelli è la maggiore complessità del processo di traduzione degli indirizzi virtuali.



**Figura 5.29** Il RISC-V utilizza quattro livelli di tabelle per tradurre un indirizzo virtuale di 48 bit in un indirizzo fisico di 40 bit. Invece di richiedere 64 miliardi di elementi in una singola tabella come in Figura 5.27, l'approccio gerarchico ne richiede solamente una piccola frazione. Ciascun passo della traduzione utilizza 9 bit dell'indirizzo virtuale per identificare la tabella del livello successivo, fino quando i 36 bit più significativi dell'indirizzo virtuale non sono stati mappati sull'indirizzo fisico della pagina desiderata di 4 KiB. Ciascun elemento di una tabella della pagina RISC-V è di 8 byte, per cui i 512 elementi di una tabella riempiono una singola pagina di 4 KiB. Il registro base della tabella delle pagine del supervisore (SPTBR, Supervisor Page Table Base Register) fornisce l'indirizzo di partenza della prima tabella delle pagine.

## Che cosa succede in scrittura?

La differenza tra il tempo di accesso a una cache e il tempo di accesso alla memoria principale è dell'ordine delle decine o centinaia di cicli di clock; risulta quindi possibile utilizzare uno schema di scrittura di tipo write-through per le cache, anche se occorre introdurre un buffer di scrittura per nascondere al processore la latenza della scrittura. In un sistema di memoria virtuale, la scrittura al livello inferiore della gerarchia, il disco, richiede milioni di cicli di clock; perciò, utilizzare un buffer di scrittura per consentire al processore di scrivere direttamente sul disco sarebbe del tutto inutile. I sistemi di memoria virtuale devono invece utilizzare lo schema write-back, in cui il singolo dato viene scritto nella pagina della memoria principale; l'intera pagina viene poi copiata nella memoria secondaria quando questa deve essere sostituita in memoria principale.

## Interfaccia hardware/software

Lo schema write-back ha un altro grande vantaggio per i sistemi di memoria virtuale. Dato che il tempo di trasferimento su disco è piccolo rispetto al tempo di accesso al disco, copiare un'intera pagina su disco è molto più efficiente che non scrivere le singole parole. Tuttavia, un'operazione di scrittura su disco è comunque costosa anche utilizzando lo schema di scrittura write-back, che è più efficiente del trasferimento di parole singole. Perciò è opportuno sapere se sia effettivamente *necessario* copiare una pagina sul disco quando si sceglie di sostituirla. Per sapere se il sistema abbia scritto su una certa pagina a partire dall'istante in cui la pagina è stata caricata in memoria, si aggiunge un bit, detto *dirty bit* (bit sporco), alla tabella delle pagine. Il dirty bit viene impostato a 1 nel momento in cui una qualsiasi parola della pagina viene scritta. Quando il sistema operativo sceglie di sostituire la pagina, il dirty bit indica se è necessario scrivere la pagina sul disco prima di rilasciare lo spazio che occupava in memoria per destinarlo a una nuova pagina. Una pagina il cui contenuto è stato modificato viene chiamata *dirty page* (pagina sporca).

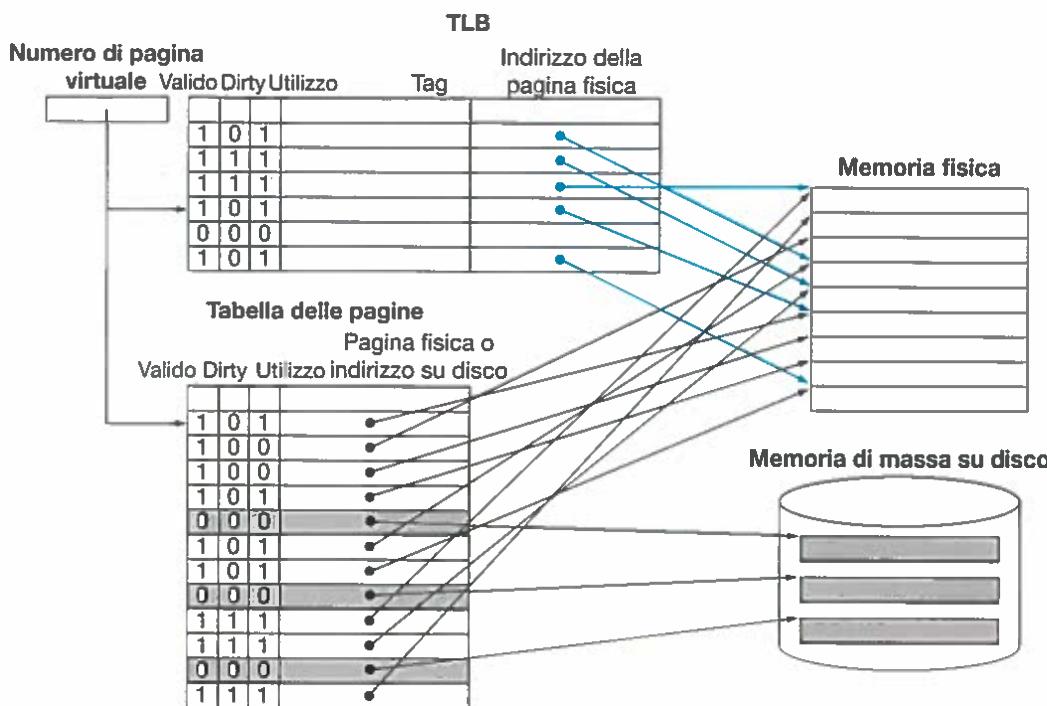
## Come rendere più veloce la traduzione degli indirizzi: il TLB

Dato che la tabella delle pagine deve essere salvata nella memoria principale, ogni accesso alla memoria da parte di un programma richiede almeno il doppio del tempo, dovendo compiere un accesso alla memoria per ottenere l'indirizzo fisico e un secondo accesso per avere il dato richiesto. La chiave per migliorare le prestazioni in termini di tempo di accesso è lo sfruttamento della località degli accessi alla tabella delle pagine. Quando si utilizza la traduzione del numero di una certa pagina virtuale, è molto probabile che la stessa traduzione sia richiesta di nuovo poco dopo, dato che le parole contenute in una pagina sono caratterizzate da località sia spaziale sia temporale.

Basandosi su queste considerazioni, molti processori moderni contengono una cache speciale che tiene traccia delle traduzioni utilizzate più di recente. Questa cache speciale, che specifica la traduzione degli indirizzi virtuali utilizzati più di recente, è chiamata **translation lookaside buffer (TLB)**, letteralmente "buffer di traduzione consultato a parte", sebbene sia più corretto chiamarla *cache di traduzione*. Tornando all'analogia con la biblioteca, il TLB corrisponde a quel pezzettino di carta che utilizziamo per scrivere la posizione dei libri che abbiamo cercato nel catalogo: invece di consultare ogni volta l'intero schedario, annotiamo la posizione dei diversi libri e utilizziamo i nostri appunti come se fossero una cache del catalogo della biblioteca.

La Figura 5.30 mostra che il campo tag del TLB contiene una parte del numero della pagina virtuale, mentre il campo dati contiene il numero di una

**Translation lookaside buffer (TLB):** una cache che tiene traccia degli indirizzi virtuali tradotti più di recente per evitare l'accesso alla tabella delle pagine.



**Figura 5.30** Il TLB si comporta come una cache della tabella delle pagine solo per le pagine mappate su una pagina fisica. Il TLB contiene un sottoinsieme delle corrispondenze tra gli indirizzi virtuali e gli indirizzi fisici che sono memorizzati nella tabella delle pagine. Le corrispondenze definite nel TLB sono disegnate in blu. Dato che il TLB è una cache, deve contenere il campo tag. Se nessun elemento del TLB corrisponde alla pagina richiesta, si rende necessario esaminare la tabella delle pagine: questa può fornire il numero della corrispondente pagina fisica (che può essere inserito assieme alla pagina virtuale corrispondente nel TLB), oppure indica che la pagina si trova sul disco (nel qual caso si genera un page fault). Poiché la tabella delle pagine contiene un elemento per ogni pagina virtuale, essa non richiede alcun campo etichetta. In altre parole, a differenza del TLB, la tabella delle pagine *non è* una memoria cache.

pagina fisica. Poiché non c'è più bisogno di accedere alla tabella delle pagine a ogni richiesta di accesso alla memoria (dato che stiamo utilizzando il TLB), nelle linee del TLB devono essere presenti altri bit di stato: il bit di accesso e il dirty bit. Anche se la Figura 5.30 mostra una tabella delle pagine singola, i TLB funzionano bene anche con le tabelle delle pagine multilivello. Il TLB semplicemente carica l'indirizzo fisico e i tag di protezione associati alla tabella delle pagine di primo livello.

Quando viene richiesta una pagina, per prima cosa si cerca il numero della pagina virtuale nel TLB. Se la ricerca ha avuto successo, il numero di pagina fisica corrispondente viene utilizzato per costruire l'indirizzo e il bit di utilizzo viene impostato a 1; se il processore sta effettuando una scrittura, anche il dirty bit viene impostato a 1. Se l'accesso al TLB genera una miss, occorre verificare se si tratti di un vero page fault, oppure soltanto di una miss del TLB. Se la pagina è presente nella memoria principale, allora la miss del TLB indica solamente che manca la traduzione. In tal caso, il processore può gestire la miss del TLB caricando dalla tabella delle pagine (dell'ultimo livello) nel TLB la traduzione dell'indirizzo virtuale, per poi ripetere l'accesso al TLB. Se invece la pagina non si trova nella memoria principale, allora la miss del TLB indica che si è verificato un vero e proprio page fault. In questo caso, il processore richiede l'intervento del sistema operativo, sollevando un'eccezione. Dato che il TLB contiene molti meno elementi rispetto al numero di pagine della memoria principale, le miss del TLB saranno molto più frequenti dei page fault effettivi.

Le miss del TLB possono essere gestite via hardware oppure via software. La differenza di prestazioni tra le due soluzioni è minima, poiché le operazioni fondamentali che devono essere eseguite sono le stesse in entrambi i casi.

Quando si verifica una miss del TLB e la traduzione mancante è stata recuperata dalla tabella delle pagine, occorre individuare quale elemento del TLB debba essere sostituito. Dato che il bit di utilizzo e il dirty bit fanno parte di ogni elemento del TLB, quando un elemento del TLB viene sostituito occorre copiare questi bit nel blocco corrispondente della tabella delle pagine. Questi bit rappresentano l'unica parte di un elemento del TLB che può essere modificata. Utilizzare una strategia di write-back, ossia copiare questi bit quando si verifica una miss invece di farlo ogni volta che vengono modificati, è molto efficiente, poiché ci aspettiamo che la frequenza di miss del TLB sia bassa. Alcuni sistemi utilizzano altre tecniche per impostare in modo approssimato il bit di utilizzo e il dirty bit, eliminando la necessità di scrivere nel TLB se non quando viene caricato un nuovo elemento a seguito di una miss.

Alcuni tipici valori di un TLB possono essere:

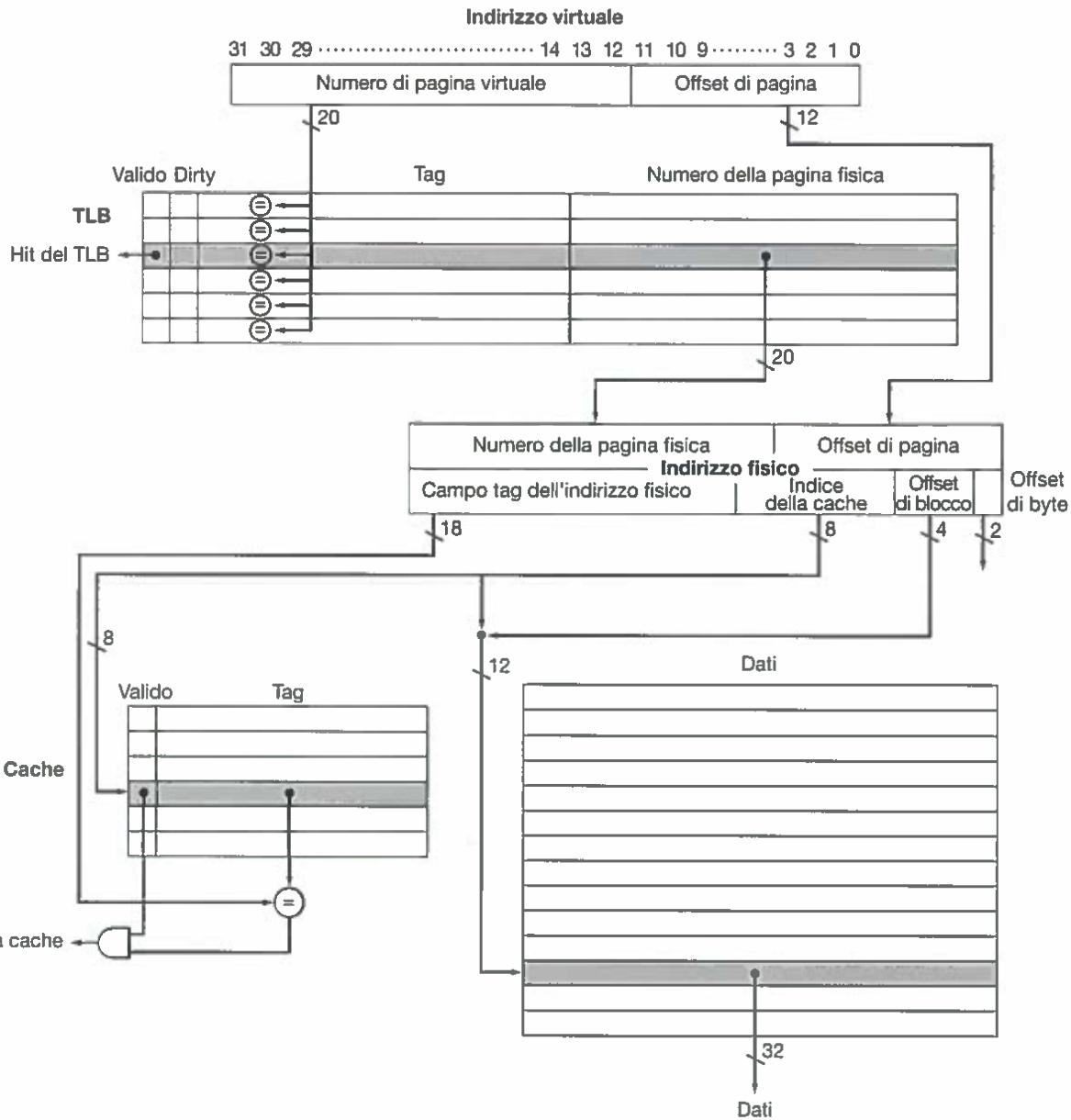
- dimensione del TLB: 16-512 elementi;
- dimensione del blocco: 1-2 elementi della tabella delle pagine (tipicamente 4-8 byte ciascuno);
- tempo di hit: da 0,5 a 1 ciclo di clock;
- penalità di miss: 10-100 cicli di clock;
- frequenza di miss: 0,01%-1%.

I progettisti hanno utilizzato una vasta gamma di gradi di associatività nei TLB. Molti sistemi utilizzano piccoli TLB completamente associativi, perché una corrispondenza di questo tipo consente una frequenza di miss più bassa; inoltre, viste le ridotte dimensioni del TLB, il costo della mappatura completamente associativa non è troppo elevato. Altri sistemi utilizzano TLB di grandi dimensioni, spesso con basso grado di associatività. Con una mappatura completamente associativa, scegliere l'elemento da sostituire diventa complicato, dal momento che l'implementazione hardware di uno schema LRU diventa troppo costoso. Inoltre, essendo le miss del TLB molto più frequenti dei page fault e dovendo quindi essere gestite in maniera più efficiente, non si può accettare l'utilizzo di un algoritmo software pesante da un punto di vista computazionale, come invece è possibile fare per la sostituzione delle pagine a seguito di un page fault. Di conseguenza, molti sistemi prevedono di sostituire un elemento del TLB scelto a caso. Esamineremo gli schemi di sostituzione con maggior dettaglio nel paragrafo 5.8.

### TLB del processore FastMATH Intrinsity

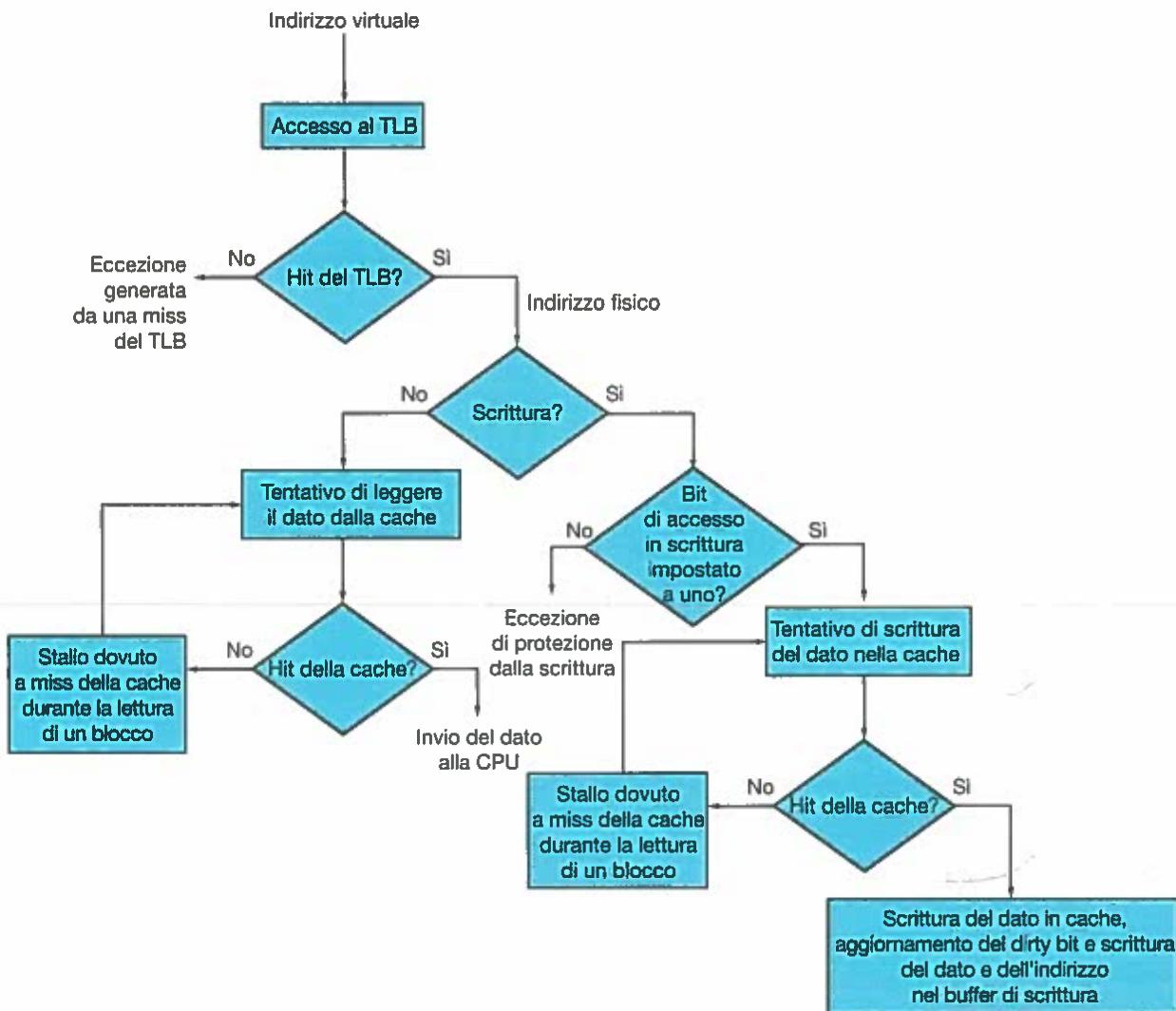
Per vedere come vengono applicati questi concetti a un processore reale, esaminiamo più da vicino il TLB del processore FastMATH Intrinsity. Il sistema di memoria utilizza pagine di 4 KiB e uno spazio di indirizzamento a soli 32 bit; quindi, il numero della pagina virtuale viene codificato su 20 bit. L'indirizzo fisico ha la stessa dimensione dell'indirizzo virtuale. Il TLB contiene 16 elementi, è completamente associativo e contiene sia dati che istruzioni; ogni elemento del TLB è di 64 bit e comprende un campo tag di 20 bit (contenente il numero della pagina virtuale associata all'elemento del TLB), il numero della pagina fisica associata (anch'esso su 20 bit), un bit di validità, un dirty bit e altri bit richiesti per la gestione della pagina virtuale. Come la maggior parte delle architetture MIPS, le miss del TLB vengono gestite via software.

La Figura 5.31 mostra il TLB e una delle cache, mentre la Figura 5.32 illustra i passi necessari per soddisfare una richiesta di lettura o scrittura. Quando avviene una miss del TLB, l'hardware provvede a salvare il numero di pagina associata all'indirizzo richiesto in un registro speciale e solleva un'eccezione; l'eccezione attiva il sistema operativo, il quale gestisce la miss via software. Per trovare l'indirizzo fisico, la procedura software di risposta all'eccezione di miss



**Figura 5.31 Il TLB e la cache implementano, nel processore FastMATH Intrinsity, il meccanismo che consente di trasformare un indirizzo virtuale nel dato corrispondente.** Questa figura mostra l'organizzazione del TLB e della cache dati supponendo una dimensione della pagina di 4 KiB. Si noti che gli indirizzi per questo computer sono ampi solo 32 bit. Lo schema si focalizza sulla lettura; la Figura 5.32 descriverà come vengono gestite le scritture. Si noti che, diversamente dalla Figura 5.12, la memoria che contiene i campi tag e quella che contiene i dati sono separate. Indirizzando la memoria dati (RAM), lunga e stretta, concatenando i campi dell'indice e dell'offset di blocco, viene selezionata la parola del blocco desiderata senza aver bisogno di un multiplexer 16:1. Mentre la cache è a mappatura diretta, il TLB è completamente associativo. Implementare un TLB completamente associativo richiede che ogni tag del TLB venga confrontato con il numero della pagina virtuale, dato che l'elemento di interesse può essere ovunque nel TLB (vedi la sezione *Approfondimento* di pagina 353). Se il bit di validità dell'elemento identificato è impostato a 1, l'accesso si traduce in una hit del TLB e i bit del numero della pagina fisica insieme ai bit provenienti dall'offset della pagina concorrono a formare l'indice che viene utilizzato per accedere alla cache.

indicizza la tabella delle pagine utilizzando il numero di pagina dell'indirizzo virtuale e il registro della tabella delle pagine che specifica l'indirizzo di partenza della tabella delle pagine associata al processo attivo. Attraverso uno speciale insieme di istruzioni di sistema che aggiornano il TLB, il sistema operativo scrive nel TLB l'indirizzo fisico preso dalla tabella delle pagine. Una miss del TLB può essere gestita in circa 13 cicli di clock, supponendo che il codice di risposta all'eccezione e l'elemento della tabella delle pagine siano contenuti



**Figura 5.32** Gestione di una lettura o di una scrittura di tipo write-through nel TLB e nella cache nel processore FastMATH Intrinsity. Se l'accesso al TLB produce una hit, l'accesso alla cache avviene con l'indirizzo fisico ottenuto. In lettura, la cache può generare una hit o una miss e quindi fornire il dato richiesto o causare uno stallo che dura il tempo necessario perché il dato sia prelevato dalla memoria principale. In scrittura, se l'accesso genera una hit, una parte dell'elemento della cache selezionato viene sovrascritto e il dato viene inviato al buffer di scrittura, poiché utilizziamo un approccio di tipo write-through. Una miss in scrittura è simile a una miss in lettura, ad eccezione del fatto che il blocco viene modificato dopo che è stato letto dalla memoria principale. La politica write-back richiederebbe di impostare a 1 (a ogni scrittura) il dirty bit associato al blocco della cache e di caricare nel buffer di scrittura l'intero blocco solo quando si verifica una miss in lettura o scrittura e il blocco che deve essere sostituito è stato modificato. Si noti inoltre che una hit del TLB e una hit della cache sono eventi indipendenti, ma una hit nella cache può verificarsi solo se è avvenuta una hit del TLB, il che significa che il dato deve essere presente nella memoria principale. La relazione tra miss del TLB e miss della cache verrà esaminata nel prossimo esempio e negli esercizi alla fine del capitolo. Si noti che gli indirizzi per questo computer sono ampi solo 32 bit.

rispettivamente nella cache istruzioni e nella cache dati. Un vero e proprio page fault si verifica se il blocco individuato nella tabella delle pagine non contiene un indirizzo fisico valido. L'hardware mantiene un indice che segnala quale elemento del TLB sia opportuno sostituire; in questo caso l'elemento da sostituire è scelto a caso.

Esiste una complicazione aggiuntiva per le richieste di scrittura, dovuta al fatto che occorre controllare il bit di accesso in scrittura del TLB. Questo bit impedisce a un programma di scrivere nelle pagine alle quali può accedere solamente in lettura. Se il programma cerca di scrivere ma il bit di accesso in scrittura è impostato a 0, viene generata un'eccezione. Il bit di accesso in scrittura fa parte del meccanismo di protezione che discuteremo tra breve.

## Integrazione della memoria virtuale, dei TLB e delle cache

La memoria virtuale e le cache lavorano insieme, in un'unica gerarchia, cosicché un dato non può trovarsi all'interno della cache a meno che non sia presente anche nella memoria principale. Il sistema operativo gioca un ruolo importante nel sostenere questa gerarchia, scaricando il contenuto delle pagine presenti in cache quando decide di spostare queste pagine sul disco. Allo stesso tempo, il sistema operativo modifica la tabella delle pagine e il TLB, cosicché un tentativo di accesso a un qualsiasi dato all'interno di una pagina spostata su disco genera un errore di page fault.

Nella migliore delle ipotesi, un indirizzo virtuale viene tradotto dal TLB e inviato alla cache, dove il dato cercato viene individuato, letto e inviato al processore. Nel caso peggiore l'accesso si traduce in una miss di tutti e tre i componenti della gerarchia delle memorie: il TLB, la tabella delle pagine e la cache. Il seguente esempio illustra più in dettaglio l'interazione fra i tre componenti.

### Complesso delle operazioni di una gerarchia delle memorie

In una gerarchia delle memorie come quella di Figura 5.31, che comprende un TLB e una cache, un accesso alla memoria può generare tre differenti tipi di miss: una miss del TLB, un page fault e una miss della cache. Considerare tutte le possibili combinazioni di questi tre tipi di miss, tenendo presente che una miss può essere generata da uno o più componenti, per un totale di sette combinazioni diverse. Per ogni combinazione, si dica se l'evento può realmente accadere e in quali situazioni.

La Figura 5.33 mostra tutte le combinazioni possibili e se queste possano presentarsi o meno nella realtà.

ESEMPIO

SOLUZIONE

TLB	Tabella delle pagine	Cache	Possibile? Se sì, in quale situazione?
Hit	Hit	Miss	Possibile, sebbene la tabella delle pagine non venga mai realmente controllata se si verifica una hit del TLB
Miss	Hit	Hit	Miss del TLB, ma l'elemento si trova nella tabella delle pagine; al secondo tentativo, il dato viene trovato nella cache
Miss	Hit	Miss	Miss del TLB, ma l'elemento si trova nella tabella delle pagine; al secondo tentativo, si verifica però una miss della cache
Miss	Miss	Miss	Miss del TLB, seguita da un page fault; al secondo tentativo, l'accesso al dato deve provocare una miss della cache
Hit	Miss	Miss	Impossibile: il TLB non può fornire la traduzione di una pagina che non è presente in memoria
Hit	Miss	Hit	Impossibile: il TLB non può fornire la traduzione di una pagina che non è presente in memoria
Miss	Miss	Hit	Impossibile: i dati non possono trovarsi nella cache se la pagina non è presente in memoria

**Figura 5.33** Possibili combinazioni di eventi nel TLB, nel sistema di memoria virtuale e nella cache. Tre di queste combinazioni sono impossibili; una combinazione è teoricamente possibile, ma in pratica non si verifica mai: hit del TLB, hit della memoria virtuale e miss della cache.



**Cache indirizzata virtualmente:** una cache a cui si accede tramite l'indirizzo virtuale invece che tramite l'indirizzo fisico.

**Aliasing:** una situazione in cui si accede allo stesso oggetto attraverso due indirizzi; può verificarsi nella memoria virtuale quando esistono due indirizzi virtuali per la stessa pagina fisica.

**Cache indirizzata fisicamente:** una cache che viene indirizzata utilizzando l'indirizzo fisico.

**Approfondimento.** In Figura 5.33 si presuppone che tutti gli indirizzi di memoria siano tradotti in indirizzi fisici prima che avvenga l'accesso alla cache. In questa organizzazione, la cache è *indicizzata fisicamente* e contiene *tag fisici* (il campo tag e il campo indice sono riferiti all'indirizzo fisico e non a quello virtuale). In tale schema, la quantità di tempo richiesta per accedere alla memoria, in caso di hit della cache, deve essere sufficiente per accedere sia al TLB sia alla cache; ovviamente questi due accessi possono essere eseguiti in pipeline.

In alternativa, il processore può indirizzare la cache utilizzando un indirizzo completamente o parzialmente virtuale. In questo caso si parla di **cache indirizzata virtualmente** e contiene campi tag associati a indirizzi virtuali; quindi, una cache di questo tipo viene *indicizzata virtualmente* e contiene *tag virtuali*. In questo schema, la traduzione hardware dell'indirizzo virtuale in indirizzo fisico non viene utilizzata nei normali accessi alla cache, dato che gli accessi avvengono attraverso l'indirizzo virtuale che non è ancora stato tradotto in indirizzo fisico. Questo porta il TLB al di fuori del cammino critico, riducendo così la latenza della cache. Tuttavia, quando si verifica una miss della cache, il processore ha comunque bisogno di tradurre l'indirizzo virtuale in indirizzo fisico, in modo tale da poter prelevare dalla memoria principale il blocco da scrivere nella cache.

Quando si accede alla cache con un indirizzo virtuale e alcune pagine fisiche sono condivise tra diversi programmi (i quali possono accedere a queste pagine utilizzando indirizzi virtuali diversi), esiste la possibilità che si verifichi **aliasing** (sovraposizione). L'aliasing si verifica quando lo stesso oggetto ha due nomi diversi; in questo caso due indirizzi virtuali diversi per la stessa pagina. Questa ambiguità diventa critica, perché una parola appartenente a una pagina di questo tipo può essere caricata in due diverse locazioni della cache, corrispondenti a due indirizzi virtuali diversi, consentendo quindi a un programma di scrivere nella pagina senza che l'altro programma possa accorgersi che i dati vengono modificati. Le cache a indirizzamento virtuale richiedono che siano imposti dei vincoli di progetto sulla cache e sul TLB per ridurre l'aliasing, oppure richiedono che il sistema operativo, ed eventualmente l'utente, facciano in modo che l'aliasing non si possa mai verificare.

Un compromesso spesso utilizzato, che si trova a cavallo fra le due scelte di progetto estreme, è quello di utilizzare cache che contengono nei campi tag indirizzi fisici ma che vengono indicizzate virtualmente, utilizzando a volte solo l'offset interno alla pagina dell'indirizzo virtuale (che è anche l'offset dell'indirizzo fisico, poiché non viene tradotto). Queste cache, che sono quindi *indicizzate virtualmente ma contengono nei campi tag indirizzi fisici*, hanno l'obiettivo di raggiungere le prestazioni delle cache indicate virtualmente sfruttando la semplicità architettonica delle **cache indirizzate fisicamente**. In questo tipo di cache, per esempio, il problema dell'aliasing non esiste. In Figura 5.31 si presuppone una dimensione della pagina di 4 KiB, ma in realtà la dimensione è di 16 KiB, per cui il FastMATH Intrinsics può implementare questo stratagemma; tuttavia, perché funzioni, ci deve essere un attento coordinamento tra la dimensione minima della pagina, la dimensione della cache e il grado di associatività. Il RISC-V richiede che la cache si comporti come se avesse indici e tag fisici, ma non richiede che sia implementata in questo modo. Per esempio, cache con indici virtuali e tag fisici potrebbero utilizzare della logica aggiuntiva per garantire che il software non si accorga della differenza.

## Meccanismi di protezione basati sulla memoria virtuale

Oggigiorno, la funzione più importante della memoria virtuale è forse quella di consentire a più processi di condividere un'unica memoria principale, fornendo allo stesso tempo un meccanismo di protezione della memoria ai diversi processi e al sistema operativo. Il meccanismo di protezione deve garantire che, anche se più processi condividono la stessa memoria principale, non sia possibile per un processo scrivere nello spazio d'indirizzamento di un altro processo utente o nello spazio riservato al sistema operativo, di fatto non rispettando le regole in modo intenzionale o accidentale. Il bit di accesso in scrittura del TLB permette di proteggere la pagina associata, evitando che questa possa essere scritta. Senza questo livello di protezione, i virus dei calcolatori sarebbero ancora più diffusi.

Per permettere al sistema operativo di implementare i meccanismi di protezione nel sistema di memoria virtuale, l'hardware deve fornire almeno le tre seguenti funzionalità di base. Si noti che le prime due sono le stesse richieste per le macchine virtuali (par. 5.6).

1. Supportare almeno due modalità di funzionamento: la prima associata all'esecuzione di processi utente e la seconda all'esecuzione di processi del sistema operativo; questa seconda modalità viene chiamata anche **modalità kernel**, **modalità supervisore** o **modalità executive**.
2. Mettere a disposizione dei processi utente una parte dello stato del processore che i processi potranno leggere ma non scrivere. Questa parte dello stato comprende il bit che determina la modalità di funzionamento, utente o kernel, il puntatore alla tabella delle pagine e il TLB. Per scrivere in questi componenti il processore utilizza delle istruzioni speciali che sono disponibili solo in modalità kernel.
3. Fornire dei meccanismi che permettano al processore di passare dalla modalità utente alla modalità kernel e viceversa. Il primo passaggio viene solitamente realizzato mediante un'eccezione di **chiamata di sistema** (*system call*). Questa eccezione viene sollevata da un'istruzione speciale, che nell'insieme di istruzioni RISC-V è l'istruzione `ecall`; essa trasferisce il controllo a un indirizzo pre-determinato nello spazio di indirizzamento della modalità supervisore. Come per tutte le altre eccezioni, il contenuto del program counter associato all'istruzione `ecall` viene salvato nel program counter delle eccezioni del supervisore, SEPC (*Supervisor Exception Program Counter*), e il processore inizia a lavorare in modalità supervisore. Per ritornare alla modalità utente, al termine della risposta all'eccezione, viene utilizzata l'istruzione di *ritorno da eccezione del supervisore* (`sret`, *Supervisor exception RETurn*), che ripristina la modalità utente e continua l'esecuzione a partire dall'indirizzo memorizzato nel SEPC.

Mediante l'utilizzo di questi meccanismi e memorizzando la tabella delle pagine nello spazio di indirizzamento del sistema operativo, il sistema operativo può modificare le tabelle delle pagine e, allo stesso tempo, vietare ai processi utente di modificarle, garantendo quindi che un processo utente possa solo accedere allo spazio di memoria ad esso assegnato dal sistema operativo.

Occorre anche evitare che un processo legga i dati di un altro processo. Per esempio, non vogliamo che il programma di uno studente legga i voti di un esame quando questi si trovano nella memoria del processore. Se si permette di condividere la memoria principale, occorre fornire a ogni processo la possibilità di proteggere i propri dati dalla lettura e scrittura da parte di altri processi; in caso contrario, la condivisione della memoria si potrebbe rivelare una maledizione!

Ricordiamo che ciascun processo ha il proprio spazio di indirizzamento virtuale; perciò se il sistema operativo organizza le tabelle delle pagine in modo che pagine virtuali indipendenti corrispondano a pagine fisiche differenti, un processo non sarà in grado di accedere ai dati di un altro processo. Ovviamente, questo richiede anche che un processo utente non possa alterare la traduzione delle pagine virtuali contenuta nella tabella delle pagine. Il sistema operativo può garantire la protezione se può impedire ai processi utente di modificare la propria tabella delle pagine; tuttavia, il sistema operativo deve poter modificare le tabelle delle pagine. Inserendo le tabelle delle pagine nello spazio di indirizzamento protetto del sistema operativo, si possono soddisfare entrambi i requisiti.

Quando i processi vogliono condividere delle informazioni, devono essere assistiti dal sistema operativo, dato che accedere alle informazioni di un altro processo implica cambiare la tabella delle pagine per il processo che effettua l'accesso. Si può utilizzare il bit d'accesso in scrittura per limitare la condivisione alla sola lettura; tale bit, come il resto della tabella delle pagine, può essere modificato solo dal sistema operativo. Per consentire a un processo, per esem-

## Interfaccia hardware/software

**Modalità kernel:** detta anche **modalità supervisore**, è una modalità di funzionamento nella quale il processo in esecuzione è un processo del sistema operativo.

**Chiamata di sistema:** è un'istruzione speciale che trasferisce il controllo dalla modalità utente a un indirizzo predeterminato nello spazio dei processi di supervisione, sollevando un'eccezione.

pio P1, di leggere una pagina del processo P2, P2 dovrebbe chiedere al sistema operativo di creare un nuovo elemento della tabella delle pagine contenente una pagina virtuale nello spazio di indirizzamento di P1; questo nuovo elemento punterà alla pagina fisica che P2 vuole condividere. Il sistema operativo può usare il bit di protezione in scrittura per impedire a P1 di scrivere dei dati in questa pagina fisica, se questo è il desiderio di P2. Tutti i bit che determinano i diritti di accesso a una pagina devono essere contenuti sia nella tabella delle pagine sia nel TLB, dato che la tabella delle pagine viene consultata solo in corrispondenza delle miss del TLB.

**Cambio di contesto:** la modifica dello stato interno del processore che consente a un processo di subentrare a un altro processo nell'utilizzo del processore; comprende il salvataggio dello stato, necessario a ritornare in seguito al processo che era precedentemente in esecuzione.

**Approfondimento.** L'azione che compie il sistema operativo quando decide di cambiare il processo in esecuzione, da P1 a P2, è chiamata **cambio di contesto** (*context switch*), o anche **cambio di processo** (*process switch*). Il cambio di contesto deve garantire che P2 non possa accedere alla tabella delle pagine di P1, perché ciò potrebbe compromettere la protezione dei dati di P1. Se non c'è un TLB, è sufficiente modificare il registro della tabella delle pagine in modo che punti alla tabella delle pagine di P2, anziché a quella di P1. Se invece è presente un TLB, occorre invalidare gli elementi del TLB, i quali corrispondono a P1, sia per proteggere i dati di P1 sia per forzare il TLB a caricare gli elementi di P2. Se la frequenza con cui avvengono i cambi di contesto fosse elevata, questo meccanismo potrebbe rivelarsi facilmente inefficiente. Per esempio, P2 potrebbe fare in tempo a caricare solo pochi elementi del TLB prima che il sistema operativo restituisca di nuovo il controllo a P1; P1 a quel punto si troverebbe nella situazione in cui tutti gli elementi del TLB sono stati cancellati e dovrebbe ricaricarli di nuovo attraverso una sequenza di miss del TLB. Questo problema sorge perché gli indirizzi virtuali utilizzati da P1 e da P2 sono gli stessi e occorre, quindi, svuotare il TLB per evitare confusione tra essi.

Un'alternativa comunemente utilizzata è quella di estendere lo spazio di indirizzamento virtuale aggiungendo un *identificatore di processo*, o *identificatore di task*. Per esempio, il FastMATH Intrinsicity prevede a questo scopo un campo per l'identificazione dello spazio di indirizzamento su 8 bit, chiamato ASID (*Address Space Identifier*). Questo piccolo campo identifica il processo correntemente in esecuzione ed è contenuto in un registro apposito, che viene aggiornato dal sistema operativo ogni volta che cambia il processo in esecuzione. Il RISC-V offre anche un ASID per ridurre i flush del TLB quando si verifica una commutazione di contesto. L'identificatore di processo viene concatenato con il campo tag del TLB, in modo tale che una hit del TLB può verificarsi solo nel caso in cui sia il numero di pagina sia l'identificatore del processo corrispondono a quelli della pagina cercata. Questa soluzione permette di eliminare la necessità di cancellare gli elementi del TLB, che saranno cancellati solamente in rare occasioni. Possono nascere dei problemi di natura simile anche per la cache, poiché quando cambia il processo in esecuzione la cache contiene i dati del processo precedentemente in esecuzione. I problemi per le cache indirizzate fisicamente sono diversi da quelli delle cache indirizzate virtualmente e diverse soluzioni, come quella degli identificatori di processo, possono essere utilizzate per garantire che ciascun processo possa accedere ai propri dati.

## Gestione delle miss del TLB e dei page fault

Mentre la traduzione degli indirizzi virtuali in indirizzi fisici attraverso il TLB è semplice da capire quando l'accesso al TLB produce una hit, la gestione delle miss del TLB e dei page fault è più complessa. Una miss del TLB si verifica quando nessuno degli elementi del TLB corrisponde all'indirizzo virtuale e può essere dovuta a una delle due seguenti cause:

1. la pagina è presente in memoria e bisogna soltanto inserire nel TLB l'elemento corrispondente;
2. la pagina non è presente in memoria ed è quindi necessario trasferire il controllo al sistema operativo per gestire il page fault che viene generato.

Per gestire una miss del TLB o un page fault si utilizza il meccanismo della risposta alle eccezioni: viene interrotto il processo attivo e si trasferisce il controllo

al sistema operativo, per poi riprendere l'esecuzione del processo interrotto. Il page fault viene identificato all'interno del ciclo di clock in cui si accede alla memoria. Per riprendere l'esecuzione dell'istruzione dopo avere gestito l'errore di page fault, è necessario salvare il contenuto del program counter dell'istruzione che ha provocato il page fault. Il registro *SEPC*, *program counter delle eccezioni del supervisore*, viene utilizzato per mantenere questo valore.

Inoltre una miss del TLB, o un'eccezione di page fault, deve essere sollevata prima della fine del ciclo di clock in cui si verifica l'accesso alla memoria, in modo tale che il programma di risposta all'eccezione possa iniziare già nel ciclo di clock immediatamente successivo, invece di continuare la normale esecuzione delle istruzioni. Se il page fault fosse riconosciuto più tardi, un'eventuale istruzione di load potrebbe modificare il contenuto di un registro con conseguenze disastrose nel momento in cui il processo viene fatto ripartire dall'istruzione interrotta. Per esempio, si consideri l'istruzione  $1b \times 10, 0(x10)$ : il calcolatore deve evitare la scrittura nella pipeline del dato letto dalla memoria, altrimenti non sarebbe più possibile far ripartire correttamente il programma, poiché il contenuto originale del registro  $x10$  verrebbe distrutto. Un problema simile si verifica con le istruzioni di store: occorre evitare che la scrittura in memoria venga effettivamente completata quando si verifica un page fault; di solito ciò viene ottenuto impostando a 0 il segnale di controllo della scrittura della memoria.

Nel lasso di tempo che intercorre tra l'istante in cui il sistema operativo inizia l'esecuzione del programma di risposta alle eccezioni e quello in cui termina di salvare tutto lo stato del processo, il sistema operativo è particolarmente vulnerabile. Per esempio, se si dovesse verificare una seconda eccezione mentre il sistema operativo sta ancora gestendo la prima, l'unità di controllo sovrascriverebbe il program counter delle eccezioni rendendo impossibile il ritorno all'istruzione che ha causato l'errore di page fault. È possibile evitare questa situazione disastrosa fornendo un meccanismo di **disabilitazione e abilitazione delle eccezioni**. Quando si verifica la prima eccezione, il processore impone a 1 un bit che disabilita tutte le altre eccezioni: questo può avvenire nello stesso istante in cui il processore impone a 1 il bit relativo alla modalità supervisore. Il sistema operativo salverà quindi quella parte dello stato sufficiente a far ripartire l'esecuzione anche se dovesse arrivare un'altra eccezione, vale a dire il program counter delle eccezioni del supervisore (SEPC) e il registro *causa del supervisore* (SCAUSA), che, come abbiamo visto nel Capitolo 4, memorizza la causa dell'eccezione. Nel RISC-V, SEPC e SCAUSA sono due degli speciali registri di controllo che vengono utilizzati nella gestione delle eccezioni, delle miss del TLB e dei page fault. Al termine del salvataggio dello stato, il sistema operativo può riabilitare le eccezioni. Questo meccanismo fa in modo che le eccezioni non causino la perdita dello stato di nessun processo, cosa che renderebbe impossibile il riavvio dell'esecuzione a partire dall'istruzione interrotta.

## Interfaccia hardware/software

**Abilitazione delle eccezioni:** o abilitazione degli interrupt, è un segnale o un'azione che controlla se il processo debba o meno rispondere a un'eccezione; è necessario per evitare che si verifichino eccezioni prima che il processo abbia salvato in modo sicuro lo stato richiesto per riprendere l'esecuzione.

Dopo che il sistema operativo ha individuato l'indirizzo virtuale che ha provocato il page fault, deve eseguire i seguenti tre passi:

1. esaminare l'elemento della tabella delle pagine individuato tramite l'indirizzo virtuale e identificare la posizione nella memoria secondaria della pagina desiderata;
2. scegliere la pagina fisica da sostituire; se la pagina scelta è stata modificata, cioè il dirty bit è impostato a 1, la pagina deve essere salvata nella memoria secondaria prima di copiare una nuova pagina virtuale sulla pagina fisica;
3. avviare la lettura dalla memoria secondaria e trasferire la pagina richiesta alla memoria principale nella pagina fisica individuata.

Ovviamente quest'ultimo passo richiede milioni di cicli di clock del processore, come pure il secondo nel caso in cui la pagina da sostituire sia stata modificata; di conseguenza il sistema operativo, di solito, sceglie un altro processo da fare eseguire al processore finché l'accesso al disco non è stato completato. Dato che il sistema operativo ha salvato l'intero stato del processo corrente, può tranquillamente assegnare il controllo del processore a un altro processo.

Quando la lettura della pagina dalla memoria secondaria è stata completata, il sistema operativo può ripristinare lo stato del processo che aveva generato il page fault ed eseguire l'istruzione di ritorno dall'eccezione. Questa istruzione fa tornare il processore dalla modalità di funzionamento kernel alla modalità utente e ripristina il contenuto del program counter. A questo punto, il processo utente può rieseguire l'istruzione che aveva provocato il page fault, riuscendo ad accedere alla pagina desiderata, e può quindi continuare la sua esecuzione.

La risposta alle eccezioni di page fault causate dagli accessi ai dati è difficile da implementare correttamente in un processore, per la combinazione delle seguenti tre caratteristiche:

1. queste eccezioni sono generate durante l'esecuzione di un'istruzione, a differenza dei page fault sulle istruzioni;
2. l'istruzione non può essere completata prima di aver gestito l'eccezione che è stata generata;
3. dopo avere gestito l'eccezione, l'istruzione deve riprendere come se nulla fosse accaduto.

**Istruzione riavviabile:** un'istruzione che può riprendere l'esecuzione dopo che l'eccezione è stata risolta, senza che ciò influenzi il risultato dell'istruzione.

Fare in modo che un'**istruzione** sia **riavviabile** (*restartable instruction*) vuol dire poter continuare l'esecuzione dopo avere gestito l'eccezione; ciò è relativamente semplice in un'architettura come il RISC-V. Dato che ogni istruzione scrive solo un dato e che la scrittura del dato avviene alla fine del ciclo di esecuzione dell'istruzione, è sufficiente evitare che l'istruzione venga completata, impedendo la scrittura, per poi riprendere l'esecuzione dell'istruzione dall'inizio.

**Approfondimento.** Per i calcolatori dotati di istruzioni più complesse, che possono richiedere l'accesso a molte celle di memoria e scrivere molti dati, è molto più difficile fare in modo che le istruzioni siano riavviable. Una singola istruzione può generare diversi page fault nel corso della sua esecuzione. Per esempio, i calcolatori x86 sono dotati di istruzioni di spostamento di blocchi di memoria che spostano dati costituiti da migliaia di parole. In questi processori spesso non è possibile far ripartire un'istruzione dall'inizio, come si fa per le istruzioni RISC-V; l'istruzione deve essere interrotta e poi fatta ripartire dallo stadio di esecuzione nel quale era stata fermata. Riprendere l'esecuzione di un'istruzione interrotta a metà di solito richiede che venga salvato uno stato speciale, che venga risolta l'eccezione e che venga poi ripristinato questo stato speciale. Per fare in modo che tutto questo funzioni correttamente è necessario garantire un preciso coordinamento tra l'hardware e quella parte di codice del sistema operativo che si occupa della gestione delle eccezioni.

**Approfondimento.** Piuttosto che pagare il costo di un passo in più di indirizzamento a ogni accesso a memoria, il monitor della macchina virtuale (par. 5.6) mantiene una *tabella delle pagine ombra* (*shadow page table*) che mappa direttamente gli indirizzi dallo spazio virtuale dei sistemi ospitati negli indirizzi fisici dell'hardware. Intercettando tutte le modifiche alla tabella delle pagine dei sistemi ospitati, il VMM è in grado di assicurare che il contenuto della tabella delle pagine ombra utilizzato dall'hardware per la mappatura corrisponda al contenuto delle pagine del SO ospitato, con la differenza che il numero della pagina reale è sostituito dal numero di pagina fisica nelle tabelle delle pagine del SO ospitato. Quindi, il VMM deve sollevare un'eccezione di trap a ogni tentativo di un SO ospitato di modificare la sua tabella delle pagine o di accedere al puntatore alla tabella delle

pagine. Questo viene di solito ottenuto proteggendo dalla scrittura le tabelle delle pagine dei sistemi ospitati e impedendo ai SO ospitati l'accesso al puntatore alla tabella delle pagine. Come notato in precedenza, la seconda condizione è ottenuta naturalmente se l'accesso al puntatore della tabella delle pagine è un'operazione privilegiata.

**Approfondimento.** L'ultima parte dell'architettura da virtualizzare è il sistema di I/O. Questa è di gran lunga la parte di sistema più difficile da virtualizzare per il numero crescente di dispositivi di I/O che si possono collegare al calcolatore e per la loro diversità. Un'altra difficoltà è rappresentata dalla condivisione dei dispositivi reali tra più VM, e un'altra ancora proviene dalla necessità di supportare la miriade di driver dei dispositivi presenti, specialmente se i diversi SO ospitati devono essere supportati dalla stessa VM. L'illusione di una VM può essere mantenuta assegnando a ciascuna VM una generica versione del driver di ogni tipo di dispositivo di I/O, per poi lasciare al VMM il compito di gestire le reali operazioni di I/O.

**Approfondimento.** Oltre alla virtualizzazione dell'insieme di istruzioni, per una macchina virtuale un'altra sfida è rappresentata dalla virtualizzazione della memoria virtuale, dato che ciascun SO ospitato in ogni VM gestisce le proprie tabelle delle pagine. Per fare funzionare il tutto, il VMM distingue tra *memoria reale* e *memoria fisica*, anche se queste sono spesso trattate come se fossero la stessa cosa, e rende la memoria reale un livello intermedio separato tra la memoria virtuale e la memoria fisica; alcuni chiamano questi tre livelli: *memoria virtuale*, *memoria fisica* e *memoria macchina*. Il SO ospitato mappa la memoria virtuale sulla memoria reale attraverso le proprie tabelle delle pagine e il VMM mappa la memoria reale del sistema ospitato sulla memoria fisica attraverso la tabella delle pagine del VMM. L'architettura della memoria virtuale può essere definita attraverso la tabella delle pagine, come nell'IBM VM/370, nell'x86 e nel RISC-V.

## Riepilogo

“Memoria virtuale” è il nome che viene dato al livello della gerarchia delle memorie che si occupa del trasferimento dei dati tra la memoria principale e il disco. La memoria virtuale permette ai singoli programmi di espandere il proprio spazio di indirizzamento oltre i limiti della memoria principale. Caratteristica ancora più importante è il fatto che la memoria virtuale supporta la condivisione della memoria principale tra più processi attivi simultaneamente, in modo protetto.

La gestione della coppia di memorie costituita da memoria principale e disco, all'interno della gerarchia delle memorie è molto complessa per l'alto costo dei page fault. Per ridurre la frequenza delle miss sono state introdotte alcune tecniche:

1. vengono definite pagine di grandi dimensioni per trarre vantaggio dalla località spaziale e per ridurre così la frequenza delle miss;
2. la mappatura tra indirizzi virtuali e indirizzi fisici, contenuta nella tabella delle pagine, viene resa completamente associativa, in modo che una pagina virtuale possa essere scritta in una pagina qualsiasi della memoria principale;
3. il sistema operativo utilizza tecniche come LRU e il bit di utilizzo per scegliere quale pagina sostituire.

Le scritture su disco sono costose; per questo motivo la memoria virtuale utilizza uno schema write-back e tiene anche traccia delle pagine che non sono state modificate, attraverso il dirty bit, per evitare di salvarle su disco.

Il meccanismo della memoria virtuale fornisce la traduzione degli indirizzi dallo spazio di indirizzamento virtuale utilizzato dal programma allo spazio degli indirizzi fisici utilizzati per accedere alla memoria principale. Tale traduzione permette la condivisione protetta della memoria principale e offre altri vantaggi, tra cui la semplificazione dell'allocazione della memoria. Per

garantire che i processi siano reciprocamente protetti, è necessario che solo il sistema operativo possa modificare la traduzione degli indirizzi. Ciò viene ottenuto impedendo ai programmi utente di modificare le tabelle delle pagine. Si può realizzare una condivisione controllata delle pagine tra diversi processi con l'aiuto del sistema operativo e dei bit di accesso contenuti nella tabella delle pagine, i quali indicano se il programma utente abbia diritto di lettura o scrittura su una pagina.

Se il processore dovesse accedere a una tabella delle pagine residente in memoria per tradurre ogni accesso, il carico di lavoro introdotto dalla memoria virtuale sarebbe troppo elevato e l'utilizzo delle cache risulterebbe inutile. Per questo motivo si utilizza un TLB, che agisce come una cache e memorizza alcune delle mappature contenute nella tabella delle pagine. Gli indirizzi vengono quindi tradotti da virtuali a fisici dal TLB.

Le cache, la memoria virtuale e il TLB sono caratterizzati dallo stesso insieme di principi e modalità di gestione; nel prossimo paragrafo descriveremo tale visione unificata.

## Capire le prestazioni dei programmi

Sebbene la memoria virtuale sia stata inventata per consentire a una memoria di piccole dimensioni di comportarsi come una memoria più grande, la differenza di prestazioni tra disco e memoria implica che, se un programma accede spesso alla memoria virtuale invece che alla memoria fisica, verrà eseguito molto lentamente. Tale programma scambierà continuamente le pagine tra memoria e disco, generando un fenomeno che prende il nome di *thrashing*. Il thrashing comporta prestazioni disastrose, ma per fortuna si verifica molto raramente. Se un programma incorre nel thrashing, la soluzione più semplice è eseguirlo su un calcolatore dotato di molta memoria oppure acquistare della memoria aggiuntiva. Una soluzione più complessa consiste nel riesaminare l'algoritmo e le strutture dati del programma per vedere se è possibile modificare la definizione dei dati e delle funzioni con lo scopo di aumentare la località e ridurre quindi il numero di pagine che vengono utilizzate simultaneamente dal programma. Questo insieme di pagine viene informalmente chiamato *insieme di lavoro (working set)*.

Un problema più comune che affligge le prestazioni è quello delle miss del TLB. Dato che il TLB può gestire solo 32-64 pagine alla volta, un programma può facilmente avere una frequenza di miss del TLB elevata, dato che il processore può accedere direttamente al massimo a un quarto di mebibyte ( $64 \times 4 \text{ KiB} = 0,25 \text{ MiB}$ ). Per esempio, le miss nel TLB costituiscono un grosso problema per l'algoritmo Radix Sort. Per cercare di risolverlo, molte architetture di calcolatori supportano una dimensione della pagina variabile. Per esempio, in aggiunta alla dimensione minima di pagina di 4 KiB, l'hardware RISC-V supporta pagine di dimensioni di 2 MiB e 1 GiB. Quindi, se un programma utilizza una dimensione di pagina più grande, può accedere direttamente a una porzione di memoria maggiore senza incorrere in una miss del TLB.

La sfida, dal punto di vista pratico, è fare in modo che il sistema operativo riesca a far selezionare ai programmi pagine di dimensioni variabili. Ancora una volta, la soluzione più complessa per ridurre le miss del TLB consiste nel riesaminare l'algoritmo e le strutture dati del programma, con lo scopo di ridurre il numero di pagine che costituiscono l'insieme di lavoro. Data l'importanza delle prestazioni degli accessi a memoria e della frequenza delle miss del TLB, alcuni programmi caratterizzati da insiemi di lavoro di grandi dimensioni sono stati ristrutturati per raggiungere questo obiettivo.

**Approfondimento.** Il RISC-V supporta pagine di dimensioni più ampie attraverso la tabella delle pagine multilivello di Figura 5.29. Inoltre, per puntare alla tabella delle pagine successiva nei livelli 1 e 2, utilizza una traduzione "superpage" per mappare l'indirizzo virtuale su un indirizzo fisico di 1 GiB (se la traduzione del blocco è al livello 1) o su un indirizzo fisico di 2 MiB (se la traduzione è al livello 2).

## 5.8 Schema comune per le gerarchie delle memorie

Finora abbiamo visto come i diversi tipi di gerarchie delle memorie abbiano molte caratteristiche in comune. Sebbene molti aspetti differiscano dal punto di vista quantitativo, molte delle caratteristiche che definiscono il funzionamento della gerarchia sono simili da un punto di vista qualitativo. La Figura 5.34 mostra gli intervalli di variazione di alcuni parametri caratteristici delle gerarchie. In questo paragrafo esamineremo le soluzioni operative più comuni adottate per le gerarchie delle memorie e come queste ne determinino il comportamento. Esamineremo queste strategie attraverso quattro domande che possono essere formulate per ogni coppia di livelli di una gerarchia delle memorie, anche se, per semplicità, adotteremo principalmente la terminologia utilizzata per le cache.

### Domanda 1: dove può essere posizionato un blocco?

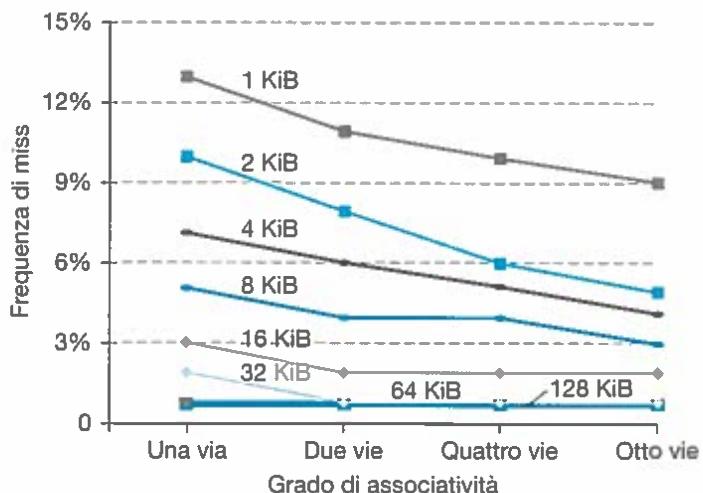
Abbiamo visto che si può utilizzare una vasta gamma di schemi per identificare un blocco nel livello superiore della gerarchia: dalla mappatura diretta a quella set-associativa, a quella completamente associativa. Come si è detto in precedenza, i diversi schemi possono essere visti come una variazione dello schema set-associativo, nel quale il numero di linee e il numero di blocchi contenuti in una linea variano:

Nome dello schema	Numero di linee	Blocchi per linea
Mappatura diretta	Numero di blocchi della cache	1
Set-associativo	Numero di blocchi della cache Grado di associatività	Grado di associatività (tipicamente 2-16)
Completamente associativo	1	Numero di blocchi della cache

Il vantaggio che deriva dall'aumento del grado di associatività è che generalmente diminuisce la frequenza delle miss, perché si riduce il numero di miss dovute agli accessi che competono per la stessa locazione. Esamineremo questo aspetto in maggior dettaglio tra breve. Ci focalizziamo ora sul guadagno dovuto all'aumento di associatività: la Figura 5.35 mostra la frequenza di miss al variare del grado di associatività, da quello di una cache a mappatura diretta a quello di una cache set-associativa a otto vie. Il guadagno maggiore si ottiene passando da una cache a mappatura diretta a una cache set-associativa a due vie: questo passaggio produce una riduzione della frequenza di miss compresa tra il 20 e il 30%. Al crescere della dimensione della cache, il miglioramento relativo

Caratteristica	Valori tipici per una cache L1	Valori tipici per una cache L2	Valori tipici per una memoria a pagine	Valori tipici per un TLB
Dimensione totale in blocchi	250-2000	2500-25 000	16 000-250 000	40-1024
Dimensione totale in kilobyte	16-64	125-2000	1 000 000-1 000 000 000	0,25-16
Dimensione del blocco in byte	16-64	64-128	4000-64 000	4-32
Penalità di miss in cicli di clock	10-25	100-1000	10 000 000-100 000 000	10-1000
Frequenza di miss (globale per L2)	2-5%	0,1-2%	0,00001-0,0001%	0,01-2%

**Figura 5.34** Intervallo dei parametri caratterizzanti i principali livelli di memoria di una gerarchia delle memorie di un calcolatore. Questi sono i tipici valori dei parametri riferiti all'anno 2012. Si noti come il campo di variazione dei parametri sia ampio; ciò è dovuto in parte al fatto che molti dei parametri, che sono cambiati nel tempo, sono tra loro correlati; per esempio, al crescere delle dimensioni delle cache, per evitare penalità di miss più alte, anche la dimensione dei blocchi deve aumentare. Anche se non mostrato in tabella, i microprocessori utilizzati nei server hanno anche una cache L3, di dimensioni variabili tra 2 e 8 MiB, che può contenere molti più blocchi della cache L1 o L2. Le cache L3 abbassano la penalità di miss a 30-40 cicli di clock.



**Figura 5.35** La frequenza delle miss sui dati di una cache per ognuna delle otto diverse dimensioni considerate diminuisce al crescere del grado di associatività. Mentre il beneficio ottenuto passando da una cache set-associativa a una via (a mappatura diretta) a una cache a due vie è significativo, il miglioramento ottenuto con un incremento ulteriore del grado di associatività è ridotto. Per esempio, il miglioramento nel passaggio da due a quattro vie è compreso tra l'1% e il 10%, contro il 20-30% di miglioramento nel passaggio da una a due vie. Il miglioramento è ancora inferiore passando da quattro a otto vie, dove ci si avvicina molto alla frequenza di miss di una cache completamente associativa. Cache di dimensioni ridotte ottengono un miglioramento della frequenza di miss assoluta maggiore all'aumentare del grado di associatività, perché la frequenza di base delle miss di una cache piccola è più elevata. Nella didascalia di Figura 5.16 si spiega in che modo sono stati raccolti questi dati.

aumenta di poco con l'aumento del grado di associatività della cache: dato che è minore la frequenza globale delle miss in una cache di dimensioni maggiori, l'opportunità di migliorare ulteriormente la frequenza di miss diminuisce e il miglioramento assoluto della frequenza di miss dovuto al grado di associatività si riduce significativamente. I potenziali svantaggi dell'aumento di associatività, come già menzionato, sono l'incremento del costo e dei tempi di accesso.

### Domanda 2: come si individua un blocco?

La scelta di come individuare un blocco dipende dallo schema utilizzato per la disposizione dei blocchi, poiché esso definisce il numero di possibili locazioni. Possiamo riassumere gli schemi come segue:

Associatività	Metodo di ricerca	Numero di confronti richiesti
Mappatura diretta	Indice	1
Set-associativa	Indice per la linea, ricerca tra gli elementi della linea	Grado di associatività
Completamente associativa	Ricerca in tutti gli elementi all'interno della cache	Dimensione della cache
	Tramite una tabella di lookup separata	0

In ogni gerarchia delle memorie la scelta tra mappatura diretta, set-associativa o completamente associativa dipende dal confronto tra il costo delle miss e il costo dell'implementazione dell'associatività, sia in termini di tempo sia in termini di hardware aggiuntivo. Includere la cache L2 all'interno del chip del processore permette di utilizzare un grado di associatività molto maggiore, poiché il tempo di hit non è più così critico e il progettista non deve più utilizzare i chip standard di SRAM come blocchi di base. Le cache completamente

associative sono proibitive eccetto che per piccole dimensioni: in questo caso, il costo dei comparatori non è così determinante e il miglioramento della frequenza assoluta di miss è massimo.

Nei sistemi di memoria virtuale si utilizza una tabella separata, detta tabella delle pagine, per indicizzare la memoria. Oltre allo spazio richiesto per memorizzarla, la tabella delle pagine richiede un doppio accesso alla memoria. La scelta di associatività totale per la tabella delle pagine e l'utilizzo di una tabella aggiuntiva sono motivati da questi tre aspetti:

1. l'associatività totale è un beneficio, essendo le miss molto dispendiose;
2. l'associatività totale permette al software di utilizzare schemi di sostituzione molto raffinati, progettati al fine di ridurre la frequenza di miss;
3. l'intera tabella può essere facilmente indicizzata senza hardware aggiuntivo e non serve implementare alcuna ricerca al suo interno.

Per questi motivi, i sistemi di memoria virtuale utilizzano di solito una disposizione delle pagine completamente associativa.

La disposizione set-associativa viene spesso utilizzata per le cache e per i TLB, nei quali l'accesso richiede l'indicizzazione unita alla ricerca all'interno di un piccolo insieme. Alcuni sistemi impiegano cache a mappatura diretta per i vantaggi che ne derivano in termini di tempo d'accesso e semplicità. Il vantaggio relativo al tempo d'accesso è dovuto al fatto che per trovare il blocco richiesto non c'è bisogno di alcun confronto. Tali scelte di progetto dipendono da molti dettagli dell'implementazione hardware, per esempio dal fatto che la cache sia contenuta nel chip del processore o meno, dalla tecnologia utilizzata per implementare la cache e dal ruolo critico che ha il tempo d'accesso alla cache nel determinare il periodo di clock del processore.

### **Domanda 3: quale blocco deve essere sostituito in caso di miss della cache?**

Quando si verifica una miss in una cache associativa, bisogna decidere quale blocco sostituire. In una cache completamente associativa ogni blocco è un potenziale candidato alla sostituzione, mentre se la cache è set-associativa bisogna scegliere tra i blocchi contenuti sulla stessa linea. Ovviamente la scelta diventa facile in una cache a mappatura diretta, perché c'è solo un candidato.

Sono due le principali strategie utilizzate per la scelta del blocco da sostituire nelle cache set-associative o completamente associative:

- *casuale*: il blocco è scelto a caso, eventualmente utilizzando il supporto di alcuni componenti hardware.
- *blocco utilizzato meno di recente* (LRU): il blocco prescelto è quello che è rimasto inutilizzato da più tempo.

Nella pratica la politica LRU è troppo costosa da implementare per le memorie che hanno un grado di associatività non basso (tipicamente da due a quattro), poiché è costoso tenere traccia dell'utilizzo dei diversi blocchi. Anche negli schemi set-associativi a quattro vie la strategia LRU viene spesso approssimata, per esempio suddividendo le linee in due coppie di blocchi, tenendo traccia della coppia che è stata utilizzata meno di recente (serve un solo bit) e memorizzando poi il blocco nella coppia che è stato a sua volta utilizzato meno di recente (anche questa operazione richiede un solo bit per ogni coppia).

Quando il grado di associatività cresce, la strategia LRU viene sempre approssimata, oppure si preferisce una scelta casuale del blocco. Nelle cache l'algoritmo di sostituzione è implementato in hardware, il che significa che l'algoritmo deve

poter essere realizzato in maniera semplice. La sostituzione casuale è facile da realizzare in hardware, e in una cache set-associativa a due vie la sostituzione casuale produce una frequenza di miss pari a circa 1,1 volte quella ottenuta con la strategia LRU. Al crescere delle dimensioni della cache, la frequenza delle miss diminuisce con entrambi gli approcci e la differenza in termini assoluti diventa minima; la sostituzione casuale può certe volte dare risultati migliori delle semplici approssimazioni della politica LRU che possono essere facilmente implementate in hardware.

Per la memoria virtuale si adotta sempre una politica di tipo LRU approssimata, dato che anche una minima riduzione della frequenza di miss può essere importante (essendo così alto il costo di una miss). Il bit di utilizzo o modalità equivalenti di verifica dell'uso di una pagina vengono spesso impiegati per aiutare il sistema operativo a tenere traccia delle pagine richieste meno di recente. Poiché le miss sono così costose e relativamente infrequenti, approssimare queste informazioni principalmente via software è considerato accettabile.

#### Domanda 4: come vengono gestite le scritture?

Una caratteristica fondamentale di tutte le gerarchie delle memorie è la gestione della scrittura. Abbiamo già visto le due opzioni di base:

- *write-through*: l'informazione viene scritta sia nel blocco della cache sia nel blocco del livello inferiore della gerarchia delle memorie, che per la cache è la memoria principale (le cache presentate nel paragrafo 5.3 adottavano questo schema);
- *write-back*: l'informazione viene scritta solo nel blocco della cache; i blocchi che sono stati modificati vengono poi scritti nel livello inferiore della gerarchia solo quando devono essere sostituiti. I sistemi di memoria virtuale utilizzano sempre lo schema write-back per le ragioni illustrate nel paragrafo 5.7.

Entrambi gli schemi hanno vantaggi e svantaggi. I vantaggi principali della modalità write-back sono i seguenti:

- le singole parole possono essere scritte dal processore alla frequenza con cui la cache, e non la memoria principale, è in grado di scriverle;
- scritture multiple all'interno dello stesso blocco richiedono una sola scrittura nel livello inferiore della gerarchia;
- quando i blocchi vengono scritti nel livello inferiore della gerarchia, il sistema può sfruttare appieno la banda molto larga del trasferimento, dato che viene trasferito un intero blocco.

I vantaggi principali della modalità write-through sono i seguenti:

- le miss sono più semplici da gestire e meno costose, perché non richiedono mai la scrittura di un blocco intero nel livello inferiore;
- lo schema write-through è più facile da implementare dello schema di write-back anche se, per essere efficace, una cache con write-through deve essere dotata anche di un buffer di scrittura.

Nei sistemi di memoria virtuale l'unica modalità praticabile è quella di write-back a causa della lunga latenza richiesta dall'operazione di scrittura nel livello inferiore (disco). La frequenza con cui vengono generate le scritture da parte del processore, in genere, supera la frequenza con cui il sistema di memoria può gestirle, anche quando la memoria è più larga sia nella dimensione fisica che in quella logica, e viene utilizzata la modalità di trasferimento a raffica. Di conseguenza, attualmente le cache del livello inferiore della gerarchia utilizzano solitamente la strategia write-back.

## QUADRO D'INSIEME

Anche se cache, TLB e memoria virtuale possono inizialmente apparire fra loro estremamente diversi, si basano sui medesimi due principi di località. Il loro funzionamento si può comprendere rispondendo alle seguenti quattro domande:

**Domanda 1:** Dove può essere posizionato un blocco?

**Risposta:** In una sola posizione (mappatura diretta), in poche posizioni (set-associatività) o in una qualsiasi posizione (associatività completa).

**Domanda 2:** Come si individua un blocco?

**Risposta:** Sono disponibili quattro metodi: con indice (come in una cache a mappatura diretta), con una ricerca limitata (come in una cache set-associativa), con una ricerca completa (come in una cache completamente associativa) o tramite una tabella di lookup separata (come nella tabella delle pagine).

**Domanda 3:** Quale blocco deve essere sostituito in caso di miss?

**Risposta:** Di solito il blocco utilizzato meno di recente oppure un blocco scelto a caso.

**Domanda 4:** Come vengono gestite le scritture?

**Risposta:** Ogni livello della gerarchia può utilizzare lo schema write-through oppure lo schema write-back. ■

## Le tre C: un modello intuitivo per comprendere il comportamento delle gerarchie delle memorie

In questo paragrafo esamineremo un modello che aiuta a comprendere la causa delle miss in una gerarchia delle memorie e le conseguenze sulle miss di eventuali modifiche alla gerarchia. Spiegheremo i diversi concetti facendo riferimento alle cache, ma gli stessi concetti si possono estendere a tutti gli altri livelli della gerarchia. In questo modello tutte le miss vengono classificate in base alle tre seguenti categorie (**modello delle tre C**):

- **miss obbligate** (*compulsory misses*): queste miss sono causate dal primo accesso a un blocco della cache che non era mai stato caricato in precedenza; sono dette anche **miss da partenza a freddo** (*cold start misses*);
- **miss di capacità** (*capacity misses*): queste miss sono generate quando la cache non è in grado di contenere tutti i blocchi di memoria necessari durante l'esecuzione di un programma; si verificano quando un blocco viene sostituito per essere ricaricato più tardi;
- **miss di conflitto** (*conflict misses*): queste miss sono generate nelle cache a mappatura diretta o set-associative quando più blocchi competono per la stessa linea; possono essere eliminate utilizzando una cache completamente associativa della stessa dimensione. Queste miss sono anche dette **miss di collisione** (*collision misses*).

La Figura 5.36 mostra che le miss possono essere suddivise in base a tre cause principali; l'incidenza di ciascuna causa può essere ridotta modificando alcuni aspetti del progetto della cache. Siccome le miss di conflitto nascono direttamente dalla competizione per lo stesso blocco della cache, incrementare il grado

**Modello delle tre C:** un modello di cache in cui tutte le miss vengono classificate in base alle seguenti tre categorie: miss obbligate, miss di capacità e miss di conflitto.

**Miss obbligate:** dette anche **miss da partenza a freddo**, sono le miss della cache causate dal primo accesso a un blocco che non era mai stato caricato in cache.

**Miss di capacità:** sono le miss che si verificano perché la cache, anche se fosse completamente associativa, non potrebbe contenere tutti i blocchi necessari a soddisfare le richieste.

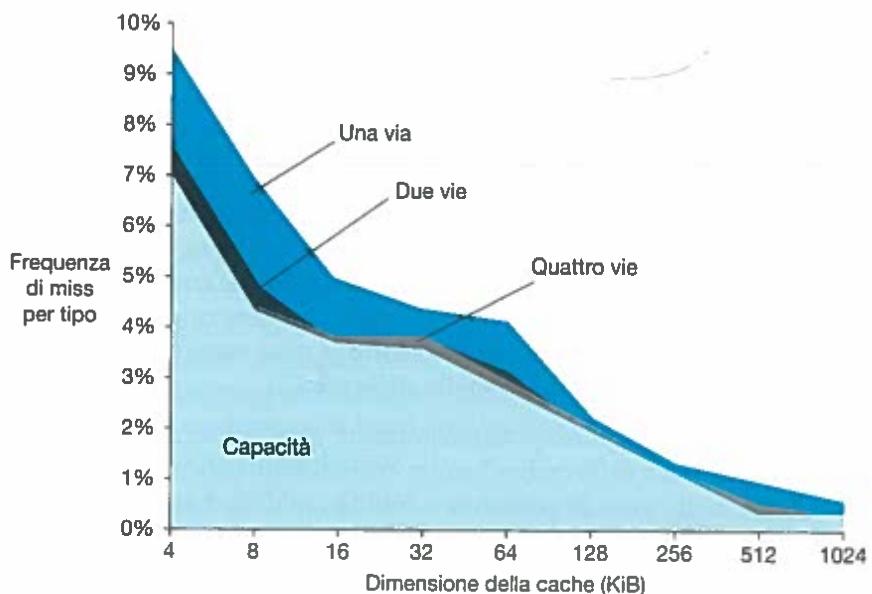
**Miss di conflitto o miss di collisione:** indicano le miss che avvengono in una cache set-associativa o a mappatura diretta quando più blocchi competono per la stessa linea. Le miss di conflitto possono essere eliminate utilizzando una cache completamente associativa della stessa dimensione.

di associatività ne riduce l'incidenza. Tuttavia, l'aumento dell'associatività può causare un aumento del tempo di accesso, provocando una diminuzione delle prestazioni complessive.

Le miss di capacità si possono facilmente ridurre aumentando la dimensione della cache; infatti, la capacità delle cache di secondo livello negli ultimi anni è cresciuta con un tasso costante. Ovviamente, quando si progettano cache più grandi, occorre anche prestare attenzione all'aumento del tempo d'accesso, che può provocare una diminuzione delle prestazioni globali; per questo motivo la dimensione delle cache di primo livello è cresciuta con un tasso molto inferiore, o addirittura è rimasta invariata.

Dato che le miss obbligate sono provocate dal primo accesso a un blocco di memoria, il sistema migliore che si può adottare per ridurre questo tipo di miss consiste nell'aumentare la dimensione del blocco. Questo accorgimento permette di ridurre il numero di accessi ai blocchi di memoria che contengono un programma, visto che il programma sarà formato da un minor numero di blocchi. Come spiegato in precedenza, però, aumentare troppo la dimensione dei blocchi può avere un effetto negativo sulle prestazioni, perché determina anche un aumento della penalità di miss.

La suddivisione delle miss secondo il modello delle tre C è utile dal punto di vista qualitativo. Nella progettazione delle cache reali, molte delle scelte progettuali interagiscono tra loro e la modifica di una caratteristica della cache influenza spesso le altre caratteristiche, che a loro volta hanno conseguenze sulla frequenza delle miss. Nonostante questi limiti, il modello è utile per prevedere le prestazioni che potrà avere una cache.



**Figura 5.36** La frequenza delle miss può essere scomposta in tre diverse componenti. Questo grafico mostra la frequenza complessiva delle miss e le sue componenti in funzione della dimensione della cache. Questi dati sono relativi ai benchmark SPEC CPU2000 per numeri interi e in virgola mobile, e provengono dalla stessa fonte dei dati di Figura 5.35. Le miss obbligate sono pari allo 0,006% e non sono visibili in questo grafico. La seconda componente è rappresentata dalle miss di capacità, le quali dipendono dalla dimensione della cache. La frequenza delle miss di conflitto, che dipende sia dal grado di associatività sia dalle dimensioni della cache, viene mostrata per un grado di associatività che varia da una a otto vie: la parte del grafico etichettata corrisponde all'incremento della frequenza di miss quando il grado di associatività passa dal livello superiore al livello indicato dalla scritta. Per esempio, la parte del grafico più scura (*due vie*) indica la frequenza di miss in più che si ottiene quando la cache passa da un grado di associatività pari a quattro a un grado di associatività pari a due. Quindi, la differenza della frequenza di miss di una cache a mappatura diretta rispetto a quella di una cache completamente associativa della stessa dimensione è data dalla somma delle aree indicate come *quattro vie*, *due vie* e *una via*. La differenza tra otto vie e quattro vie è così piccola che è difficilmente visibile nel grafico.

Cambiamento nel progetto	Effetto sulla frequenza di miss	Eventuale effetto negativo sulle prestazioni
Aumento della dimensione della cache	Diminuiscono le miss di capacità	Può aumentare il tempo di accesso
Aumento del grado di associatività	Diminuisce la frequenza delle miss causate da conflitti	Può aumentare il tempo di accesso
Aumento della dimensione del blocco	Diminuisce la frequenza delle miss per un'ampia gamma di dimensioni dei blocchi a causa della località spaziale	Aumenta la penalità di miss. Blocchi molto grandi potrebbero aumentare la frequenza di miss

Figura 5.37 Problemi da affrontare nella progettazione delle memorie.

## QUADRO D'INSIEME

La sfida nella progettazione di una gerarchia delle memorie consiste nel fatto che ogni modifica che può migliorare la frequenza di miss può anche influenzare in maniera negativa le prestazioni globali, come è qui riassunto in Figura 5.37. Questa combinazione di effetti positivi e negativi rende interessante la progettazione di una gerarchia delle memorie. ■

### Autovalutazione

Quali delle seguenti affermazioni (se ce ne sono) sono generalmente vere?

1. Non c'è modo di ridurre le miss obbligate.
2. Le cache completamente associative non hanno miss di conflitto.
3. Il grado di associatività è più importante della capacità se si vuole ridurre il numero di miss.

## 5.9 | Come utilizzare una macchina a stati finiti per controllare una cache semplificata

Possiamo implementare ora l'unità di controllo di una cache, esattamente come abbiamo fatto nel Capitolo 4 per l'unità di controllo dell'unità di elaborazione a singolo ciclo di quella con pipeline. Questo paragrafo inizia con la definizione di una cache semplificata, per poi descrivere le macchine a stati finiti (FSM, *Finite-State Machines*) e terminare con la sintesi della FSM che controlla questa cache. Questo argomento viene approfondito nel paragrafo 5.12 , in cui viene mostrata la cache e la sua unità di controllo in un nuovo linguaggio di descrizione dell'hardware.

### Una cache semplificata

Ci accingiamo a progettare l'unità di controllo di una cache semplificata che ha le seguenti caratteristiche:

- cache a mappatura diretta;
- write-back con modalità write-allocate;
- dimensione del blocco di 4 parole (16 byte, ossia 128 bit);

- dimensione della cache di 16 KiB; la cache contiene quindi 1024 blocchi;
- indirizzi su 32 bit;
- presenza del bit di validità e del dirty bit per ogni blocco.

Da quanto visto nel paragrafo 5.3, possiamo ora calcolare la dimensione dei campi che compongono l'indirizzo della cache:

- l'indice della cache è di 10 bit;
- l'offset di blocco è di 4 bit;
- la dimensione del tag è quindi di  $32 - (10 + 4) = 18$  bit.

I segnali tra il processore e la cache sono:

- segnale di lettura e scrittura su 1 bit ciascuno;
- segnale di validità su 1 bit: indica se si sta eseguendo un'operazione sulla cache o meno;
- indirizzi su 32 bit;
- dati trasferiti fra il processore e la cache su 32 bit;
- dati trasferiti fra la cache e il processore su 32 bit;
- segnale di cache ready (pronto) su 1 bit: segnala quando la cache ha terminato l'attività.

L'interfaccia tra memoria principale e cache ha gli stessi campi dell'interfaccia esistente tra cache e processore, fatta eccezione per il campo dati che qui è largo 128 bit. Questa maggiore larghezza della memoria si trova facilmente negli attuali microprocessori, che utilizzano parole di 32 o 64 bit nel processore, mentre i controllori delle DRAM gestiscono spesso dati su 128 bit. Rendere il blocco della cache largo quanto la larghezza della DRAM semplifica il progetto. Questi sono i segnali utilizzati tra cache e memoria principale:

- segnale di lettura e scrittura su 1 bit ciascuno;
- segnale di validità su 1 bit: indica se si sta eseguendo un'operazione sulla memoria o meno;
- indirizzi su 32 bit;
- dati trasferiti da cache a memoria su 128 bit;
- dati trasferiti da memoria a cache su 128 bit;
- segnale di cache ready (cache pronta) su 1 bit: indica se la memoria ha terminato il suo lavoro.

Si noti che l'interfaccia con la memoria non richiede un numero fisso di cicli di clock, ma si suppone che il controllore della memoria notifichi alla cache tramite il segnale di ready quando è terminata la lettura o la scrittura della memoria.

Prima di descrivere il controllore della cache, dobbiamo rivedere i concetti su cui si basano le macchine a stati finiti, che permettono di controllare operazioni che possono durare più cicli di clock.

## Macchine a stati finiti

**Macchina a stati finiti:** una macchina astratta che implementa una logica sequenziale. Essa è definita da un insieme di ingressi, di uscite e di stati, da una funzione stato prossimo che mappa gli ingressi e lo stato corrente in un nuovo stato e da una funzione di uscita che mappa lo stato corrente (ed eventualmente gli ingressi) in un insieme di segnali in uscita.

Per costruire l'unità di controllo dell'unità di elaborazione a singolo ciclo, abbiamo utilizzato una serie di tabelle della verità che specificano il valore che dovevano assumere i segnali di controllo in funzione del tipo di istruzione. Nel caso di una cache l'unità di controllo è più complessa, perché le operazioni di lettura e scrittura sono costituite da una successione di passi elementari. L'unità di controllo della cache deve sia impostare i segnali di controllo che vengono utilizzati a ogni passo sia determinare quale sarà il passo successivo.

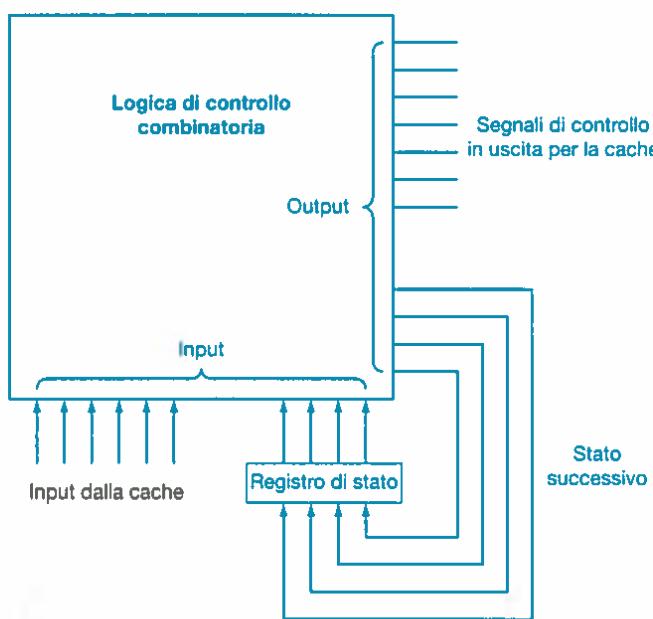
Il metodo di controllo più diffuso per operazioni che si sviluppano su più passi è basato sulle **macchine a stati finiti**. Queste vengono di solito rappresentate mediante grafi: per definire una macchina a stati finiti occorre specificare un

insieme di stati e per ogni stato bisogna dire quale sarà lo stato successivo. La scelta dello stato successivo viene effettuata da una particolare funzione, detta **funzione stato prossimo**, che mappa lo stato corrente e gli ingressi in quello successivo. Quando si utilizza una macchina a stati finiti per realizzare l'unità di controllo, occorre specificare anche l'insieme delle uscite che devono essere asserite in corrispondenza di ciascuno stato. L'implementazione circuitale di una macchina a stati finiti viene realizzata supponendo che quando un'uscita non viene esplicitamente asserita deve essere deasserita. In modo analogo, il funzionamento corretto dell'unità di elaborazione implica che un segnale non esplicitamente asserito deve essere deasserito.

Il segnale di controllo di un multiplexer funziona in modo leggermente diverso, poiché esso seleziona uno degli ingressi a seconda che assuma il valore 0 o 1. Perciò, in una macchina a stati finiti dobbiamo sempre impostare il segnale di controllo dei multiplexer che ci servono. Quando realizziamo una macchina a stati finiti con la logica digitale, impostare un segnale di controllo a 0 può essere la scelta di default e può quindi non richiedere porte logiche. Un semplice esempio di macchina a stati finiti è riportato nell'Appendice A; se non avete familiarità con i concetti su cui si basano le macchine a stati finiti, potete leggere l'Appendice A prima di proseguire con la lettura di questo paragrafo.

Una macchina a stati finiti può essere implementata con un registro temporaneo che mantiene lo stato corrente e un blocco di logica combinatoria che determina sia i segnali dell'unità di elaborazione che devono essere asseriti sia lo stato successivo. La Figura 5.38 mostra lo schema di una macchina a stati finiti e l'Appendice B descrive in dettaglio come si possa realizzare una macchina a stati finiti siffatta. Nel paragrafo A.3 la logica di controllo combinatoria di una macchina a stati finiti viene implementata utilizzando una ROM (*Read-Only Memory*, memoria a sola lettura) o una PLA (*Programmable Logic Array*, matrice logica programmabile). Vedi l'Appendice A anche per una descrizione di questi componenti logici.

**Funzione stato prossimo:** una funzione di logica combinatoria che determina lo stato successivo della macchina a partire dagli ingressi e dallo stato corrente.



**Figura 5.38** I controllori basati su macchine a stati finiti vengono tipicamente implementati utilizzando un blocco di logica combinatoria e un registro per salvare lo stato corrente (modello di Huffman). Le uscite della logica combinatoria sono lo stato successivo e i segnali di controllo che devono essere asseriti per ogni stato. Gli ingressi alla logica combinatoria sono lo stato corrente e tutti gli ingressi necessari per determinare lo stato successivo. Si noti che nella macchina a stati finiti utilizzata in questo capitolo le uscite dipendono solamente dallo stato corrente, non dagli ingressi. La sezione *Approfondimento* spiega la macchina a stati finiti più in dettaglio.

**Approfondimento.** Si noti che la cache associata alla Figura 5.39 è una *cache bloccante*, cioè il processore deve attendere che la cache termini la sua richiesta. Il paragrafo 5.12  descrive l'alternativa, rappresentata dalle *cache non bloccanti*.

**Approfondimento.** Il tipo di macchina a stati finiti utilizzato in questo libro è chiamato macchina di Moore, dal nome del suo inventore, Edward Moore. La sua caratteristica distintiva è che le uscite dipendono solamente dallo stato corrente. In una macchina di Moore il riquadro etichettato "Logica di controllo combinatoria" può essere suddiviso in due parti: una parte riceve solamente lo stato presente in ingresso e fornisce i segnali di controllo in uscita, mentre l'altra riceve in ingresso lo stato presente e gli ingressi e fornisce in uscita solamente lo stato successivo.

Un tipo alternativo di macchina a stati finiti è la macchina di Mealy, dal nome del suo inventore, George Mealy. La macchina di Mealy consente di utilizzare sia lo stato corrente sia gli ingressi per determinare le uscite. Le macchine di Moore hanno un vantaggio teorico in termini di velocità e dimensioni dell'unità di controllo: la maggiore velocità è dovuta al fatto che i segnali di controllo forniti in uscita, che devono essere utilizzati dall'unità di elaborazione all'inizio del ciclo di clock, non dipendono dagli ingressi ma solamente dallo stato corrente. Nell'Appendice A , dove lo schema della macchina a stati finiti viene tradotto in un circuito di porte logiche, il vantaggio in termini di dimensioni diventa chiaramente visibile. Lo svantaggio teorico di una macchina di Moore è dovuto al fatto che può richiedere un numero maggiore di stati. Per esempio, nel caso in cui due sequenze di stati differiscano per un solo stato, la macchina di Mealy può unificare gli stati rendendo le uscite dipendenti dagli ingressi.

### FSM per il controllore semplificato della cache

La Figura 5.39 mostra i quattro stati del controllore della nostra cache semplificata:

- *Inattivo (idle)*: questo è lo stato in cui si trova il controllore quando è in attesa di una richiesta valida di lettura o scrittura dal processore. Quando questa richiesta arriva, lo stato passa a *confronto dei tag*.

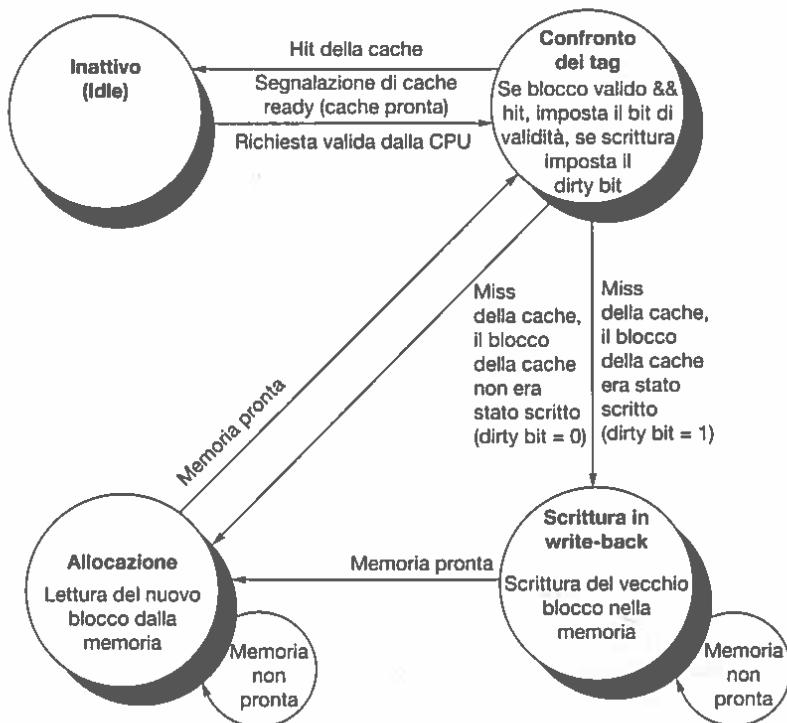


Figura 5.39 I quattro stati del controllore semplificato della cache.

- **Confronto dei tag:** come suggerito dal nome, in questo stato il controllore verifica se la richiesta di lettura o scrittura si risolve in una hit o in una miss: utilizza la parte dell'indirizzo riservata all'indice per selezionare la linea della cache che occorre analizzare. Se il bit di validità della linea è impostato a 1 e il campo tag della linea coincide con quello dell'indirizzo, si verifica una hit. In questo caso la parola contenente il dato viene letta se si tratta di una load o scritta se si tratta di una store e il bit *cache ready* viene impostato a 1. Se si tratta di una scrittura, occorre impostare a 1 anche il dirty bit. Si noti che in una hit in scrittura occorre anche impostare a 1 il bit di validità e scrivere il contenuto del campo tag; anche se a prima vista ciò sembrerebbe non necessario, occorre farlo perché i bit accessori vengono gestiti come un'unica struttura, per cui se si modifica il contenuto del dirty bit verrà modificato anche il contenuto del campo tag e del bit di validità. Se si verifica una hit e il blocco è valido, la FSM ritorna allo stato *inattivo*. Se invece si verifica una miss, per prima cosa si aggiorna il campo tag e quindi si passa o allo stato *scrittura in write-back* (se il blocco indirizzato ha il dirty bit impostato a 1), o allo stato *allocazione* (se il blocco indirizzato ha il dirty bit uguale a 0).
- **Scrittura in write-back:** in questo stato viene trasferito alla memoria il blocco di 128 bit della cache selezionato utilizzando l'indirizzo composto dal campo tag e dall'indice. Si rimane in questo stato in attesa del segnale di ready dalla memoria. Quando la memoria ha terminato la scrittura, la FSM passa allo stato *allocazione*.
- **Allocazione:** il nuovo blocco viene prelevato dalla memoria e si rimane in questo stato in attesa del segnale ready dalla memoria. Quando la memoria ha terminato la lettura, la FSM passa allo stato *confronto dei tag*. Per completare l'operazione saremmo potuti andare in uno stato diverso da *confronto dei tag*; tuttavia le operazioni da eseguire sono in gran parte le stesse, compreso l'aggiornamento della parola del blocco corretta, quando l'accesso viene effettuato in scrittura.

Questo semplice modello può essere facilmente esteso introducendo altri stati per cercare di migliorare le prestazioni. Per esempio, nello stato *confronto dei tag* il controllore confronta i campi tag e legge o scrive il dato in cache, tutto in un singolo ciclo di clock. Tuttavia, spesso, confronto e accesso ai dati della cache vengono eseguiti in stati separati per poter ridurre la durata del periodo del clock. Un'altra ottimizzazione consiste nell'aggiunta di un buffer di scrittura nel quale salvare il blocco se era stato modificato, prima di leggere il nuovo blocco; in questo modo, il processore non deve aspettare due accessi alla memoria quando si verifica una miss su un blocco modificato, ma la cache scriverà il dato nel buffer di scrittura mentre il processore continuerà con l'elaborazione dei dati che aveva richiesto.

Il paragrafo 5.12  contiene ulteriori dettagli sulle FSM, mostrando il controllore completo definito in un linguaggio di descrizione dell'hardware e uno schema a blocchi di questa cache semplificata.

## 5.10 Parallelismo e gerarchie delle memorie: coerenza delle cache

Dato che i multiprocessori multicore contengono più processori sullo stesso chip, è verosimile che questi processori condividano lo stesso spazio di indirizzamento fisico. La scrittura in cache di dati condivisi introduce un nuovo problema: poiché ciascun processore vede la memoria attraverso la propria cache, senza opportune precauzioni due processori finirebbero per “vedere” allo stesso indirizzo due dati diversi. La Figura 5.40 illustra questo problema

Passo	Evento	Contenuto della cache della CPU A	Contenuto della cache della CPU B	Contenuto della memoria nella locazione X
0				0
1	CPU A legge X	0		0
2	CPU B legge X	0	0	0
3	CPU A scrive 1 in X	1	0	1

**Figura 5.40** Il problema della coerenza delle cache, illustrato con una locazione di memoria (X), che viene letta e scritta da due processori (A e B). Supponiamo inizialmente che nessuna delle due cache contenga dati e che la locazione X della memoria principale contenga il valore 0. Supponiamo anche che la cache sia di tipo write-through; una cache di tipo write-back aggiunge alcuni problemi, ma della stessa natura. Dopo che la CPU A ha scritto il valore 1 nella locazione X della memoria, la memoria e la cache A contendono entrambe il valore aggiornato, mentre la cache B contiene ancora 0: se B richiedesse il contenuto della locazione X, leggerebbe quindi il valore 0!

e mostra come due processori diversi possano avere dati diversi nella stessa posizione di memoria. Questo problema è chiamato *problema della coerenza della cache*.

In modo informale si può affermare che un sistema di memoria è coerente se la lettura di un dato restituisce il valore che è stato scritto più di recente in quella cella di memoria. Questa definizione, sebbene possa apparire intuitivamente chiara, è vaga e semplicistica; la realtà è molto più complessa. Nella definizione di coerenza sono contenuti due diversi aspetti del comportamento di un sistema di memoria, che sono entrambi critici per la scrittura in una memoria condivisa. Il primo aspetto, chiamato *coerenza*, definisce *quale valore* debba essere restituito da una lettura. Il secondo aspetto, chiamato *consistenza*, determina *quando* un valore che viene scritto in memoria possa essere restituito da una lettura.

Esaminiamo anzitutto la coerenza. Un sistema di memoria risulta coerente quando:

1. un processore P legge una locazione X della memoria dopo che P stesso ha scritto nella locazione X, senza che altri processori abbiano scritto in X tra la scrittura e la lettura di P, e il valore ottenuto è quello che è stato scritto da P; perciò in Figura 5.40 la CPU A, dopo il passo 3, deve leggere da X il valore 1;
2. un processore legge una locazione X della memoria dopo che un altro processore ha scritto in quella locazione X, e il valore letto corrisponde al valore scritto dall'altro processore, se la lettura e la scrittura sono sufficientemente distanti nel tempo e non intervengono altre scritture nella locazione X tra l'accesso in scrittura e quello in lettura; perciò in Figura 5.40 avremmo bisogno di un meccanismo che scriva 1 nella cache della CPU B al posto del valore 0 quando la CPU A scrive il valore 1 all'indirizzo X della memoria al passo 3;
3. le scritture nella stessa locazione di memoria vengono *serializzate*, ossia due scritture nella stessa locazione fatte da due processori diversi devono essere viste nella stessa sequenza anche dagli altri processori; per esempio, se la CPU B scrivesse 2 nella locazione X della memoria dopo il passo 3, i processori non dovrebbero mai leggere dalla locazione X prima il valore 2 e poi il valore 1.

La prima proprietà semplicemente preserva l'ordine delle istruzioni e ci aspettiamo che sia soddisfatta nei processori a singolo core. La seconda proprietà definisce che cosa significhi avere una visione coerente della memoria: se un processore continuasse a leggere il vecchio contenuto della memoria, potremmo affermare con sicurezza che la memoria è incoerente.

La necessità di *serializzare le scritture* è più sottile, ma ugualmente importante. Supponiamo di non avere serializzato le scritture e che un processore P1 scriva nella locazione X e subito dopo un altro processore P2 scriva anch'esso nella locazione X. La serializzazione delle scritture garantisce che, prima o poi, ciascun processore possa vedere il dato scritto in memoria da P2. Se non serializzassimo le scritture, si potrebbe verificare la situazione in cui uno dei processori vede il dato scritto da P2 prima di vedere quello scritto da P1; in tal caso manterebbe indefinitamente nella sua cache il dato scritto da P1. Il modo più semplice per evitare questo problema è assicurare che tutte le scritture nella stessa locazione di memoria siano viste nello stesso ordine; questa proprietà viene chiamata *serializzazione delle scritture*.

## Schemi di base per garantire la coerenza

In un multiprocessore in cui viene garantita la coerenza, le cache devono consentire la *migrazione* e la *replicazione* dei dati condivisi.

- *Migrazione*: un dato può essere trasferito in una cache locale e utilizzato in quella cache in modo trasparente. La migrazione riduce sia la latenza di accesso ai dati condivisi allocati nella memoria remota sia la larghezza di banda richiesta dalla memoria condivisa.
- *Replicazione*: quando un dato condiviso viene letto simultaneamente da più processori, le cache fanno una copia del dato nella cache locale. La replicazione riduce sia la latenza di accesso sia la competizione nella lettura dei dati condivisi.

Il supporto alla migrazione e alla replicazione è critico per ottenere buone prestazioni nell'accesso ai dati condivisi, per cui molti multiprocessori introducono un protocollo hardware per mantenere la coerenza delle cache; questo protocollo viene chiamato *protocollo di coerenza delle cache*. L'elemento chiave per l'implementazione di questi protocolli è il tracciamento dello stato della condivisione di tutti i blocchi.

Il protocollo più diffuso che viene utilizzato per mantenere la coerenza della cache è chiamato *protocollo di snooping* (letteralmente, "spiare"). Ogni cache che contiene una copia dei dati contenuti in un blocco della memoria fisica ha anche una copia dello stato della condivisione di quel blocco; non viene mantenuta, invece, una copia centralizzata dello stato di condivisione del blocco. Le cache sono tutte accessibili attraverso un dispositivo di collegamento, un bus o una rete, e i controllori di tutte le cache controllano o "spiano" il traffico sul dispositivo di collegamento per determinare se uno dei propri blocchi sia richiesto da un altro dispositivo collegato al bus o alla rete.

Nel prossimo paragrafo spiegheremo come la coerenza delle cache venga implementata su un bus condiviso utilizzando lo snooping, ma qualsiasi altro dispositivo di comunicazione che trasmetta le miss di una cache a tutti i processori può essere utilizzato per implementare uno schema di coerenza basato su snooping. La trasmissione a tutte le cache rende semplici da implementare i protocolli di snooping, ma limita anche la loro scalabilità.

## Protocolli di snooping

Un metodo per garantire la coerenza consiste nell'assicurare che il processore abbia accesso esclusivo a un dato prima di scriverlo. Questo tipo di protocollo viene chiamato *write invalidate protocol* (protocollo di invalidazione in scrittura), perché la copia di un dato contenuto nelle altre cache viene invalidata quando si verifica una scrittura di quel dato. L'accesso esclusivo assicura che non ci siano copie del dato disponibili né in lettura né in scrittura nel momento

Attività del processore	Attività del bus	Contenuto della cache della CPU A	Contenuto della cache della CPU B	Contenuto della locazione X della memoria
				0
CPU A legge X	Miss cache per X	0		0
CPU B legge X	Miss cache per X	0	0	0
CPU A scrive 1 in X	Invalidazione di X	1		0
CPU B legge X	Miss cache per X	1	1	1

**Figura 5.41** Esempio di protocollo di invalidazione su un bus con snooping per l'accesso a un singolo blocco della cache, corrispondente alla locazione X della memoria, con una cache di tipo write-back. Si suppone che all'inizio nessuna delle due cache contenga il dato della locazione X della memoria principale e che X contenga inizialmente il valore 0. Il contenuto della CPU e della memoria viene mostrato dopo che il processore e il bus hanno entrambi terminato le loro attività. Lo spazio vuoto indica che non si è verificata alcuna attività o che non è stata scritta in cache alcuna copia del dato. Quando si verifica la seconda miss, relativa alla cache B, la CPU A risponde alla richiesta fornendo il dato e cancellando la richiesta della memoria. Inoltre, il contenuto sia della cache B sia della locazione X della memoria viene aggiornato. Questa modalità di aggiornamento della memoria, che si verifica quando un blocco viene condiviso, semplifica il protocollo, è comunque possibile tenere traccia del processore che aggiorna i dati e forzare la scrittura in memoria solamente quando il blocco deve essere sostituito. Ciò richiede l'aggiunta di un altro stato, chiamato *owner* (proprietario), che indica che un blocco può essere condiviso, ma che il processore proprietario è responsabile dell'aggiornamento della cache degli altri processori e della memoria quando modifica il contenuto del blocco o lo deve sostituire.

in cui il dato viene scritto; tutte le altre copie presenti nelle altre cache vengono invalidate.

La Figura 5.41 mostra un esempio di protocollo di invalidazione che funziona su un bus con snooping, con una cache di tipo write-back. Per vedere come questo protocollo assicuri la coerenza, si consideri la scrittura di un dato seguita dalla sua lettura da parte di un altro processore: poiché la scrittura richiede l'accesso al dato in modo esclusivo, deve essere invalidata la copia del dato contenuto nella cache dell'altro processore che deve effettuare la lettura del dato, da cui il nome del protocollo. Perciò, quando l'altro processore cerca di leggere il dato, viene generata una miss e la cache sarà costretta a prelevare la nuova copia dalla memoria principale. In caso di scrittura occorre che il processore abbia l'accesso esclusivo al dato e bisogna evitare che gli altri processori cerchino di scrivere lo stesso dato simultaneamente; solo uno dei processori può vincere l'eventuale competizione per la scrittura del dato (*data race*) e invaliderà la copia del dato contenuta nella cache degli altri processori. Perché possano completare la loro scrittura, gli altri processori devono ottenere prima una nuova copia del dato, che conterrà il valore aggiornato. Questo protocollo forza perciò anche la sequenzialità delle scritture.

## Interfaccia hardware/software

**Falsa condivisione:** quando due variabili condivise, non in relazione tra loro, sono situate nello stesso blocco e l'intero blocco viene scambiato tra due processori, anche se i processori, in realtà, accedono a variabili diverse.

Si può intuire che la dimensione dei blocchi della cache gioca un ruolo fondamentale nel meccanismo di coerenza. Per esempio, consideriamo il caso di snooping su una cache con blocchi della dimensione di otto parole e di due processori che alternativamente leggono e scrivono la stessa parola. Quasi tutti i protocolli scambiano interi blocchi tra processori e richiedono quindi una notevole larghezza di banda per mantenere la coerenza.

I blocchi di grandi dimensioni possono anche causare quella che viene definita **falsa condivisione:** quando due variabili condivise non in relazione tra loro sono situate nello stesso blocco della cache, l'intero blocco viene scambiato tra i processori anche se i processori, in realtà, stanno lavorando su variabili diverse. I programmati e i compilatori devono stare molto attenti a come dispongono i dati per evitare false condivisioni.

**Approfondimento.** Sebbene le tre proprietà descritte all'inizio del paragrafo 5.10 siano sufficienti ad assicurare la coerenza, sapere quando il valore di un dato scritto in cache verrà visto dagli altri processori è un aspetto altrettanto importante. Per capirne il motivo, notiamo che non si può pensare che il dato scritto in una locazione X da un processore sia immediatamente disponibile per la lettura (Figura 5.40). Per esempio, se un processore scrive nella locazione X poco prima che un altro processore legga in quella stessa locazione, può diventare impossibile assicurare che la lettura restituiscia il valore corretto, perché al momento della lettura il dato può non aver lasciato ancora il processore. I modelli di consistenza della memoria si occupano di definire esattamente quando un valore scritto può essere visto e letto correttamente.

Formuliamo le due ipotesi seguenti. Innanzitutto, una scrittura non può essere considerata completata (consentendo quindi la scrittura successiva) finché tutti i processori non vedono il suo effetto. In secondo luogo, un processore non deve cambiare l'ordine delle scritture rispetto alla sequenza degli accessi alla memoria. Queste due condizioni implicano che se un processore scrive prima nella locazione X e poi nella locazione Y, i processori che vedono il nuovo valore di Y devono vedere anche il nuovo valore di X. Questi vincoli consentono ai processori di riordinare le letture, ma forzano i processori a completare le scritture nel giusto ordine.

**Approfondimento.** Dato che l'input può modificare la memoria che sta sotto le cache e che l'output può avere bisogno del valore più recente contenuto in una cache di tipo write-back, ci imbattiamo in un problema di coerenza della memoria anche tra l'I/O e la cache, anche nel caso di processore singolo, esattamente come tra le cache di un multiprocessore. Il problema della coerenza delle cache per i multiprocessori e per l'I/O (vedi Cap. 6), sebbene abbia la stessa origine, ha caratteristiche differenti, che fanno sì che la soluzione più appropriata per risolverlo possa essere diversa. A differenza dell'I/O, in cui la replicazione di uno stesso dato è un evento raro (da evitare il più possibile), un programma che viene eseguito su un sistema multiprocessore dispone normalmente di una copia degli stessi dati nelle diverse cache.

**Approfondimento.** Oltre al protocollo di coerenza delle cache basato su snooping, in cui lo stato dei blocchi condivisi è distribuito, esiste un protocollo di coerenza delle cache basato su directory, il quale mantiene lo stato della condivisione di ogni blocco della memoria fisica in un unico posto, chiamato directory. L'implementazione della coerenza basata su directory ha un costo leggermente superiore allo snooping, ma può ridurre il traffico tra le cache e può quindi essere utilizzata anche su sistemi costituiti da molti processori.

## 5.11 Parallelismo e gerarchie delle memorie: i dischi RAID

Questo paragrafo, disponibile online , descrive come utilizzando molti dischi assieme si possa ottenere un throughput maggiore, che è stata la motivazione che ha portato all'invenzione dei RAID (*Redundant Arrays of Inexpensive Disks*, insieme ridondante di dischi economici). La popolarità dei RAID, tuttavia, è dovuta all'affidabilità molto maggiore ottenuta introducendo un piccolo numero di dischi aggiuntivi. Questo paragrafo mostra le differenze in prestazioni, costi e affidabilità tra i livelli RAID.



## 5.12 Argomenti avanzati: come implementare i controllori delle cache

Questo paragrafo, disponibile online , mostra come implementare il controllore di una cache in modo del tutto simile a quanto abbiamo fatto nel Capitolo 4 per il controllore delle unità di elaborazione a singolo ciclo e con pipeline. Il

paragrafo inizia con la descrizione delle macchine a stati finiti e del controllore di una cache semplificata, e include la descrizione del controllore mediante un linguaggio di descrizione dell'hardware. Viene inoltre esaminato in dettaglio un esempio di protocollo che mantiene la coerenza delle cache, analizzando i problemi che possono sorgere nella sua implementazione.

## 5.13 Due casi reali: la gerarchia delle memorie del Cortex-A53 ARM e del Core i7 Intel

In questo paragrafo prenderemo in considerazione la gerarchia delle memorie degli stessi due microprocessori descritti nel Capitolo 4: il Cortex-A53 ARM e il Core i7 Intel. Questo paragrafo è basato sul paragrafo 2.6 della quinta edizione del testo *Computer Architecture: A Quantitative Approach* degli stessi autori.

La Figura 5.42 riporta la dimensione degli indirizzi e del TLB dei due processori. Si noti che l'A53 ha due micro-TLB di 10 elementi completamente associativi supportati da un TLB principale di 512 elementi set-associativo a quattro vie con uno spazio di indirizzamento virtuale di 48 bit e uno spazio di indirizzamento fisico di 40 bit. Il Core i7 ha tre TLB con un indirizzo virtuale su 48 bit e un indirizzo fisico su 44 bit. Sebbene i registri a 64 bit del Core i7 possano contenere un indirizzo virtuale più ampio, il software non ne ha bisogno e un indirizzo su 48 bit consente di ridurre la dimensione della memoria per la tabella delle pagine e dell'hardware del TLB.

Caratteristica	Cortex-A53 ARM	Core i7 Intel
Indirizzo virtuale	48 bit	48 bit
Indirizzo fisico	40 bit	44 bit
Dimensione pagina	Variabile: 4, 16, 64 KiB, 1, 2 MiB, 1 GiB	Variabile: 4 KiB, 2/4 MiB
Organizzazione TLB	1 TLB per le istruzioni e 1 TLB per i dati	1 TLB per le istruzioni e 1 TLB per i dati per core
	Entrambi i micro-TLB sono completamente associativi, con 10 elementi, sostituzione di tipo round robin (circolare)	Entrambi i TLB di L1 sono set-associativi a quattro vie, con sostituzione LRU
	I TLB sono set-associativi a quattro vie e contengono 64 elementi	
	Le miss dei TLB sono gestite in hardware	Il TLB Istruzioni di L1 ha 128 elementi per pagine piccole, 7 per ogni thread per pagine grandi
		Il TLB Dati ha 64 elementi per pagine piccole, 32 elementi per pagine grandi
		Il TLB di L2 è set associativo a quattro vie, con sostituzione LRU
		Il TLB L2 ha 512 elementi
		Le miss dei TLB sono gestite in hardware

**Figura 5.42 Traduzione degli indirizzi e struttura hardware dei TLB del Cortex-A53 ARM e del Core i7 Intel.** Entrambi i processori forniscono il supporto a pagine di grandi dimensioni, impiegate per esempio nella gestione del sistema operativo o nella mappatura della memoria video. Le pagine grandi evitano di utilizzare un elevato numero di elementi per mappare un singolo oggetto che è sempre presente in memoria.

La Figura 5.43 mostra le cache di questi processori. Il Cortex-A53 contiene da uno a quattro processori o core mentre il Core i7 ne ha quattro. Il Cortex-A53 ha una cache istruzioni L1 da 16 a 64 KiB, set-associativa a due vie (per core) mentre il Core i7 ha una cache istruzioni L1 da 32 KiB, set-associativa, a quattro vie (per core). Entrambe le cache hanno blocchi da 64 byte. Il Cortex-A53 aumenta l'associatività a quattro vie per la cache dati mentre gli altri parametri rimangono invariati. Analogamente, il Core i7 mantiene tutti i parametri invariati tranne il grado di associatività, che aumenta a otto. Il Core i7 utilizza una cache L2 unificata set-associativa a otto vie con blocchi da 64 byte. Il Cortex-A53 invece è dotato di una cache L2 che viene condivisa da uno o quattro core. Questa cache è set-associativa a 16 vie con blocchi da 64 byte e ha una dimensione compresa tra 128 KiB e 2 MiB. Dato che i Core i7 vengono utilizzati per i server, offrono anche una cache L3 condivisa da tutti i core del chip. La sua dimensione varia a seconda del numero di core. Con quattro core, come in questo caso, la dimensione è di 8 MiB.

Una sfida importante che devono affrontare i progettisti delle cache è il supporto a processori come il Cortex-A53 e il Core i7 che possono eseguire più di un'istruzione sulla memoria per ciclo di clock. Una tecnica molto utilizzata consiste nel suddividere la cache in banchi per consentire accessi multipli, indipendenti, in parallelo a banchi diversi. Questa tecnica è simile ai banchi di DRAM con interleaving (par. 5.2).

Il Cortex-A53 e il Core i7 contengono delle ottimizzazioni ulteriori che consentono di ridurre la penalità di miss. La prima consiste nel restituire la parola



PARALLELISMO

Caratteristica	Cortex-A53 ARM	Core i7 Intel
Organizzazione cache L1	Cache separate per dati e istruzioni	Cache separate per dati e istruzioni
Dimensione cache L1	Configurabile, da 16 a 64 KiB ciascuna per istruzioni/dati	32 KiB per istruzioni/dati per ogni core
Grado associatività cache L1	Set-associativa a 2 vie (I), a 4 vie (D)	Set-associativa a 4 vie (I), a 8 vie (D)
Modalità sostituzione L1	Random	LRU approssimato
Dimensione blocco L1	64 byte	64 byte
Politica scrittura L1	Write-back, diverse politiche di allocazione (il default è write-allocate)	Write-back, no-write-allocate
Tempo hit L1 (utilizzo load)	Due cicli di clock	4 cicli di clock (pipeline)
Organizzazione cache L2	Un'unica cache per istruzioni e dati	Un'unica cache per istruzioni e dati
Dimensione cache L2	Da 128 KiB a 2 MiB	256 KiB (0,25 MiB)
Grado associatività cache L2	Set-associativa a 16 vie	Set-associativa a 8 vie
Modalità sostituzione L2	LRU approssimato	LRU approssimato
Dimensione blocco L2	64 byte	64 byte
Politica scrittura L2	Write-back, write-allocate	Write-back, write-allocate
Tempo hit L2	12 cicli di clock	10 cicli di clock
Organizzazione cache L3	—	Un'unica cache per istruzioni e dati
Dimensione cache L3	—	8 MiB, condivisa dai core
Grado associatività cache L3	—	Set-associativa a 16 vie
Modalità sostituzione L3	—	LRU approssimato
Dimensione blocco L3	—	64 byte
Politica scrittura L3	—	Write-back, write-allocate
Tempo hit L3	—	35 cicli di clock

Figura 5.43 Cache del Cortex-A53 ARM e del Core i7 Intel.

**Cache non bloccante:** una cache che permette al processore di effettuare accessi anche quando la cache sta gestendo una miss verificatasi precedentemente.

che ha causato la miss prima di completare il trasferimento del blocco di memoria dal livello inferiore. Inoltre, può continuare a eseguire altre istruzioni anche se accedono alla cache dati mentre la gestione della miss è ancora in corso. Una cache di questo tipo, chiamata **cache non bloccante**, viene comunemente utilizzata dai progettisti per nascondere la latenza dovuta alle miss della cache attraverso l'esecuzione fuori ordine. Essi implementano due tipi di tecniche non bloccanti. La *hit sotto una miss* consente di ottenere altre hit quando la gestione della miss è in corso, mentre la *miss sotto una miss* consente di gestire miss multiple. Lo scopo della prima tecnica è di mascherare parte della latenza dovuta a una miss permettendo al processore di svolgere altro lavoro, mentre lo scopo della seconda è di sovrapporre temporalmente le latenze di due diverse miss.

Per sovrapporre le latenze di più miss, occorre disporre di un sistema di memoria con una larghezza di banda elevata, capace di gestire miss multiple in parallelo. Nei dispositivi mobili, il sistema di memoria sottostante è spesso in grado di organizzare gli accessi in pipeline, unire, riordinare e assegnare priorità diverse alle richieste in modo appropriato. I grandi server e i sistemi multiprocessori hanno spesso un sistema di memoria capace di gestire più di una miss in parallelo.

Il Cortex A-53 e il Core i7 possiedono al loro interno un meccanismo hardware per il caricamento anticipato dei dati. Questo meccanismo osserva la successione delle miss sui dati e utilizza questa informazione per cercare di predire il prossimo indirizzo da cui caricare il dato prima che la miss successiva venga generata. Tale sistema funziona generalmente in modo ottimale quando si accede a vettori all'interno di cicli.

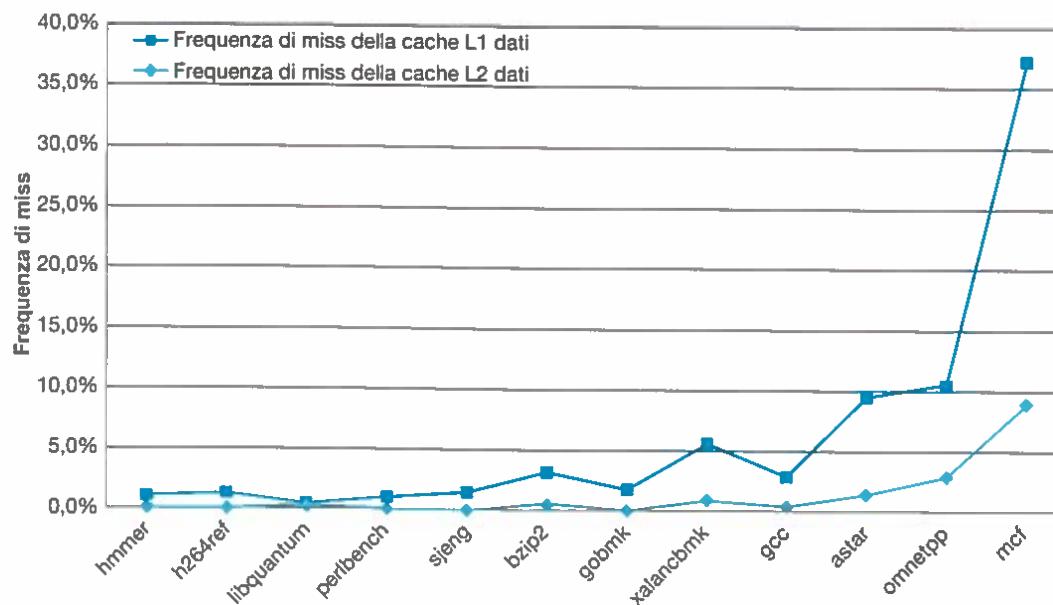
Le sofisticate gerarchie delle memorie di questi chip e la larga porzione di silicio dedicato alle cache e ai TLB illustrano chiaramente il grande sforzo dei progettisti per cercare di colmare il divario tra durata del ciclo di clock del processore e latenza della memoria.

### Prestazioni della gerarchia delle memorie del Cortex-A53 e del Core i7

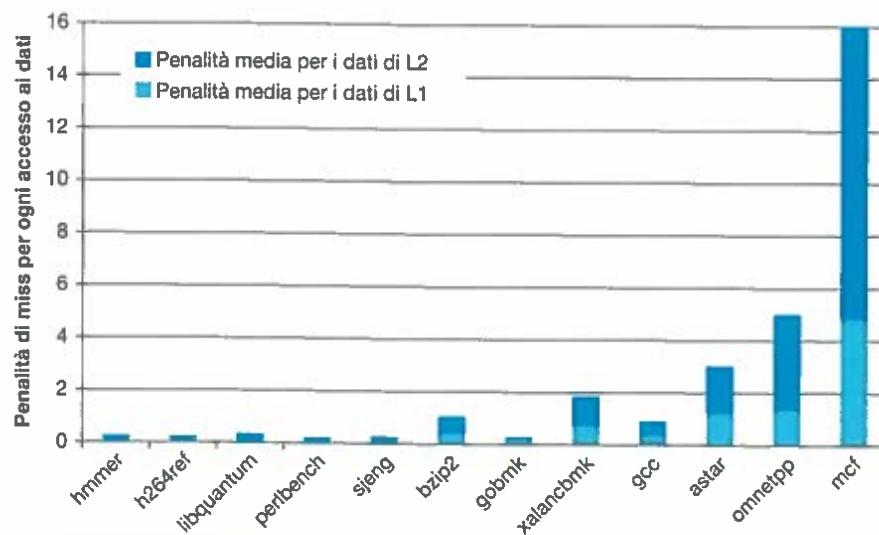
La gerarchia delle memorie del Cortex-A53 è stata valutata utilizzando una cache istruzioni L1 set-associativa a due vie da 32 KiB, una cache dati L1 set-associativa a quattro vie da 32 KiB, e una cache L2 set-associativa a 16 vie. La valutazione è stata fatta sul benchmark SPEC2006 su interi.

La frequenza di miss della cache istruzioni del Cortex-A53 per questi benchmark è molto bassa. La Figura 5.44 mostra i risultati per la cache dati del Cortex-A53 che mostra frequenze di miss L1 e L2 significative. La frequenza di miss della cache dati L1 va dallo 0,5 al 37,3%, con un valore medio del 6,4% e un valore mediano del 2,4%. La frequenza di miss (globale) della cache L2 varia tra lo 0,1 e il 9%, con una valore medio dell'1,3% e un valore mediano dello 0,3%. La penalità di miss della cache L1 per un Cortex-A53 da 1 GHz è di 12 cicli di clock, mentre la penalità di miss della cache L2 è di 124 cicli di clock. Utilizzando questi valori di penalità di miss, la Figura 5.45 mostra il valore medio della penalità di miss per l'accesso ai dati. Quando questi valori piccoli di frequenza di miss vengono moltiplicati per le penalità di miss associate, potete vedere che rappresentano una frazione significativa del CPI per 5 dei 12 programmi che costituiscono i programmi dello SPEC2006.

La Figura 5.46 mostra la frequenza di miss delle cache per il Core i7 utilizzando i benchmark dello SPEC2006. La frequenza di miss della cache L1 delle istruzioni varia da 0,1 a 1,8%, con una media appena superiore a 0,4%. Questo dato è in linea con quelli riportati in altri studi sul comportamento delle cache istruzioni nei benchmark SPEC CPU2006, che mostrano una frequenza di miss bassa. Data una frequenza di miss della cache L1 dei dati compresa tra il 5 e il 10%, e a volte anche più alta, l'importanza delle cache L2 e L3 è evidente. Dato



**Figura 5.44** Frequenza di miss della cache dati per un Cortex-A53 ARM nell'esecuzione dei benchmark SPEC2006int. Applicazioni che richiedono una maggiore occupazione di memoria tendono ad avere una frequenza di miss maggiore nelle cache L1 e L2. Si noti che la frequenza di miss nella cache L2 è una frequenza di miss globale, cioè si considerano tutti gli accessi che hanno successo compresi quelli alla cache L1 (vedi la sezione *Approfondimento* del par. 5.4). Il programma mcf è noto come *cache buster* ("acchiappa cache"). Si noti che questa figura è basata sugli stessi sistemi e sugli stessi benchmark di Figura 4.74.



**Figura 5.45** Penalità media di accesso alla memoria misurata in numero di cicli di clock per accesso alla memoria richiesti dalle cache L1 e L2 di un processore ARM nell'esecuzione dei benchmark SPEC2006int. Sebbene la frequenza di miss di L1 sia significativamente più alta, la penalità di miss della cache L2, che è oltre cinque volte maggiore, significa che le miss della cache L2 possono dare un contributo significativo.

che il costo di una miss della memoria è superiore ai 100 cicli di clock e la frequenza media di miss della cache L2 è del 4%, la cache L3 è chiaramente critica. Supponendo che metà delle istruzioni siano load o store, senza la cache L3 le miss nella cache L2 aggiungerebbero due cicli di clock per istruzione al CPI! Confrontando i risultati, la frequenza di miss media della cache L3 è dell'1%, ancora significativa, ma quattro volte inferiore alla frequenza di miss della cache L2 e sei volte inferiore alla frequenza di miss della cache L1.

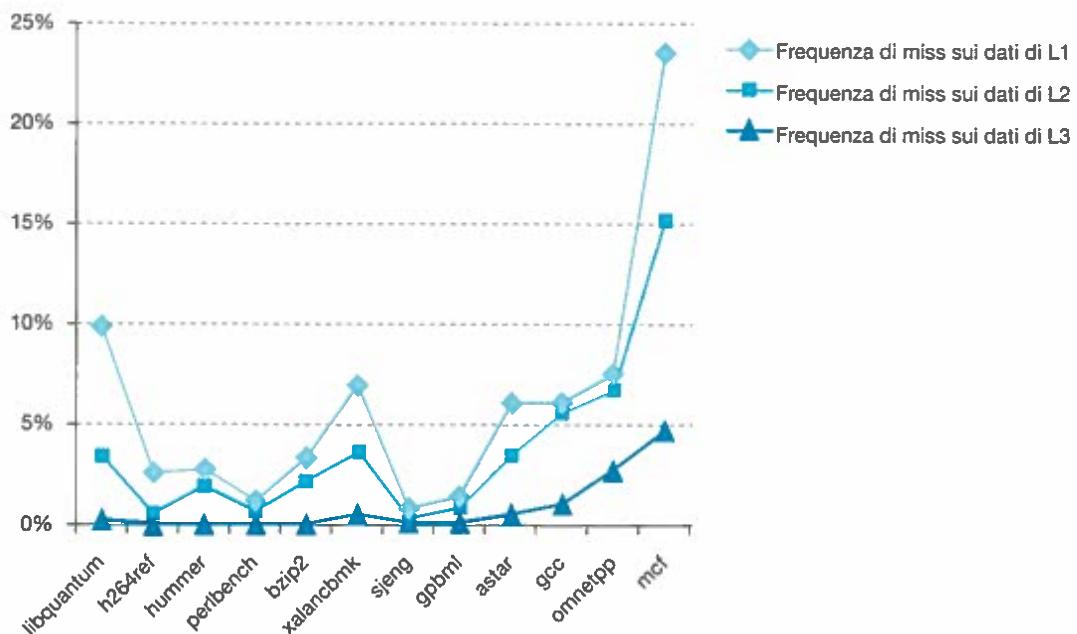


Figura 5.46 Frequenza di miss della cache L1, L2 e L3 per un Core i7 920 Intel durante l'esecuzione dei benchmark interi completi dello SPEC2006.

**Approfondimento.** Dato che la speculazione può dare a volte risultati sbagliati (vedi Cap. 4), ci sono degli accessi alla cache L1 che non corrispondono a load o store utili all'esecuzione. I dati in Figura 5.45 sono misurati rispetto a tutti gli accessi alla cache compresi quelli che vengono cancellati. La frequenza di miss misurata solamente rispetto agli accessi utili diventa maggiore di 1,6 volte (9,5% in media, contro il 5,9%, per la miss della cache dati L1).

## 5.14 Un caso reale: il resto del sistema RISC-V e le istruzioni speciali

La Figura 5.47 riporta le rimanenti istruzioni RISC-V classificate nella categoria delle istruzioni generiche e di sistema.

Le istruzioni FENCE (letteralmente siepe) forniscono delle barriere di sincronizzazione per le istruzioni (fence.i), per i dati (fence) e per la traduzione degli indirizzi (sfence.vma). La prima, fence.i, informa il processore che il software ha modificato la memoria istruzioni, in modo tale che il fetch delle istruzioni rifletta l'aggiornamento delle istruzioni. La seconda, fence, influenza l'ordine con cui accedono ai dati le elaborazioni multiple e l'I/O. La terza, sfence.vma, informa il processore che il software ha modificato le tabelle delle pagine, in modo tale che possa garantire che la traduzione degli indirizzi rifletta questo aggiornamento.

Le sei istruzioni che accedono ai registri di controllo e di stato (CSR, *Control and Status Register*) spostano i dati tra i registri a uso generale e i CSR. L'istruzione csrrwi (lettura/scrittura immediata di un CSR) copia un CSR in un registro e poi sovrscrive il CSR con una costante. csrrsi (lettura/impostazione immediata di un CSR) copia un CSR in un registro intero e sovrscrive il CSR con una costante. csrrci (leggi/azzera un CSR) si comporta come la csrrsi ma imposta i bit a zero invece di scrivere una costante. Le istruzioni csrrw, csrrs, e csrrc utilizzano un operando su registro invece di una costante, ma fanno le stesse cose.

Nome simbolico	Nome esteso
FENCE.I	Fence Istruzioni
FENCE	Fence
SFENCE.VMA	Fence sulla Traduzione degli Indirizzi
CSRRWI	Lettura/scrittura immediata CSR
CSRRSI	Lettura/impostazione immediata CSR
CSRRCI	Lettura/azzeramento immediato CSR
CSRRW	Lettura/scrittura CSR
CSRRS	Lettura/impostazione CSR
CSRRC	Lettura/azzeramento CSR
ECALL	Chiamata di ambiente
EBREAK	Breakpoint di ambiente
SRET	Ritorno da eccezione del supervisore
WFI	Attesa di interrupt

**Figura 5.47** Elenco delle istruzioni in linguaggio assembler per le funzioni di sistema e funzioni speciali dell'insieme completo di istruzioni RISC-V.

L'unico compito di due istruzioni è generare un'eccezione: ecall genera un'eccezione di chiamata di ambiente per invocare il sistema operativo, ed ebreak genera un'eccezione di breakpoint per invocare il debugger. L'istruzione del supervisore di ritorno da eccezione (sret), naturalmente, consente a un programma di ritornare al punto di interruzione dopo l'esecuzione della procedura di gestione dell'eccezione.

Infine, l'istruzione di attesa degli interrupt, wfi, informa il processore che può entrare in uno stato di attesa (*idle*) fino a quando non arriva un interrupt.

## 5.15 Come andare più veloci: blocchi di cache e moltiplicazione tra matrici

Il passo successivo per aumentare le prestazioni della funzione DGEMM adattandola progressivamente all'hardware a disposizione consiste nell'utilizzo di blocchi di cache adeguati, da aggiungere alla parallelizzazione a livello di parola e di istruzioni, viste nei Capitoli 3 e 4. La Figura 5.48 mostra la versione di DGEMM riportata in Figura 4.77 modificata per sfruttare la gestione a blocchi della cache. Le modifiche sono le stesse apportate in precedenza quando siamo passati dalla versione DGEMM non ottimizzata di Figura 3.21 alla versione a blocchi di Figura 5.21. Questa volta prendiamo come base la versione di DGEMM utilizzata nel Capitolo 4 in cui i cicli sono stati espansi e la richiamiamo più volte sulle sottomatrici di A, B e C. Le linee 28-34 e 7-8 di Figura 5.48 sono uguali alle linee 14-20 e 5-6 di Figura 5.21, ad eccezione dell'incremento dell'indice di fine ciclo alla linea 7, pari al numero dell'espansione del ciclo.

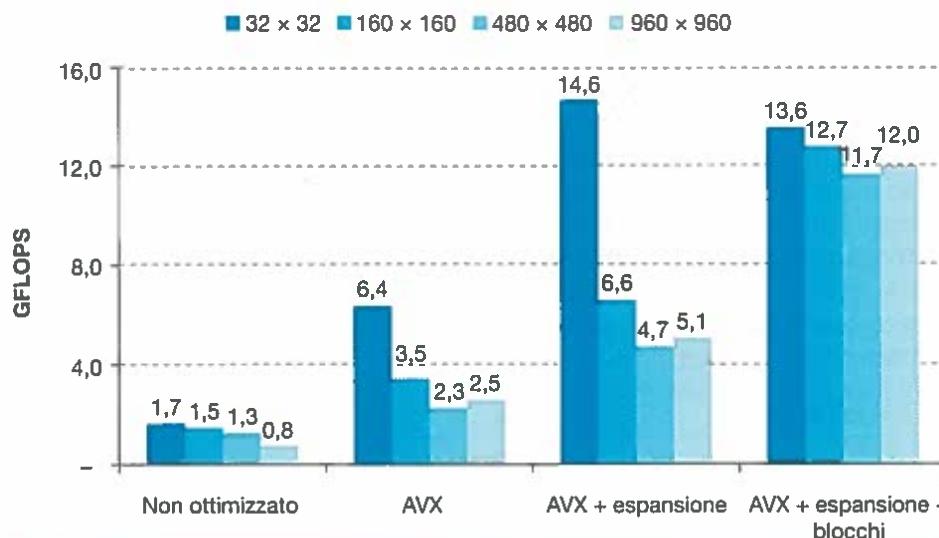
Diversamente dai paragrafi precedenti, in questo caso non mostreremo il codice x86 risultante, poiché il codice del ciclo interno è praticamente identico al codice di Figura 4.78, dato che l'utilizzo dei blocchi non influenza il calcolo ma solamente l'ordine con cui si accede ai dati in memoria. Quello che cambia sono le istruzioni intere utilizzate per implementare il ciclo: queste passano

```

1 #include <x86intrin.h>
2 #define ESPANDI (4)
3 #define DIM_BLOCCO 32
4 void esegui_blocco (int n, int si, int sj, int sk,
5                      double *A, double *B, double *C)
6 {
7     for (int i = si; i < si+DIM_BLOCCO; i+=ESPANDI*4)
8         for (int j = sj; j < sj+DIM_BLOCCO; j++) {
9             __m256d c[4];
10            for (int x = 0; x < ESPANDI; x++)
11                c[x] = _mm256_load_pd(C+i+x*4+j*n);
12 /* c[x] = C[i][j] */
13            for( int k = sk; k < sk+DIM_BLOCCO; k++)
14            {
15                __m256d b = _mm256_broadcast_sd(B+k+j*n);
16 /* b = B[k][j] */
17                for (int x = 0; x < ESPANDI; x++)
18                    c[x] = _mm256_add_pd(c[x], /* c[x]+=A[i][k]*b */
19                               _mm256_mul_pd(_mm256_load_pd(A+n*k+x*4+i), b));
20            }
21
22
23            for (int x = 0; x < ESPANDI; x++)
24                _mm256_store_pd(C+i+x*4+j*n, c[x]);
25 /* C[i][j] = c[x] */
26        }
27
28 void dgemm (int n, double* A, double* B, double* C)
29 {
30     for (int sj = 0; sj < n; sj += DIM_BLOCCO)
31         for (int si = 0; si < n; si += DIM_BLOCCO)
32             for (int sk = 0; sk < n; sk += DIM_BLOCCO)
33                 esegui_blocco(n, si, sj, sk, A, B, C);
34 }
```

**Figura 5.48** Versione di DGEMM di Figura 4.78 ottimizzata utilizzando blocchi di cache. Queste modifiche sono le stesse di Figura 5.21. Il linguaggio assembler prodotto dal compilatore per la funzione `esegui_blocco` è quasi identico a quello di Figura 4.79. Anche qui non c'è costo addizionale nella chiamata della funzione `esegui_blocco`, perché la funzione viene inserita direttamente nel codice (modalità inline).

da 14 prima del ciclo interno e 8 dopo il ciclo interno (vedi Figura 4.78) a 40 e 28 rispettivamente nel codice riportato in Figura 5.48. Tuttavia, il numero di istruzioni aggiunte è del tutto trascurabile in confronto all'aumento delle prestazioni ottenuto riducendo le miss della cache. La Figura 5.49 mostra i risultati del confronto tra la versione di DGEMM non ottimizzata, la versione ottimizzata sfruttando il parallelismo a livello di parola, quella ottimizzata sfruttando anche il parallelismo a livello di parola a livello di istruzioni, e quella ottimizzata sfruttando anche la gestione a blocchi della cache. Solamente l'ultima ottimizzazione consente di aumentare le prestazioni su codice AVX nel quale i cicli sono stati espansi di un fattore compreso tra 2 e 2,5 su matrici di grandi dimensioni. Confrontando la versione non ottimizzata con la versione contenente tutte e quattro le ottimizzazioni si osserva un miglioramento delle prestazioni di un fattore compreso tra 8 e 15, con il miglioramento più ampio ottenuto per le matrici più grandi.



**Figura 5.49** Prestazioni delle quattro versioni di DGEMM su matrici di dimensioni variabili da  $32 \times 32$  a  $960 \times 960$ . Il codice contenente tutte le ottimizzazioni è quasi 15 volte più veloce del codice non ottimizzato di Figura 3.22.

**Approfondimento.** Come nella sezione *Approfondimento* del paragrafo 3.8, questi risultati sono stati ottenuti con la modalità Turbo disattivata. Come abbiamo visto nei Capitoli 3 e 4, l'attivazione della modalità Turbo consente di aumentare le prestazioni di un fattore pari al rapporto tra i clock:  $3,3/2,6 = 1,27$ . La modalità Turbo funziona particolarmente bene in questo caso, perché si utilizza un solo core degli otto disponibili sul chip. Per migliorare ulteriormente le prestazioni, dovremmo utilizzare tutti gli otto core; vedremo come nel Capitolo 6.

## 5.16 Errori e trabocchetti

Anche nell'uso e nella progettazione delle gerarchie delle memorie sono stati compiuti molti errori (i cui effetti si sono trascinati per anni) e si è incappati in svariati trabocchetti, alcuni dei quali hanno avuto ricadute estremamente negative. Iniziamo questo paragrafo esaminando due trabocchetti in cui cadono spesso gli studenti quando svolgono gli esercizi.

**Trabocchetto:** ignorare il comportamento del sistema di memoria quando si scrivono i programmi o quando si genera il codice tramite compilatore.

Questo trabocchetto può condurre facilmente all'errore: «I programmati possono ignorare le gerarchie delle memorie quando scrivono i loro programmi». La valutazione del programma di ordinamento di Figura 5.19 e dell'utilizzo dei blocchi della cache nel paragrafo 5.14 dimostrano chiaramente che i programmati possono facilmente raddoppiare le prestazioni dei loro programmi se considerano il comportamento del sistema di memoria quando progettano i loro algoritmi.

**Trabocchetto:** dimenticare di tenere conto dell'indirizzamento al byte o della dimensione dei blocchi quando si simula una cache.

Quando si simula una cache, manualmente o con il calcolatore, bisogna tenere in considerazione l'indirizzamento al byte e il numero di parole di cui sono

costituiti i blocchi quando si determina il blocco della cache corrispondente a un certo indirizzo. Per esempio, se abbiamo una cache a mappatura diretta di 32 byte con blocchi di 4 byte, l'indirizzo del byte numero 36 corrisponderà al blocco 1 della cache, dato che l'indirizzo del byte 36 viene mappato sul blocco avente indirizzo  $9: 9 \text{ modulo } 8 = 1$ . D'altra parte, se l'indirizzo 36 fosse un indirizzo di parola, allora corrisponderebbe al blocco:  $36 \text{ modulo } 8 = 4$ . Verificate sempre l'unità di misura di base rispetto alla quale viene espresso l'indirizzo.

In maniera analoga occorre tenere in considerazione la dimensione dei blocchi: si ipotizzi di avere una cache di 256 byte con blocchi di 32 byte. In quale blocco cade l'indirizzo 300 espresso in byte? Se scomponiamo l'indirizzo 300 nei diversi campi, la risposta è immediata:

63	62	61	...	...	...	11	10	9	8	7	6	5	4	3	2	1	0
0	0	0	...	...	...	0	0	0	1	0	0	1	0	1	1	0	0
Numero del blocco della cache										Offset del blocco							
Indirizzo del blocco																	

L'indirizzo 300 in byte corrisponde al blocco:

$$\left[ \frac{300}{32} \right] = 9$$

Il numero di blocchi nella cache è pari a:

$$\left[ \frac{256}{32} \right] = 8$$

Il blocco numero 9 corrisponde quindi al blocco della cache numero:  $9 \text{ modulo } 8 = 1$ .

Molte persone commettono questo errore, compresi gli autori di questo libro (che conteneva l'errore nelle bozze!) e i docenti che si dimenticano di specificare se gli indirizzi debbano essere intesi come indirizzi di parole, di byte o di blocco. Ricordatevi di questo trabocchetto quando svolgete gli esercizi.

**Trabocchetto:** avere un grado di associatività di una cache condivisa inferiore al numero di core o di thread che condividono quella cache.



Senza particolari attenzioni, un programma parallelo che viene eseguito su  $2^n$  processori o che è costituito da  $2^n$  thread potrebbe facilmente allocare le strutture dati a indirizzi che vengono mappati sulla stessa linea di una cache L2 condivisa. Se la cache ha un grado di associatività pari ad almeno  $2^n$ , gli eventuali conflitti possono essere nascosti al programma dall'hardware; altrimenti, i programmati si troverebbero a dovere capire l'origine di un calo delle prestazioni apparentemente inspiegabile, ma in realtà dovuto alle miss di conflitto della cache L2. Questo calo di prestazioni, per esempio, può insorgere quando si porta il codice da un'architettura a 16 core su un'architettura a 32 core mantenendo però una cache associativa L2 a 16 vie.

**Trabocchetto:** utilizzare il tempo di accesso medio alla memoria per valutare la gerarchia delle memorie di un processore fuori ordine.

Se un processore va in stallo durante una miss della cache, è possibile calcolare separatamente il tempo di stallo dovuto alla memoria e il tempo di esecuzione

del processore, e quindi valutare la gerarchia delle memorie indipendentemente utilizzando il tempo medio di accesso alla memoria (par. 5.4).

Se durante la gestione di una miss il processore continua a eseguire altre istruzioni e può anche gestire altre miss della cache, l'unica valutazione accurata della gerarchia delle memorie consiste nella simulazione del processore fuori ordine e della sua gerarchia delle memorie.

**Trabocchetto:** estendere lo spazio di indirizzamento aggiungendo altri segmenti a uno spazio di indirizzamento non segmentato.

Durante gli anni '70 molti programmi crebbero a tal punto che non tutte le istruzioni e i dati potevano essere indirizzati utilizzando solamente i 16 bit a disposizione per l'indirizzo.

I calcolatori furono riprogettati in modo da offrire indirizzi su 32 bit attraverso uno spazio di indirizzamento non segmentato, detto anche *spazio di indirizzamento piatto (flat address space)*, oppure aggiungendo agli indirizzi esistenti (su 16 bit) altri 16 bit per specificare un segmento. Dal punto di vista del mercato, l'aggiunta dei segmenti – che erano visibili ai programmati e che obbligavano programmati e compilatori a scomporre i programmi in segmenti – poteva risolvere il problema dell'indirizzamento. Si verificava però un problema ogni volta che un programma doveva indirizzare dati oltre la dimensione di un singolo segmento, per esempio quando bisognava scorrere mediante un indice gli elementi di un vettore di grandi dimensioni, oppure quando si utilizzavano puntatori senza restrizioni o senza indirizzare parametri. Inoltre, l'aggiunta dei segmenti può spezzare tutti gli indirizzi in due parole, una per il segmento e una per l'offset all'interno del segmento, causando dei problemi nell'utilizzo degli indirizzi che sono contenuti nei registri.

**Errore:** la frequenza dei guasti di un disco utilizzato in un calcolatore è uguale a quella riportata nelle specifiche.

Due studi recenti hanno preso in considerazione alcune grandi collezioni di dischi per valutare la relazione tra i risultati ottenuti sul campo e i dati riportati nelle specifiche. Il primo studio riguardava quasi 100 000 dischi che riportavano un MTTF compreso tra 1 000 000 e 1 500 000 ore, ossia un AFR compreso tra lo 0,6 e lo 0,8%; i ricercatori, invece, hanno trovato che l'AFR era compreso tra il 2 e il 4%, valori spesso da tre a cinque volte superiori a quelli riportati nelle specifiche [Schroeder e Gibson, 2007].

Un secondo studio su più di 100 000 dischi di Google, le cui specifiche riportavano un AFR dell'1,5%, ha misurato un AFR dell'1,7% per le unità disco nel loro primo anno di vita, che saliva all'8,6% nel loro terzo anno di vita, un valore da cinque a sei volte superiore a quello riportato nelle specifiche [Pinheiro, Weber e Barroso, 2007].

**Errore:** i sistemi operativi sono il posto migliore per programmare gli accessi a disco.

Come accennato nel paragrafo 5.2, le interfacce ad alto livello dei dischi consentono al sistema operativo di indirizzare blocchi logici. Essendo presente questo alto livello di astrazione, la cosa migliore che il SO può fare per cercare di migliorare le prestazioni è ordinare gli indirizzi dei blocchi logici in ordine crescente. Tuttavia, poiché è il disco a sapere quale sia la mappatura degli indirizzi logici sui diversi componenti fisici del disco (settori, tracce e piatti), è il controllore del disco che può ridurre la latenza di rotazione e di ricerca riorganizzando gli accessi.

Per esempio, si supponga di avere un carico di lavoro di quattro letture [Anderson, 2003]:

Operazione	Indirizzo di partenza del blocco	Lunghezza
Lettura	724	8
Lettura	100	16
Lettura	9987	1
Lettura	26	128

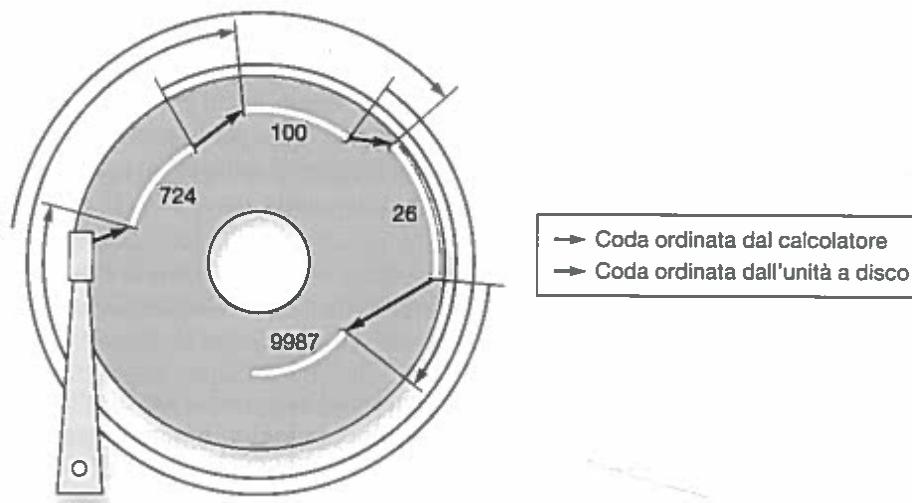
Il calcolatore potrebbe riordinare le quattro letture a seconda dell'indirizzo logico dei blocchi:

Operazione	Indirizzo di partenza del blocco	Lunghezza
Lettura	26	128
Lettura	100	16
Lettura	724	8
Lettura	9987	1

In base alla posizione relativa dei dati sul disco, la riorganizzazione degli accessi potrebbe determinare un peggioramento delle prestazioni, come mostrato in Figura 5.50. Le letture riorganizzate dal disco richiederebbero solamente tre quarti di un giro completo del disco, mentre la riorganizzazione fatta dal sistema operativo richiederebbe tre giri.

*Trabocchetto: implementare il monitor di una macchina virtuale per un'architettura dell'insieme delle istruzioni che non è stata progettata per essere virtualizzabile.*

Molti progettisti negli anni '70 e '80 non consideravano con la dovuta attenzione il fatto che tutte le istruzioni che leggevano o scrivevano informazioni colle-



**Figura 5.50** Esempio che mostra la differenza tra gli accessi ordinati dal SO e gli accessi ordinati dall'unità a disco. Nel primo caso sono richiesti tre giri del disco per completare le quattro letture, mentre nel secondo caso le letture possono essere effettuate utilizzando solamente tre quarti del tempo di una singola rotazione [da Anderson, 2003].

gate alle risorse dell'hardware dovessero essere istruzioni privilegiate. Questo approccio "lassista" ha causato svariati problemi nella realizzazione dei VMM per molte architetture, compresa la x86 che utilizziamo qui come esempio. La **Figura 5.51** descrive le 18 istruzioni che sono problematiche per la virtualizzazione [Robin e Irvine, 2000]; esse possono essere suddivise in due ampie classi:

- istruzioni di lettura dei registri di controllo in modalità utente, che rivelano che il sistema operativo ospitato sta girando su una macchina virtuale, come l'istruzione POPF, già introdotta in questo capitolo;
- istruzioni di controllo delle protezioni (richieste dalle architetture che utilizzano la segmentazione) che suppongono che il sistema operativo sia in esecuzione con i privilegi più elevati.

Per semplificare le implementazioni dei VMM sugli x86, sia AMD sia Intel hanno introdotto alcune estensioni dell'architettura aggiungendo una nuova modalità di funzionamento. L'estensione Intel VT-x fornisce una nuova modalità di esecuzione da utilizzare per eseguire le VM, una definizione interna all'architettura dello stato di una VM, delle istruzioni per sostituire velocemente le VM e un ampio insieme di parametri per scegliere in quali circostanze occorra richiamare il VMM. Complessivamente, l'estensione VT-x aggiunge 11 nuove istruzioni all'insieme dell'x86 e l'estensione Pacifica di AMD utilizza soluzioni simili.

Un'alternativa alla modifica dell'hardware è quella di operare delle piccole modifiche al sistema operativo per evitare di utilizzare quelle parti dell'architettura che possono creare problemi. Questa tecnica si chiama *paravirtualizzazione* e il software open source Xen VMM ne è un buon esempio. Lo Xen VMM fornisce un SO ospitato contenente un'astrazione di macchina virtuale che utilizza solamente le parti facilmente "virtualizzabili" dell'hardware fisico dell'x86 sul quale eseguire il VMM.

Categoria di problemi	Istruzioni x86 che causano problemi
Registri sensibili all'accesso, che non generano un'eccezione di trap quando l'esecuzione è in modalità utente	Memorizzazione del registro della tabella dei descrittori globali (SGDT) Memorizzazione del registro della tabella dei descrittori locali (SLDT) Memorizzazione del registro della tabella dei descrittori degli interrupt (SIDT) Memorizzazione della parola contenente lo stato della macchina (SMSW) Push dei flag (PUSHF, PUSHFD) Pop dei flag (POPF, POPFD)
Istruzioni che non superano il controllo di protezione dell'x86 quando si accede ai meccanismi della memoria virtuale in modalità utente	Caricamento dei diritti di accesso dal descrittore dei segmenti (LAR) Caricamento dei limiti dei segmenti dal descrittore dei segmenti (LSL) Verifica della leggibilità del descrittore dei segmenti (VERR) Verifica della scrivibilità del descrittore dei segmenti (VERW) Copia da stack dei registri di segmento (POP CS, POP SS, ...) Scrittura dei registri di segmento in stack (PUSH CS, PUSH SS, ...) Chiamata di tipo far a istruzioni con un diverso livello di privilegio (CALL) Ritorno da procedura di tipo far a istruzioni con un diverso livello di privilegio (RET) Salto di tipo far a istruzioni con un diverso livello di privilegio (JMP) Interrupt software (INT) Memorizzazione del registro di selezione di segmento (STR) Spostamento ai/dai registri di segmento (MOVE)

**Figura 5.51** Elenco delle 18 istruzioni dell'x86 che creano problemi alla virtualizzazione [Robin e Irvine, 2000]. Le prime cinque istruzioni del primo gruppo consentono a un programma in modalità utente di leggere un registro di controllo, per esempio il registro di una tabella di descrizione, senza sollevare un'eccezione di trap. L'istruzione pop dei flag modifica un registro di controllo con informazioni sensibili, ma fallisce quando viene eseguita in modalità utente senza darne avviso. Il controllo delle protezioni dell'architettura segmentata dell'x86, nel secondo gruppo in basso, si è rivelato particolarmente disastroso, dato che ciascuna di queste istruzioni controlla il livello di privilegio in modo implicito durante l'esecuzione dell'istruzione stessa, leggendo uno dei registri di controllo. Il controllo presuppone che il SO si trovi al livello più alto di privilegio, cosa che non succede in una VM ospitata. Solamente l'istruzione di spostamento a uno dei registri di segmento cerca di modificare lo stato di controllo, ma il controllo delle protezioni fa comunque fallire il tentativo di modifica.

## 5.17 Note conclusive

La difficoltà nel costruire un sistema di memoria che possa tenere il passo di CPU sempre più veloci è dovuta al fatto che il componente di base utilizzato per la memoria principale, il chip di DRAM, è praticamente lo stesso nei calcolatori più veloci e in quelli più lenti ed economici.

È il principio di località che ci dà la possibilità di superare il limite imposto dai lunghi tempi di latenza degli accessi alla memoria, e l'efficacia di questa strategia è dimostrata a tutti i livelli della gerarchia delle memorie. Sebbene i diversi livelli di una gerarchia sembrino piuttosto diversi dal punto di vista quantitativo, essi funzionano mettendo in pratica strategie simili e sfruttano le stesse proprietà di località.

Le cache multilivello rendono possibile utilizzare più facilmente diverse ottimizzazioni della cache, per due ragioni.

In primo luogo, i parametri di progetto di una cache di livello più basso nella gerarchia sono differenti da quelli di una cache di primo livello. Per esempio, poiché una cache di livello inferiore conterrà parole molto più ampie, sarà possibile utilizzare blocchi di dimensioni maggiori. In secondo luogo, una cache di livello più basso non è costantemente utilizzata dal processore, come avviene invece per una cache di primo livello. Questo permette di far fare qualcosa di utile a una cache di livello inferiore mentre è inattiva, per prevenire miss future.

Un'altra strategia utilizzata è quella di "farsi aiutare" dal software: gestire in modo efficace una gerarchia delle memorie ricorrendo a trasformazioni dei programmi o a particolari componenti hardware è uno degli argomenti più caldi nello sviluppo dei compilatori. In questo ambito vengono esplorate due idee principali.

La prima consiste nella riorganizzazione del codice per aumentarne la località spaziale e temporale. Questo approccio si focalizza sui programmi che fanno un uso intensivo dei cicli e che utilizzano vettori di grandi dimensioni come strutture principali (le operazioni di algebra lineare, come DGEMM, ne sono un esempio). Riorganizzando le istruzioni contenute nei cicli che accedono agli elementi dei vettori, si può migliorare significativamente la località e, quindi, le prestazioni delle cache.

Un altro approccio è rappresentato dal **caricamento anticipato** (*prefetching*), in cui un blocco di dati viene portato all'interno della cache prima ancora che venga effettivamente indirizzato. Molti microprocessori utilizzano un caricamento anticipato guidato dall'hardware, cercando di predire gli accessi ai dati che sarebbe difficile identificare via software.

Un terzo approccio è rappresentato da istruzioni speciali che tengono conto della struttura della cache per ottimizzare il trasferimento dalla memoria. Per esempio, i microprocessori descritti nel paragrafo 6.10 adottano un'ottimizzazione per la quale quando un programma si accinge a scrivere in cache e si verifica una miss in scrittura, non viene caricato in cache il contenuto del blocco della memoria principale quando il programma si accinge a scrivere l'intero blocco.

Questa ottimizzazione riduce significativamente il traffico con la memoria a carico del processore.

Come vedremo nel Capitolo 6, i sistemi di memoria costituiscono un argomento chiave nella progettazione dei processori paralleli. L'importanza crescente delle gerarchie delle memorie nel determinare le prestazioni dell'intero sistema farà in modo che questo aspetto continuerà a ricevere, nei prossimi anni, una particolare attenzione sia da parte dei progettisti sia da parte dei ricercatori.



**Caricamento anticipato**  
**(prefetching):** una tecnica nella quale i blocchi di dati che serviranno in futuro vengono scritti nella cache in anticipo, utilizzando delle istruzioni speciali che specificano l'indirizzo di questi blocchi.



## 5.18 Inquadramento storico e approfondimenti

Questo paragrafo, disponibile online , fornisce una visione di insieme delle tecnologie delle memorie che si sono susseguite negli anni, dalle linee di ritardo al mercurio fino alle DRAM. Esso tratta l'invenzione della gerarchia delle memorie, i meccanismi di protezione e le macchine virtuali, e termina con una breve storia dei sistemi operativi, tra cui CTSS, MULTICS, UNIX, UNIX BSD, MS-DOS, Windows e Linux.

## 5.19 Esercizi

Si supponga che la memoria sia indirizzabile al byte e che le parole siano di 64 bit, a meno che non sia specificato diversamente.

**5.1** In questo esercizio esamineremo le proprietà di località della memoria mediante il calcolo matriciale. Il seguente frammento di codice è scritto in C; gli elementi delle matrici posti sulla stessa riga sono contenuti in celle contigue della memoria. Si supponga che ciascuna parola sia costituita da interi su 64 bit.

```
for (I=0; I<8; I++)
    for (J=0; J<8000; J++)
        A[I][J] = B[I][0] + A[J][I];
```

**5.1.1** [5] <5.1> Quanti interi su 64 bit possono essere contenuti in una linea di cache di 16 byte?

**5.1.2** [5] <5.1> L'accesso a quali variabili gode della località temporale?

**5.1.3** [5] <5.1> L'accesso a quali variabili gode della località spaziale?

La località è influenzata sia dall'ordine degli accessi sia dalla struttura dei dati. Lo stesso calcolo matriciale può essere scritto in Matlab come riportato sotto; il codice Matlab differisce dal codice C perché sono gli elementi delle colonne delle matrici a essere memorizzati in locazioni adiacenti.

```
for I=1:8
    for J=1:8000
        A(I,J) = B(I,0) + A(J,I);
    end
end
```

**5.1.4** [5] <5.1> L'accesso a quali variabili gode della località temporale?

**5.1.5** [5] <5.1> L'accesso a quali variabili gode della località spaziale?

**5.1.6** [15] <5.1> Quante linee di cache contenenti 16 byte sono necessarie per memorizzare tutti gli elementi della matrice di 64 bit a cui il programma accede, utilizzando la modalità di memorizzazione di Matlab? Quante utilizzando la modalità di memorizzazione del C? (Si supponga che ciascuna linea contenga più di un elemento.)

**5.2** Le cache sono importanti per fornire ai processori una gerarchia delle memorie ad alte prestazioni. Di seguito viene riportata una sequenza di indirizzi della memoria su 64 bit, espressi come indirizzi di parole.

0x03, 0xb4, 0x2b, 0x02, 0xbf, 0x58,  
0xbe, 0x0e, 0xb5, 0x2c, 0xba, 0xfd

**5.2.1** [10] <5.3> Per ciascuno di questi indirizzi di memoria, identificare il corrispondente indirizzo binario, il tag e l'indice per una cache a mappatura diretta contenente 16 blocchi di una parola. Riportare anche se ciascun accesso si risolve in una hit o una miss, supponendo che la cache sia inizialmente vuota.

**5.2.2** [10] <5.3> Per ciascuno di questi indirizzi di memoria, identificare il corrispondente indirizzo binario, il tag e l'indice per una cache a mappatura diretta contenente otto blocchi di due parole ciascuno. Stabilire anche se ogni accesso si risolve in una hit o una miss, supponendo che la cache sia inizialmente vuota.

**5.2.3** [20] <5.3, 5.4> Si deve ottimizzare la struttura della cache per la sequenza di accessi alla memoria riportata sopra. Si può scegliere tra tre tipi di cache a mappatura diretta, tutte contenenti dati per un totale di otto parole:

- la cache C1 ha blocchi di una parola,
- C2 ha blocchi di due parole
- C3 ha blocchi di quattro parole.

**5.3** Per convenzione, la dimensione di una cache indicata si riferisce alla quantità di dati che contiene (per es. una cache da 4 KiB può contenere 4 KiB di dati); tuttavia, le cache richiedono memoria SRAM anche per memorizzare dei metadati come i tag e i bit di validità. In questo esercizio, si esaminerà come la configurazione di una cache influenzi la quantità totale di SRAM richiesta per implementarla e le prestazioni della cache. Per tutti i problemi si supponga che le cache siano indirizzabili al byte e che gli indirizzi e le parole siano di 64 bit.

**5.3.1** [10] <5.3> Calcolare il numero totale di bit richiesti per implementare una cache da 32 KiB con blocchi di due parole.

**5.3.2** [10] <5.3> Calcolare il numero totale di bit richiesti per implementare una cache da 64 KiB con blocchi da 16 parole. Di quanto è più grande questa cache rispetto alla cache da 32 KiB descritta nell'Esercizio 5.3.1? (Si noti che cambiando la dimensione del blocco, si raddoppia la quantità di dati senza raddoppiare la dimensione totale della cache.)

**5.3.3** [5] <5-3> Spiegare perché questa cache da 64 KiB, nonostante la dimensione maggiore dei dati, potrebbe fornire prestazioni più lente della prima cache.

**5.3.4** [10] <5.3, 5.4> Generare una serie di richieste di lettura che abbiano una frequenza di miss inferiore in una cache set-associativa a due vie da 32 KiB che nella cache descritta nell'Esercizio 5.3.1.

**5.4** [15] <5.3> Il paragrafo 5.3 mostra il tipico metodo per indicizzare una cache a mappatura diretta, e precisamente (indirizzo di blocco) modulo (numero di blocchi nella cache). Supponendo indirizzi su 64 bit e che la cache contenga 1024 blocchi, consideriamo una differente funzione di indicizzazione, precisamente (indirizzo di blocco[63:54]) XOR (indirizzo di blocco[53:44]). È possibile utilizzare questa funzione per indicizzare una cache a mappatura diretta? Se è possibile, spiegare perché e discutere tutte le modifiche necessarie da apportare alla cache. Se non è possibile spiegare perché.

**5.5** I seguenti gruppi di bit degli indirizzi su 64 bit, suddivisi nei diversi campi, vengono utilizzati per l'accesso a una cache a mappatura diretta.

Tag	Indice	Offset
63-10	9-5	4-0

**5.5.1** [5] <5.3> Qual è la dimensione della linea della cache (in numero di parole)?

**5.5.2** [5] <5.3> Quanti elementi contiene la cache?

**5.5.3** [5] <5.3> Qual è il rapporto tra il numero totale di bit contenuti in questa cache e il numero di bit utilizzati per memorizzare i dati?

A partire dall'accensione del calcolatore, si verifica una sequenza di accessi alla cache con i seguenti indirizzi, riferiti al byte (vedi tabella riportata in fondo alla pagina).

**5.5.4** [20] <5.3> Per ciascun indirizzo, riportare (1) il suo tag, indice e offset, (2) se l'accesso risulta in una hit o miss, e (3) quali byte vengono sostituiti (in caso di sostituzione).

**5.5.5** [5] <5.3> Qual è la frequenza di hit?

**5.5.6** [5] <5.3> Riportare lo stato finale della cache, indicando (per ogni blocco che contiene un dato valido) indice, tag e dato contenuto. Per esempio:

<0, 3, Mem[0xC00]-Mem[0xC1F]>

**5.6** Si ricordi che sono due le possibili politiche di scrittura e di allocazione in scrittura, e che diverse loro combinazioni possono essere implementate nella cache L1 o L2. Supponiamo che siano state effettuate queste scelte:

L1	L2
Write-through, No-write-allocate	Write-back, Write-allocate

**5.6.1** [5] <5.3, 5.8> I buffer vengono utilizzati per ridurre la latenza di accesso tra livelli di memoria adiacenti. Per le strategie riportate in tabella, elencare i buffer necessari tra la cache L1 e la cache L2 e tra la cache L2 e la memoria.

**5.6.2** [20] <5.3, 5.8> Descrivere la procedura per la gestione delle miss in scrittura di L1, considerando il componente coinvolto e la possibilità di sostituire i blocchi che sono stati modificati (quelli per i quali è impostato a 1 il dirty bit).

**5.6.3** [20] <5.3, 5.8> Per una cache multilivello di tipo esclusivo, cioè nella quale un blocco può essere contenuto o nella cache L1 o nella cache L2, descrivere la

Indirizzo												
Hex	00	04	10	84	E8	A0	400	1E	8C	C1C	B4	884
Dec	0	4	16	132	232	160	1024	30	140	3100	180	2180

procedura per la gestione delle miss in scrittura di L1, considerando i componenti coinvolti e la possibilità di sostituire i blocchi che sono stati modificati.

**5.7** Si consideri il programma e il comportamento della cache con le caratteristiche riportate nella tabella in fondo alla pagina.

**5.7.1** [10] <5.3, 5.8> Si supponga che una CPU con una cache di tipo write-through e write-allocate, raggiunga un CPI di 2. Qual è la larghezza di banda in lettura e scrittura, misurata in numero di byte per ciclo di clock, tra la RAM e la cache? (Si supponga che ciascuna miss generi una richiesta per ciascun blocco.)

**5.7.2** [10] <5.3, 5.8> Per una cache di tipo write-back e write-allocate, qual è la minima larghezza di banda in lettura e scrittura, richiesta per ottenere un CPI di 2, supponendo che il 30% dei blocchi di dati sostituiti siano blocchi modificati?

**5.8** Le applicazioni multimediali che riproducono file audio o video, fanno parte di una famiglia di carichi di lavoro detti *streaming*, letteralmente “flusso continuo”, perché acquisiscono dall'esterno una grande quantità di dati ma ne riutilizzano pochi. Si consideri il carico di lavoro di uno streaming video di 512 KiB che accede sequenzialmente ai seguenti indirizzi:

0, 1, 2, 3, 4, 5, 6, 7, 8, 9, ...

**5.8.1** [10] <5.4, 5.8> Si supponga di avere a disposizione una cache a mappatura diretta di 64 KiB con linee di 32 byte. Quale sarà la frequenza di miss per la sequenza di indirizzi riportata in precedenza? Quanto è sensibile la frequenza di miss alla dimensione della cache o dello stream? Come si classificano le miss di questo carico di lavoro secondo il modello delle 3C?

**5.8.2** [5] <5.1, 5.8> Ricalcolare la frequenza di miss quando la linea della cache diventa di 16 byte, di 64 byte o di 128 byte. Quale tipo di località viene sfruttata da questo carico di lavoro?

**5.8.3** [10] <5.13> L'anticipazione del caricamento (*prefetching*) è una tecnica che fa leva sulla predicitività delle sequenze di indirizzi per caricare in anticipo i dati nelle linee della cache, ogni volta che si accede a una linea della cache stessa. Un esempio di caricamento anticipato è rappresentato dal buffer di

stream: quando vengono caricati i dati in una linea della cache, viene caricato anche il contenuto delle linee successive in un buffer separato. Se i dati si trovano nel buffer di prefetch, viene generata una hit, i dati vengono copiati nella cache e la linea successiva della cache viene caricata nel buffer di prefetch. Si supponga di avere un buffer di stream a due elementi e che la latenza della cache sia tale per cui una linea di cache può essere caricata prima che l'elaborazione del contenuto della linea precedente sia stata completata. Qual è la frequenza di miss per la sequenza di indirizzi riportata?

**5.9** La dimensione del blocco della cache, che chiameremo B, influenza sia la frequenza di miss sia la latenza di miss. Data una macchina con un CPI pari a 1 e una media di 1,35 accessi cumulativi per le istruzioni e i dati, individuare la dimensione ottimale del blocco della cache, date le frequenze di miss riportate nella tabella seguente in funzione della dimensione dei blocchi.

8: 4%	16: 3%	32: 2%	64: 1,5%	128: 1%
-------	--------	--------	----------	---------

**5.9.1** [10] <5.3> Qual è la dimensione ottimale del blocco per una latenza di miss di  $20 \times B$  cicli di clock?

**5.9.2** [10] <5.2> Qual è la dimensione ottimale del blocco per una latenza di miss di  $24 + B$  cicli di clock?

**5.9.3** [10] <5.2> Data una certa latenza di miss, costante, qual è la dimensione ottimale del blocco?

**5.10** In questo esercizio esamineremo i diversi modi in cui la capacità di una cache influenza le prestazioni complessive. In generale, il tempo di accesso di una cache è proporzionale alla sua capacità. Si supponga che gli accessi alla memoria principale richiedano 70 ns e che costituiscano il 36% di tutte le istruzioni. La tabella seguente mostra i dati relativi a una cache L1 collegata a due processori diversi: P1 e P2.

	Dimensione L1	Frequenza di miss L1	Tempo di hit L1
P1	2 KiB	8,0%	0,66 ns
P2	4 KiB	6,0%	0,90 ns

**5.10.1** [5] <5.4> Supponendo che il tempo di hit di L1 determini la durata del periodo del clock di P1 e P2, quale sarà la loro frequenza di clock?

Lettura dati ogni 1000 istruzioni	Scrittura dati ogni 1000 istruzioni	Frequenza di miss sulla cache istruzioni	Frequenza di miss sulla cache dati	Dimensione del blocco (byte)
250	100	0,30%	2%	64

**5.10.2 [5] <5.4>** Qual è l'*Average Memory Access Time* (AMAT – tempo medio di accesso alla memoria) di P1 e di P2 (in cicli di clock)?

**5.10.3 [5] <5.4>** Supponendo un CPI di base di 1,0 senza stalli della memoria, quale sarà il CPI totale per P1 e P2? Quale dei due processori è più veloce? (Quando si dichiara un “CPI di base pari a 1,0”, intendiamo che le istruzioni vengono completate in un ciclo di clock, a meno che l’accesso ai dati in memoria da parte dell’istruzione non causi una miss.)

Per i tre problemi seguenti, si aggiunga una cache L2 a P1, per esempio per compensare la limitata capacità della cache L1. Utilizzare la capacità e il tempo di hit della cache L1 riportati nella precedente tabella per risolvere questi problemi. La frequenza di miss di L2 riportata nella tabella seguente è la frequenza di miss locale.

Dimensione L2	Frequenza di miss L2	Tempo di hit L2
1 MiB	95%	5,62 ns

**5.10.4 [10] <5.4>** Qual è l’AMAT di P1 quando si aggiunge la cache L2? L’AMAT è migliore o peggiore con la cache L2?

**5.10.5 [5] <5.4>** Supponendo un CPI di base di 1,0 senza stalli della memoria, quale sarà il CPI globale per P1 con l’aggiunta della cache L2?

**5.10.6 [10] <5.4>** Quale dovrebbe essere la frequenza di miss di L2 perché il processore P1 con una cache L2 risulti più veloce del processore P1 senza una cache L2?

**5.10.7 [15] <5.4>** Quale dovrebbe essere la frequenza di miss di L2 perché P1 con una cache L2 risulti più veloce del processore P2 senza una cache L2?

**5.11** Questo esercizio esamina l’impatto di diversi schemi di cache; in particolare verranno confrontate le cache associative con le cache a mappatura diretta descritte nel paragrafo 5.4. Per questi esercizi si faccia riferimento alla sequenza di indirizzi seguente:

0x03, 0xb4, 0x2b, 0x02, 0xbe, 058, 0xbf,  
0x0e, 0x1f, 0xb5, 0xbf, 0xba, 0x2e, 0xce

**5.11.1 [10] <5.4>** Delineare la struttura di una cache set-associativa a tre vie con blocchi di due parole e una dimensione totale di 48 parole. La struttura dovrà essere simile a quella riportata in Figura 5.18, ma dovrà essere indicata chiaramente l’ampiezza del campo tag e dei dati.

**5.11.2 [10] <5.4>** Riportare il comportamento della cache dell’Esercizio 5.11.1. Si supponga una politica di sostituzione completamente LRU. Per ciascun accesso a cache identificare:

sostituzione completamente LRU. Per ciascun accesso a cache identificare:

- l’indirizzo binario della parola;
- il tag;
- l’indice;
- l’offset;
- se l’accesso si risolve in una hit o in una miss;
- quale tag è contenuto in ciascuna delle tre vie della cache quando è terminata la gestione dell’accesso alla cache.

**5.11.3 [5] <5.4>** Delineare la struttura di una cache completamente associativa con blocchi di una parola e una dimensione totale di otto parole. La struttura dovrà essere simile a quella riportata in Figura 5.18, ma dovranno essere indicati chiaramente l’ampiezza del campo tag e dei dati.

**5.11.4 [10] <5.4>** Riportare il comportamento della cache dell’Esercizio 5.11.3. Si supponga una politica di sostituzione completamente LRU. Per ciascun accesso a cache identificare:

- l’indirizzo binario della parola;
- il tag;
- l’indice;
- l’offset;
- se l’accesso si risolve in una hit o in una miss;
- il contenuto della cache quando è terminata la gestione di ciascun accesso alla cache.

**5.11.5 [5] <5.4>** Delineare la struttura di una cache completamente associativa con blocchi di due parole e una dimensione totale di otto parole. La struttura dovrà essere simile a quella riportata in Figura 5.18, ma dovranno essere indicati chiaramente l’ampiezza del campo tag e dei dati.

**5.11.6 [10] <5.4>** Riportare il comportamento della cache dell’Esercizio 5.11.5. Si supponga una politica di sostituzione completamente LRU. Per ciascun accesso a cache identificare:

- l’indirizzo binario della parola;
- il tag;
- l’indice;
- l’offset;
- se l’accesso si risolve in una hit o in una miss;
- il contenuto della cache quando è terminata la gestione di ciascun accesso alla cache.

**5.11.7 [10] <5.4>** Ripetere l’Esercizio 5.11.6 utilizzando la politica di sostituzione MRU (*Most Recently Used* – utilizzato più di recente).

**5.11.8 [15] <5.4>** Ripetere l’Esercizio 5.11.6 utilizzando la politica ottima di sostituzione (cioè quella che fornisce la frequenza di miss più bassa).

**5.12** Le cache multilivello rappresentano una tecnica importante per superare il problema della dimensione limitata dello spazio che una cache di primo livello può mettere a disposizione senza perdere in velocità. Si consideri un processore con i parametri indicati nella tabella che segue.

CPI di base, nessuno stallo della memoria	Velocità del processore	Tempo di accesso alla memoria principale	Frequenza di miss per istruzione della cache di primo livello*	Velocità della cache di secondo livello a mappatura diretta	Frequenza di miss globale con una cache di secondo livello a mappatura diretta	Velocità della cache di secondo livello, set-associativa a otto vie	Frequenza di miss globale con una cache di secondo livello, set-associativa a otto vie
1,5	2 GHz	100 ns	7%	12 cicli	3,5%	28 cicli	1,5%

\* Frequenza di miss per istruzione per la cache di primo livello. Si supponga che il numero totale di miss della cache L1 (istruzioni e dati) sia uguale al 7% del numero di istruzioni.

**5.12.1** [10] <5.4> Calcolare il CPI dei due processori descritti in tabella, utilizzando: (1) solamente una cache di primo livello, (2) una cache di secondo livello a mappatura diretta e (3) una cache di secondo livello set-associativa a otto vie. Come varia il CPI se il tempo di accesso alla memoria principale raddoppia? E se viene dimezzato? (Fornire ogni valore come CPI assoluto e percentuale di variazione.) Si noti fino a che punto una cache L2 può nascondere gli effetti di una memoria lenta.

**5.12.2** [10] <5.4> È possibile avere una gerarchia delle memorie su un numero di livelli maggiore di due? Un progettista vorrebbe aggiungere una cache di terzo livello al processore riportato nella precedente tabella, dotato di una cache di secondo livello a mappatura diretta; questa cache di terzo livello richiederebbe 50 cicli di clock per l'accesso e produrrebbe una frequenza di miss globale all'1,3%. Questa modifica fornirebbe prestazioni migliori? In generale, quali sono i vantaggi e gli svantaggi nell'aggiungere un terzo livello di cache?

**5.12.3** [20] <5.4> Nei vecchi processori, quali il Pentium di Intel o l'Alpha 21264, la cache di secondo livello era esterna al processore, contenuta in un chip separato da quello contenente il processore e la cache di primo livello. Questo consentiva di avere cache di secondo livello di grandi dimensioni, ma la latenza per accedere a queste cache era molto maggiore e la larghezza di banda tipicamente inferiore, perché la cache di secondo livello lavorava a una frequenza inferiore. Si supponga di avere a disposizione una cache di

secondo livello di 512 KiB su un chip separato con una frequenza di miss del 4%. Se ogni 512 KiB aggiunti alla cache riducesse dello 0,7% la frequenza di miss e la cache secondaria avesse un tempo di accesso totale di 50 cicli di clock, quanto dovrebbe essere grande la cache per potere avere le stesse prestazioni della cache di secondo livello a mappatura diretta riportata nella tabella?

**5.13** Il tempo medio tra due guasti (MTBF, *Mean Time Between Failures*), il tempo medio di sostituzione (MTTR, *Mean Time To Replacement*) e il tempo di funzionamento medio prima che si verifichi un guasto (MTTF, *Mean Time To Failure*) sono misure utili per valutare l'affidabilità e la disponibilità di una risorsa di memoria di massa. Esploreremo queste misure risolvendo alcuni problemi riguardanti un dispositivo che ha le seguenti specifiche:

MTTF	MTTR
3 anni	1 giorno

**5.13.1** [5] <5.5> Calcolare l'MTBF per questo dispositivo.

**5.13.2** [5] <5.5> Calcolare la disponibilità di questo dispositivo.

**5.13.3** [5] <5.5> Che cosa succede alla disponibilità se l'MTTR si avvicina a 0? È una situazione realistica?

**5.13.4** [5] <5.5> Che cosa succede alla disponibilità se l'MTTR diventa molto elevato, per esempio perché il dispositivo è difficile da riparare? Un MTTR molto elevato implica che il dispositivo abbia una disponibilità bassa?

**5.14** Questo esercizio esamina il codice di Hamming per la *correzione dell'errore singolo* e il *rilevamento dell'errore doppio* (SEC/DED, *Single Error Correcting/Double Error Detection*).

**5.14.1** [5] <5.5> Qual è il numero minimo di bit di parità richiesti per proteggere una parola di 128 bit utilizzando il codice SEC/DED?

**5.14.2** [5] <5.5> Il paragrafo 5.5 riporta che i moduli di memoria (DIMM) nei server moderni utilizzano il codice ECC SEC/DED per proteggere dati su 64 bit con 8 bit di parità. Determinare il rapporto costo/prestazioni di questo codice e del codice utilizzato nell'Esercizio 5.14.1. In questo caso, il costo è rappresentato dal numero relativo di bit di parità richiesti, mentre le prestazioni sono rappresentate dal numero relativo di errori che possono essere corretti. Quale dei due è migliore?

Anno	Costo DRAM (\$/Mib)	Dimensione della pagina (Kib)	Costo del disco (\$)	Frequenza di accesso a disco (accessi/s)
1987	5000	1	15000	15
1997	15	8	2000	64
2007	0,05	64	80	83

**5.14.3 [5] <5.5>** Si consideri un codice SEC che protegge parole di 8 bit con 4 bit di parità. Se leggiamo il numero 0x375 si è verificato un errore? Se sì, correggere l'errore.

**5.15** Per un sistema ad alte prestazioni quale l'individuazione mediante B-albero (*B-tree*), o albero bilanciato, di un database, la dimensione della pagina è determinata principalmente dalla dimensione dei dati e dalle prestazioni del disco. Si supponga che, in media, una pagina contenente l'albero bilanciato degli indici sia piena al 70% e che contenga elementi di dimensione prefissata. L'utilità di una pagina sta nel livello di profondità del B-albero associato ed è pari al  $\log_2$  del numero di elementi. La tabella seguente mostra che la dimensione ottimale della pagina è di 16 KiB, per dati su 16 byte su un disco vecchio di 10 anni con una latenza di 10 ms e una velocità di trasferimento di 10 MB/s.

Dimensione pagina (KiB)	Utilità di pagina o profondità del B-albero (n. di accessi a disco salvati)	Costo di accesso a indice di pagina (ms)	Utilità/Costo
2	6,49 [ $0 \log_2(2048/16 \times 0,7)$ ]	10,2	0,64
4	7,49	10,4	0,72
8	8,49	10,8	0,79
16	9,49	11,6	0,82
32	10,49	13,2	0,79
64	11,49	16,4	0,70
128	12,49	22,8	0,55
256	13,49	35,6	0,38

**5.15.1 [10] <5.7>** Quale sarebbe la dimensione ottimale della pagina se gli elementi diventassero di 128 bit?

**5.15.2 [10] <5.7>** Seguendo il ragionamento dell'Esercizio 5.15.1, quale sarebbe la dimensione ottimale della pagina se le pagine fossero piene a metà?

**5.15.3 [20] <5.7>** Seguendo il ragionamento dell'Esercizio 5.15.2, quale sarebbe la dimensione ottimale della pagina se si utilizzasse un disco più recente con una latenza di 3 ms e una velocità di trasferimento di 100 MB/s? Spiegare perché è verosimile che i server futuri abbiano pagine di dimensioni maggiori.

Mantenere nella DRAM le pagine utilizzate più di frequente, dette anche "pagine calde" (*hot page*), può ridurre il numero di accessi a disco; ma come possiamo determinare le hot page per un certo sistema? I progettisti utilizzano il rapporto tra il costo d'accesso alla DRAM e quello d'accesso a disco per determinare una soglia sul tempo di riutilizzo di una pagina; sotto questa soglia una pagina si può considerare "calda". Il costo di un accesso a disco viene misurato come costo del disco/ accessi al secondo, mentre il costo per mantenere una pagina in DRAM viene misurato come costo per MiB della DRAM/dimensione della pagina. Il costo tipico delle DRAM e dei dischi e le dimensioni caratteristiche delle pagine dei database riferiti a tre anni particolari sono riportati nella tabella in alto alla pagina.

**5.15.4 [20] <5.7>** Quali altri fattori si possono modificare per continuare a mantenere la stessa dimensione della pagina, evitando perciò di dover riscrivere il software? Discutere quanto siano adatti alla tecnologia corrente considerando il trend e i costi.

**5.16** Come descritto nel paragrafo 5.7, la memoria virtuale utilizza una tabella delle pagine per tracciare la traduzione degli indirizzi virtuali in indirizzi fisici. Questo esercizio mostra come questa tabella debba essere aggiornata quando gli indirizzi vengono utilizzati. La tabella in fondo alla pagina e quelle che seguono nella pagina accanto contengono una sequenza di indirizzi virtuali visti da un sistema. Si supponga che la pagina sia di 4 KiB, il TLB sia completamente associativo e contenga 4 elementi, e che venga utilizzata una strategia di sostituzione LRU non approssimata. Se una pagina deve essere copiata dal disco in memoria, incrementare il numero di pagina più grande, tra quelli delle pagine che seguono.

Decimale	4669	2227	13916	34587	48870	12608	49225
Eseadecimale	0x123d	0x08b3	0x365c	0x871b	0xbbe6	0x3140	0xc049

TLB

Valido	Tag	Numero di pagina fisica	Tempo dall'ultimo accesso
1	0xb	12	4
1	0x7	4	1
1	0x3	6	3
0	0x4	9	7

Tabella delle pagine

Indice	Valida	Numero di pagina fisica o su disco
0	1	5
1	0	Disco
2	0	Disco
3	1	6
4	1	9
5	1	11
6	0	Disco
7	1	4
8	0	Disco
9	0	Disco
a	1	3
b	1	12

**5.16.1 [10] <5.7>** Per ciascuno degli accessi mostrati, riportare:

- se si verifica una hit o una miss del TLB;
- se si verifica una hit o una miss della tabella delle pagine;
- se si verifica un page fault;
- lo stato aggiornato del TLB.

**5.16.2 [15] <5.7>** Risolvere l'Esercizio 5.16.1, utilizzando una pagina da 16 KiB invece che da 4 KiB. Qual è il vantaggio di utilizzare una pagina di dimensioni maggiori? Quali sono alcuni degli svantaggi?

**5.16.3 [15] <5.7>** Risolvere l'Esercizio 5.16.1, utilizzando una pagina da 4 KiB e un TLB set-associativo a due vie.

**5.16.4 [15] <5.7>** Risolvere l'Esercizio 5.16.1, utilizzando una pagina da 4 KiB e un TLB a mappatura diretta.

**5.16.5 [10] <5.4, 5.7>** Spiegare perché una CPU debba avere un TLB per fornire prestazioni elevate. Come verrebbero gestiti gli accessi alla memoria virtuale se non ci fosse un TLB?

**5.17** Diversi sono i parametri che influenzano la dimensione totale della tabella delle pagine; alcuni

parametri chiave sono riportati nella tabella seguente:

Dimensione dell'indirizzo virtuale	Dimensione della pagina	Dimensione dell'elemento della pagina
32 bit	8 KiB	4 byte

**5.17.1 [5] <5.7>** Dati i parametri riportati nella tabella precedente, calcolare la dimensione totale della tabella delle pagine per un sistema che sta eseguendo 5 applicazioni che utilizzano metà della memoria disponibile.

**5.17.2 [10] <5.7>** Dati i parametri riportati nella tabella precedente, calcolare la dimensione totale della tabella delle pagine per un sistema che sta eseguendo 5 applicazioni che utilizzano metà della memoria disponibile; si consideri un sistema con una tabella delle pagine su due livelli contenente 256 elementi. Si supponga che ciascun elemento della tabella delle pagine principale sia di 6 byte. Calcolare la minima e a massima quantità di memoria richiesta.

**5.17.3 [10] <5.7>** Un progettista vuole aumentare la dimensione di una memoria cache, attualmente di 4 KiB, contenente tag fisici e indicizzata virtualmente. Considerando le diverse dimensioni di pagina riportate nella tabella precedente, è possibile realizzare una cache a mappatura diretta da 16 KiB, contenente parole di 64 bit per ogni blocco? Che cosa dovrebbe fare il progettista per aumentare la quantità di memoria della cache riservata ai dati?

**5.18** In questo esercizio esamineremo l'ottimizzazione delle tabelle delle pagine in termini di capacità e velocità. La tabella seguente mostra i parametri di un sistema di memoria virtuale.

Indirizzi virtuali (bit)	Memoria fisica DRAM, installata	Dimensione della pagina	Dimensione del PTE (byte)
43	16 GiB	4 KiB	4

**5.18.1 [10] <5.7>** Quanti elementi della tabella delle pagine (PTE, *Page Table Element*) sono contenuti in una tabella delle pagine con un singolo livello? Quanta memoria fisica è richiesta per memorizzare la tabella delle pagine?

**5.18.2 [10] <5.7>** Utilizzando una tabella delle pagine su più livelli si può ridurre la quantità di memoria fisica consumata dalla tabella delle pagine, mantenendo in memoria solamente i PTE che sono attivi. Quanti livelli deve avere la tabella delle pagine se le tabelle di segmento (le tabelle delle pagine del livello superiore)

ID elemento	Valido	Ind Virt Pagina	Modificata	Protezione	Ind Fisico Pagina
1	1	140	1	RW	30
2	0	40	0	RX	34
3	1	200	1	RO	32
4	1	280	0	RW	31

possono avere dimensioni illimitate? Quanti accessi alla memoria sono richiesti per tradurre un indirizzo che non si trova nel TLB?

**5.18.3 [15] <5.7>** Si supponga che i segmenti siano limitati a una dimensione della pagina di 4 KiB (in modo tale che possano essere paginati). Elementi di 4 byte sono sufficientemente ampi per tutte le tabelle delle pagine comprese quelle dei segmenti?

**5.18.4 [10] <5.7>** Quanti livelli di tabelle delle pagine sono richiesti se i segmenti sono limitati a una dimensione di pagina di 4 KiB?

**5.18.5 [15] <5.7>** Una tabella delle pagine invertita può essere utilizzata per ottimizzare ulteriormente lo spazio richiesto e la velocità. Quanti PTE sono necessari per memorizzare una tabella delle pagine? Supponendo un'implementazione mediante tabella *hash*, quale sarà il numero di accessi alla memoria richiesto per gestire una miss del TLB nel caso più comune e nel caso peggiore?

**5.19** La tabella in alto alla pagina mostra il contenuto di un TLB con quattro elementi.

**5.19.1 [5] <5.7>** In quale situazione verrebbe impostato a 0 il bit di validità associato all'elemento 3?

**5.19.2 [5] <5.7>** Che cosa succede quando un'istruzione scrive nella pagina virtuale il cui indirizzo è 30? In quali casi un TLB gestito via software si può rivelare più veloce di un TLB gestito via hardware?

**5.19.3 [5] <5.7>** Che cosa succede quando un'istruzione scrive nella pagina con indirizzo virtuale 200?

**5.20** In questo esercizio esamineremo l'impatto sulla frequenza di miss delle diverse politiche di sostituzione. Si supponga una cache set-associativa a 2 vie contenente 4 blocchi. Considerare la seguente sequenza di indirizzi di parola:

0, 1, 2, 3, 4, 2, 3, 4, 5, 6, 7, 0, 1, 2, 3, 4, 5, 6, 7, 0

**5.20.1 [5] <5.4, 5.8>** Supponendo che venga utilizzata una strategia di sostituzione di tipo LRU, quante hit si verificano in questa sequenza di accessi?

**5.20.2 [5] <5.4, 5.8>** Supponendo che venga utilizzata una strategia di sostituzione di tipo MRU (blocco utilizzato più di recente), quante hit si verificano in questa sequenza di accessi?

**5.20.3 [5] <5.4, 5.8>** Simulare una politica di sostituzione casuale ottenuta lanciando in aria una moneta; per esempio, se esce "testa" viene eliminato il primo blocco di un insieme e se esce "croce" viene eliminato il secondo blocco di un insieme. Quante hit si verificano in questa sequenza di accessi?

**5.20.4 [10] <5.4, 5.8>** Quale blocco deve essere eliminato a ogni sostituzione in modo da massimizzare il numero di hit? Quante hit si verificherebbero in questa sequenza di accessi se si seguisse questa strategia ottimale?

**5.20.5 [10] <5.4, 5.8>** Spiegare perché è difficile implementare una politica di sostituzione dei blocchi di una cache che risulti ottimale per tutte le possibili sequenze di accessi.

**5.20.6 [10] <5.4, 5.8>** Si supponga di poter decidere per ogni accesso alla memoria se inserire in cache il dato richiesto o meno. Quale impatto avrebbe ciò sulla frequenza di miss?

**5.21** Uno dei maggiori ostacoli alla diffusione delle macchine virtuali è la diminuzione delle prestazioni dovute all'esecuzione della macchina virtuale stessa. Nella tabella in fondo alla pagina sono riportati alcuni parametri di valutazione delle prestazioni e alcune caratteristiche delle applicazioni.

**5.21.1 [10] <5.6>** Calcolare il CPI per il sistema riportato nella tabella sottostante supponendo che non ci siano accessi al sistema di I/O. Qual è il CPI se raddoppia l'impatto sulle prestazioni del VMM? E se viene

CPI di base	Accessi a istruzioni privilegiate del SO ogni 10 000 istruzioni	Impatto sulle prestazioni delle eccezioni di trap sul SO ospitato	Impatto sulle prestazioni delle eccezioni di trap sul VMM	Accesso all'I/O ogni 10 000 istruzioni	Tempo di accesso dell'I/O, comprendente il tempo per la gestione delle eccezioni di trap per il SO ospitato
1,5	120	15 cicli di clock	175 cicli di clock	30	1100 cicli di clock

ridotto della metà? Se una società volesse produrre software per macchine virtuali che venga eseguito con una diminuzione delle prestazioni del 10%, quale sarebbe la massima durata della gestione dell'eccezione di trap per richiamare il VMM?

**5.21.2 [15] <5.6>** Gli accessi all'I/O hanno spesso un grande impatto sulle prestazioni complessive del sistema. Calcolare il CPI di una macchina utilizzando i parametri di prestazione riportati nella tabella precedente, supponendo un sistema non virtualizzato. Calcolare il CPI per un sistema virtualizzato. Come cambierebbe il CPI se il sistema effettuasse la metà degli accessi di I/O? Spiegare perché la virtualizzazione ha un impatto minore sulle applicazioni legate all'I/O. Spiegare perché applicazioni limitate dall'I/O soffrono di un impatto inferiore della virtualizzazione.

**5.22 [15] <5.6, 5.7>** Confrontare e discutere i concetti su cui sono basate la memoria virtuale e le macchine virtuali. Confrontare gli obiettivi dei due sistemi. Quali sono i vantaggi e gli svantaggi di ciascuno dei due sistemi? Riportare alcuni casi nei quali è desiderabile avere a disposizione una memoria virtuale e altri nei quali è desiderabile avere a disposizione una macchina virtuale.

**5.23 [20] <5.6>** Il paragrafo 5.6 tratta la virtualizzazione sotto l'ipotesi che il sistema virtualizzato sia basato sulla stessa ISA dell'hardware su cui viene eseguito. Tuttavia, uno dei possibili utilizzi della virtualizzazione è l'emulazione delle ISA non native. Un esempio è il QEMU, che è un emulatore di diverse ISA, tra cui MIPS, SPARC e PowerPC. Quali sono alcuni dei problemi che si incontrano in questo tipo di virtualizzazione? È possibile che un sistema emulato possa girare più velocemente che sull'ISA nativo?

**5.24** In questo esercizio esamineremo l'unità di controllo di una cache per un processore dotato di buffer di scrittura. Utilizzare la Figura 5.39 come punto di partenza per progettare la macchina a stati finiti. Si supponga che l'unità di controllo della cache debba gestire la semplice cache a mappatura diretta discussa nel paragrafo 5.9 con l'aggiunta di un buffer di scrittura delle dimensioni di un blocco.

Si ricordi che lo scopo del buffer di scrittura è fornire un supporto alla memorizzazione temporanea dei blocchi in modo tale che il processore non debba aspettare due accessi alla memoria quando si verifica una miss su un blocco di cache che è stato modificato. Invece di scrivere questo blocco direttamente nella memoria principale, il processore scrive nel buffer di scrittura il blocco da scaricare dalla cache in modo tale da iniziare subito a leggere il nuovo blocco dalla

memoria principale. Il blocco scaricato verrà poi trasferito dal buffer di scrittura nella memoria principale mentre il processore sta lavorando.

**5.24.1 [10] <5.8, 5.9>** Che cosa succederebbe se il processore inviasse una richiesta che genera una hit della cache mentre si sta trasferendo un blocco dal buffer di scrittura alla memoria principale?

**5.24.2 [10] <5.8, 5.9>** Che cosa succederebbe se il processore inviasse una richiesta che genera una miss della cache mentre si sta trasferendo un blocco dal buffer di scrittura alla memoria principale?

**5.24.3 [30] <5.8, 5.9>** Progettare una macchina a stati finiti che utilizzi un buffer di scrittura.

**5.25** La coerenza della cache riguarda il modo in cui i blocchi di una cache vengono visti da diversi processori. La tabella seguente mostra due processori e le loro operazioni di lettura/scrittura su due diverse parole del blocco X della cache. Si supponga inizialmente  $X[0] = X[1] = 0$ .

P1	P2
$X[0] \quad ++; \quad X[1] = 3;$	$X[0] = 5; \quad X[1] \quad += 2;$

**5.25.1 [15] <5.10>** Elencare i valori che può assumere il blocco della cache considerato con un'implementazione corretta di un protocollo che forzi la coerenza. Elencare almeno uno dei valori che può assumere il blocco considerato se il protocollo non garantisce la coerenza.

**5.25.2 [15] <5.10>** Per un protocollo di snooping, riportare una sequenza valida di operazioni su ciascuna coppia processore/cache per completare le operazioni di lettura/scrittura riportate nella tabella precedente.

**5.25.3 [10] <5.10>** Qual è il numero minimo e massimo di miss della cache richiesto per completare le istruzioni di lettura/scrittura riportate sopra?

La consistenza della memoria riguarda il modo in cui vengono visti dati multipli. La tabella seguente mostra due processori e le loro operazioni di lettura/scrittura su diversi blocchi di cache; si supponga che A e B contengano inizialmente 0.

P1	P2
$A = 1; \quad B = 2; \quad A+=2; \quad B++;$	$C = B; \quad D = A;$

**5.25.4 [15] <5.10>** Elencare i valori che possono assumere C e D per un'implementazione che soddisfa i requisiti di coerenza della memoria elencati all'inizio del paragrafo 5.10.

**5.25.5 [15] <5.10>** Riportare una o più coppie di valori di C e D per le quali quei requisiti non sono soddisfatti.

**5.25.6 [15] <5.3, 5.10>** Quale combinazione di strategie di scrittura e di allocazione in scrittura (write-allocate) rende più semplice l'implementazione del protocollo?

**5.26** I multiprocessori su singolo chip (CMP, *Chip MultiProcessors*) contengono più core, ciascuno con la propria cache montata sullo stesso chip. Sui CMP, la cache L2 è il risultato di un interessante compromesso: la tabella seguente mostra la frequenza di miss e la latenza di hit per due benchmark quando si utilizza una cache L2 privata oppure una cache L2 condivisa. Si supponga che la cache L1 abbia una frequenza di miss del 3% e un tempo di accesso di 1 ciclo di clock.

	Privata	Condivisa
Benchmark A, miss per istruzione	10%	4%
Benchmark B, miss per istruzione	2%	1%

La tabella seguente mostra le latenze di hit:

Cache privata	Cache condivisa	Memoria principale
5	20	180

**5.26.1 [15] <5.13>** Qual è il modello migliore di cache per ciascuno di questi due benchmark? Utilizzare dei dati per supportare le proprie conclusioni.

**5.26.2 [15] <5.13>** La larghezza di banda per le comunicazioni tra il chip CMP e gli altri componenti dell'architettura diventa il collo di bottiglia con l'aumentare del numero di core del CMP. Come vengono influenzate le cache private e le cache condivise da questo collo di bottiglia? Scegliere il modello migliore quando la latenza della memoria esterna al chip raddoppia.

**5.26.3 [10] <5.13>** Discutere i vantaggi e gli svantaggi di una cache L2 privata e di una cache L2 condivisa, per un carico di lavoro a singolo thread, multi-thread e multiprogramma; ripetere il ragionamento supponendo di avere a disposizione sul chip del processore anche una cache L3.

**5.26.4 [10] <5.13>** Una cache L2 non bloccante fornisce un miglioramento superiore su un CMP con una cache L2 condivisa o una cache L2 privata? Perché?

**5.26.5 [10] <5.13>** Si supponga che le nuove generazioni di processori raddoppino il numero di core ogni 18 mesi. Di quanto dovrebbe aumentare la larghezza di banda della memoria esterna al chip nei processori che verranno prodotti fra tre anni per mantenere lo stesso livello di prestazioni per singolo core?

**5.26.6 [15] <5.13>** Si consideri l'intera gerarchia delle memorie. Quali tipi di ottimizzazione possono ridurre il numero di miss concorrenti?

**5.27** In questo esercizio mostreremo la definizione del log (registrazione degli eventi) di un web server ed esamineremo le ottimizzazioni del codice per migliorare la velocità di elaborazione del log. La struttura dati utilizzata per il log è definita come segue:

```
struct entry {
    int srcIP; // indirizzo IP remoto
    char URL[128]; // URL richiesto, per esempio
                    // "GET index.html"
    long long refTime; // tempo di riferimento
    int status; // stato della connessione
    char browser[64]; // nome del browser
                    // del client
} log[NUM_ELEM];
```

Si supponga la seguente funzione di elaborazione del log:

```
topK_sourceIP(int ora);
```

**5.27.1 [5] <5.15>** A quali campi di una variabile contenente il log occorre accedere per eseguire la funzione riportata sopra? Supponendo che ci siano blocchi di cache di 64 byte e che non ci sia caricamento anticipato, quante miss della cache per elemento, in media, sono causate da questa funzione?

**5.27.2 [5] <5.15>** Come può essere riorganizzata la struttura dati per migliorare l'utilizzazione della cache e la località degli accessi?

**5.27.3 [10] <5.15>** Fornire un altro esempio di funzione di elaborazione del log che utilizzi preferibilmente una diversa struttura dati. Se entrambe le funzioni fossero importanti, come si potrebbe riscrivere il programma per migliorare le prestazioni globali? Supportare la risposta con frammenti di codice e dati.

**5.28** Per i problemi riportati di seguito, utilizzare i dati pubblicati in "Cache Performance for SPEC CPU2000 Benchmarks" ([www.cs.wisc.edu/multifacet/misc/spec2000cache-data/](http://www.cs.wisc.edu/multifacet/misc/spec2000cache-data/)) e, in particolare, i dati dei benchmark seguenti:

a.	Mesa/gcc
b.	mcf/swim

**5.28.1 [10] <5.15>** Per una cache contenente 64 KiB di dati con diversi gradi di associatività, qual è la frequenza di miss per i diversi tipi di miss (a freddo, di capacità e di conflitto) per ciascuno dei due benchmark?

**5.28.2 [10] <5.15>** Scegliere il grado di associatività da utilizzare per una cache dati L1 da 64 KiB, utilizzata in entrambi i benchmark. Se la cache L1 deve essere a mappatura diretta, scegliere il grado di associatività per la cache L2 da 1 MiB.

Miss del TLB ogni 1000 istruzioni	Latenza delle miss del TLB per la tabella delle pagine annidata	Page fault ogni 1000 istruzioni	Incremento del tempo di gestione dei page fault nella paginazione ombra
0,2	200 cicli di clock	0,001	30000 cicli di clock

**5.28.3 [20] <5.15>** Fornire un esempio nel quale un aumento del grado di associatività comporta una maggiore frequenza di miss. Costruire una configurazione di cache e una sequenza di indirizzi che dimostri questa affermazione.

**5.29** Per supportare macchine virtuali multiple, sono necessari due livelli di virtualizzazione della memoria. Ciascuna macchina virtuale continua a controllare la mappatura degli *indirizzi virtuali* (VA, *Virtual Address*) negli *indirizzi fisici* (PA, *Physical Address*), mentre il supervisore mappa gli *indirizzi fisici* (PA) di ciascuna macchina virtuale negli *indirizzi della macchina* (MA, *Machine Address*) effettivi del calcolatore utilizzato. Per accelerare queste mappature, un approccio software chiamato *shadow paging* (paginazione ombra) duplica la tabella delle pagine di ciascuna macchina virtuale nel supervisore e intercetta le modifiche della mappatura da VA a PA in modo da mantenere consistenti entrambe le copie. Per eliminare la complessità delle tabelle delle pagine ombra, un approccio hardware chiamato *tabella delle pagine annidata* (NPT, *Nested Page Table*) supporta in modo esplicito due tipi di tabelle delle pagine: la prima trasforma i VA in PA e la seconda trasforma i PA in MA; queste tabelle possono essere interrogate dall'hardware.

Si consideri la seguente sequenza di operazioni: (1) creazione di un processo; (2) miss del TLB; (3) page fault; (4) cambio di contesto.

**5.29.1 [10] <5.6, 5.7>** Per questa sequenza di operazioni, che cosa succede nella tabella ombra delle pagine e nella tabella delle pagine annidata?

**5.29.2 [10] <5.6, 5.7>** Supponendo che entrambe le tabelle delle pagine dell'x86, sia quella annidata che quella ospitata dalla macchina virtuale, siano su quattro livelli, quanti accessi alla memoria servono per gestire una miss del TLB per la tabella delle pagine native e per quella annidata?

**5.29.3 [15] <5.6, 5.7>** Quali sono le misure più importanti per le tabelle delle pagine ombra tra: frequen-

za di miss del TLB, latenza di miss del TLB, frequenza di page fault e latenza per la gestione dei page fault? Quali sono le più importanti per le tabelle delle pagine annidate?

La tabella in alto alla pagina mostra i parametri di un sistema di paginazione ombra.

**5.29.4 [10] <5.6>** Quale diventa il CPI per un benchmark con un CPI di esecuzione nativo di 1, se si utilizzano tabelle delle pagine ombra o tabelle delle pagine annidate, supponendo che l'incremento del tempo di esecuzione sia solamente dovuto alla virtualizzazione della tabella delle pagine?

**5.29.5 [10] <5.6>** Quali tecniche possono essere utilizzate per ridurre l'incremento del tempo di esecuzione dovuto alle tabelle delle pagine ombra?

**5.29.6 [10] <5.6>** Quali tecniche si possono utilizzare per ridurre l'incremento del tempo di esecuzione dovuto all'utilizzo delle tabelle delle pagine annidate?

## Risposte alle domande di autovalutazione

**Paragrafo 5.1, pagina 325 – 1 e 4.** (La 3 è falsa perché il costo della gerarchia delle memorie varia a seconda del calcolatore, ma nel 2013 il costo più elevato era solitamente dovuto alla DRAM.)

**Paragrafo 5.3, pagina 344 – 1 e 4:** una penalità di miss inferiore autorizza blocchi di dimensioni inferiori, dato che non c'è una grande latenza da ammortizzare; tuttavia, una larghezza di banda della memoria più ampia è associata solitamente a blocchi più ampi, dato che la penalità di miss diventa solo di poco più grande.

**Paragrafo 5.4, pagina 362 – 1.**

**Paragrafo 5.8, pagina 401 – 2.** (Sia i blocchi di grandi dimensioni sia il prefetching possono ridurre le miss obbligate, perciò l'affermazione 1 è falsa.)

# 6

## Processori paralleli: dai client al cloud

*Cerco di colpire forte con qualsiasi cosa abbia in mano.  
Posso riuscire o fallire. Mi piace vivere più forte che posso.*

Babe Ruth, giocatore di baseball americano

### 6.1 | Introduzione

*"Oltre le Montagne della Luna,  
lungo la Vallata delle Ombre,  
cavalca, cavalca coraggioso"  
rispose l'ombra. "Se cerchi  
l'Eldorado!"*

Edgar Allan Poe, *El Dorado*,  
stanza 4, 1849

**Multiprocessore:** un'architettura contenente almeno due processori. I multiprocessori si contrappongono alle architetture monoprocessoressi, che contengono un singolo processore e che oggi sono sempre più difficili da trovare.

I progettisti dei calcolatori hanno cercato a lungo l'Eldorado della progettazione dei computer, ossia un modo per creare potenti sistemi di elaborazione semplicemente connettendo tra loro tanti piccoli calcolatori. Questa ricerca ha prodotto i **multiprocessori**. Idealmente, gli utenti possono richiedere tanti processori quanti se ne possono permettere e realizzano sistemi con prestazioni proporzionali al numero di processori acquistati. Perciò, il software per i multiprocessori deve essere progettato per poter essere eseguito su un numero variabile di processori. Come accennato nel Capitolo 1, il consumo di energia elettrica è diventato il problema principale sia per i centri di calcolo sia per i microprocessori. La sostituzione di un singolo grande processore, inefficiente dal punto di vista energetico, con tanti piccoli processori che siano più efficienti può produrre prestazioni più elevate per unità di energia sia su larga sia su piccola scala, purché il software riesca a utilizzare i processori in modo efficiente. Nel caso dei multiprocessori, quindi, il miglioramento del rendimento energetico va di pari passo con l'aumento delle prestazioni.

Dato che il software per i multiprocessori deve essere scalabile, alcune architetture possono garantire l'esecuzione dei programmi anche in presenza di guasti dell'hardware. Ossia, se in un'architettura multiprocessoressi uno degli  $n$  processori si guasta, l'architettura continuerà a garantire le funzionalità utilizzando  $n - 1$  processori. Perciò, i multiprocessori possono anche aumentare la disponibilità (vedi Cap. 5).

Elevate prestazioni possono significare un alto throughput per attività tra loro indipendenti; questo tipo di funzionamento è chiamato generalmente **parallelismo a livello di attività** o **parallelismo a livello di processo**. Le attività parallele sono applicazioni a singolo thread e rappresentano uno degli impieghi più importanti e diffusi dei calcolatori paralleli. Un tale approccio è diverso dall'esecuzione di una singola attività su più processori. Utilizzeremo il termine **programma a esecuzione parallela** (o **software parallelo**) per riferirci a un singolo programma che viene eseguito su più processori simultaneamente.

Storicamente le applicazioni in ambito scientifico hanno avuto bisogno di calcolatori sempre più veloci e potenti, e questa classe di programmi ha giustificato, nei decenni passati, lo sviluppo di nuovi calcolatori paralleli. Alcune applicazioni scientifiche possono funzionare semplicemente disponendo di **cluster**, ossia di insiemi di microprocessori contenuti in tanti server indipendenti (par. 6.7). I cluster possono servire anche per applicazioni non di ambito strettamente scientifico che richiedono un utilizzo intensivo di più macchine, come i motori di ricerca, i server web, i server di posta elettronica e i database.

Come si è detto nel Capitolo 1, l'impulso che ha ricevuto lo sviluppo dei multiprocessori è dovuto alle problematiche legate al consumo energetico, poiché è evidente che l'aumento delle prestazioni in futuro dipenderà verosimilmente dal numero di processori presenti sullo stesso chip anziché da frequenze di clock più elevate o miglioramenti nel CPI. Come abbiamo spiegato nel Capitolo 1, queste architetture sono chiamate **microprocessori multicore** (e non microprocessori multiprocessore, presumibilmente per evitare cacofonie). Di conseguenza, i processori presenti su uno stesso chip vengono solitamente chiamati *core*. Si ipotizza che il numero di core per chip possa aumentare seguendo la **Legge di Moore**. I multicore sono quasi sempre **multiprocessori a memoria condivisa (SMP)**, dato che, di solito, condividono lo stesso spazio di indirizzamento fisico. Esamineremo gli SMP nel paragrafo 6.5.

La tecnologia attuale richiede che i programmati che si preoccupano delle prestazioni diventino programmati paralleli, poiché i programmi sequenziali sono inevitabilmente più lenti.

La grande sfida che deve affrontare l'industria è la creazione di hardware e software che rendano facile scrivere un programma a esecuzione parallela corretto, in modo che questo programma sia eseguito in maniera efficiente in termini di prestazioni e di energia all'aumentare del numero di core presenti in un singolo chip.

Questa improvvisa trasformazione delle architetture dei microprocessori ha colto di sorpresa gli addetti ai lavori, per cui assistiamo a una certa confusione nel significato della terminologia comunemente utilizzata. La tabella riportata in Figura 6.1 cerca di fare un po' di chiarezza, riportando il significato dei termini *seriale*, *parallelo*, *sequenziale* e *concorrente*. Le colonne della tabella rappresentano il software (che è sequenziale o concorrente), mentre le righe rappresentano l'hardware (seriale o parallelo). Per esempio, i programmati dei compilatori considerano sequenziale il loro software, poiché i passi successivi di una traduzione sono composti da: analisi lessicale, segmentazione del codice,



## PARALLELISMO

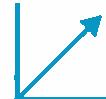
**Parallelismo a livello di attività o parallelismo a livello di processo:** utilizzo di più processori per eseguire simultaneamente diversi programmi indipendenti.

**Programma a esecuzione parallela o software parallelo:** un singolo programma che viene eseguito su più processori simultaneamente.

**Cluster:** un insieme di calcolatori connessi tramite una rete locale (LAN, Local Area Network) che funziona come un unico grande multiprocessore.

**Microprocessore multicore:** un microprocessore che contiene più processori, detti *core*, all'interno di un unico circuito integrato. Oggi, praticamente tutti i microprocessori contenuti nei desktop e nei server sono multicore.

**Multiprocessore a memoria condivisa (SMP):** un processore parallelo con un unico spazio di indirizzamento fisico.



## LEGGE DI MOORE

		Software	
		Sequenziale	Concorrente
Hardware	Serial	Moltiplicazione di matrici in MATLAB in esecuzione su un Pentium 4 Intel	Sistema operativo Windows Vista in esecuzione su un Pentium 4 Intel
	Parallel	Moltiplicazione di matrici in MATLAB in esecuzione su un Core i7 Intel	Sistema operativo Windows Vista in esecuzione su un Core i7 Intel

Figura 6.1 Classificazione hardware/software ed esempi di applicazioni concorrenti e di tipi di hardware.

generazione della traduzione del codice, ottimizzazione e così via. I programmati dei sistemi operativi, invece, generalmente considerano concorrenti i loro programmi, essendo costituiti da processi cooperativi che gestiscono eventi di I/O attraverso attività indipendenti in esecuzione sul calcolatore.

Dalla tabella in Figura 6.1 si evince che il software concorrente può essere eseguito su hardware seriale, come il sistema operativo del Pentium 4 di Intel monoprocesso, oppure su hardware parallelo, come il sistema operativo della più recente architettura Core i7 di Intel. Lo stesso discorso vale per il software sequenziale. Per esempio, un programmatore MATLAB scrive un programma per la moltiplicazione di matrici ragionando in modo sequenziale, ma il programma può essere eseguito sequenzialmente sul Pentium 4 o in parallelo sul Core i7.

Si può intuire che la principale sfida della rivoluzione del parallelismo non sia solo di trovare il modo per ottenere prestazioni elevate su hardware parallelo per software che è naturalmente sequenziale, ma anche di rendere concorrenti i programmi e ottenere prestazioni elevate sui multiprocessori con il crescere del numero dei processori. Dopo aver chiarito questa distinzione, nel seguito del capitolo utilizzeremo il termine *programma a esecuzione parallela* o *software parallelo* per indicare software sequenziale o concorrente eseguito su hardware parallelo. Il prossimo paragrafo illustra il motivo per cui è difficile costruire programmi a esecuzione parallela efficienti.

Prima di procedere nella descrizione del parallelismo, elenchiamo i paragrafi dei capitoli precedenti in cui sono stati trattati argomenti correlati:

- Capitolo 2, paragrafo 2.11: Parallelismo e istruzioni: la sincronizzazione.
- Capitolo 3, paragrafo 3.6: Parallelismo e aritmetica dei calcolatori: parallelismo a livello di parola.
- Capitolo 4, paragrafo 4.10: Parallelismo a livello di istruzioni.
- Capitolo 5, paragrafo 5.10: Parallelismo e gerarchie delle memorie: coerenza delle cache.

### Autovalutazione

Vero o falso? Per sfruttare un multiprocessore, un'applicazione deve essere concorrente.

## 6.2 Le difficoltà nel creare programmi a esecuzione parallela

Il problema principale del parallelismo non è l'hardware, bensì il fatto che pochi programmi applicativi importanti sono stati riscritti per essere eseguiti più velocemente sui multiprocessori. È difficile scrivere software che siano in grado di svolgere un compito più rapidamente utilizzando più processori, e questa difficoltà aumenta al crescere del numero dei processori.

Come mai? Perché i programmi a esecuzione parallela sono più difficili da scrivere dei programmi sequenziali?

Il primo motivo è che si devono ottenere prestazioni più elevate e una maggiore efficienza quando i programmi sono eseguiti in parallelo sui multiprocessori, altrimenti converrebbe utilizzare programmi sequenziali su sistemi monoprocesso, la cui programmazione risulta di gran lunga più semplice. Si

tenga presente che le tecniche avanzate di progettazione dei monoproessori paralleli (struttura superscalare ed esecuzione fuori ordine) sfruttano già i vantaggi del parallelismo a livello di istruzioni (vedi Cap. 4), di solito senza che il programmatore ne sia coinvolto. Queste innovazioni hanno ridotto il bisogno di riscrivere i programmi per i sistemi multiprocessore, poiché, senza fare nulla, i programmi sequenziali sono comunque più veloci sui nuovi calcolatori.

In questo modo, però, non abbiamo risposto alla domanda. Perché è difficile scrivere programmi paralleli che siano veloci, soprattutto se aumenta il numero di processori? Nel Capitolo 1 abbiamo utilizzato l'analogia degli otto giornalisti che lavorano tutti a uno stesso articolo con la speranza che l'articolo sia scritto otto volte più in fretta. Per raggiungere l'obiettivo, il compito deve essere separato in otto parti di uguale dimensione, altrimenti alcuni giornalisti finirebbero per non avere nulla da fare se non aspettare che gli altri abbiano terminato di scrivere la loro parte, più consistente. Un altro problema per le prestazioni potrebbe verificarsi se i giornalisti passassero troppo tempo a comunicare tra loro invece di scrivere la loro parte. Questo esempio ci fa capire che le sfide sono rappresentate dalla programmazione delle attività (*scheduling*), dal bilanciamento dei carichi di lavoro, dal tempo richiesto per la sincronizzazione e dal carico di lavoro aggiuntivo richiesto dalla comunicazione tra le diverse unità. Il problema diventa ancora più complesso mano a mano che aumenta il numero di giornalisti (analogamente, mano a mano che aumentano i processori).

Nel Capitolo 1 abbiamo messo in evidenza l'esistenza di un altro problema, rappresentato dalla Legge di Amdahl. Questa legge ci ricorda che anche le più piccole parti di un programma devono essere rese parallele se si vuole che il programma venga eseguito in modo efficiente su più core.

## Sfida della velocità

Supponiamo di voler ottenere un aumento di velocità di 90 volte con 100 processori. In quale percentuale le operazioni originali possono rimanere sequenziali?

Secondo la Legge di Amdahl (vedi Cap. 1) si ha che:

**ESEMPIO**

**SOLUZIONE**

Tempo di esecuzione dopo il miglioramento =

$$\frac{\text{Tempo di esecuzione influenzato dal miglioramento}}{\text{Quantità di miglioramento}} + \text{Tempo di esecuzione non influenzato dal miglioramento}$$

Possiamo riformulare la Legge di Amdahl in termini di incremento di velocità (*speed-up*) rispetto al tempo di esecuzione iniziale:

$$\text{Speed-up} = \frac{\text{Tempo di esecuzione prima}}{(\text{Tempo di esecuzione prima} - \text{Tempo di esecuzione influenzato}) + \frac{\text{Tempo di esecuzione influenzato}}{\text{Quantità di miglioramento}}}$$

Questa formula di solito viene riscritta supponendo che il tempo di esecuzione prima del miglioramento sia uguale a 1, in una data unità di misura, e che il tempo modificato dal miglioramento sia espresso in termini di una frazione del tempo originario di esecuzione:

(continua)

(continua)

$$\text{Speed-up} = \frac{1}{(1 - \text{Frazione di tempo influenzata}) + \frac{\text{Frazione di tempo influenzata}}{\text{Quantità di miglioramento}}}$$

Inserendo il valore dell'incremento di velocità che vogliamo ottenere (90 volte) e la quantità di miglioramento (100) nella formula precedente, si ha:

$$90 = \frac{1}{(1 - \text{Frazione di tempo influenzata}) + \frac{\text{Frazione di tempo influenzata}}{100}}$$

Semplificando la formula e risolvendo per la frazione di tempo influenzata dal miglioramento, si ottiene:

$$\begin{aligned} 90 \times (1 - 0,99 \times \text{Frazione di tempo influenzata}) &= 1 \\ 90 - (90 \times 0,99 \times \text{Frazione di tempo influenzata}) &= 1 \\ 90 - 1 &= 90 \times 0,99 \times \text{Frazione di tempo influenzata} \\ \text{Frazione di tempo influenzata} &= 89/89,1 = 0,999 \end{aligned}$$

Perciò, per ottenere un incremento di velocità di 90 volte con 100 processori, la percentuale di elaborazione che può rimanere sequenziale non può superare lo 0,1%.

Esistono, però, applicazioni che hanno un sostanziale grado di parallelismo, come vedremo tra poco.

### Sfida dell'aumento di velocità: problemi di grandi dimensioni

#### ESEMPIO

Supponiamo di voler eseguire due somme: la prima è la somma di 10 variabili scalari, la seconda è la somma di due matrici bidimensionali  $10 \times 10$ . Per ora supponiamo che sia parallelizzabile solamente la somma di matrici (vedremo presto come parallelizzare le somme scalari). Che incremento di velocità si ottiene passando da 10 a 40 processori? Calcolare l'aumento di velocità supponendo che la matrice diventi  $20 \times 20$ .

#### SOLUZIONE

Supponendo che le prestazioni siano funzione del tempo relativo all'esecuzione di una singola addizione,  $t$ , le 10 addizioni delle variabili scalari non otterranno alcun beneficio dalla presenza di più processori paralleli, mentre le 100 addizioni relative alle matrici sì. Se il tempo di esecuzione su un singolo processore è di  $110t$ , il tempo di esecuzione su 10 processori sarà:

Tempo di esecuzione dopo miglioramento =

$$= \frac{\text{Tempo di esecuzione influenzato dal miglioramento}}{\text{Quantità di miglioramento}} + \text{Tempo di esecuzione non modificato dal miglioramento}$$

(continua)

(continua)

$$\text{Tempo di esecuzione dopo miglioramento} = \frac{100t}{10} + 10t = 20t$$

per cui l'incremento di velocità utilizzando 10 processori sarà  $110t/20t = 5,5$ . Il tempo di esecuzione con 40 processori è:

$$\text{Tempo di esecuzione dopo miglioramento} = \frac{100t}{40} + 10t = 12,5t$$

per cui l'incremento di velocità quando si passa a 40 processori è quindi  $110t/12,5t = 8,8$ . Perciò, per problemi di questo ordine di grandezza, otteniamo circa il 55% dell'incremento di velocità potenziale utilizzando 10 processori, ma solamente il 22% con 40 processori.

Vediamo adesso che cosa succede nel momento in cui aumentano le dimensioni della matrice. Il programma sequenziale richiede in questo caso un tempo pari a  $10t + 400t = 410t$ . Il tempo di esecuzione su 10 processori è:

$$\text{Tempo di esecuzione dopo miglioramento} = \frac{400t}{10} + 10t = 50t$$

per cui l'incremento di velocità con 10 processori è  $410t/50t = 8,2$ . Il tempo di esecuzione su 40 processori è:

$$\text{Tempo di esecuzione dopo miglioramento} = \frac{400t}{40} + 10t = 20t$$

per cui l'aumento di velocità con 40 processori è di  $410t/20t = 20,5$ . Perciò, nel caso di problemi di grandi dimensioni, otteniamo un aumento di velocità dell'82% rispetto all'aumento potenziale con 10 processori e del 51% con 40 processori.

Questi esempi mostrano che aumentare la velocità di esecuzione sui multiprocessori mantenendo fisse le dimensioni del problema è più difficile che ottenere un miglioramento delle prestazioni incrementando le dimensioni del problema. Si tratta di un comportamento che ci permette di introdurre due termini che descrivono due diverse modalità con cui si valuta l'aumento delle prestazioni: la **scalabilità forte** (*strong scaling*) misura l'incremento di velocità quando si mantiene fissa la dimensione del problema, mentre la **scalabilità debole** (*weak scaling*) si riferisce all'incremento di velocità quando le dimensioni del problema crescono proporzionalmente al numero dei processori. Supponiamo che la dimensione del problema,  $M$ , sia pari allo spazio di lavoro allocato nella memoria principale e che abbiamo a disposizione  $P$  processori; la memoria per ciascun processore, quindi, sarà all'incirca pari a  $M/P$  in caso di scalabilità forte e pari a  $M$  nel caso di scalabilità debole.

Si noti che la **gerarchia delle memorie** fa sì che l'intuizione che la scalabilità debole si possa ottenere più facilmente della scalabilità forte possa rivelarsi falsa. Per esempio, se i dati, a seguito dell'incremento delle dimensioni dovuto alla scalabilità debole, non possono più essere tutti contenuti nella cache di livello

**Scalabilità forte:** incremento di velocità che si ottiene in un multiprocessore senza aumentare la dimensione del problema.

**Scalabilità debole:** incremento di velocità che si ottiene in un multiprocessore quando le dimensioni del problema aumentano proporzionalmente al numero dei processori.



più elevato di un microprocessore multicore, le prestazioni potrebbero essere molto peggiori di quelle ottenute con la scalabilità forte.

A seconda del tipo di applicazione con cui abbiamo a che fare, si possono trovare argomenti a favore di entrambi gli approcci. Per esempio, il benchmark del database delle transazioni debito-credito, TPC-C, richiede che il numero degli account dei clienti cresca con il numero di transazioni nell'unità di tempo. Il motivo è che non ha senso pensare che un cliente tipico inizi improvvisamente a utilizzare il bancomat 100 volte al giorno solamente perché la banca ha velocizzato la sua rete di calcolatori; dovendo dimostrare che un sistema è in grado di eseguire un numero di transazioni nell'unità di tempo 100 volte superiore, è necessario disporre di un numero di clienti che è 100 volte superiore. Di solito problemi più grandi richiedono più dati, che è uno degli argomenti in favore della scalabilità debole.

L'esempio seguente illustra l'importanza del bilanciamento del carico.

### Sfida dell'aumento di velocità: bilanciamento del carico

#### ESEMPIO

Per ottenere un aumento di velocità di 20,5 volte per il problema di grandi dimensioni considerato nell'esempio precedente (con 40 processori), abbiamo supposto che il carico fosse perfettamente bilanciato. Cioè, ognuno dei 40 processori svolgeva il 2,5% del lavoro totale. Mostriamo ora l'impatto sulle prestazioni della situazione in cui il carico di uno dei processori è maggiore di quello degli altri. Svolgere i calcoli quando un processore ha un carico pari al doppio (5%) o a cinque volte (12,5%) il carico totale. Quanto vengono utilizzati gli altri processori?

#### SOLUZIONE

Se un processore riceve il 5% del carico parallelo, allora deve eseguire il 5% di 400 addizioni, ossia 20 addizioni, mentre le rimanenti 380 addizioni saranno suddivise fra gli altri 39 processori. Dato che i processori lavorano simultaneamente, possiamo calcolare il tempo di esecuzione globale come il massimo tempo di esecuzione tra tutti i processori:

$$\text{Tempo di esecuzione dopo miglioramento} = \text{Max}\left(\frac{380t}{39}, \frac{20t}{1}\right) + 10t = 30t$$

L'incremento di velocità scende quindi da 20,5 a  $410t/30t = 14$ . Gli altri 39 processori vengono utilizzati per meno della metà del tempo: mentre aspettano 20t che il processore con il carico maggiore termini, devono effettuare calcoli per soli  $380t/39 = 9,7t$ .

Se un processore avesse il 12,5% del carico, eseguirebbe 50 addizioni. La formula è:

$$\text{Tempo di esecuzione dopo miglioramento} = \text{Max}\left(\frac{350t}{39}, \frac{50t}{1}\right) + 10t = 60t$$

L'incremento di velocità in questo caso scende ulteriormente a  $(410t/60t = 7)$ . Gli altri processori vengono utilizzati per meno del 20% del tempo ( $9t/50t$ ). Questo esempio dimostra l'importanza di un carico bilanciato, dato che un singolo processore con il doppio del carico degli altri riduce l'aumento delle prestazioni di 1/3 e con un carico di 5 volte di quasi un fattore 3.

Ora che abbiamo chiarito gli obiettivi e le sfide della programmazione parallela, forniamo una panoramica sul resto del capitolo. Nel paragrafo 6.3 descriveremo uno schema di classificazione molto più vecchio di quello di Figura 6.1 e due stili di architetture di insiemi di istruzioni che supportano l'esecuzione di applicazioni sequenziali su hardware parallelo: le architetture *SIMD* e le architetture *vettoriali*. Nel paragrafo 6.4 descriveremo il *multithreading*, un termine spesso confuso con multiprocesssing, in parte perché è basato su concetti di concorrenza tra i programmi simili. Il paragrafo 6.5 descrive architetture che implementano la condivisione di un unico spazio di indirizzamento fisico tra tutti i processori. Alternativamente, lo spazio di indirizzamento fisico può essere non condiviso. Le due versioni più popolari di queste due alternative sono chiamate *processori a memoria condivisa* (SMP, *Shared Memory Multiprocessors*) e *cluster*; questo paragrafo descrive i processori SMP. Il paragrafo 6.6 descrive un tipo di calcolatore relativamente recente, proveniente dalla comunità dell'hardware per la grafica, chiamato *GPU* (*Graphics-Processing Unit*) che utilizza anch'esso un unico spazio di indirizzamento fisico (l'Appendice C descrive le GPU ancora più in dettaglio). Il paragrafo 6.8 mostra tipiche topologie per connettere assieme molti processori: nodi server nei cluster o core nei microprocessori. Il paragrafo 6.9 descrive l'hardware e il software necessari per fare comunicare i nodi di un cluster utilizzando Ethernet e mostra come ottimizzare le prestazioni utilizzando software e hardware commerciale. Discuteremo quindi la difficoltà nel reperire benchmark paralleli nel paragrafo 6.10. Questo paragrafo comprende anche un semplice ma potente modello di valutazione delle prestazioni che aiuta nella progettazione delle applicazioni e delle architetture. Utilizzeremo questo modello e benchmark paralleli nel paragrafo 6.11 per comparare un calcolatore multicore con una GPU. Il paragrafo 6.12 descrive l'ultimo e più grande passo del nostro percorso volto ad accelerare la moltiplicazione tra matrici. Per le matrici che non possono essere contenute interamente in cache, l'elaborazione parallela sfrutta 16 core per aumentare le prestazioni di un fattore 14. Chiuderemo il capitolo con il paragrafo su errori e trabocchetti e con le nostre osservazioni sul parallelismo.

Nel prossimo paragrafo introduciamo alcuni acronimi che avete probabilmente già sentito per identificare i diversi tipi di calcolatori paralleli.

## Autovalutazione

Vero o falso? La scalabilità forte non è limitata dalla Legge di Amdahl.

## 6.3 | SISD, MIMD, SIMD, SPMD e processori vettoriali

Un tipo di classificazione dell'hardware parallelo proposto negli anni '60 (e utilizzato ancora oggi) è basato sul numero dei flussi di istruzioni e di dati. La Figura 6.2 mostra le diverse categorie. Secondo questa classificazione, un monoprocessoore convenzionale conterrà un unico flusso di istruzioni e un unico flusso di dati, mentre un multiprocessoare convenzionale conterrà flussi multipli di istruzioni e di dati. Queste due categorie sono chiamate rispettivamente **SISD** e **MIMD**.

Dato un problema di grandi dimensioni, si potrebbero scrivere diversi programmi separati che possono essere eseguiti in parallelo sui diversi processori di un calcolatore MIMD, in modo da cooperare per risolvere il problema dato. Tuttavia, i programmati preferiscono di solito scrivere un unico program-

**SISD:** *Single Instruction stream, Single Data stream* (flusso singolo di istruzioni, flusso singolo di dati). Un monoprocessoore.

**MIMD:** *Multiple Instruction streams, Multiple Data streams* (flussi multipli di istruzioni, flussi multipli di dati). Un multiprocessoare.

		Flussi di dati	
		Singoli	Multipli
Flussi di istruzioni	Singoli	SISD: Pentium 4 Intel	SIMD: Istruzioni SSE dell'x86
	Multipli	MISD: Nessun esempio a oggi	MIMD: Core i7 Intel

**Figura 6.2** Classificazione dell'hardware ed esempi basati sul numero di flussi di istruzioni e di dati: SISD, SIMD, MISD e MIMD.

**SPMD:** *Single Program, Multiple Data streams*. Il modello di programmazione MIMD convenzionale, in cui un singolo programma viene eseguito su più processori.

**SIMD:** *Single Instruction stream, Multiple Data streams*. La stessa istruzione viene applicata a diversi flussi di dati, come nei processori vettoriali.

ma che viene eseguito su tutti i processori del calcolatore MIMD, basandosi sulle istruzioni di salto condizionato per fare eseguire ai diversi processori parti diverse del codice. Questo metodo di programmazione è chiamato **SPMD** (*Single Program Multiple Data*) e costituisce il modo usuale di programmare i calcolatori MIMD.

Mentre è difficile fornire esempi di calcolatori caratterizzati da flussi multipli di istruzioni e flussi singoli di dati (MISD), la situazione opposta (flussi singoli di istruzioni e multipli di dati) è molto più diffusa. I processori **SIMD** (*Single Instruction stream, Multiple Data streams*) operano su vettori di dati. Per esempio, una singola istruzione SIMD potrebbe sommare tra loro 64 coppie di numeri, inviandoli a 64 ALU per ottenere 64 somme, in un singolo ciclo di clock. Le istruzioni che operano sulla parola che abbiamo visto nei paragrafi 3.6 e 3.7 sono un esempio di istruzioni SIMD: la seconda S dell'acronimo SSE utilizzato da Intel sta proprio per SIMD.

Il vantaggio dei processori SIMD è che tutte le unità di elaborazione parallele sono sincronizzate tra loro e rispondono a una singola istruzione che proviene da un unico program counter (PC). Dal punto di vista del programmatore, questo approccio è simile alle architetture SISD già incontrate. Sebbene ciascuna unità funzionale esegua la stessa operazione, le unità funzionali hanno i loro registri degli indirizzi, per cui ciascuna unità può contenere degli indirizzi diversi dei dati. Riferendoci alla Figura 6.1, un'applicazione sequenziale potrebbe essere compilata per essere eseguita su hardware seriale organizzato come un'architettura SISD o su hardware parallelo organizzato come una SIMD.

La motivazione originale che ha portato alla realizzazione delle architetture SIMD è l'ammortamento del costo dell'unità di controllo, che viene distribuito su dozzine di unità funzionali. Un altro vantaggio è rappresentato dalle dimensioni ridotte della memoria dei programmi, poiché le architetture SIMD richiedono solo una copia delle istruzioni, le quali vengono eseguite simultaneamente, mentre le architetture MIMD basate su scambi di messaggi possono richiedere una copia delle istruzioni per ogni processore e le MIMD basate sulla condivisione della memoria richiedono cache istruzioni multiple.

Le architetture SIMD lavorano a un livello ottimale quando devono gestire vettori all'interno di cicli *for*; quindi, affinché il parallelismo in tali architetture funzioni, occorre che i dati siano strutturati in modo simile, per cui si parla di **parallelismo a livello di dati**. Il punto debole delle architetture SIMD è l'esecuzione dei costrutti *case* o *switch*, dove ciascuna unità funzionale deve eseguire un'operazione diversa sui propri dati, differente a seconda del tipo di dato a disposizione. Le unità funzionali abbinate a dati sbagliati nelle SIMD vengono disabilitate, in modo tale che le unità che hanno ricevuto i dati corretti possano continuare l'elaborazione. Questa organizzazione fa sì che il codice venga eseguito a  $1/n$  delle prestazioni di picco, dove  $n$  è il numero dei casi nei costrutti *case* o *switch*.

I cosiddetti processori vettoriali che hanno ispirato le architetture SIMD sono ormai diventati reperti storici (par. 6.15 ), ma le due interpretazioni correnti delle architetture SIMD vengono utilizzate ancora oggi.

**Parallelismo a livello di dati:** parallelismo ottenuto eseguendo la stessa operazione su dati fra loro indipendenti.

## SIMD negli x86: le estensioni multimediali

Come descritto nel Capitolo 3, il parallelismo a livello di parola per dati interi stretti è stato l'ispirazione originale che ha portato all'estensione multimediale (MMX) dell'x86 nel 1996. Continuando ad applicare la Legge di Moore, sono state aggiunte via via nuove istruzioni contenute nell'estensione SSE (*Streaming SIMD Extension*) e più recentemente nell'estensione vettoriale AVX (*Advanced Vector Extension*). AVX supporta l'esecuzione simultanea di quattro operazioni su numeri a 64 bit in virgola mobile. L'ampiezza dell'operazione e dei registri è codificata nel codice operativo di queste istruzioni multimediali. Con l'aumento dell'ampiezza dei registri e del numero di operazioni, il numero dei codici operativi delle istruzioni multimediali è esploso: oggi abbiamo a disposizione centinaia di istruzioni SSE e AVX (vedi Cap. 3).



## Architetture vettoriali

Un'interpretazione più vecchia ma più elegante delle architetture SIMD è rappresentata dalle cosiddette *architetture vettoriali*, spesso identificate con i calcolatori progettati da Seymour Cray a partire dagli anni '70. Queste architetture erano particolarmente adatte ai problemi caratterizzati da un parallelismo spinto a livello dei dati. Invece di avere 64 ALU che eseguono 64 addizioni simultaneamente, come nei vecchi processori vettoriali, le architetture vettoriali mettono in pipeline una singola ALU per ottenere buone prestazioni a un costo inferiore. La filosofia di base delle architetture vettoriali consiste nel leggere i dati dalla memoria, metterli in ordine scrivendoli in un gran numero di registri, elaborarli su una unità di elaborazione dotata di **pipeline**, scrivere il risultato nei registri e infine salvare il risultato in memoria. La caratteristica principale delle architetture vettoriali è quindi il vasto insieme di registri vettoriali a disposizione: un'architettura vettoriale può averne 32, ciascuno contenente 64 elementi da 64 bit.



## Confronto tra codice convenzionale e codice vettoriale

Si supponga di volere estendere l'architettura dell'insieme delle istruzioni RISC-V con istruzioni e registri vettoriali. Le operazioni vettoriali avranno lo stesso nome delle operazioni RISC-V, ma con la lettera "V" aggiunta alla fine del nome dell'istruzione. Per esempio, l'istruzione fadd.d.v somma due vettori in doppia precisione. Aggiungiamo anche 32 registri vettoriali, v0-v31, ciascuno con 64 elementi da 64 bit. Le istruzioni vettoriali accettano come input una coppia di registri vettoriali V (fadd.d.v), oppure un registro scalare e un registro vettoriale (fadd.d.vs). In quest'ultimo caso, il contenuto del registro scalare viene utilizzato come input di tutte le operazioni: l'operazione fadd.d.vs sommerà il contenuto del registro scalare a tutti gli elementi del registro vettoriale. I nomi fld.v e fsd.v indicano le istruzioni di lettura e scrittura di un vettore in memoria e trasferiscono a un intero vettore di dati in doppia precisione. Uno dei due operandi è il registro vettoriale, il cui contenuto deve essere letto o scritto in memoria, mentre l'altro operando, che è contenuto in uno dei registri RISC-V di utilizzo generale, contiene la posizione di partenza del vettore in memoria.

### ESEMPIO

Sulla base di questa breve descrizione, scrivere il codice per un processore RISC-V convenzionale e quello per la versione vettoriale che calcoli

(continua)

(continua)

la seguente espressione:

$$Y = a \times X + Y$$

dove  $X$  e  $Y$  sono due vettori contenenti 64 numeri in doppia precisione, inizialmente residenti in memoria, e  $a$  è una variabile scalare in doppia precisione. Questo tipo di calcolo è chiamato ciclo DAXPY e costituisce il ciclo interno dei benchmark Linpack; DAXPY sta per “doppia precisione  $a \times X$  più  $Y$ ”. Si supponga che l’indirizzo di partenza di  $X$  e  $Y$  sia contenuto rispettivamente in  $x19$  e  $x20$ .

**SOLUZIONE**

Questo è il codice scritto per un RISC-V convenzionale che esegue il ciclo DAXPY:

```

fld    f0, a(x3)      // caricamento dello scalare a
addi   x5, x19, 512   // indirizzo ultimo elemento da caricare
ciclo:fld  f1, 0(x19) // caricamento di x[i]
        fmul.d f1, f1, f0 // a * x[i]
        fld    f2, 0(x20) // caricamento di y[i]
        fadd.d f2, f2, f1 // a * x[i] + y[i]
        fsd    f2, 0(x20) // scrittura in y[i]
        addi   x19, x19, 8  // incremento dell’indice per x
        addi   x20, x20, 8  // incremento dell’indice per y
        bltu  x19, x5, ciclo // controllo sulla fine del ciclo

```

Questo è il codice del ciclo DAXPY scritto per il RISC-V vettoriale:

```

fld      f0, a(x3)  // caricamento dello scalare a
fld.v   v0, 0(x19) // caricamento del vettore X
fmul.d.vs v0, v0, f0 // moltiplicazione vettore-scalare
fld.v   v1, 0(x20) // caricamento del vettore Y
fadd.d.v v1, v1, v0 // somma di Y al prodotto
fsd.v   v1, 0(x20) // memorizzazione del vettore Y

```

Ci sono alcune differenze interessanti tra i due frammenti di codice di questo esempio. Quella più eclatante è che il processore vettoriale riduce enormemente la larghezza di banda richiesta dalle istruzioni, poiché esegue solamente sei istruzioni rispetto alle oltre 500 eseguite dal RISC-V convenzionale. Questa riduzione è dovuta al fatto che le operazioni vettoriali lavorano su 64 elementi alla volta e che le istruzioni di servizio all’interno del ciclo, che costituiscono circa la metà delle istruzioni totali, sono assenti nella versione vettoriale del codice. Come ci si può aspettare, questa riduzione del numero di istruzioni permette anche di risparmiare energia.

Un’altra importante differenza è relativa alla frequenza degli hazard della pipeline (*vedi Cap. 4*). Nel codice RISC-V convenzionale ogni istruzione fadd.d deve attendere il completamento della fmul.d, ogni fsd deve attendere il completamento della fadd.d e tutte le fadd.d e fmul.d devono attendere il completamento della fld. Nel processore vettoriale, invece, ogni istruzione genererà uno stallo solamente per il primo elemento del vettore, mentre gli elementi successivi scorreranno regolarmente nella pipeline. Perciò, lo stallo delle pipeline viene richiesto una sola volta per ogni *operazione* vettoriale e non per ogni *elemento* del vettore. In questo esempio la frequenza di stallo del RISC-V convenzionale sarebbe circa 64 volte superiore a quella del RISC-V vettoriale. Si potrebbero ridurre gli stalli nella pipeline convenzionale utilizzando un’architettura di pipeline più avanzata, ma ciò aumenterebbe il costo del chip.



zando la tecnica dell'espansione dei cicli (*vedi Cap. 4*). In ogni caso, la grande differenza nella larghezza di banda per il caricamento delle istruzioni non può essere colmata.

Dato che gli elementi del vettore sono indipendenti, possiamo operare su di essi in parallelo, come nel caso del parallelismo a livello di parola delle istruzioni AVX. Tutti i moderni calcolatori vettoriali hanno unità funzionali vettoriali con pipeline parallele multiple (chiamate *vector lanes* o *corsie vettoriali*; Figure 6.2 e 6.3) che possono produrre due o più risultati per ogni ciclo di clock.

**Approfondimento.** Il ciclo dell'esempio precedente operava su dati della lunghezza esatta dei vettori. Quando i cicli sono su un numero minore di elementi, le architetture vettoriali utilizzano un registro per ridurre il numero di elementi (del vettore) su cui effettuare le operazioni. Quando i cicli sono su più elementi, occorre aggiungere delle istruzioni di prenotazione per poter iterare le operazioni su tutti gli elementi del vettore e su quelli che rimangono. Questo secondo processo viene chiamato *strip mining* (estrarre strisce).

## Confronto tra architetture vettoriali e scalari

Le istruzioni vettoriali godono delle seguenti importanti proprietà rispetto alle istruzioni delle architetture convenzionali (che in questo contesto chiameremo *architetture scalari*).

- Una singola istruzione vettoriale specifica una grande quantità di lavoro: è equivalente a un intero ciclo. La larghezza di banda richiesta per caricare e decodificare le istruzioni risulta enormemente ridotta.
- Utilizzando istruzioni vettoriali, il compilatore o il programmatore possono specificare che i calcoli eseguiti sugli elementi omologhi di due o più vettori siano indipendenti dai calcoli eseguiti sugli altri elementi dello stesso vettore, per cui l'hardware non deve preoccuparsi di individuare hazard sui dati all'interno della stessa istruzione vettoriale.
- Le architetture vettoriali e i compilatori rendono molto più facile scrivere applicazioni efficienti quando queste contengono un elevato livello di parallelismo a livello dei dati, rispetto ai multiprocessori MIMD.
- L'hardware deve controllare gli hazard sui dati per le istruzioni vettoriali solamente una volta, ossia per il primo elemento del vettore e non per ciascun elemento del vettore. Anche la riduzione del numero dei controlli fa risparmiare energia oltre che tempo.
- Le istruzioni vettoriali che accedono alla memoria utilizzano una sequenza di indirizzi prevedibile. Se gli elementi dei vettori sono tutti adiacenti, il caricamento dei dati del vettore da un insieme di banchi di memoria fortemente interallacciati funziona molto bene. Perciò, il costo della latenza della memoria principale si manifesta una sola volta per l'intero vettore invece che per ogni elemento del vettore.
- Dato che un intero ciclo è sostituito da un'istruzione vettoriale, il cui comportamento è predeterminato, vengono eliminati gli hazard sul controllo originati normalmente dalle istruzioni di salto condizionato che si utilizzano per il controllo di fine ciclo.
- Il risparmio relativo alla larghezza di banda per le istruzioni e al controllo degli hazard, oltre all'utilizzo efficiente della larghezza di banda, fornisce alle architetture vettoriali enormi vantaggi rispetto alle architetture scalari in termini di potenza assorbita ed energia di funzionamento.

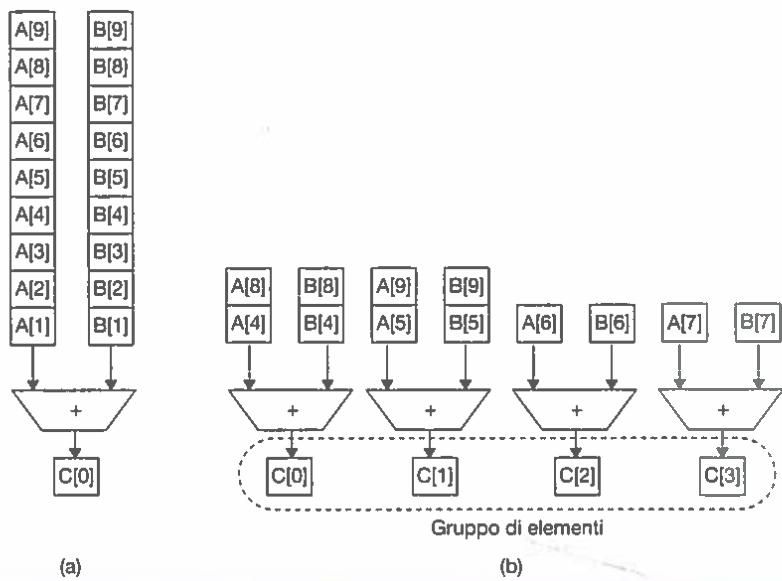
Per questi motivi, le operazioni vettoriali possono essere rese più veloci di una sequenza equivalente di operazioni scalari sugli stessi dati; i progettisti sono quindi indotti a includere unità vettoriali se il dominio di applicazione può farne un uso frequente.

## Processori vettoriali ed estensioni multimediali

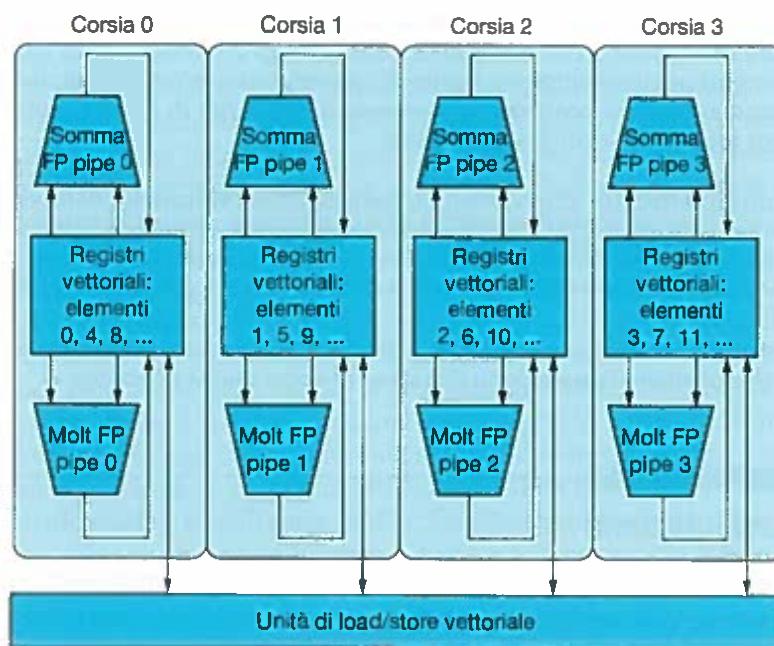
Come le estensioni multimediali delle istruzioni AVX dell'x86, un'istruzione vettoriale specifica operazioni multiple. Tuttavia, un'istruzione appartenente alle estensioni multimediali, tipicamente, specifica poche operazioni, mentre un'istruzione vettoriale può specificare decine di operazioni. A differenza delle estensioni multimediali, il numero di elementi in un'operazione vettoriale non è contenuto nel codice operativo, ma in un registro separato. Questo significa che si possono implementare diverse versioni di un'architettura vettoriale, ciascuna in grado di eseguire le operazioni su un numero diverso di elementi semplicemente modificando il contenuto di quel registro (mantenendo quindi la compatibilità binaria). Invece, nelle estensioni multimediali dell'architettura x86 (MMX, SSE, SSE2, AVX, AVX2 ecc.), ogni volta che la lunghezza del "vettore" viene modificata, occorre anche aggiungere un gran numero di codici operativi.

Inoltre, a differenza delle estensioni multimediali, i dati da trasferire non devono essere necessariamente contigui. Le architetture vettoriali supportano sia gli accessi estesi, nei quali l'hardware carica dalla memoria un elemento ogni  $n$  elementi, sia gli accessi indicizzati, nei quali l'hardware carica dalla memoria gli elementi il cui indirizzo è contenuto in un registro vettoriale. Gli accessi indicizzati sono anche chiamati *gather-scatter* (raggruppa-spargi), perché le operazioni di caricamento, tramite gli indici, raggruppano gli elementi presi dalla memoria principale in elementi contigui del vettore e le operazioni di store, tramite gli indici, spargono gli elementi del vettore nella memoria principale.

Come le estensioni multimediali, le istruzioni vettoriali gestiscono facilmente la variabilità sull'ampiezza dei dati, per cui è facile eseguire operazioni che lavorano su 32 elementi di 64 bit, su 64 elementi di 32 bit, su 128 elementi di 16 bit, o ancora su 256 elementi di 8 bit. La semantica parallela delle istruzioni vettoriali consente di implementare l'esecuzione di queste operazioni mediante un'unità funzionale dotata di pipeline profonda, un insieme di unità funzionali parallele o una combinazione di unità funzionali parallele e in pipeline. La Figura 6.3 illustra come migliorare le prestazioni sui vettori utilizzando pipeline parallele per eseguire la somma degli elementi di due vettori.



**Figura 6.3 Utilizzo di unità funzionali multiple per aumentare le prestazioni di una singola istruzione vettoriale di somma:  $C = A + B$ . Il processore vettoriale sulla sinistra (a) ha una singola pipeline di somma e può eseguire una sola addizione per ciclo di clock. Il processore vettoriale sulla destra (b) ha quattro pipeline di somma, dette anche corsie, e può eseguire quattro somme per ciclo di clock. Gli elementi di ogni vettore vengono interallacciati nelle quattro corsie.**



**Figura 6.4** Struttura di una unità vettoriale contenente quattro corsie. Il registro che memorizza i dati è suddiviso tra le diverse corsie: ogni corsia contiene un elemento ogni quattro del registro vettoriale. La figura mostra tre unità funzionali vettoriali: un sommatore FP, un moltiplicatore FP e un'unità di load/store. Ciascuna unità aritmetica contiene quattro pipeline di esecuzione, una per corsia, e agisce in modo coordinato con le altre per completare le singole istruzioni vettoriali. Si noti come ciascuna porzione del register file vettoriale debba solo fornire porte in numero sufficiente per la lettura e scrittura (vedi Cap. 4) dei dati per le unità funzionali della corsia a cui è assegnato.

Le istruzioni aritmetiche sui vettori solitamente utilizzano gli N elementi di un registro vettoriale in operazioni che coinvolgono gli N elementi di un secondo registro vettoriale. Questo consente di semplificare enormemente la costruzione di un'unità vettoriale altamente parallela, che può essere strutturata come un insieme di più **corsie vettoriali**. Come per il traffico sulle autostrade, possiamo aumentare il throughput di un'unità vettoriale aggiungendo più corsie. La Figura 6.4 mostra la struttura di un'unità vettoriale a quattro corsie. Perciò, passando da una a quattro corsie, il numero di cicli di clock per istruzione vettoriale viene ridotto approssimativamente di un fattore 4. Perché un'architettura dotata di quattro corsie sia vantaggiosa, sia le applicazioni sia l'architettura devono supportare vettori lunghi. Altrimenti, le operazioni verrebbero eseguite così velocemente che si rimarrebbe in fretta senza istruzioni. Vengono quindi richieste tecniche di **parallelismo** a livello di istruzioni come quelle viste nel Capitolo 4 per fornire un numero sufficiente di istruzioni vettoriali.

In generale, le architetture vettoriali sono molto efficienti nell'eseguire programmi che consentono un'esecuzione parallela dei dati, supportano meglio la tecnologia dei compilatori e sono più facili da sviluppare nel tempo rispetto alle estensioni multimediali dell'architettura x86.

A partire da queste classiche categorie, vediamo ora come si possano sfruttare i flussi di elaborazione paralleli di istruzioni per migliorare le prestazioni di un **singolo** processore, che riutilizzeremo poi all'interno dei processori multipli.

**Approfondimento.** Considerati i vantaggi delle architetture vettoriali, come mai queste sono assai diffuse solo nell'ambito del calcolo ad alte prestazioni? Uno dei motivi è l'aumento del tempo dovuto al cambio di contesto; questo, infatti, dipende dalla dimensione dello stato, che in questo caso deve contenere anche registri di grandi dimensioni, adatti a gestire i dati vettoriali. Un altro motivo è l'aumento della difficoltà nel gestire gli errori di page fault che si possono verificare durante

**Corsia vettoriale:** una o più unità funzionali vettoriali abbinate a una porzione del register file. Il nome è ispirato alle corsie delle autostrade che consentono di aumentare la velocità del traffico: corsie multiple consentono di eseguire operazioni vettoriali simultaneamente.



PARALLELISMO



il caricamento e la scrittura dei vettori di dati; le architetture SIMD hanno comunque alcuni dei vantaggi offerti dalle istruzioni vettoriali. Inoltre, fino a quando i miglioramenti nel parallelismo a livello di istruzioni consentiranno di mantenere il miglioramento delle prestazioni promesso dalla Legge di Moore, non ci sono motivi per modificare tipo di architettura.

**Approfondimento.** Un altro vantaggio delle istruzioni vettoriali e delle estensioni multimediali consiste nel fatto che è relativamente facile estendere un'architettura di un insieme di istruzioni scalari attraverso le istruzioni vettoriali, con lo scopo di migliorare le prestazioni delle operazioni che lavorano in parallelo sui dati.

**Approfondimento.** I processori x86 di Intel della generazione Haswell supportano l'AVX2, che contiene l'operazione di gather ma non quella di scatter.

### Autovalutazione

Vero o falso? Come abbiamo visto nell'x86, le estensioni multimediali possono essere viste come un'architettura vettoriale con vettori corti che supportano solo trasferimenti di dati vettoriali contigui.

## 6.4 Multithreading hardware

**Multithreading hardware:** aumento dello sfruttamento di un processore realizzato passando all'esecuzione di un altro thread quando il thread in esecuzione viene messo in stallo.

**Thread:** un thread comprende il program counter, il registro di stato e lo stack. È un processo leggero; i thread di solito condividono lo spazio di indirizzamento, mentre i processi non lo condividono.

**Processo:** un processo comprende uno o più thread, lo spazio di indirizzamento e lo stato del sistema operativo. Quindi, il cambio di processo di solito richiede una chiamata al sistema operativo, ma non un cambio di thread.

**Multithreading a grana fine:** un tipo di multithreading hardware in cui si passa da un thread all'altro a ogni istruzione.

**Multithreading a grana grossa:** un tipo di multithreading hardware in cui si passa da un thread all'altro solo quando si verificano eventi particolari, per esempio una miss dell'ultimo livello di cache.

Un concetto collegato alle architetture MIMD, specialmente dal punto di vista del programmatore, è il **multithreading hardware**. Mentre le architetture MIMD sono basate su **processi** o **thread** multipli, per cercare di mantenere impegnati tutti i processori, il multithreading hardware permette a più thread di condividere le unità funzionali di un *singolo* processore in modo sovrapposto, per cercare di ottimizzare l'utilizzo delle risorse hardware. Per ottenere una tale condivisione delle risorse, il processore deve duplicare lo stato di ciascun thread. Per esempio, ogni thread avrà una propria copia del register file e del PC. La memoria stessa può essere condivisa attraverso il meccanismo della memoria virtuale, che già supporta la multiprogrammazione. In aggiunta, l'hardware deve essere in grado di supportare il passaggio rapido dall'esecuzione di un thread all'altro. In particolare, questo passaggio deve essere molto più efficiente dei cambi di processo, che tipicamente richiedono centinaia o migliaia di cicli di clock per essere portati a termine (il cambio di thread deve essere istantaneo).

Esistono due approcci principali al multithreading: nel **multithreading a grana fine** a ogni istruzione si passa da un thread all'altro, producendo in pratica un'esecuzione intrecciata dei vari thread. Questa esecuzione intrecciata viene spesso realizzata utilizzando una modalità di tipo *round robin* (a rotazione), passando da un thread all'altro a turno e saltando i thread che si trovano in stallo in quel momento. Per rendere possibile il multithreading a grana fine, il processore deve essere capace di cambiare il thread in esecuzione a ogni ciclo di clock. Un vantaggio chiave del multithreading a grana fine è che può nascondere perdite di throughput che nascono da stalli sia lunghi sia corti, poiché quando un thread è in stallo si possono eseguire le istruzioni degli altri thread. Lo svantaggio principale del multithreading a grana fine è che rallenta l'esecuzione dei singoli thread, dal momento che il thread pronto per essere eseguito e che non richiede stalli viene comunque ritardato dall'esecuzione delle istruzioni degli altri thread.

Il **multithreading a grana grossa** rappresenta l'alternativa. In questo secondo approccio il thread in esecuzione cambia solo quando si presenta uno stallo

“costoso”, come quelli associati alle miss dell’ultimo livello di cache. Inoltre, si rimuove il vincolo che il cambio di thread sia essenzialmente a costo zero e si rende molto meno probabile il rallentamento dell’esecuzione del singolo thread, poiché le istruzioni degli altri thread non vengono avviate in esecuzione fino a quando non si incontra uno stall. Nonostante ciò, il multithreading a grana grossa ha un grande svantaggio: la sua capacità di mascherare le perdite di throughput è limitata, specialmente in presenza degli stalli più brevi. Questa limitazione è dovuta all’alto costo di riavvio della pipeline. Dato che una CPU con multithreading a grana grossa carica le istruzioni da un singolo thread, quando si verifica uno stall la pipeline deve essere svuotata o bloccata. Il nuovo thread che inizia l’esecuzione subito dopo lo stall deve riempire la pipeline prima che l’esecuzione delle sue istruzioni possa terminare. A causa di questo vincolo sulla pipeline, il multithreading a grana grossa è molto più utile per ridurre la penalità degli stalli gravi, cioè quando il tempo speso per ricaricare la pipeline è trascurabile se confrontato con la durata dello stall.

Il **multithreading simultaneo (SMT)**, *Simultaneous Multithreading*) è una variazione del multithreading hardware e utilizza le risorse già presenti in un processore con pipeline a parallelizzazione dinamica dell’esecuzione e riordinamento dinamico del codice per realizzare il multithreading mentre implementa l’esecuzione parallela delle istruzioni (vedi Cap. 4). L’intuizione chiave che sta alla base dell’SMT è che i processori dotati di parallelizzazione dinamica dell’esecuzione hanno spesso più unità funzionali che lavorano in parallelo di quante ne possa generalmente utilizzare un singolo thread. Inoltre, grazie alla ridenominazione dei registri e alla riorganizzazione dinamica del codice, più istruzioni appartenenti a thread indipendenti possono essere caricate senza curarsi delle reciproche dipendenze; la risoluzione delle dipendenze può essere gestita dal controllore della riorganizzazione dinamica del codice.

Poiché in questo caso ci si affida ai meccanismi esistenti di riorganizzazione dinamica del codice, l’SMT non deve riassegnare le risorse a ogni ciclo di clock, ma esegue *continuamente* istruzioni appartenenti a thread diversi, lasciando all’hardware il compito di associare ai relativi thread sia le diverse istruzioni lanciate contemporaneamente in esecuzione sia i diversi registri ridenominati.

La Figura 6.5 illustra le differenze tra i modi in cui un processore sfrutta le risorse di un’architettura superscalare in base alle seguenti configurazioni: nella parte superiore viene mostrato come quattro thread verrebbero eseguiti in maniera indipendente su un processore superscalare senza supporto del multithreading; nella parte inferiore si mostra invece come i quattro thread potrebbero essere combinati per essere eseguiti in modo più efficiente dal processore utilizzando tre diverse opzioni di multithreading:

- architettura superscalare con multithreading a grana grossa;
- architettura superscalare con multithreading a grana fine;
- architettura superscalare con multithreading simultaneo.

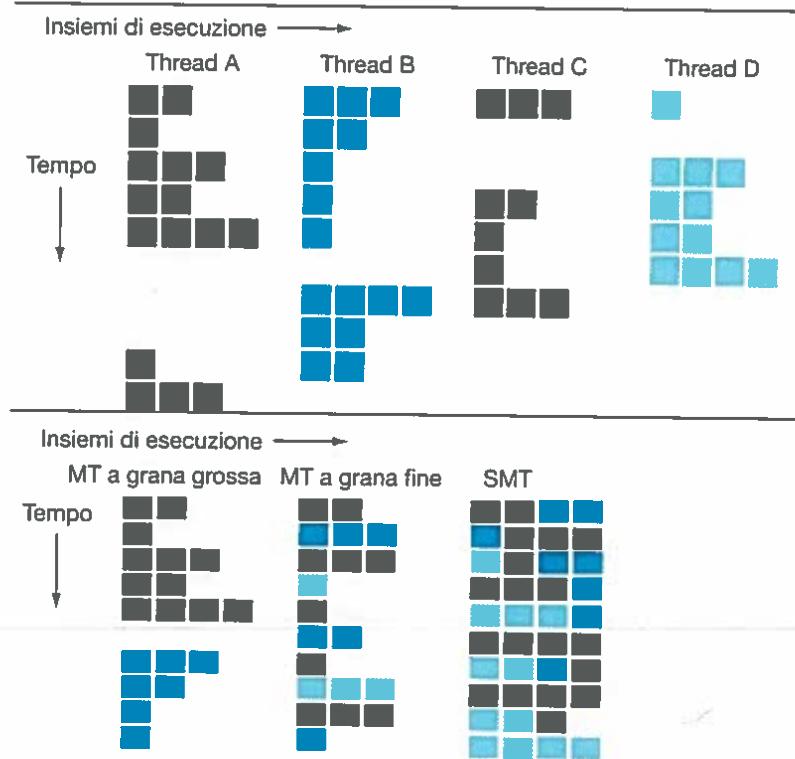
Nel processore superscalare senza supporto del multithreading, il completo sfruttamento del parallelismo è limitato da una carenza di **parallelismo a livello di istruzioni (ILP)**. Inoltre, uno stall importante, come una miss della cache istruzioni, può lasciare inattivo l’intero processore.

Nelle architetture superscalari dotate di multithreading a grana grossa, gli stalli lunghi vengono parzialmente mascherati, sostituendo il thread in esecuzione con un altro thread che utilizza le risorse del processore. Sebbene questo riduca il numero di cicli di clock di completa inattività, il sovraccarico di lavoro richiesto per riavviare la pipeline produce comunque dei cicli di clock di inattività; inoltre, le limitazioni sull’ILP fanno sì che non tutti gli insiemi di istruzioni lanciate in esecuzione possano essere completamente riempiti. Nel caso del multithreading a grana fine, l’intreccio dei thread elimina la maggior parte



**Multithreading simultaneo (SMT):** un tipo di multithreading hardware che abbassa il costo del multithreading sfruttando le risorse utilizzate dalle microarchitetture per la parallelizzazione dell’esecuzione e il riordinamento dinamico del codice.



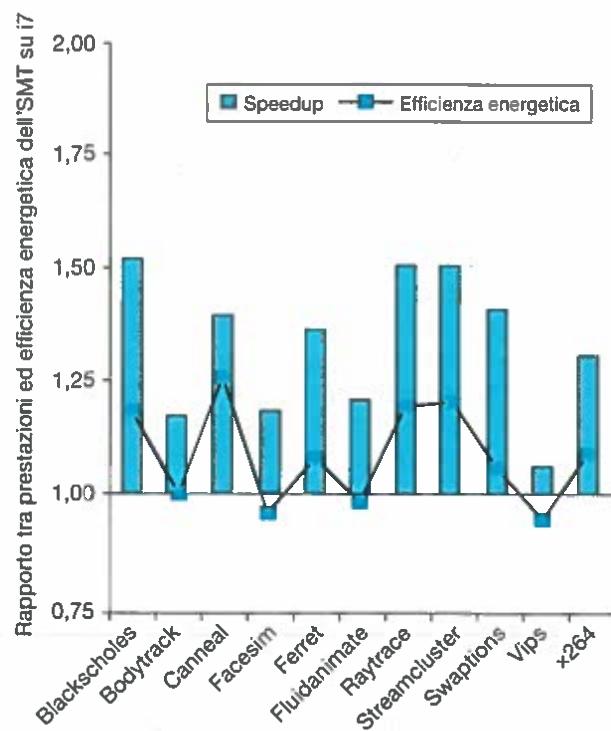


**Figura 6.5** Modalità con cui vengono riempite le quattro istruzioni dell'insieme di esecuzione di un processore superscalare secondo i diversi approcci. I quattro thread in alto mostrano come ogni thread verrebbe eseguito su un processore superscalare tradizionale senza supporto del multithreading. I tre esempi in basso mostrano come i thread verrebbero eseguiti in parallelo con i tre diversi tipi di multithreading. Sull'asse orizzontale è riportato il numero di istruzioni lanciate in esecuzione contemporaneamente a ogni ciclo di clock, mentre sull'asse verticale compare la sequenza dei cicli di clock. Un quadratino vuoto (bianco) indica che non viene inserita alcuna istruzione in quella posizione dell'insieme di esecuzione. I quadratini grigi o colorati corrispondono a quattro thread differenti. L'effetto del riavvio della pipeline richiesto dal multithreading a grana grossa, che non è mostrato in figura, porta a un'ulteriore perdita di prestazioni in termini di throughput.

degli insiemi di esecuzione completamente vuoti; poiché un solo thread lancia in esecuzione le istruzioni in ogni ciclo di clock, le limitazioni sul parallelismo a livello di istruzioni possono ancora produrre qualche insieme di esecuzione vuoto in alcuni cicli di clock.

Nel caso dell'approccio SMT al multithreading, il parallelismo a livello di thread e quello a livello di istruzioni vengono sfruttati contemporaneamente: le istruzioni di più thread sono inserite nello stesso insieme di istruzioni lanciato in esecuzione nello stesso ciclo di clock e quindi, in teoria, il riempimento degli insiemi di esecuzione è limitato solamente dalla discrepanza tra le risorse richieste dai vari thread e le risorse disponibili. In pratica, altri fattori possono ridurre il numero di istruzioni che possono essere inserite in tali insiemi. Sebbene in Figura 6.5 il modo di funzionamento di questi processori sia molto semplificato, la figura illustra i potenziali vantaggi in termini di prestazioni del multithreading in generale e dell'SMT in particolare.

La Figura 6.6 mostra le prestazioni e i benefici in termini di energia del multithreading su un singolo processore del Core i7 960 di Intel, che ha il supporto hardware per due thread. L'aumento di prestazioni è di 1,31 volte, che non è male data la ridotta quantità di risorse aggiuntive richiesta per il multithreading hardware. Il miglioramento medio nell'efficienza energetica è di 1,07, che è eccellente. In generale, possiamo essere contenti quando si ha un aumento di efficienza senza richiedere un maggiore consumo di energia.



**Figura 6.6** L'aumento di prestazioni ottenuto utilizzando il multithreading su un processore i7 è in media pari a 1,31 per i benchmark PARSEC (par. 6.9) e l'aumento dell'efficienza energetica è pari a 1,07. Questi dati sono stati raccolti e analizzati da Esmaeilzadeh *et al.* [2011].

Ora che abbiamo visto come più thread possano utilizzare efficacemente le risorse di un singolo processore, vediamo come utilizzare queste tecniche nei processori multipli.

### Autovalutazione

1. Vero o falso? Sia il multithreading sia le architetture multicore si basano sul parallelismo per ottenere dai chip un'esecuzione più efficiente.
2. Vero o falso? Il *multithreading simultaneo* (SMT) utilizza i thread per migliorare lo sfruttamento delle risorse di un processore con riordinamento dinamico del codice ed esecuzione fuori ordine.

## 6.5 I multicore e gli altri multiprocessori a memoria condivisa

Mentre il multithreading hardware migliora l'efficienza dei processori al prezzo di un costo modesto, la grande sfida di questo decennio è ottenere l'aumento di prestazioni previsto dalla **Legge di Moore** programmando in modo efficiente i diversi processori, sempre più numerosi, sullo stesso chip.

Data la difficoltà nel riscrivere i vecchi programmi per eseguirli in modo efficiente su hardware parallelo, dobbiamo chiederci che cosa possano fare i progettisti dei calcolatori per semplificare questo compito. Una possibile risposta è

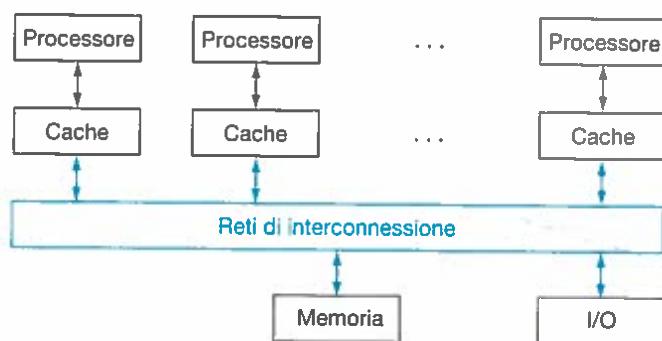


fornire un unico spazio di indirizzamento fisico che possa essere condiviso da tutti i processori, così che i programmi non debbano sapere su quale processore vengono eseguiti, ma soltanto che possono essere eseguiti in parallelo. In questo approccio, tutte le variabili di un programma possono essere rese disponibili, in ogni momento, a un qualsiasi processore. L'alternativa è avere uno spazio di indirizzamento separato per ogni processore, una strategia che richiede che la condivisione sia resa esplicita; discuteremo questa seconda possibilità nel paragrafo 6.7. Quando lo spazio di indirizzamento fisico è comune, l'hardware tipicamente deve provvedere alla coerenza delle cache per fornire una visione coerente della memoria condivisa (*vedi* par. 5.8).

Come si è detto precedentemente, un *multiprocessore a memoria condivisa* (SMP) è un processore che offre ai programmati un *unico spazio di indirizzamento fisico* per i diversi processori, che è quasi sempre il caso dei chip multicore; il loro nome più appropriato sarebbe *multiprocessori a indirizzi condivisi*. I processori comunicano attraverso variabili condivise contenute in memoria, poiché tutti i processori possono accedere a una qualsiasi locazione attraverso le istruzioni di load e store. La Figura 6.7 mostra la classica organizzazione di un processore SMP. Si noti che questi processori possono eseguire programmi indipendenti nei loro propri spazi di indirizzamento virtuali, anche quando condividono lo stesso spazio di indirizzamento fisico.

I processori con un unico spazio di indirizzamento sono suddivisi in due categorie. Alla prima appartengono i processori che impiegano lo stesso tempo per accedere alla memoria principale indipendentemente dalla parola richiesta e dal processore che ha effettuato la richiesta; tali architetture sono denominate **multiprocessori con memoria ad accesso uniforme (UMA, Uniform Memory Access)**. Nella seconda categoria alcuni accessi alla memoria possono essere molto più rapidi di altri, a seconda del dato richiesto e del processore che ha inoltrato la richiesta. Tali macchine prendono il nome di **multiprocessori con memoria ad accesso non uniforme (NUMA, Nonuniform Memory Access)**. Come ci si potrebbe aspettare, ci sono più problemi a programmare i multiprocessori NUMA che i multiprocessori UMA, ma il numero di multiprocessori che possono essere contenuti in un'architettura NUMA può essere maggiore di quello delle architetture UMA; inoltre, le architetture NUMA possono avere latenze inferiori per accessi ad aree di memoria vicine.

Poiché i processori che operano in parallelo generalmente condividono i dati, devono anche coordinarsi quando operano su dati condivisi. Altrimenti, un processore potrebbe iniziare a lavorare su un dato prima che un altro processore abbia terminato di modificarlo. Questo coordinamento viene chiamato **sincronizzazione** (*vedi* Cap. 2): quando la condivisione viene supportata da uno spazio di indirizzamento unico, ci deve essere un meccanismo di sincronizzazione. Un approccio per implementare la sincronizzazione è basato sul **lock** (blocco) delle



**Figura 6.7** Organizzazione classica di un multiprocessore a memoria condivisa.

variabili condivise: un solo processore alla volta può bloccare una cella di memoria contenente una variabile condivisa; gli altri dovranno attendere fino a quando quel processore non avrà sbloccato la variabile. Le istruzioni del RISC-V per bloccare e sbloccare una cella di memoria sono state descritte nel paragrafo 2.11.

### Un semplice programma a elaborazione parallela per uno spazio di indirizzamento condiviso

Si supponga di dover sommare 64 000 numeri su un calcolatore multiprocessore a memoria condivisa con memoria ad accesso uniforme. Si supponga di avere a disposizione 64 processori.

Il primo passo per ottenere un carico bilanciato tra i diversi processori è suddividere l'insieme dei numeri da sommare in sottoinsiemi di uguale dimensione. Non allochiamo i diversi sottoinsiemi in spazi di memoria diversi, poiché c'è un unico spazio di memoria per l'intera macchina, ma assegniamo semplicemente diversi indirizzi di partenza ai diversi processori. Chiameremo  $P_n$  il numero che identifica ciascun processore, compreso tra 0 e 63. Tutti i processori iniziano a eseguire il loro programma con un ciclo che somma tutti i numeri del sottoinsieme a loro assegnato:

```
somma[Pn] = 0;
for (i = 1000*Pn; i < 1000*(Pn+1); i += 1)
    somma[Pn] += A[i]; /* somma gli elementi del sottoinsieme
                        assegnato */
```

(Si noti che la notazione  $C\ i\ +=\ 1$  è un modo abbreviato di scrivere  $i = i + 1$ )

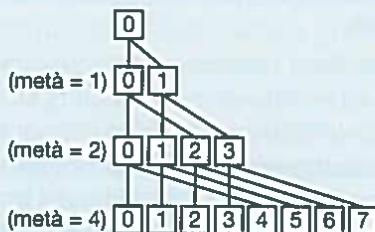
Il passo successivo consiste nel sommare tra loro le 64 somme parziali ottenute. Questa fase è chiamata **riduzione** (*reduction*) e applica una strategia di tipo “*divide et impera*”: metà dei processori sommerà coppie di somme parziali, un quarto dei processori sommerà le nuove somme parziali ottenute e così via, fino a quando non si avrà una singola somma finale. La Figura 6.8 illustra la natura gerarchica di questa riduzione.

In questo esempio, i due processori devono sincronizzarsi prima che il processore “consumatore” possa leggere la somma parziale dalla locazione della memoria in cui la somma è stata scritta dal processore “produttore”, altrimenti il consumatore potrebbe leggere il vecchio valore. Vogliamo, inoltre, che ciascun processore abbia la propria versione privata del contatore di ciclo,  $i$ , per cui dobbiamo specificare che si tratta

### ESEMPIO

### SOLUZIONE

**Riduzione:** una funzione che elabora una struttura dati e restituisce un singolo valore.



**Figura 6.8 Ultimi quattro livelli della riduzione relativa alla somma dei risultati di ciascun processore, dal basso verso l'alto.** In tutti i processori il cui numero d'ordine,  $i$ , è inferiore a metà, viene sommata la somma parziale prodotta dal processore  $(i + \text{metà})$  a quella del processore  $i$  stesso.

(continua)

(continua)

di una variabile “privata”. Ecco il codice risultante (anche la variabile metà è privata):

```
metà = 64; /* 64 processori nel multiprocessore */
do
    synch(); /* attesa completamento somme parziali */
    if (metà % 2 != 0 && Pn == 0)
        somma[0] += somma[metà-1];
    /* Quando metà è dispari, P0 prende l'elemento spaiato */
    metà = metà/2; /* punto che separa quali processori
                     eseguono la somma */
    if (Pn < metà) somma[Pn] += somma[Pn+metà];
while (metà > 1); /* uscita; la somma finale è in somma[0] */
```

## Interfaccia hardware/software

**OpenMP:** un'API per il multiprocessing con memoria condivisa in C, C++ o Fortran che viene eseguita sulle piattaforme UNIX e Microsoft. OpenMP comprende le direttive di compilazione, una libreria e le direttive di esecuzione.

Dato l'interesse per la programmazione parallela, ci sono stati centinaia di tentativi di costruire sistemi che la supportino. Un esempio limitato ma molto diffuso è **OpenMP**: è un'API (*Application Programmer Interface*) associata a un insieme di direttive di compilazione, variabili d'ambiente e funzioni di libreria run-time che possono estendere i linguaggi di programmazione standard. OpenMP offre un semplice modello, portabile e scalabile, per multiprocessori a memoria condivisa. L'obiettivo principale è parallelizzare i cicli ed eseguire le riduzioni.

La maggior parte dei compilatori C già supporta OpenMP. Il comando per utilizzare le API OpenMP in C con i compilatori UNIX è semplicemente:

```
cc -fopenmp foo.c
```

OpenMP estende il C utilizzando i *pragma*, che sono dei comandi al pre-processore C delle macro, come `#define` e `#include`. Per impostare a 64 il numero dei processori che vogliamo utilizzare, come per l'esempio sopra, utilizziamo il comando:

```
# define P 64 /* definiamo una costante che utilizzeremo altre
               volte */
# pragma omp parallel num_threads(P)
```

Questa coppia di macro definisce che le librerie run-time dovranno utilizzare 64 thread paralleli.

Per trasformare il ciclo sequenziale in un ciclo parallelo che suddivida il lavoro tra tutti i thread che abbiamo definito, possiamo scrivere (supponendo che somma sia inizializzata a 0):

```
#pragma omp parallel for
for (Pn = 0; Pn < P; Pn += 1)
    for (i = 0; 1000*Pn; i < 1000*(Pn+1); i += 1)
        somma[Pn] += A[i]; /* somma le aree assegnate */
```

Per operare la riduzione, possiamo utilizzare un altro comando che dice a OpenMP qual è l'operatore di riduzione da utilizzare e in quale variabile inserire il risultato dell'operazione di riduzione.

```
#pragma omp parallel for reduction(+ : SommaFinale)
for (I = 0; I < P; I += 1)
    SommaFinale += somma[i]; /* riduzione a un singolo numero */
```

Si noti che è ora il turno della libreria OpenMP di trovare un codice efficiente per sommare i 64 numeri in modo efficiente utilizzando 64 processori.

Se OpenMP rende facile scrivere del codice parallelo semplice, non è molto utile nella ricerca degli errori, per cui molti programmati paralleli utilizzano sistemi più sofisticati di OpenMP per la programmazione parallela, proprio come molti programmati moderni utilizzano linguaggi di programmazione più produttivi del C.

Dopo l'introduzione delle architetture MIMD hardware e software, il prossimo argomento sarà la presentazione di un tipo diverso di architetture MIMD, con un background e una prospettiva diversa sulle sfide della programmazione parallela.

### Autovalutazione

Vero o falso? I multiprocessori a memoria condivisa non possono sfruttare il parallelismo a livello di task.

**Approfondimento.** Alcuni autori hanno proposto di cambiare la traduzione dell'acronimo SMP in *Symmetric MultiProcessor* (multiprocessore simmetrico) per indicare che la latenza tra il singolo processore e la memoria è uguale per tutti i processori. Questo per rimarcare la differenza con i multiprocessori NUMA a grande scala, dato che entrambi utilizzano un singolo spazio di indirizzamento. Dato che i cluster hanno avuto una diffusione molto maggiore dei multiprocessori NUMA, in questo libro utilizzeremo SMP con il suo significato originale, e lo utilizzeremo per rimarcare la differenza con le architetture che utilizzano spazi di indirizzamento multipli quali i cluster.

**Approfondimento.** Un'alternativa alla condivisione dello spazio di indirizzamento fisico sarebbe disporre di spazi di indirizzamento fisico separati ma con lo stesso spazio di indirizzamento virtuale, lasciando al sistema operativo l'onere di gestire la comunicazione. Questo approccio è stato sperimentato, ma il sovraccarico di lavoro si è rivelato troppo elevato perché possa costituire un valido modello di astrazione della memoria condivisa per il programmatore.

## 6.6 | Introduzione alle unità di elaborazione grafica

Uno dei motivi principali che sta alla base dell'aggiunta di istruzioni SIMD alle architetture esistenti è che la maggior parte dei microprocessori era collegata a un terminale grafico, per cui una frazione sempre più grande del tempo di elaborazione era utilizzata per gestire la grafica. Inoltre, poiché il numero di transistor disponibili su un processore ha continuato a crescere in linea con la **Legge di Moore**, è risultato sensato migliorare l'elaborazione grafica.

L'industria dei videogiochi, sviluppati sia per PC sia per console dedicate (come la PlayStation di Sony), ha rappresentato una delle maggiori forze propulsive per il miglioramento dell'elaborazione grafica. La rapida crescita del mercato dei videogiochi ha incoraggiato molte società a investire ingenti risorse per sviluppare hardware grafico sempre più veloce, e questo meccanismo ha fatto sì che le architetture per l'elaborazione grafica si sviluppassero a una velocità maggiore rispetto ai microprocessori tradizionali a utilizzo generale.

Poiché l'industria della grafica e dei videogiochi ha obiettivi diversi da quelli delle aziende che sviluppano i microprocessori, in quell'ambito sono nati un tipo di elaborazione e una terminologia differenti. A causa dell'aumento della loro potenza di elaborazione, inoltre, i processori grafici si sono guadagnati il nome di *Graphics Processing Units* (GPU, unità di elaborazione grafica), per distinguerli dalle CPU.

Chiunque oggi, con poche centinaia di dollari, può acquistare una GPU contenente centinaia di unità di elaborazione in virgola mobile, che rendono il calcolo ad alte prestazioni più accessibile. L'interesse per il calcolo tramite



LEGGE DI MOORE

GPU è fiorito quando il suo potenziale è stato combinato con un linguaggio di programmazione che consentisse di programmare le GPU più facilmente: oggi i programmatore di applicazioni scientifiche e multimediali valutano attentamente se utilizzare le GPU o le CPU.

(Questo paragrafo è focalizzato su come utilizzare le GPU per il calcolo. Per vedere come il calcolo su GPU sia legato al ruolo tradizionale di acceleratore grafico, *vedi l'Appendice C*.)

Queste sono alcune delle caratteristiche principali che rendono le GPU diverse dalle CPU:

- Le GPU sono acceleratori che complementano la CPU, per cui non hanno bisogno di eseguire tutti i compiti che sono eseguiti normalmente da una CPU. Questo diverso ruolo delle GPU consente loro di dedicare tutte le risorse alla grafica. È accettabile che una GPU esegua alcuni compiti in maniera non efficiente o che non li esegua affatto, poiché in un sistema che contiene sia una GPU sia una CPU sarà la CPU a eseguire questi compiti.
- Le dimensioni dei problemi risolti dalla GPU sono tipicamente dell'ordine delle centinaia di Megabyte o dei Gigabyte, ma non arrivano alle centinaia di Gigabyte o al Terabyte.

Queste differenze hanno portato allo sviluppo di architetture di tipo diverso:

- Forse la differenza maggiore è che le GPU non si basano su cache multilivello, come le CPU, per superare la latenza elevata della memoria, ma cercano di utilizzare molti thread per nascondere la latenza della memoria (par. 6.4). Ossia, dall'istante in cui viene generata una richiesta di lettura in memoria all'istante in cui il dato è disponibile, la GPU esegue centinaia o migliaia di thread che sono indipendenti dalla richiesta.
- La memoria principale della GPU viene perciò ottimizzata per massimizzare la larghezza di banda di trasferimento invece che per minimizzare la latenza. Nelle GPU, inoltre, vengono anche montati dei chip di DRAM separati, più ampi e con una larghezza di banda maggiore dei chip di DRAM presenti nelle CPU. Inoltre, la memoria principale delle GPU è più piccola di quella dei microprocessori convenzionali: nel 2013 una GPU in media era dotata di una memoria da 4 a 6 GiB, mentre una CPU aveva una memoria principale compresa tra i 32 e i 256 GiB. Infine, bisogna ricordare che per elaborazioni diverse da quelle tipiche della grafica occorre considerare anche il tempo necessario per trasferire i dati tra la memoria della CPU e la memoria della GPU, poiché in questo caso la GPU funziona da coprocessore.
- Poiché si utilizzano molti thread per ottenere una buona larghezza di banda, le GPU possono contenere sia molti processori paralleli (MIMD) sia molti thread. Quindi, ciascuna GPU è un processore multithread, con un numero di thread più elevato di quello delle CPU; inoltre possiede un numero più elevato di processori.

## Interfaccia hardware/software

Sebbene le GPU inizialmente siano state progettate per un numero ristretto di applicazioni, alcuni programmatore si chiesero se potevano tradurre le loro applicazioni in una forma che consentisse di utilizzare le elevate prestazioni delle GPU. I programmatore iniziarono presto a stancarsi di dovere utilizzare le API grafiche e il linguaggio di shading per definire i loro problemi e scrivere il codice per risolverli, per cui svilupparono un linguaggio di programmazione, ispirato al C, che consentisse di scrivere i programmi direttamente per la GPU. Un esempio è il linguaggio di programmazione di NVIDIA, chiamato CUDA (*Compute Unified Device Architecture*), che consente ai programmatore di scrivere in C i programmi da eseguire sulle GPU, anche se con alcune restrizioni. Alcuni esempi di codice CUDA sono riportati nell'Appendice C. Un altro esempio è OpenCL, un'iniziativa di più società di sviluppare un linguaggio di programmazione portabile che fornisca molti dei benefici offerti da CUDA.

NVIDIA ha deciso che il tema unificatore di tutte queste forme di parallelismo può essere il *thread CUDA*. Utilizzando questo livello più basso di parallelismo come primitiva di programmazione, il compilatore e l'hardware possono raggruppare migliaia di thread CUDA per sfruttare i diversi tipi di parallelismo contenuti all'interno della GPU: multithreading, MIMD, SIMD e parallelismo a livello di istruzioni. Questi thread vengono raggruppati a ogni istante in blocchi di 32 thread ed eseguiti in parallelo. Un processore multithread di una GPU esegue questi blocchi di thread e una GPU è costituita da un numero variabile da 8 a 32 di questi processori a thread multipli.

## Introduzione alle architetture GPU di NVIDIA

Utilizziamo i sistemi NVIDIA come esempio rappresentativo delle architetture GPU: seguiremo qui la terminologia adottata dal linguaggio di programmazione parallelo CUDA e utilizzeremo l'architettura Fermi come esempio.

Come le architetture vettoriali, le GPU lavorano bene solamente con i problemi che esibiscono parallelismo a livello dei dati. Entrambe queste architetture sono dotate delle funzioni di trasferimento dei dati gather/scatter, ma le GPU hanno un numero maggiore di registri rispetto ai processori vettoriali. A differenza di molte architetture vettoriali, le GPU si basano sul multithreading hardware all'interno di ogni processore SIMD a thread multipli per nascondere la latenza della memoria (par. 6.4).

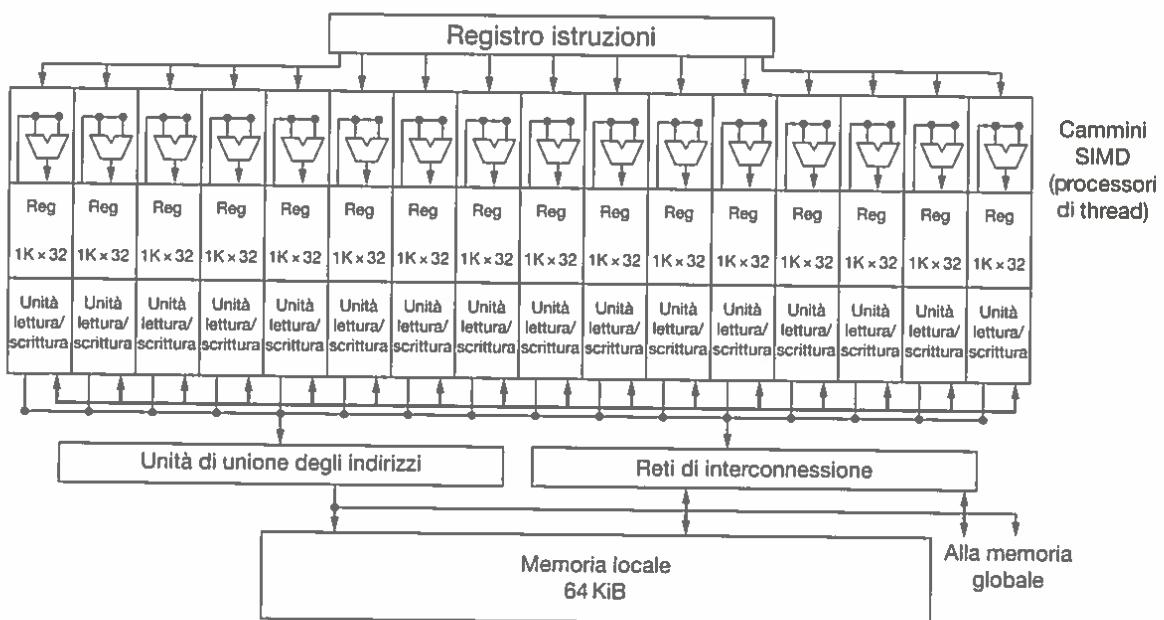
Un processore SIMD multithread è simile a un processore vettoriale, con la differenza che il primo ha molte unità funzionali che lavorano in parallelo, mentre il secondo ne ha poche dotate di una pipeline profonda.

Come abbiamo accennato in precedenza, una GPU contiene un insieme di processori SIMD multithread, ovvero una GPU è un'architettura MIMD composta da processori SIMD multithread. Per esempio, NVIDIA ha proposto quattro versioni dell'architettura Fermi, con costi diversi e un numero diverso di processori SIMD multithread: 7, 11, 14 o 15. Per consentire la scalabilità in modo trasparente tra modelli di GPU con un diverso numero di processori SIMD multithread, il componente hardware denominato *thread block scheduler* (scheduler di blocchi di thread) assegna blocchi di thread ai diversi processori SIMD multithread. Il diagramma a blocchi semplificato di un processore SIMD multithread è mostrato in Figura 6.9.

Scendendo più nel dettaglio, l'oggetto che viene creato, gestito ed eseguito dall'hardware è il *thread di istruzioni SIMD*, detto anche *thread SIMD*. È un thread tradizionale che contiene solo istruzioni SIMD. Questi thread hanno il proprio program counter e vengono eseguiti su un processore SIMD multithread. Lo *scheduler di thread SIMD* contiene un controllore che informa l'architettura su quali thread di istruzioni SIMD siano pronti per l'esecuzione e un'unità che li invia a esecuzione sul processore SIMD multithread. È identico a uno scheduler di thread hardware di un tradizionale processore multithread (par. 6.4), tranne che invia a esecuzione thread di istruzioni SIMD. Perciò le GPU hanno uno scheduler hardware su due livelli:

1. lo *scheduler di blocchi di thread* che assegna blocchi di thread ai diversi processori SIMD multithread;
2. lo scheduler di thread *all'interno* di ciascun processore SIMD, che decide quando un thread SIMD debba essere eseguito.

Le istruzioni SIMD di questi thread sono ampie 32 bit, quindi ciascun thread delle istruzioni SIMD effettua i calcoli su 32 elementi. Poiché un thread consiste in istruzioni SIMD, il processore SIMD deve avere unità funzionali parallele per eseguire le operazioni. Chiamiamo queste unità *cammini di elaborazione SIMD* (*SIMD Lanes*), molto simili ai cammini vettoriali (*Vector Lanes*) del paragrafo 6.3.



**Figura 6.9** Schema a blocchi semplificato del cammino di elaborazione dei dati di un processore SIMD multithread. Il processore ha 16 cammini di elaborazione SIMD. Lo scheduler dei thread SIMD gestisce molti thread SIMD che vengono scelti per l'esecuzione su questo processore.

**Approfondimento.** Il numero di cammini di elaborazione in ciascun processore SIMD delle GPU varia molto da generazione a generazione. Con le architetture Fermi, ciascun thread di istruzioni SIMD ampio 32 elementi viene mappato su 16 cammini di elaborazione SIMD, così che ciascuna istruzione SIMD, appartenente a un thread di istruzioni SIMD, richiede due cicli di clock per essere completata. Ciascun thread di istruzioni SIMD viene eseguito in un tempo fissato. Continuando nell'analogia tra processori SIMD e architetture vettoriali, possiamo affermare che abbiamo 16 cammini di elaborazione e che la lunghezza del vettore da elaborare è di 32 elementi. L'elaborazione uguale di un insieme ampio di dati è il motivo per cui viene utilizzato il termine più intuitivo processore SIMD invece di processore vettoriale.

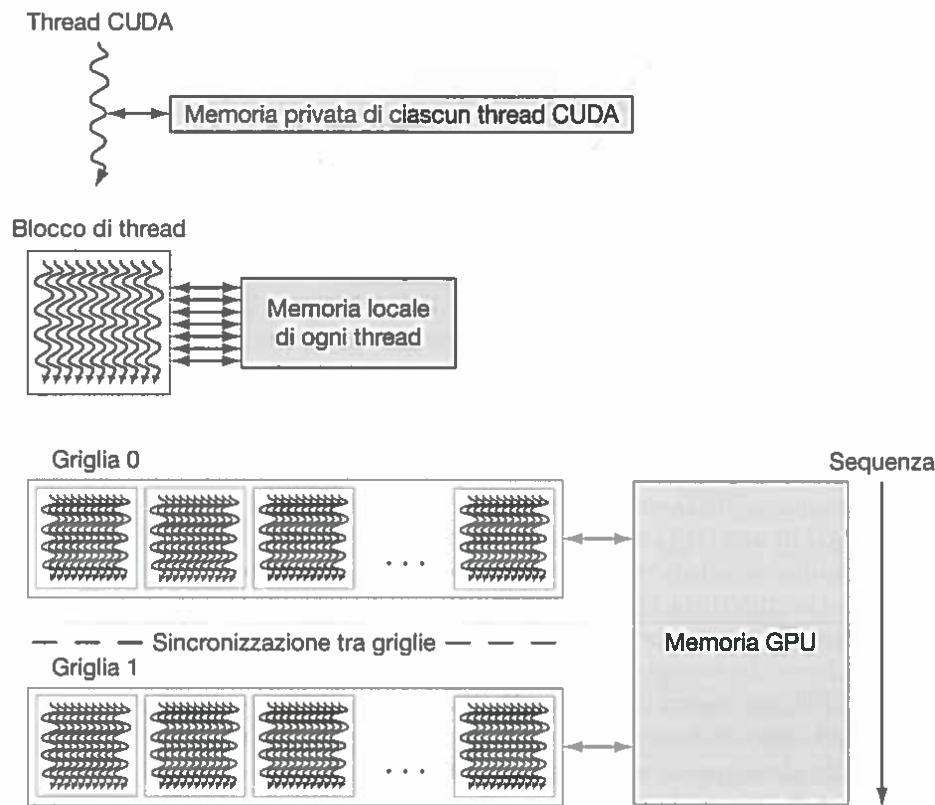
Dato che, per definizione, i thread delle istruzioni SIMD sono indipendenti, lo scheduler dei thread SIMD può prelevare uno qualsiasi dei thread di istruzioni SIMD pronti, e non deve per forza prendere l'istruzione SIMD successiva all'interno dello stesso thread. Cioè, utilizzando la terminologia del paragrafo 6.4, un processore SIMD utilizza un multithread a grana fine.

Per tenere in memoria gli elementi, un processore SIMD Fermi contiene un numero impressionante di registri a 32 bit: 32 768. Esattamente come in un processore vettoriale, questi registri vengono suddivisi in modo logico tra i diversi cammini di elaborazione, o, in questo caso, tra i diversi cammini SIMD. Ciascun thread SIMD non può occupare più di 64 registri e avrà quindi a disposizione un massimo di 64 registri vettoriali, ciascuno contenente 32 elementi, ciascuno ampio 32 bit.

Dato che l'architettura Fermi ha 16 cammini di elaborazione SIMD, ciascun cammino contiene 2048 registri. Ciascun thread CUDA legge un elemento da ciascuno dei registri vettoriali. Si noti che un thread CUDA è una sezione verticale di un thread di istruzioni SIMD, corrispondente a un elemento eseguito da uno dei cammini SIMD. Fate attenzione che i thread CUDA sono molto diversi dai thread POSIX: non si possono effettuare chiamate arbitrarie al sistema operativo o sincronizzarli in modo arbitrario.

## Strutture di memoria delle GPU NVIDIA

La Figura 6.10 mostra le strutture di memoria di una GPU NVIDIA. La memoria a disposizione sul chip, locale a ciascun processore SIMD multithread, si chiama *memoria locale*. Viene condivisa tra i diversi cammini di elabora-



**Figura 6.10 Struttura della memoria delle GPU.** La memoria delle GPU viene condivisa tra i cicli vettoriali. Tutti i thread di istruzioni SIMD all'interno dello stesso blocco di thread condividono la stessa memoria locale.

zione SIMD all'interno dello stesso processore SIMD multithread, ma non viene condivisa anche con gli altri processori SIMD multithread. La memoria DRAM esterna al chip, condivisa da tutta la GPU e da tutti i blocchi di thread, è chiamata *memoria della GPU*.

Invece di basarsi su cache di grandi dimensioni per contenere l'intero insieme di lavoro di un'applicazione, le GPU utilizzano cache più piccole per inviare i dati in stream e sfruttano molto il multithread dei thread delle istruzioni SIMD per nascondere la grande latenza della DRAM, dato che i loro insiemi di lavoro possono raggiungere le centinaia di Megabyte e non potrebbero essere contenuti nell'ultimo livello di cache di un microprocessore multicore. Dato che viene utilizzato il multithread hardware per nascondere la latenza della DRAM, l'area del chip utilizzata per le cache nei processori viene utilizzata nelle GPU per risorse di calcolo e registri aggiuntivi per memorizzare lo stato dei molti thread delle istruzioni SIMD.

**Approfondimento.** Anche se nascondere la latenza della memoria è la filosofia delle GPU, nelle GPU e nei processori vettoriali di ultima generazione sono state aggiunte le memorie cache. Per esempio, la recente architettura Fermi contiene le cache, ma queste sono state pensate come filtri per ridurre la larghezza di banda verso la memoria della GPU o come acceleratore per poche variabili la cui latenza non può essere nascosta dal multithread. La memoria locale per i frame di stack, le chiamate a funzione e il trasferimento dei dati dai registri sono particolarmente adatti alle cache, dato che la latenza è importante quando si chiama una funzione. Le cache consentono anche di risparmiare energia, poiché l'accesso a una cache presente sul chip consente di risparmiare molta energia rispetto all'accesso ai molti chip esterni di DRAM.

## La prospettiva delle GPU

Ad alto livello, i calcolatori multicore con le estensioni contenenti istruzioni SIMD hanno diverse somiglianze con le GPU. La Figura 6.11 riassume le somiglianze e le differenze tra le due architetture. Sono entrambe architetture MIMD nelle quali i processori utilizzano più cammini di elaborazione SIMD, ma le GPU contengono più processori e un numero molto maggiore di cammini. Entrambe utilizzano un multithreading hardware per migliorare l'utilizzazione del processore, ma le GPU forniscono il supporto hardware per un gran numero di thread. Entrambe utilizzano le cache: le GPU utilizzano delle cache più piccole per inviare dati in streaming, i calcolatori multicore utilizzano grandi cache multilivello che possano contenere l'intero insieme di lavoro. Entrambe hanno uno spazio di indirizzamento su 64 bit, ma la memoria principale delle GPU è molto più piccola. Anche se le GPU supportano la protezione della memoria a livello di pagine, non supportano ancora la paginazione intensa.

I processori SIMD sono simili anche ai processori vettoriali. I molti processori SIMD in una GPU si comportano come core MIMD indipendenti, proprio come molti calcolatori vettoriali hanno più processori vettoriali. Secondo questa visione, l'architettura Fermi GTX 580 sarebbe una macchina a 16 core con il supporto hardware per il multithreading, dove ogni core ha 16 cammini di elaborazione. La maggiore differenza è il multithreading, che è fondamentale per la GPU ma manca nella maggior parte dei processori vettoriali.

Le GPU e le CPU non hanno un comune progenitore: non c'è un'architettura passata che le possa spiegare entrambe. Il risultato di ciò è che le GPU non utilizzano la terminologia comune alla comunità delle architetture dei calcolatori, cosa che ha creato un po' di confusione sulle GPU e sulla loro modalità di funzionamento. Per aiutarvi a chiarire questa confusione, i termini più utilizzati in questo paragrafo vengono riportati in Figura 6.12, assieme al termine più simile utilizzato nel mondo dei calcolatori, il termine ufficiale utilizzato nelle GPU NVIDIA, nel caso foste interessati, e una breve descrizione del termine (da sinistra a destra). Questa "Stele di Rosetta delle GPU" vi aiuterà a mettere in relazione questo paragrafo e il suo contenuto con le descrizioni convenzionali delle GPU, come quella riportata nell'Appendice C.

Mentre le GPU si stanno muovendo verso il calcolo, non possono abbandonare il ruolo nel quale eccellono: la grafica. Per questo motivo, la domanda che i progettisti di GPU dovrebbero porre è come l'hardware sviluppato per la grafica possa essere utilizzato per aumentare le prestazioni in un insieme più ampio di applicazioni.

Dopo avere descritto due diversi tipi di architetture MIMD che hanno lo stesso spazio di indirizzamento condiviso, introduciamo ora dei processori paralleli nei quali ciascun processo ha il suo spazio di indirizzamento privato,

Caratteristica	Multicore con SIMD	GPU
Processori SIMD	Da 4 a 8	Da 8 a 16
Cammini SIMD per processore	Da 2 a 4	Da 8 a 16
Supporto hardware al multithreading per i thread SIMD	Da 2 a 4	Da 16 a 32
Dimensione della cache più grande	8 MiB	0,75 MiB
Aampiezza degli indirizzi di memoria	64 bit	64 bit
Dimensione della memoria principale	Da 8 GiB a 256 GiB	Da 4 GiB a 6 GiB
Protezione della memoria a livello di pagina	Sì	Sì
Paginazione	Sì	No
Coerenza della cache	Sì	No

Figura 6.11 Similarità e differenze tra multicore con estensioni multimediali SIMD e GPU recenti.

Tipo	Nome descrittivo	Termine più simile al di fuori delle GPU	Termine ufficiale CUDA / GPU NVIDIA	Definizione utilizzata nel libro
Astrazioni nei programmi	Cicli vettorizzabili	Cicli vettorizzabili	Grid (griglia)	Un ciclo vettorizzabile, eseguito su GPU, costituito da uno o più blocchi di thread (corpo del ciclo vettorizzato) che possono essere eseguiti in parallelo
	Corpo dei cicli vettorizzati	Corpo dei cicli vettorizzati (strip-mined)	Blocco di thread	Un ciclo vettorizzato, eseguito su un processore SIMD multithread, costituito da uno o più blocchi di thread di istruzioni SIMD, che comunicano attraverso la memoria locale
	Sequenza di operazioni nel cammino elaborazione SIMD	Un'iterazione di un ciclo scalare	Thread CUDA	Un taglio verticale di un thread di istruzioni SIMD corrispondenti a un elemento eseguito da un cammino SIMD. Il risultato viene memorizzato a seconda del valore di una maschera e dei registri predicatori
Oggetti delle architetture	Un thread di istruzioni SIMD	Thread di istruzioni vettoriali	Warp	Un thread tradizionale che contiene solo istruzioni SIMD che vengono eseguite su un processore SIMD multithread. I risultati vengono memorizzati a seconda del contenuto di una maschera definita per ogni elemento
	Istruzione SIMD	Istruzione vettoriale	Istruzione PTX	Una singola istruzione SIMD eseguita in parallelo attraverso diversi cammini di elaborazione SIMD
Hardware di elaborazione	Processore SIMD multithread	Processore vettoriale (multithread)	Multiprocessore streaming	Un processore SIMD multithread esegue thread di istruzioni SIMD, indipendenti da altri processori SIMD
	Scheduler dei blocchi di thread	Processore scalare	Motore di Giga thread	Assegna blocchi multipli di thread (corpi di cicli vettorizzati) a processori SIMD multithread
	Scheduler di thread SIMD	Scheduler di thread in CPU multithread	Scheduler di warp	Componente hardware che invia a esecuzione thread di istruzioni SIMD quando sono pronte per essere eseguite. Memorizza la traccia dell'esecuzione del thread
Hardware della memoria	Cammino SIMD	Cammino vettoriale	Processore di thread	Un'unità hardware che programma e invia a esecuzione thread di istruzioni SIMD su un singolo elemento. I risultati vengono memorizzati a seconda del valore che assume una maschera
	Memoria GPU	Memoria principale	Memoria globale (Global Memory)	Memoria DRAM accessibile a tutti i processori SIMD multithread di una GPU
	Memoria locale	Memoria locale	Memoria condivisa (Shared Memory)	Memoria veloce SRAM locale per un processore multithread SIMD, non disponibile agli altri processori SIMD
Registri dei cammini SIMD		Registri dei cammini vettoriali	Registri di processore di thread	Registri utilizzati da un singolo cammino SIMD allocati all'interno di un intero blocco di thread (corpo dei cicli vettorizzati)

**Figura 6.12 Guida rapida ai termini utilizzati per le GPU.** Riportiamo nella prima colonna i termini hardware, raggruppati in quattro insiemi. Dall'alto al basso: astrazioni nei programmi, oggetti delle architetture, hardware di elaborazione e hardware della memoria.

cosa che rende molto più facile costruire sistemi di dimensioni molto maggiori. I servizi Internet che utilizzate ogni giorno dipendono da questi sistemi a grandissima scala.

**Approfondimento.** Sebbene quando le GPU sono state introdotte avessero una memoria separata dalla CPU, sia AMD sia Intel hanno annunciato prodotti "integriti" che combinano GPU e CPU in modo da utilizzare una memoria condivisa. La sfida è mantenere l'elevata larghezza di banda della memoria propria della GPU anche in queste architetture integrate.

## Autovalutazione

Vero o falso? Le GPU sono basate su chip di DRAM grafica per ridurre la latenza della memoria e quindi aumentare le prestazioni nelle applicazioni grafiche.

## 6.7 Cluster, calcolatori per centri di calcolo e altri multiprocessori a scambio di messaggi

### Scambio di messaggi:

comunicazione tra diversi processori ottenuta inviando e ricevendo informazioni in modo esplicito.

### Procedura di invio dei messaggi:

una procedura utilizzata da un processore in un calcolatore dotato di memoria privata per passare un messaggio a un altro processore.

### Procedura di ricezione dei messaggi:

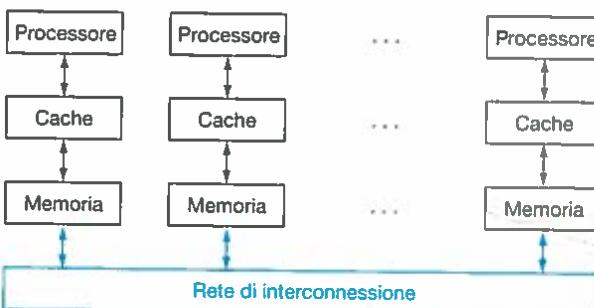
una procedura utilizzata da un processore in un calcolatore dotato di memoria privata per ricevere un messaggio da un altro processore.

Il modello alternativo utilizzato per la condivisione dello spazio di indirizzamento prevede che ciascun processore abbia il proprio spazio privato di indirizzamento fisico. La Figura 6.13 illustra l'organizzazione classica di un multiprocessore con spazi privati multipli di indirizzamento. Questo multiprocessore comunica attraverso uno **scambio di messaggi** (*message passing*) esplicito, che dà tradizionalmente il nome a questo tipo di calcolatori. Un sistema di questo tipo è dotato di **procedure di invio** (*send*) e **procedure di ricezione** (*receive*) **dei messaggi**. La coordinazione è quindi garantita dallo scambio dei messaggi, poiché ogni processore sa quando invia un messaggio e il processore ricevente sa quando il messaggio gli è arrivato. Se il processore che spedisce il messaggio deve sapere se il messaggio è arrivato a destinazione, il ricevente può inviare al mittente un messaggio di ricevuta.

Sono stati fatti alcuni tentativi di costruire calcolatori ad alte prestazioni basati su reti di scambio di messaggi. Questi calcolatori hanno permesso di raggiungere prestazioni assolute di comunicazione migliori di quelle ottenute con i cluster costruiti utilizzando reti LAN (*Local Area Networks*, reti locali di area). Oggi molti supercalcolatori utilizzano reti dedicate. Il problema è che si sono rivelati troppo costosi rispetto alle LAN basate su Ethernet e poche applicazioni, al di fuori del calcolo ad alte prestazioni, potrebbero giustificare l'aumento di velocità di trasferimento, dati i costi molto maggiori.

## Interfaccia hardware/software

I calcolatori che basano la comunicazione sullo scambio di messaggi invece che sulla memoria condivisa con coerenza delle cache sono molto più facili da progettare (vedi par. 5.8). Il vantaggio per i programmati è che la comunicazione è resa esplicita, il che significa che si possono verificare meno sorprese dal punto di vista delle prestazioni rispetto alla comunicazione implicita dei calcolatori con memoria condivisa. Il rovescio della medaglia per i programmati è che è più difficile trasportare un programma sequenziale su un'architettura basata su scambio di messaggi, poiché tutti gli scambi di informazione devono essere identificati prima dell'esecuzione (altrimenti il programma non potrà funzionare). Le memorie condivise con coerenza delle cache lasciano all'hardware il compito di identificare quali dati occorra trasmettere, cosa che rende il trasporto dei programmi più facile. Ci sono opinioni discordanti su quale sia la strategia migliore per ottenere prestazioni elevate, considerati i pro e i contro della comunicazione implicita, ma non c'è confusione su quale sia la collocazione nel mercato delle due architetture. I microprocessori multicore utilizzano la memoria fisica condivisa, mentre i nodi di un cluster comunicano con gli altri nodi utilizzando lo scambio di messaggi.



**Figura 6.13** Organizzazione classica di un multiprocessore con spazi di indirizzamento privati multipli, chiamato tradizionalmente multiprocessore a scambio di messaggi. Si noti che, a differenza dell'SMP di Figura 6.7, la rete di interconnessione non si trova tra le cache e la memoria, ma tra i diversi nodi processore-memoria.

Alcune applicazioni concorrenti vengono eseguite efficientemente su hardware parallelo sia basato su un indirizzamento condiviso della memoria sia basato su scambio di messaggi. In particolare, il parallelismo a livello di task e le applicazioni che richiedono poche comunicazioni, come la ricerca sul web, i server mail e i file server, non hanno bisogno di un indirizzamento condiviso della memoria per potere essere eseguiti efficientemente. I **cluster** sono quindi diventati l'esempio più diffuso di calcolatore parallelo basato su scambio di messaggi. Dato che posseggono memorie separate, i nodi di un cluster hanno in esecuzione una copia del sistema operativo distinta da quella in esecuzione sugli altri nodi; i core di un microprocessore, invece, sono collegati mediante una rete ad alta velocità contenuta nel loro chip e un sistema di condivisione della memoria su più chip utilizza le interconnessioni della memoria per la comunicazione. Queste interconnessioni hanno una banda maggiore e una latenza minore, e consentono quindi una comunicazione molto migliore nei multiprocessori a memoria condivisa.

L'utilizzo di memorie separate per i processi utente, che è un punto di debolezza dal punto di vista della programmazione parallela, diventa un punto di forza per l'affidabilità (vedi par. 5.5). Dato che un cluster è costituito da calcolatori indipendenti connessi attraverso una LAN, è molto più facile sostituire un calcolatore indipendente senza dovere fermare il sistema in un cluster che in un multiprocessore a memoria condivisa. Fondamentalmente, l'utilizzo di indirizzi condivisi significa che diventa difficile isolare un processore e sostituirlo senza un'intensa attività del sistema operativo e un'accurata progettazione del server. È anche facile per i cluster ridurre il numero di calcolatori in modo graduale quando un server si guasta, aumentando così l'affidabilità. Dato che il software dei cluster è uno strato che viene eseguito sopra il sistema operativo locale in esecuzione sui singoli calcolatori, è molto più facile disconnettere e sostituire un calcolatore guasto. Dato che i cluster sono costruiti a partire da interi calcolatori e da reti indipendenti e scalabili, è facile espandere il sistema senza dovere fermare l'applicazione che viene eseguita dal cluster.

Il costo inferiore, la maggiore disponibilità e la facile e scalabile estendibilità rendono i cluster attraenti per i fornitori di servizi Internet, nonostante le scarse prestazioni nella comunicazione rispetto ai multiprocessori a larga scala con memoria condivisa. I motori di ricerca che centinaia di milioni di persone utilizzano ogni giorno dipendono da questa tecnologia: Amazon, Facebook, Microsoft e altre società dispongono tutte di centri di calcolo multipli, ciascuno contenente decine di migliaia di server. Sicuramente, l'utilizzo di processori multipli da parte delle società che offrono servizi attraverso Internet si è rivelato un grandissimo successo.

**Cluster:** insieme di calcolatori connessi attraverso l'I/O a una rete standard per formare un multiprocessore a scambio di messaggi.



AFFIDABILITÀ

## Calcolatori per grandi centri di calcolo

I servizi Internet, quali quelli descritti poco fa, richiedono la costruzione di nuovi edifici per contenere, alimentare e raffreddare 100 000 server. Sebbene si possano classificare come cluster enormi, la loro architettura e il loro funzionamento sono più sofisticati: essi funzionano come un singolo calcolatore gigante e hanno un costo dell'ordine di 150 milioni di dollari per l'edificio, l'infrastruttura deputata al raffreddamento e alle interconnessioni, i server e i dispositivi di interconnessione che collegano e ospitano un numero di server tra i 50 000 e i 100 000. Consideriamo questi sistemi come un nuovo tipo di calcolatori destinati ai grossi centri di calcolo (WSC, *Warehouse-Scale Computers*).

*Chiunque può costruire una CPU veloce. Il difficile è costruire un sistema veloce.*

Seymour Cray, considerato il "padre" dei supercalcolatori

## Interfaccia hardware/software

I framework più popolari per l'elaborazione batch in un WSC sono MapReduce (mappa e riduci [Dean, 2008]) e il suo gemello open-source Hadoop. Ispirata alla funzione Lisp che porta lo stesso nome, la funzione Map, eseguita su migliaia di server, per prima cosa applica una funzione fornita dal programmatore a ciascun dato dell'input logico per produrre un risultato intermedio rappresentato da coppie di valori chiave. La funzione Reduce raccoglie gli output delle funzioni e li unisce utilizzando una seconda funzione definita dal programmatore. Con un supporto software appropriato, queste due funzioni possono essere eseguite in modo altamente parallelo e allo stesso tempo facile da capire e da utilizzare: in 30 minuti, un programmatore alle prime armi può eseguire una funzione MapReduce su migliaia di server.

Per esempio, un programma MapReduce può calcolare il numero di volte che viene ripetuta ogni parola inglese in un grande insieme di documenti. Riportiamo qui sotto una versione semplificata del programma, la quale mostra solo il ciclo più interno che cerca una sola parola tra tutte le parole possibili:

```
map(String chiave, String valore):
    // chiave: nome documento
    // valore: contenuto documento
    for ciascuna parola w in valore:
        EmitIntermediate(w, "1"); // Producì un elenco di tutte le parole
reduce(String chiave, Iterator valore):
    // chiave: una parola
    // valore: una lista di conteggi
    int risultato = 0;
    for ciascun v in valore:
        risultato += ParseInt(v); // ottieni un intero dalla coppia di valori
    Emit(AsString(risultato));
```

La funzione `EmitIntermediate()` utilizzata nella funzione Map emette, cioè restituisce, ogni occorrenza della parola cercata nel documento e il valore 1. La funzione Reduce somma 1 per tutte le occorrenze della parola trovata in tutti i documenti utilizzando `ParseInt()`. L'ambiente di esecuzione di MapReduce assegna le funzioni di mappatura e di riduzione ai diversi server di un WSC.

A questo livello di scala estrema, che richiede innovazione nella distribuzione della potenza, del raffreddamento, della supervisione e delle operazioni, un WSC è il discendente moderno dei supercalcolatori del 1970, rendendo così Seymour Cray il padrino dei progettisti delle architetture WSC. I suoi potentissimi calcolatori riuscivano a eseguire calcoli che non potevano essere eseguiti da nessun'altra macchina, ma erano così costosi che solamente poche società erano in grado di acquistarli. Oggi l'obiettivo è fornire la tecnologia dell'informazione al mondo e non un'elevata potenza di calcolo agli scienziati e agli ingegneri. Quindi i WSC di oggi hanno un ruolo sociale che i supercalcolatori di Cray non avevano.

I WSC hanno alcuni obiettivi in comune con i server, ma presentano tre principali differenze: vediamo quali sono.

1. *Ampio, facile parallelismo.* Un progettista di server si deve preoccupare che le applicazioni richieste dal mercato di riferimento abbiano abbastanza parallelismo da giustificare la quantità di hardware parallelo e che il costo dell'hardware che gestisce le comunicazioni, necessario per sfruttare il parallelismo, non sia troppo alto. Un progettista di WSC non ha queste preoccupazioni. Per prima cosa, applicazioni batch come MapReduce beneficiano del gran numero di insiemi di dati indipendenti che richiedono un'elaborazione indipendente, quali i miliardi di pagine web elaborate da un web crawler. Inoltre, applicazioni che offrono servizi Internet interattivi, note anche come **software come servizio (SaaS, Software as a Service)**, possono beneficiare

**Software come servizio (SaaS):** invece di vendere il software, che viene poi installato e fatto girare dall'acquirente sul suo calcolatore, il software viene fatto girare su un sistema remoto e messo a disposizione dei clienti attraverso Internet, tipicamente attraverso un'interfaccia web. Ai clienti viene fatturato in base a quanto utilizzano il sistema.

dei milioni di utenti indipendenti che utilizzano i servizi offerti. Le letture e le scritture raramente sono dipendenti in un SaaS, per cui i SaaS raramente necessitano di sincronizzazione. Per esempio, la ricerca utilizza un indice a sola lettura e la posta elettronica di solito legge e scrive informazioni indipendenti. Chiameremo questo tipo di parallelismo facile **parallelismo a livello di richiesta**, dato che molti sforzi possono procedere in parallelo in modo naturale con pochissima necessità di comunicazione o sincronizzazione.



2. *Conteggio dei costi operativi.* Tradizionalmente, i server vengono progettati per massimizzare le prestazioni all'interno di un certo budget e i progettisti si preoccupano dell'energia solamente per essere sicuri di non superare la capacità di raffreddamento del contenitore; in genere non si curano dei costi operativi di un server, che di solito impallidiscono di fronte al costo di acquisto. I WSC hanno un tempo di vita più lungo (l'edificio, l'infrastruttura elettrica e il raffreddamento vengono ammortizzati nell'arco di dieci anni o più), per cui i costi operativi diventano significativi: l'energia, la distribuzione dell'elettricità e il raffreddamento rappresentano più del 30% dei costi di un WSC considerati nell'arco di 10 anni.

3. *Scala, opportunità e problemi associati con la scala.* Per costruire un singolo WSC, occorre acquistare 100 000 server con la relativa infrastruttura di supporto, che comporta sconti tipici per acquisti di grandi volumi. Quindi i WSC sono così grandi all'interno che si può fare economia di scala anche se non esistono tanti WSC al mondo. Questa economia di scala ha portato al *cloud computing*: il minor costo per unità di calcolo di un WSC ha consentito alle società che lavorano con il cloud di noleggiare server ottenendo un guadagno e rimanendo comunque al di sotto del costo che altre società potrebbero praticare se offrissero i servizi direttamente. Il rovescio della medaglia delle opportunità offerte dalla grande scala è la gestione dei guasti e della loro frequenza a questa scala. Anche se un server ha un MTTF (tempo prima di un guasto) di 25 anni (200 000 ore), un progettista di WSC deve tenere conto di una media di 5 guasti di server al giorno. Il paragrafo 5.16 riporta che la frequenza di guasti per anno (AFR, *Annualized Failure Rate*) misurata da Google è del 2-4%. Se ci fossero 4 dischi per server e la loro frequenza annuale di guasto fosse del 2%, un progettista di WSC dovrebbe affrontare un guasto dei dischi ogni ora. La tolleranza ai guasti è perciò più importante per i progettisti di WSC che per i progettisti di server.

L'economia di scala ottenuta con i WSC ha realizzato il sogno di lunga data di offrire l'elaborazione come servizio. Il *cloud computing* consente a chiunque, ovunque nel mondo, che abbia una buona idea, un modello di business e una carta di credito, di caricare su migliaia di server la propria visione per diffonderla istantaneamente in tutto il mondo. Naturalmente ci sono alcuni ostacoli che limitano la crescita del *cloud computing*, quali la sicurezza, la privacy, gli standard e il tasso di crescita della banda di Internet, ma crediamo che questi ostacoli verranno superati e che i WSC e il *cloud computing* continueranno a fiorire.

Per mettere il tasso di crescita del *cloud computing* nella giusta prospettiva, i servizi web di Amazon hanno annunciato nel 2012 che stanno aggiungendo *ogni giorno* nuova capacità ai server per supportare tutta l'infrastruttura globale di Amazon dal 2003, quando fatturava 5,2 miliardi di dollari e aveva 6000 dipendenti.

Ora che abbiamo capito l'importanza dei multiprocessori a scambio di messaggi, specialmente per il *cloud computing*, esaminiamo in quali modi si possono interconnettere i nodi di un WSC. Grazie alla **Legge di Moore** e all'aumento del numero di core per chip, dobbiamo inserire delle reti di interconnessione all'interno dei chip di microprocessori, perciò queste topologie di interconnessione sono importanti sia nel piccolo sia nel grande.



**Approfondimento.** La funzione MapReduce ordina le coppie chiave-valore al termine della fase di Map per produrre gruppi che condividono la stessa chiave. Questi gruppi vengono quindi passati alla fase di Reduce.

**Approfondimento.** Un'altra forma di calcolo a larga scala è il *grid computing* (calcolo su griglia), in cui i calcolatori sono sparsi su aree geografiche di grandi dimensioni e i programmi che vengono eseguiti dai diversi calcolatori comunicano attraverso reti ad ampio raggio. La forma più popolare e originale di grid computing è rappresentata dal progetto SETI@home (SETI, Serch for ExtraTerrestrial Intelligence, ricerca di intelligenza extraterrestre). L'osservazione su cui si basava il progetto quando fu concepito era che, in ogni momento, ci sono milioni di PC accesi inattivi, che non fanno nulla; se solo qualcuno avesse sviluppato il software adatto, sarebbe diventato possibile collegare e far lavorare questi PC per una buona causa. Ogni calcolatore poteva quindi essere impegnato su una piccola frazione del problema, in maniera indipendente dagli altri. Al progetto SETI hanno partecipato più di 5 milioni di utenti di calcolatori di oltre 200 Paesi, con più del 50% degli utenti residenti fuori dagli Stati Uniti. Alla fine del 2011, le prestazioni medie della griglia di SETI@home erano di 3,5 PetaFLOPS.

### Autovalutazione

1. Vero o falso? Come gli SMP, i calcolatori basati su scambi di messaggi si basano sui lock per la sincronizzazione.
2. Vero o falso? I cluster hanno memorie separate e perciò richiedono più copie del sistema operativo.

## 6.8 | Introduzione alle topologie delle reti di calcolatori

I chip multicore hanno bisogno di reti locali sul chip per connettere i core tra loro, e i cluster richiedono una LAN per connettere i server tra di loro. In questo paragrafo esamineremo i pro e i contro delle diverse topologie di reti di interconnessione.

Il costo di una rete dipende dal numero di commutatori, detti *switch*, dal numero dei dispositivi che un singolo switch collega alla rete, dall'ampiezza (numero di bit) del collegamento e dalla lunghezza dei collegamenti sul chip (sul quale viene stampata la rete). Per esempio, alcuni core possono essere fra loro adiacenti, mentre altri possono trovarsi sul lato opposto del chip. Anche le prestazioni delle reti sono variegate: esse possono essere valutate come latenza nell'invio e ricezione dei messaggi quando la rete è scarica, come throughput misurato come numero massimo di messaggi che possono essere trasmessi in un certo intervallo di tempo, come durata dei ritardi causati dalla contesa per l'accesso a una parte della rete o come prestazioni misurate in funzione della sequenza corrente degli accessi alla rete. Un'altra caratteristica che a volte si richiede alle reti è la tolleranza ai guasti, poiché alcuni sistemi devono poter operare anche in presenza di componenti non funzionanti. Infine, in questa epoca di attenzione verso il consumo energetico, l'efficienza energetica delle diverse organizzazioni delle reti può essere più importante degli altri aspetti.

Normalmente le reti di comunicazione vengono rappresentate attraverso grafi, in cui un arco rappresenta un collegamento; un nodo processore-memoria è raffigurato da un quadratino nero e lo switch da un cerchietto colorato. In questo paragrafo tutti i collegamenti sono *bidirezionali*, cioè l'informazione può

fluire in entrambe le direzioni; come abbiamo già detto, tutte le reti sono formate da switch, a ciascuno dei quali sono collegati i nodi processore-memoria e gli altri switch. Una prima topologia di rete connette una sequenza di nodi:



Questa topologia è chiamata *topologia ad anello (ring)*; poiché non tutti i nodi sono connessi in modo diretto, alcuni messaggi devono attraversare dei nodi intermedi per arrivare alla destinazione finale. A differenza dei bus, un insieme condiviso di linee che consentono la trasmissione a tutti i dispositivi connessi, con un collegamento ad anello si possono effettuare più trasferimenti simultaneamente.

Poiché è possibile scegliere fra molte topologie di rete, è necessario definire delle metriche che misurino le prestazioni delle diverse topologie; due metriche sono particolarmente diffuse: la prima è la **larghezza di banda totale della rete**, ossia la larghezza della banda di ciascun collegamento moltiplicata per il numero di collegamenti. Questa metrica valuta la condizione ottimale di funzionamento. Per la rete ad anello con  $P$  processori mostrata sopra, la larghezza di banda totale della rete sarebbe  $P$  volte la larghezza di banda del singolo collegamento; la larghezza totale della banda di trasferimento di un bus, invece, è semplicemente la larghezza di banda del bus stesso. Questa è la condizione di funzionamento più favorevole.

Per valutare la condizione di funzionamento più sfavorevole si introduce una seconda metrica, la **larghezza di banda di bisezione**. Questa viene calcolata suddividendo la macchina in due parti, ciascuna contenente metà dei nodi, e poi sommando la larghezza della banda di trasferimento delle connessioni che attraversano questa linea immaginaria di suddivisione. La larghezza della banda di bisezione di una rete ad anello è pari a due volte la larghezza di banda della singola connessione, mentre per un bus è uguale alla larghezza della banda di trasferimento del collegamento. Se una singola connessione è veloce quanto un bus, il collegamento tramite topologia ad anello sarà due volte più veloce del bus nel caso peggiore, ma sarà  $P$  volte più veloce nel caso migliore.

Dato che alcune topologie di rete non sono simmetriche, diventa un problema stabilire dove far passare la linea di suddivisione immaginaria utilizzata per misurare la banda di bisezione. Essendo questa metrica legata al caso peggiore, si sceglie di tracciare la linea che produce la misura peggiore; in altri termini, occorre calcolare la larghezza di banda di trasferimento per tutte le possibili bisezioni e scegliere quella associata alla banda più stretta. Una visione così pessimistica è giustificata dal fatto che i programmi paralleli sono spesso limitati dall'anello più debole della catena di comunicazione.

All'estremo opposto rispetto alle reti ad anello ci sono le **reti completamente connesse**, nelle quali ciascun processore ha un collegamento bidirezionale con tutti gli altri processori. Nel caso di reti completamente connesse, la larghezza totale della banda della rete è pari a  $P \times (P - 1)/2$ , e la banda di bisezione è pari a  $(P/2)^2$ .

L'enorme miglioramento delle prestazioni di una rete completamente connessa è controbilanciato dall'aumento dei costi, che ha indotto i progettisti a inventare nuove topologie di rete con caratteristiche intermedie. La valutazione delle diverse topologie dipende in gran parte dal tipo di comunicazione necessaria a eseguire il carico di lavoro dei programmi paralleli in esecuzione sulla macchina.

È difficile tenere il conto di tutte le topologie di rete discusse nelle pubblicazioni specialistiche, ma quelle effettivamente utilizzate nei processori paralleli

#### **Larghezza di banda della rete:**

informalmente, la velocità di trasferimento di picco di una rete; può riferirsi alla velocità di un singolo collegamento o alla velocità di trasferimento complessiva di tutti i collegamenti della rete.

#### **Larghezza di banda di bisezione:**

la larghezza della banda di trasferimento tra due parti uguali di un multiprocessore; è la misura della larghezza di banda che si ottiene suddividendo in due parti il microprocessore, nel caso più sfavorevole.

#### **Rete completamente connessa:**

una rete che collega i nodi processore-memoria fornendo un canale di comunicazione dedicato a ogni nodo.

**Rete multistadio:** una rete dotata di un piccolo switch per ogni nodo.

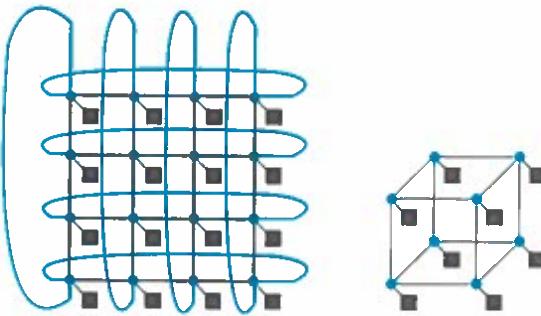
**Rete completamente connessa:** detta anche **rete crossbar**, una rete che consente a ogni nodo di comunicare con gli altri nodi con un solo passaggio.

si contano sulle dita di una mano. La Figura 6.14 illustra due delle topologie di rete più diffuse.

Anziché posizionare un processore in ciascun nodo della rete, alcuni di questi nodi possono essere occupati soltanto da uno switch; gli switch sono più piccoli dell'insieme completo costituito da processore-memoria-switch e quindi possono essere resi più densi, diminuendo le distanze e quindi migliorando le prestazioni. Le reti costruite in questo modo sono spesso chiamate **reti multistadio**, per esprimere il fatto che un messaggio deve attraversare più stadi. I tipi di reti multistadio sono altrettanto numerosi di quelli a stadio singolo: la Figura 6.15 mostra due delle reti multistadio più diffuse. Una **rete completamente connessa** o **rete crossbar** consente a ciascun nodo di comunicare con qualunque altro nodo con un solo passaggio attraverso la rete. Una **rete Omega** utilizza meno hardware di una rete crossbar ( $2n \log_2 n$  invece di  $n^2$  switch), ma si possono verificare situazioni di conflitto fra i messaggi a seconda dell'andamento della comunicazione. Per esempio, la rete Omega di Figura 6.15 non può inviare un messaggio da  $P_0$  a  $P_6$  e, contemporaneamente, un messaggio da  $P_1$  a  $P_4$ .

### Implementazione delle topologie di rete

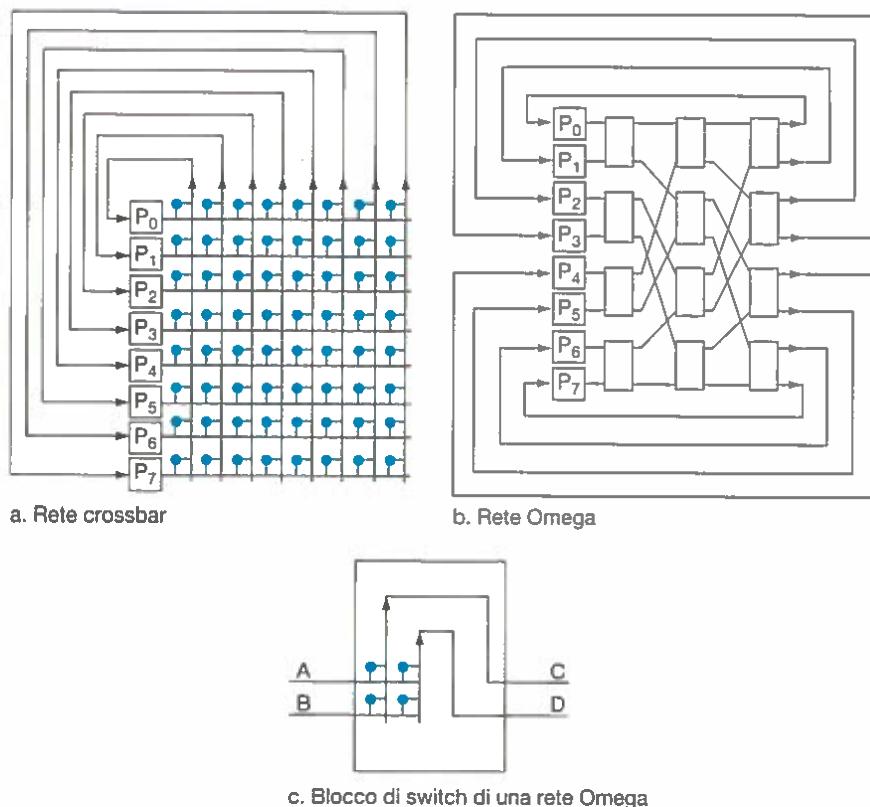
La semplice analisi delle reti fatta in questo paragrafo non considera importanti aspetti pratici di cui occorre invece tenere conto quando vengono costruite. La lunghezza di ciascun collegamento influisce sul costo della comunicazione a frequenze di clock elevate; in genere, maggiore è la lunghezza, più costoso è l'uso di un collegamento a frequenze di clock elevate. Distanze più brevi, inoltre, rendono più semplice assegnare più linee allo stesso collegamento, poiché la potenza richiesta per pilotare più linee è minore se le linee sono corte; le linee più corte, infine, sono meno costose di quelle più lunghe. Un altro vincolo pratico è che tutte le strutture tridimensionali devono essere mappate su chip che sono costituiti essenzialmente da strati bidimensionali. Un'ultima considerazione riguarda l'energia: l'assorbimento di energia può obbligare a utilizzare nei chip multicore le semplici topologie di reti a griglia. La conclusione è che topologie di rete che sembrano eleganti quando vengono disegnate alla lavagna possono rivelarsi poco funzionali quando devono essere realizzate sui wafer silicio o nei centri di elaborazione dati. Ora che abbiamo capito l'importanza dei cluster e le topologie di rete che si possono utilizzare, esaminiamo l'hardware e il software dell'interfaccia tra processore e rete.



a. Griglia 2D o rete di 16 nodi

b.  $n$ -cubo con 8 nodi ( $8 = 2^3$ , quindi  $n = 3$ )

**Figura 6.14 Topologie delle reti che vengono utilizzate nei processori paralleli in commercio.** I pallini blu rappresentano i connettori a commutazione, detti switch, mentre i quadratini neri rappresentano i nodi processore-memoria. Anche se ciascuno switch può connettere diversi canali, generalmente un solo canale alla volta è collegato attivamente al processore. La topologia booleana  $n$ -cubo è una topologia di connessione  $n$ -dimensionale con  $2^n$  nodi e  $n$  canali per ogni switch, più un canale per ogni processore, e quindi  $n$  nodi sono connessi direttamente a ogni switch. Spesso queste topologie di base vengono arricchite da archi aggiuntivi per migliorare le prestazioni e l'affidabilità.



**Figura 6.15** Topologie di rete multistadio utilizzate frequentemente nel caso di otto nodi. In questi schemi gli switch sono più semplici rispetto agli schemi precedenti, perché i collegamenti sono unidirezionali. I dati entrano a sinistra ed escono sulla destra. L'interruttore (switch) nella Figura c può collegare A a C e B a D, oppure B a C e A a D. La rete crossbar utilizza  $n^2$  switch, dove  $n$  è il numero dei processori, mentre la rete Omega utilizza  $2n \log n$  blocchi di switch, ciascuno dei quali contiene quattro switch semplici. In questo caso la rete crossbar utilizza quindi 64 switch, rispetto ai 12 blocchi di switch e ai 48 switch della rete Omega. La rete crossbar, però, può supportare qualunque combinazione di trasferimenti tra i processori, al contrario della rete Omega.

### Autovalutazione

Vero o falso? Per un anello contenente  $P$  nodi, il rapporto tra la banda totale della rete e la banda di bisezione è  $P/2$ .

## 6.9 Come comunicare con il mondo esterno: le reti dei cluster

Questo paragrafo, disponibile online , descrive l'hardware e il software utilizzato per connettere tra loro i nodi di un cluster attraverso la rete. L'esempio utilizzato è la rete Ethernet a 10 Gigabit/secondo, connessa al calcolatore attraverso il bus PCIe (*Peripheral Component Interconnect Express*). In questo paragrafo descriviamo le ottimizzazioni hardware e software utilizzate per migliorare le prestazioni della rete, compresi i messaggi a zero copia, la comunicazione nello spazio utente, l'utilizzo del polling invece degli interrupt di I/O, e il calcolo hardware dei bit di parità. Anche se l'esempio utilizzato sono le reti, le tecniche descritte in questo paragrafo si applicano anche ai controllori dei dispositivi di memoria di massa e agli altri dispositivi di I/O.

Dopo avere descritto le prestazioni di una rete a un basso livello di dettaglio, nel prossimo paragrafo mostriamo come confrontare tra loro multiprocessori di tutti i tipi mediante programmi ad alto livello.

## 6.10 Benchmark per i multiprocessori

Come abbiamo visto nel Capitolo 1, la valutazione dei sistemi mediante benchmark è un argomento particolarmente delicato, giacché costituisce un modo evidente per determinare quale sia il sistema migliore. I risultati dei benchmark non solo influenzano le vendite dei sistemi in commercio, ma determinano anche la fama dei loro progettisti. Perciò, chi progetta un calcolatore vuole ottenere le prestazioni migliori, ma vuole anche essere sicuro che, nel caso risulti migliore il sistema di qualcun altro, quel sistema abbia veramente le prestazioni più elevate. Per questo motivo sono state introdotte alcune regole per assicurare che i risultati di un benchmark non siano semplicemente frutto di qualche trucco degli ingegneri appositamente escogitato per quel particolare benchmark, bensì il risultato di miglioramenti che hanno effettivamente aumentato le prestazioni dei calcolatori.

Per evitare trucchi, una regola tipica stabilisce che i benchmark non si possono modificare. Il codice sorgente e i dati sono fissi, ed esiste un'unica misura possibile. Il mancato rispetto delle regole rende non validi i risultati ottenuti.

Molti benchmark dei multiprocessori seguono questa politica. Un'eccezione frequente è l'aumento della dimensione del problema, in modo tale da consentire la valutazione di sistemi che hanno un numero variabile di processori; ciò significa che molti benchmark scalano in modo debole, e non forte, anche se occorre fare attenzione quando si confrontano i risultati di programmi che lavorano su problemi di dimensioni diverse.

La Figura 6.16 riporta le caratteristiche di alcuni benchmark paralleli che vengono anche descritti di seguito.

- *Linpack* è una collezione di procedure di algebra lineare, e le procedure che eseguono l'eliminazione gaussiana costituiscono quello che viene chiamato benchmark Linpack. La procedura DGEMM riportata nell'esempio di pagina 185 rappresenta una piccola parte del codice sorgente del benchmark Linpack, ma utilizza la quasi totalità del tempo di esecuzione di questo benchmark. Questo benchmark scala in modo debole, dando la possibilità all'utente di scegliere problemi di qualsiasi dimensione. Inoltre, permette di riscrivere Linpack in qualsiasi altra forma e in un qualsiasi altro linguaggio, purché il risultato calcolato sia corretto ed esegua lo stesso numero di operazioni in virgola mobile per la stessa dimensione del problema. I 500 calcolatori con le prestazioni migliori, in termini di velocità, relativamente al benchmark Linpack sono pubblicati sul sito [www.top500.org](http://www.top500.org) due volte all'anno e la stampa considera il primo dell'elenco il calcolatore più veloce del mondo.
- *SPECrate* è una metrica di throughput basata sui benchmark SPEC CPU, tra cui SPEC CPU 2006 (vedi Cap. 1). Invece di riportare le prestazioni dei singoli programmi, SPECrate esegue diverse copie dello stesso programma simultaneamente, misurando così il parallelismo a livello di task, non essendoci comunicazione tra i diversi task. Si possono eseguire tante copie di un programma quante se ne vogliono, per cui si può affermare che anche questo benchmark scala in modo debole con il numero di processori.
- *SPLASH e SPLASH 2 (Stanford Parallel Applications for Shared Memory)* sono il risultato degli sforzi effettuati negli anni '90 dai ricercatori dell'Università di Stanford per raccogliere una suite di programmi paralleli di

Benchmark	Tipo di scaling	Riprogrammazione	Descrizione
Linpack	Debole	Sì	Algebra lineare con matrici dense [Dongarra, 1979]
SPECrate	Debole	No	Parallelismo di programmi indipendenti [Henning, 2007]
SPLASH 2, Stanford Parallel Applications for Shared Memory [Applicazioni parallele di Stanford per memoria condivisa] [Whoo et al., 1995]	Forte (anche se offre solo due dimensioni di problemi)	No	FFT 1D complessa Decomposizione LU a blocchi Fattorizzazione di Cholesky di matrici sparse a blocchi Ordinamento Radix Sort di interi Barnes-Hut Multipolo veloce adattativo Simulazioni oceaniche Radiosità gerarchica Ray tracing Rendering di volumi Simulazione di acqua con strutture dati spaziali Simulazione di acqua senza strutture dati spaziali
Benchmark paralleli NAS [Bailey et al., 1991]	Debole	Sì (solamente C o Fortran)	EP: applicazioni embarrassingly parallel <sup>1</sup> MG: Simplified Multigrid (griglie multiple semplificate) CG: griglie non strutturate per il calcolo del gradiente coniugato FT: soluzione delle equazioni differenziali alle derivate parziali in 3D, utilizzando la FFT IS: Large Integer Sort, ordinamento di un vettore di interi di grande dimensione
Raccolta di benchmark PARSEC [Bienia et al., 2008]	Debole	No	Blackscholes – assegnamento dei prezzi utilizzando le equazioni differenziali alle derivate parziali Black-Scholes Bodytrack – tracking del corpo di una persona Canneal – versione dell'algoritmo di simulated annealing per l'ottimizzazione dell'instradamento su rete dei pacchetti, che tiene conto della presenza della cache Dedup – algoritmo di compressione di ultima generazione, basato su deduplicazione dei dati Facesim – Simulazione della mimica di un volto umano Ferret – server di ricerca per similarità di contenuto Fluidanimate – fluidodinamica per l'animazione con il metodo SPH Freqmine – ricerca degli insiemi di elementi frequenti Streamcluster – clustering online di dati di input in stream Swaptions – assegnamento di un prezzo agli swap di un portafoglio titoli Vips – elaborazione delle immagini X264 – codifica video H.264
Moduli algoritmici di Berkeley [Asanovic et al., 2006]	Forte e debole	Sì	Macchine a stati finiti Logica combinatoria Attraversamento di grafi Griglie strutturate Matrici dense Matrici sparse Metodi spettrali (FFT) Programmazione dinamica N-corpi MapReduce Backtrack/Branch and bound Inferenza mediante modelli a grafi Griglie non strutturate

<sup>1</sup> Si dice che un'applicazione è *embarrassingly parallel* se ha bisogno di comunicare con le altre solo raramente o per nulla [N.d.T.].

Figura 6.16 Esempi di benchmark paralleli.

benchmark con un obiettivo simile a quello dei benchmark SPEC CPU. Questi benchmark contengono sia programmi kernel sia programmi applicativi, compresi molti programmi prodotti dalla comunità del calcolo ad alte prestazioni. Essi scalano in modo forte, anche se vengono forniti con due soli insiemi di dati.

- I *benchmark paralleli NAS* (*NASA Advanced Supercomputing*) costituiscono un altro tentativo degli anni '90 di costruire benchmark per i multiprocessori. Essi derivano dalla fluidodinamica computazionale e consistono di cinque kernel. Consentono una scalabilità debole, definendo un numero ridotto di insiemi di dati. Come Linpack, questi benchmark possono essere riscritti, ma le regole richiedono che il linguaggio di programmazione possa essere solamente C o Fortran.
- La recente suite di benchmark PARSEC (*Princeton Application Repository for Shared Memory Computers*) consiste in programmi multithread che utilizzano **Pthreads** (*thread POSIX*) e OpenMP (*Open MultiProcessing*). Essa è focalizzata sui mercati emergenti ed è costituita da nove applicazioni e tre kernel: otto sono basate sul parallelismo, tre sul parallelismo con pipeline e una sul parallelismo non strutturato.
- Sul fronte dei cloud, l'obiettivo del benchmark *Yahoo! Cloud Serving Benchmark* (YCSB) è comparare le prestazioni di un servizio dati erogato attraverso il cloud. Questo benchmark consente ai clienti di comparare i servizi, utilizzando i database Cassandra o HBase come esempi [Cooper, 2010].

Il rovescio della medaglia delle restrizioni imposte sui benchmark è che le innovazioni vengono principalmente limitate alle architetture e ai compilatori. Strutture dati, algoritmi e linguaggi di programmazione migliori rispetto ai precedenti spesso non vengono considerati, perché potrebbero portare a risultati fuorvianti. Un sistema potrebbe essere considerato il più veloce, per esempio, grazie all'algoritmo implementato e non all'hardware o al compilatore.

Avere delle linee guida rigide è comprensibile quando i calcolatori sono stabili, come lo erano negli anni '90 e nella prima metà di questa decade, ma diventano indesiderabili durante una rivoluzione della programmazione. Perché la rivoluzione abbia successo, occorre incentivare le innovazioni a tutti i livelli.

I ricercatori dell'Università californiana di Berkeley recentemente hanno proposto un nuovo approccio, identificando 13 moduli algoritmici che, a loro parere, faranno parte delle applicazioni future. Questi moduli sono implementati mediante framework e kernel e contengono, per esempio, matrici sparse, griglie strutturate, macchine a stati finiti, riduzione di mappe e attraversamento di grafi. Mantenendo la definizione dei problemi ad alto livello, i ricercatori contano di incentivare l'innovazione a tutti i livelli del sistema; perciò, il sistema che riesce a risolvere problemi con matrici sparse più velocemente può utilizzare una qualsiasi struttura dati, un qualsiasi algoritmo e un qualsiasi linguaggio di programmazione, oltre alle nuove architetture e ai nuovi compilatori.

## Modelli delle prestazioni

Un argomento collegato ai benchmark è quello dei modelli delle prestazioni. In questo capitolo abbiamo visto una crescente diversità delle architetture: multithread, SIMD, GPU; sarebbe perciò molto utile avere a disposizione un modello semplice che possa aiutare a valutare le prestazioni delle diverse architetture. Non è necessario che sia perfetto, è sufficiente che offra delle indicazioni.

Il modello delle 3C descritto nel Capitolo 5 è un esempio di modello delle prestazioni. Non è un modello perfetto, poiché non considera alcuni fatto-

**Pthreads:** un'API UNIX che consente di creare e di manipolare dei thread. È fornita sotto forma di libreria.

ri potenzialmente importanti, come la dimensione dei blocchi, la politica di allocazione dei blocchi e la politica di sostituzione dei blocchi; tuttavia, è un modello illuminante. Per esempio, consente di capire se una miss sia dovuta alla capacità di una particolare cache, o ai conflitti di accesso a un'altra cache delle stesse dimensioni. Il modello delle 3C viene ancora utilizzato dopo 25 anni proprio perché aiuta a capire il comportamento dei programmi, consentendo ai progettisti e ai programmatori di migliorare le loro creazioni sulla base delle indicazioni fornite da questo modello.

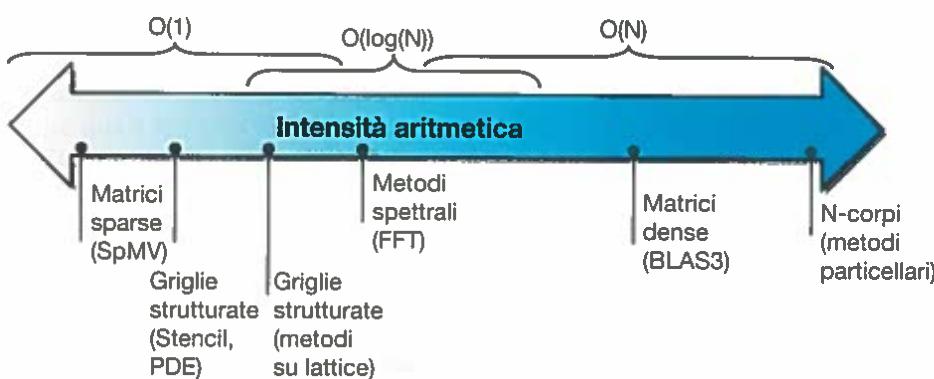
Per derivare un modello di questo tipo per le architetture parallele, iniziamo ad analizzare dei piccoli kernel come l'insieme dei 13 kernel algoritmici proposti dai ricercatori dell'Università di Berkeley e riportati nell'ultima riga della tabella di Figura 6.16. Anche se esistono versioni diverse di questi programmi kernel a seconda dei tipi di dati, i dati in virgola mobile sono i più diffusi in molte implementazioni, per cui le prestazioni di picco in virgola mobile vengono considerate il limite sulla velocità per questi kernel su un certo calcolatore. Per i chip multicore, le prestazioni di picco in virgola mobile sono una misura delle prestazioni massime dell'insieme di tutti i core sul chip. Se il sistema contiene più microprocessori, occorre moltiplicare le prestazioni di picco di un singolo chip per il numero totale di chip.

La richiesta di banda per il trasferimento da e verso la memoria può essere stimata dividendo le prestazioni in virgola mobile di picco per il numero medio di operazioni in virgola mobile che vengono effettuate per ciascun byte che viene letto:

$$\frac{\text{Operazioni in virgola mobile/s}}{\text{Operazioni in virgola mobile/byte}} = \text{Byte/s}$$

Il rapporto tra il numero di operazioni in virgola mobile e il numero di byte a cui si accede in memoria è chiamato **intensità aritmetica** e può essere calcolato dividendo il numero totale di operazioni in virgola mobile di un programma per il numero di byte dei dati che sono stati trasferiti da e verso la memoria principale durante l'esecuzione. La Figura 6.17 mostra l'intensità aritmetica di alcuni dei programmi kernel riportati nella tabella di Figura 6.16.

**Intensità aritmetica:** il rapporto tra il numero di operazioni in virgola mobile di un programma e il numero di byte di dati a cui il programma accede nella memoria principale.



**Figura 6.17** Intensità aritmetica, calcolata come numero di operazioni in virgola mobile eseguite da un programma diviso per il numero di byte a cui si accede in memoria principale [Williams, Waterman e Patterson, 2009]. Alcuni programmi kernel hanno un'intensità aritmetica che scala con la dimensione del problema, come il kernel Matrici dense, ma per molti kernel l'intensità aritmetica è indipendente dalle dimensioni del problema. I kernel del primo tipo scalano in modo debole; ciò può portare a risultati diversi, dato che richiedono meno accessi al sistema di memoria.

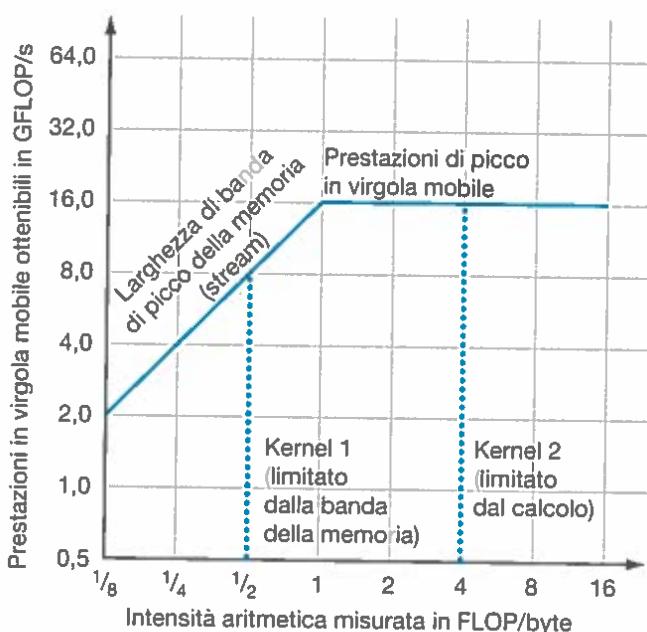
## Modello roofline

Questo semplice modello esprime le prestazioni in virgola mobile in funzione dell'intensità aritmetica e delle prestazioni della memoria mediante un grafico bidimensionale [Williams, Waterman e Patterson, 2009]. Le prestazioni in virgola mobile di picco possono essere determinate utilizzando le specifiche hardware riportate sul calcolatore. L'insieme di lavoro dei programmi kernel considerati qui non può essere interamente contenuto nelle cache, quindi occorre tenere conto delle prestazioni di picco del sistema di memoria che sta dietro le cache. Un modo per determinare le prestazioni di picco della memoria è il benchmark Stream (vedi la sezione *Approfondimento* nel paragrafo 5.2).

La Figura 6.18 illustra questo modello per un calcolatore; non verranno mostrati i risultati per i singoli kernel. Sull'asse  $y$  sono riportati i valori delle prestazioni in virgola mobile che si possono ottenere nell'intervallo 0,5-64,0 GFLOP/s. Sull'asse  $x$  è riportata l'intensità aritmetica espressa in FLOP/(byte DRAM indirizzato) nell'intervallo 1/8-16 FLOP/(byte DRAM indirizzato). Si noti che entrambi gli assi sono su scala logaritmica.

Per un dato kernel, possiamo identificare sull'asse  $x$  il punto che corrisponde alla sua intensità aritmetica. Tracciamo ora una retta verticale che passa per questo punto: le prestazioni del kernel su quel calcolatore dovranno giacere su questa retta. Possiamo quindi tracciare una linea orizzontale corrispondente alle prestazioni in virgola mobile di picco del calcolatore; ovviamente, le prestazioni in virgola mobile effettive non possono superare questa linea orizzontale, che costituisce un limite hardware.

Come possiamo identificare sul grafico le prestazioni di picco della memoria, che sono misurate in byte/secondo? Poiché l'asse  $x$  riporta i FLOP/byte e



**Figura 6.18** Modello roofline [Williams, Waterman e Patterson, 2009]. Il calcolatore di questo esempio ha prestazioni di picco in virgola mobile di 16 GFLOP/s e una larghezza di banda della memoria di 16 GB/s, misurata sul benchmark Stream. Dato che il benchmark Stream è costituito in realtà da quattro misure, questo valore è ottenuto come media delle quattro misure. La retta verticale punteggiata a sinistra rappresenta le prestazioni ottenute con il kernel 1, che ha un'intensità aritmetica di 0,5 FLOP/byte; questo kernel è limitato dalla larghezza di banda della memoria e non supera gli 8 GFLOP/s su questo Opteron X2. La retta verticale punteggiata a destra rappresenta il programma kernel 2, che ha un'intensità aritmetica di 4 FLOP/byte; questo kernel è limitato dalla capacità di calcolo, che è di 16 GFLOP/s. I dati riguardano un Opteron X2, versione F, di AMD, che utilizza processori a doppio core con frequenza di 2 GHz, in un sistema a doppio connettore.

l'asse  $y$  i FLOP/s, i byte/s saranno semplicemente dati da una retta con un'inclinazione di 45 gradi; possiamo quindi tracciare una terza linea che indica le prestazioni massime in virgola mobile che il sistema di memoria del calcolatore può supportare per ogni intensità aritmetica. Questo limite può essere espresso attraverso una formula che ci permette di tracciare la retta inclinata sul grafico di Figura 6.18:

$$\text{GFLOP/s ottenibili} = \text{Min} (\text{Banda di picco della memoria} \times \text{Intensità aritmetica}, \text{Picco nelle prestazioni in virgola mobile})$$

Le due rette (orizzontale e diagonale) danno il nome a questo semplice modello (*roofline*, linea del tetto) e forniscono i valori desiderati. La roofline fornisce il limite superiore delle prestazioni di un programma kernel e dipende dall'intensità aritmetica del programma stesso.

Se consideriamo l'intensità aritmetica come un'asta verticale che dal suolo arriva al tetto, questa potrà toccare il tetto in corrispondenza della parte piatta, e in questo caso le prestazioni sono limitate dal calcolo, oppure della parte inclinata, e in questo caso le prestazioni sono limitate dalla larghezza di banda della memoria. In Figura 6.18 il kernel 2 costituisce un esempio del primo caso, mentre il kernel 1 costituisce un esempio del secondo caso. Data la roofline di un calcolatore, si può applicare il modello ripetutamente, poiché questa linea non dipende dal kernel utilizzato per la valutazione.

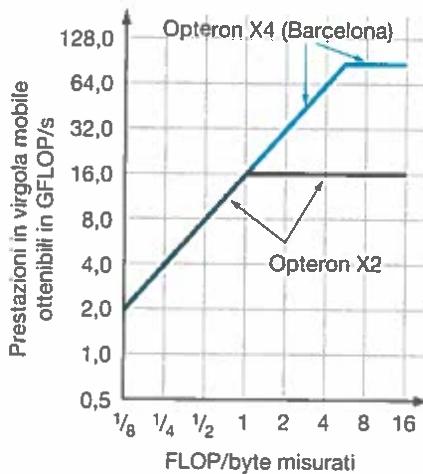
Si noti che il punto angoloso della roofline, in cui la retta orizzontale incontra quella inclinata, fornisce un'informazione interessante sul calcolatore. Se questo punto è molto spostato a destra, solamente i kernel con un'intensità aritmetica elevata potranno raggiungere le massime prestazioni; se, invece, il punto è molto spostato a sinistra, praticamente tutti i kernel potranno raggiungere le prestazioni massime. Nel prossimo paragrafo vedremo degli esempi di entrambi i casi.

## Confronto tra due generazioni di Opteron

L'Opteron X4 (Barcelona) di AMD con 4 core è il successore dell'Opteron X2 contenente due core. Per semplificare la progettazione della scheda, entrambi i processori utilizzano lo stesso connettore, e quindi hanno gli stessi canali di comunicazione con la DRAM e la stessa larghezza di banda di picco della memoria. Oltre ad aver raddoppiato il numero di core, l'Opteron X4 ha anche il doppio delle prestazioni in virgola mobile di picco per ciascun core: un core dell'Opteron X4 può inviare in esecuzione due istruzioni SSE2 in virgola mobile per ciclo di clock, mentre i core dell'Opteron X2 possono inviare in esecuzione al massimo un'istruzione. Dato che i due sistemi che vogliamo confrontare hanno quasi la stessa frequenza di clock (2,2 GHz per l'Opteron X2 e 2,3 GHz per l'Opteron X4), l'Opteron X4 ha prestazioni di picco in virgola mobile che sono più di quattro volte maggiori di quelle dell'Opteron X2, mentre ha la stessa larghezza di banda della DRAM. L'Opteron X4 contiene anche una cache L3 da 2 MiB che non è presente nell'Opteron X2.

La Figura 6.19 confronta il modello roofline per i due sistemi. Come ci si poteva aspettare, il punto angoloso si sposta dal valore 1 (per l'Opteron X2) al valore 5 (per l'Opteron X4); ciò significa che per ottenere un guadagno di prestazioni, i kernel devono avere un'intensità aritmetica superiore a 1, oppure l'insieme di lavoro dei benchmark deve essere contenuto interamente nelle cache dell'Opteron X4.

Il modello roofline fornisce il limite superiore delle prestazioni. Supponete che il vostro programma sia molto al di sotto di quel limite: quali ottimizzazioni si dovrebbero implementare e in quale ordine per aumentare le prestazioni?



**Figura 6.19 Comportamento del modello roofline per due generazioni di Opteron.** La curva del modello roofline per l'Opteron X2 (la stessa riportata in Figura 6.18) è disegnata in nero, mentre la curva relativa all'Opteron X4 è disegnata in blu. Il fatto che il punto angoloso si trovi più in alto e più a destra nell'Opteron X4 significa che i programmi kernel che erano limitati dal calcolo sull'Opteron X2 potrebbero essere limitati dalle prestazioni della memoria sull'Opteron X4.

Per ridurre i colli di bottiglia sul calcolo, le seguenti due ottimizzazioni sono utili per quasi tutti i programmi kernel.

1. *Mix di operazioni in virgola mobile.* Le prestazioni di picco in virgola mobile di un calcolatore sono ottenute di solito quando abbiamo un numero circa uguale di addizioni e moltiplicazioni che vengono eseguite contemporaneamente. Questo bilanciamento è necessario, perché il calcolatore supporta le istruzioni di moltiplicazione e di somma integrate in un'unica istruzione (*vedi* la sezione *Approfondimento* a pagina 190), oppure perché l'unità di calcolo in virgola mobile ha uno stesso numero di addizionatori e moltiplicatori in virgola mobile. Le prestazioni migliori richiedono anche che una parte significativa del mix di istruzioni sia composta da operazioni in virgola mobile e non su interi.
2. *Migliorare il parallelismo a livello di istruzioni e applicare il modello SIMD.* Nelle architetture superscalari, le prestazioni più elevate si ottengono quando è possibile caricare, eseguire e terminare da tre a quattro istruzioni per ciclo di clock (*vedi* par. 4.10). L'obiettivo, in questo caso, è fare in modo che il compilatore migliori il codice per aumentare l'ILP. Un modo per ottenere ciò è attraverso l'espansione dei cicli, come si è visto nel paragrafo 4.12. Inoltre, in un'architettura x86, una singola istruzione AVX può operare su quattro operandi in doppia precisione, per cui queste istruzioni dovrebbero essere utilizzate ogni volta che è possibile (*vedi* parr. 3.7 e 3.8).

Per ridurre i colli di bottiglia sulla memoria, invece, possono essere implementate le due seguenti ottimizzazioni.

1. *Precaricamento del software.* Di solito, per ottenere le prestazioni massime, è necessario avere subito a disposizione molte operazioni sulla memoria: questo si può ottenere più facilmente attraverso la **predizione** degli accessi e il precaricamento delle istruzioni, senza aspettare che i dati vengano richiesti dall'esecuzione.
2. *Affinità della memoria.* Attualmente la maggior parte dei microprocessori include un controllore della memoria sullo stesso chip su cui è montato il microprocessore, che migliora le prestazioni della **gerarchia delle memorie**.



PARALLELISMO



PREDIZIONE



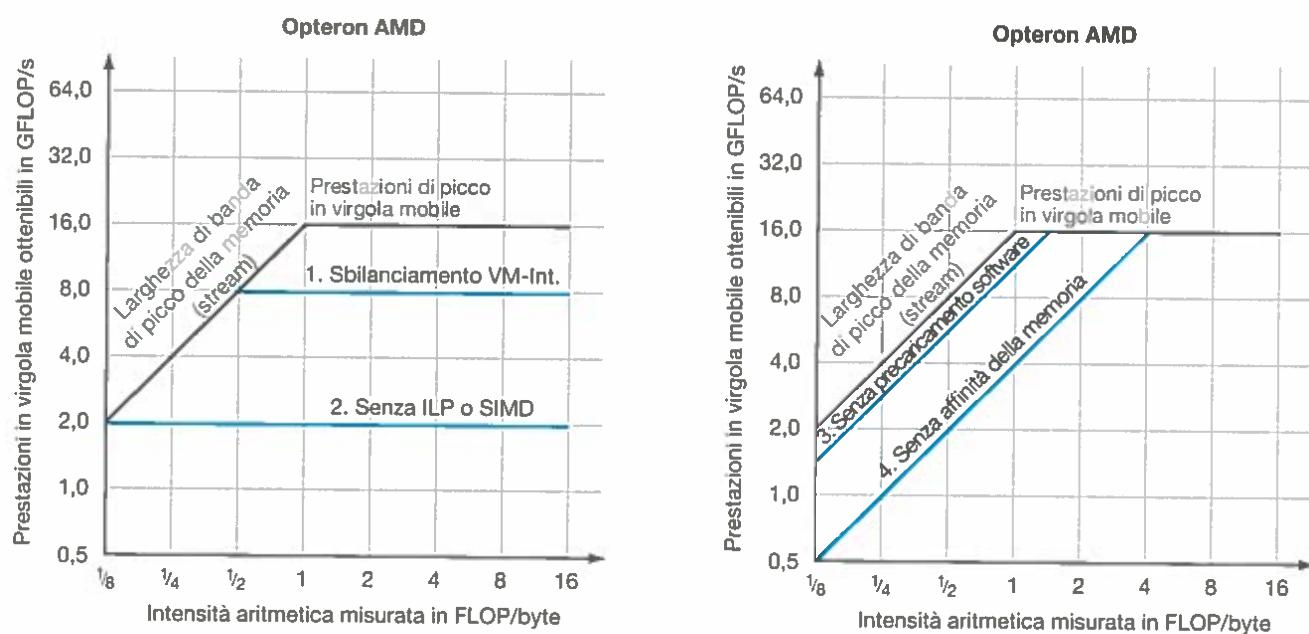
GERARCHIA

Se il sistema è costituito da chip multipli, può accadere che alcune richieste siano indirizzate alla DRAM associata al chip da cui proviene la richiesta, e che altre richieste siano inviate alla rete di interconnessione dei chip per prelevare i dati dalle DRAM associate agli altri chip. Questa suddivisione fa sì che l'accesso alla memoria non sia più uniforme, come descritto nel paragrafo 6.5. In questo secondo caso le prestazioni diminuiscono. Questa seconda ottimizzazione ha l'obiettivo di allocare i dati (e i thread relativi ai task che operano su questi dati) sulla stessa coppia processore-memoria, in modo tale che il processore debba accedere raramente alla memoria su altri chip.

Il modello roofline può aiutare a decidere quali di queste ottimizzazioni debbano essere implementate e in quale ordine. Possiamo considerare ciascuna delle problematiche affrontate da queste ottimizzazioni come un «soffitto» che si trova al di sotto della linea del tetto: si può forare il soffitto per arrivare a toccare il tetto solo implementando l'ottimizzazione associata.

La porzione della curva di roofline relativa al calcolo si può ricavare dal manuale del calcolatore, mentre la parte relativa alla memoria si può ricavare dall'esecuzione del benchmark Stream. I «soffitti» associati al calcolo, come quello relativo al bilanciamento delle operazioni in virgola mobile e intere, possono essere ricavati dal manuale del calcolatore. I «soffitti» sulla memoria, invece, implicano l'esecuzione di esperimenti sui diversi calcolatori per determinare le differenze tra i modelli. La buona notizia è che questi esperimenti possono essere effettuati una sola volta: quando si determinano questi «soffitti» per un calcolatore, il corrispondente modello roofline può essere utilizzato da chiunque per determinare le ottimizzazioni migliori per quel calcolatore.

La Figura 6.20 aggiunge i soffitti al modello roofline mostrato in Figura 6.18: nel grafico a sinistra sono evidenziati i soffitti associati al calcolo, in quello a destra i soffitti associati alla larghezza di banda della memoria. Anche se non abbiamo indicato esplicitamente che i soffitti più alti sono associati a entrambe



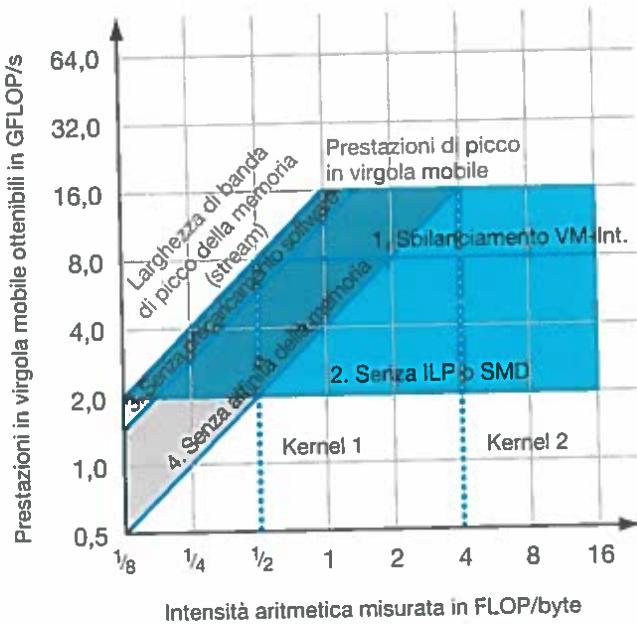
**Figura 6.20 Modello roofline con i soffitti.** Il grafico a sinistra mostra il "soffitto" sul calcolo a 8 GFLOP/s che si crea quando il mix delle operazioni in virgola mobile e intere è sbilanciato e il soffitto a 2 GFLOP/s che si crea quando sono mancati sia l'ottimizzazione dell'ILP sia dell'elaborazione SIMD. Il grafico a destra mostra il "soffitto" associato alla larghezza di banda della memoria a 11 GB/s senza precaricamento software, e quello a 4,8 GB/s se mancano anche le ottimizzazioni sull'affinità della memoria.

le ottimizzazioni, lo si evince dalle figure: per forare il soffitto superiore occorre aver forato prima i soffitti bassi.

La distanza tra due soffitti è il guadagno che si può ottenere implementando l'ottimizzazione associata. Perciò, la Figura 6.20 suggerisce che l'ottimizzazione 2, che migliora l'ILP, produca un miglioramento sensibile nel calcolo per l'Opteron, mentre l'ottimizzazione 4, che migliora l'affinità della memoria, fornisca un grande beneficio sulla larghezza di banda della memoria per quel calcolatore.

La Figura 6.21 raggruppa le soglie di Figura 6.20 in un unico grafico. L'intensità aritmetica del kernel determina la regione in cui viene implementata l'ottimizzazione, che, a sua volta, suggerisce quale altra ottimizzazione è opportuno implementare. Si noti che le ottimizzazioni sul calcolo e sulla larghezza di banda della memoria si sovrappongono per gran parte dei valori di intensità aritmetica. In Figura 6.21 vengono identificate tre regioni che si riferiscono alle diverse strategie di ottimizzazione. Per esempio, il kernel 2 cade nell'area del trapezoide blu a destra, che suggerisce di lavorare solo sulle ottimizzazioni del calcolo, mentre il kernel 1 cade nell'area del parallelogramma blu-grigio al centro, che suggerisce di implementare entrambi i tipi di ottimizzazioni; suggerisce anche di iniziare con le ottimizzazioni 2 e 4. Si noti che le rette verticali del kernel 1 cadono al di sotto delle ottimizzazioni per il bilanciamento delle operazioni, per cui l'ottimizzazione 1 può non essere necessaria. Se un kernel cadesse nel triangolo grigio all'estrema sinistra, sarebbe opportuno implementare solamente le ottimizzazioni sulla memoria.

Finora abbiamo supposto che l'intensità aritmetica fosse fissa, ma questo non è il caso reale. Innanzitutto ci sono alcuni kernel la cui intensità aritmetica aumenta con la dimensione del problema, come i kernel delle matrici dense e



**Figura 6.21** Modello roofline con i soffitti e con sovrapposte le aree tratteggiate e i due kernel di Figura 6.18. I kernel la cui intensità aritmetica si trova all'interno del trapezoide blu a destra devono puntare all'ottimizzazione computazionale, mentre i kernel la cui intensità aritmetica cade nel triangolo grigio in basso a sinistra devono focalizzarsi sulle ottimizzazioni della larghezza di banda della memoria. I kernel che cadono all'interno del parallelogramma blu-grigio al centro devono essere ottimizzati su entrambi i fronti. Dato che il kernel 1 cade nel parallelogramma al centro, occorre cercare di ottimizzare l'ILP e il livello di SIMD, l'affinità della memoria e il precaricamento software. Il kernel 2 cade nel trapezoide a destra, per cui occorre cercare di ottimizzare l'ILP e il livello di SIMD e il bilanciamento delle operazioni in virgola mobile e intere.

quello dei problemi a N corpi (Figura 6.17). Questo è uno dei motivi per cui i programmati hanno più successo quando il problema scala in modo debole anziché in modo forte. In secondo luogo, l'efficacia della gerarchia delle memorie influenza il numero di accessi che arrivano alla memoria, per cui le ottimizzazioni che migliorano le prestazioni delle cache aumentano anche l'intensità aritmetica. Un esempio di ciò è rappresentato dal miglioramento della località temporale, che si ottiene espandendo i cicli per raggruppare le istruzioni che accedono a indirizzi vicini. Molti calcolatori hanno istruzioni speciali di cache che allocano i dati in una cache, ma che non riempiono subito l'intero blocco della cache con i dati prelevati al relativo indirizzo, poiché le parole adiacenti potrebbero essere presto sovrascritte. Entrambe queste ottimizzazioni riducono il traffico con la memoria, spostando quindi la linea verticale dell'intensità aritmetica sulla destra di un fattore che può essere pari a 1,5. Questo spostamento verso destra può posizionare il programma kernel in una regione di ottimizzazione diversa.

Anche se gli esempi mostrati sopra sono stati riportati per illustrare ai programmati come migliorare le prestazioni, i progettisti hardware possono utilizzare questo modello per migliorare le prestazioni su kernel che considerano importanti.

Nel prossimo paragrafo il modello roofline verrà utilizzato per mostrare la differenza nelle prestazioni di un microprocessore multicore e di una GPU e per vedere se queste differenze si riflettono su programmi reali.

**Approfondimento.** I soffitti sono ordinati in modo tale che quelli più bassi siano più facili da "superare" attraverso le ottimizzazioni. Chiaramente, un programmatore può implementare le ottimizzazioni in un qualsiasi ordine. Seguendo l'ordine suggerito dal modello, però, si riduce la possibilità di fare sforzi inutili con un cambiamento che non produce miglioramenti a causa della presenza di altri vincoli. Come il modello delle 3C, anche il modello roofline può richiedere delle assunzioni che possono rivelarsi ottimistiche; per esempio, assume che il carico di lavoro del programma sia ripartito in modo bilanciato tra i processori. Queste ipotesi sono valide finché il modello fornisce indicazioni utili.

**Approfondimento.** Un'alternativa al benchmark Stream per stimare il flusso continuo di trasferimento massimo con la memoria è quella di utilizzare la larghezza di banda grezza della DRAM all'interno del modello roofline. La DRAM sicuramente impone un limite invalicabile alle prestazioni; tuttavia, le prestazioni effettive della memoria sono spesso così lontane da questo limite che non è utile considerarlo come se fosse un limite superiore. In altri termini, nessun programma può arrivare vicino a questo limite. L'inconveniente dell'uso del benchmark Stream è che una programmazione particolarmente attenta può superare i risultati di questo benchmark, per cui la parte della curva di roofline relativa alla memoria può risultare un limite non così invalicabile come quella relativa al calcolo. Abbiamo deciso di utilizzare il benchmark Stream perché pochi programmati sono in grado di ottenere una larghezza di banda della memoria maggiore di quella fornita da Stream.

**Approfondimento.** Il modello roofline è stato introdotto per i processori multicore, ma potrebbe essere applicato anche alle architetture monoprocessoress.



## Autovalutazione

Vero o falso? Il problema principale degli approcci convenzionali ai benchmark dei calcolatori paralleli è che le regole che garantiscono la *fairness* (equità) rallentano l'innovazione nel software.

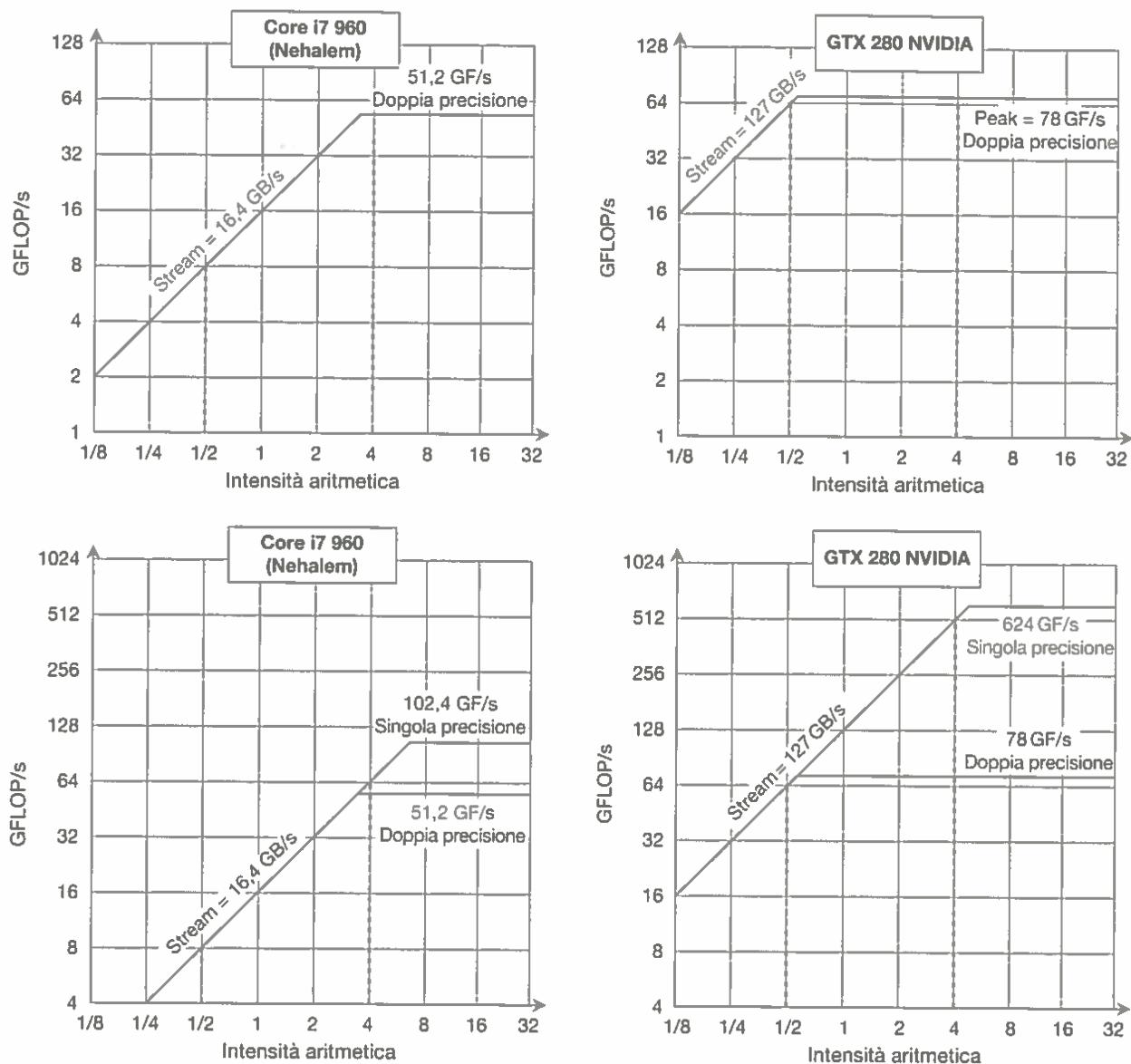
## 6.11 Un caso reale: il confronto mediante il modello roofline tra un Core i7 960 Intel e una GPU Tesla NVIDIA

Un gruppo di ricercatori Intel ha pubblicato un articolo [Lee *et al.*, 2010] che compara un Core i7 960 Intel quad-core contenente le estensioni multimediali SIMD con una GPU della generazione precedente: una Tesla GTX 280 di NVIDIA. In Figura 6.22 sono riportate le caratteristiche delle due architetture. Entrambi i prodotti sono stati acquistati nell'autunno 2009. Il Core i7 è basato sulla tecnologia Intel dei semiconduttori da 45 nanometri, mentre la GPU è basata sulla tecnologia TSMC da 65 nanometri. Anche se sarebbe più equa una comparazione effettuata da una terza parte neutrale o da entrambe le parti, lo scopo di questo paragrafo non è determinare quale e quanto uno dei due prodotti sia più veloce dell'altro, ma cercare di capire le potenzialità delle caratteristiche delle due architetture mettendole a confronto.

La roofline del Core i7 960 e del GTX 280 mostrata in Figura 6.23 illustra la differenza tra i due calcolatori. Non solo la GTX 280 ha una banda di trasferimento dalla memoria molto più elevata e prestazioni più elevate in virgola mobile in doppia precisione, ma ha anche un punto angoloso molto più

	Core i7 960	GTX 280	GTX 480	Rapporto 280/i7	Rapporto 480/i7
Numero di elementi di elaborazione (core o SM)	4	30	15	7,5	3,8
Frequenza di clock (GHz)	3,2	1,3	1,4	0,41	0,44
Dimensione del chip	263	576	520	2,2	2,0
Tecnologia	Intel 45 nm	TSMC 65 nm	TSMC 40 nm	1,6	1,0
Potenza (chip, non modulo)	130	130	167	1,0	1,3
Numero transistor	700 M	1400 M	3030 M	2,0	4,4
Larghezza di banda della memoria (GByte/s)	32	141	177	4,4	5,5
Larghezza SIMD in singola precisione	4	8	32	2,0	8,0
Larghezza SIMD in doppia precisione	2	1	16	0,5	8,0
FLOP di picco scalari in singola precisione (GFLOP/s)	26	117	63	4,6	2,5
FLOP di picco SIMD in singola precisione (GFLOP/s)	102	da 311 a 933	da 515 a 1344	3,0-9,1	6,6-13,1
(SP1 somma o moltiplicazione)	N.A.	(311)	(515)	(3,0)	(6,6)
(SP1 istruzioni moltiplicazione e addizione fuse)	N.A.	(622)	(1344)	(6,1)	(13,1)
(SP rare, pacchetti di esecuzione doppi contenenti moltiplicazioni e addizioni fuse)	N.A.	(933)	N.A.	(9,1)	-
FLOP di picco SIMD doppia precisione (GFLOP/s)	51	78	515	1,5	10,1

**Figura 6.22** Specifiche del Core i7 960 di Intel, della GTX 280 e della GTX 480 NVIDIA. Le colonne più a destra riportano il rapporto tra le prestazioni di una Tesla GTX 280 e di una Fermi GTX 480 rispetto a un Core i7. Sebbene l'obiettivo di questo lavoro sia un'analisi comparata di una Tesla 280 e di un i7, abbiamo incluso anche i valori relativi alle Fermi 480 per mostrare l'incremento di prestazioni rispetto alle Tesla 280, dato che vengono descritte in questo capitolo. Si noti che la larghezza di banda della memoria è maggiore di quella riportata in Figura 6.23, perché viene misurata sui pin della DRAM, mentre per i dati in Figura 6.23 è misurata dal processore mediante un programma di benchmark. (Dati ricavati dalla Tabella 2 di Lee *et al.*, 2010.)



**Figura 6.23 Modello roofline [Williams, Waterman e Patterson 2009].** Le linee di roofline mostrano le prestazioni in virgola mobile in doppia precisione (riga superiore) e in virgola mobile in singola precisione (riga inferiore). Per chiarezza, nella riga inferiore vengono riportate anche le prestazioni di picco in doppia precisione. Il Core i7 960 (colonna a sinistra) ha prestazioni di picco in doppia precisione di 51,2 GFLOP/s, prestazioni di picco in singola precisione di 102,4 GFLOP/s e una larghezza di banda di picco di trasferimento dalla memoria di 16,4 GByte/s. La GTX 280 NVIDIA ha prestazioni di picco in doppia precisione di 78 GFLOP/s, prestazioni di picco in singola precisione di 624 GFLOP/s e una larghezza di banda di picco di trasferimento dalla memoria di 127 GByte/s. La linea verticale tratteggiata, sulla sinistra, corrisponde a un'intensità aritmetica di 0,5 FLOP/byte, le corrispondenti prestazioni di picco sono limitate dalla memoria a 8 GFLOP/s sia in doppia sia in singola precisione sul Core i7. La linea verticale tratteggiata, sulla destra, corrisponde a un'intensità aritmetica di 4 FLOP/byte, le corrispondenti prestazioni di picco sono limitate dal calcolo a 51,2 GFLOP/s in doppia precisione e a 102,4 GFLOP/s in singola precisione sul Core i7, e a 78 GFLOP/s in singola precisione e 624 GFLOP/s in doppia precisione sulla GTX 280. Per raggiungere le prestazioni massime sul Core i7, occorre utilizzare tutti e quattro i core e istruzioni SSE che abbiano un ugual numero di moltiplicazioni e addizioni. Per quanto riguarda la GTX 280, occorre utilizzare le istruzioni di moltiplicazione e addizione fuse assieme su tutti processori SIMD multithread.

spostato a sinistra: si trova a 0,6 per la GTX, mentre per il Core i7 si trova a 3,1. Come accennato in precedenza, è più facile raggiungere le prestazioni di calcolo di picco quando il punto angoloso della roofline si trova più a sinistra. Per quanto riguarda le prestazioni in singola precisione, il punto angoloso si sposta molto a destra su entrambi i calcolatori, per cui diventa più difficile raggiungere le prestazioni di picco in singola precisione. Si noti che l'intensità aritmetica dei programmi kernel è basata sul conteggio dei byte che vengono

trasferiti nella memoria principale e non nelle memorie cache. Perciò, come accennato in precedenza, la memoria cache può modificare l'intensità aritmetica di un kernel su un certo calcolatore, se la maggior parte dei dati può essere contenuta nella memoria cache. Le situazioni in cui gli accessi alla memoria sono sparsi producono prestazioni inferiori sia sulle GTX 280 sia sui Core i7, come vedremo fra breve.

I ricercatori hanno selezionato i programmi di benchmark analizzando le caratteristiche di calcolo e di memoria di quattro recenti suite di benchmark per “formulare l'insieme di *kernel* per la stima del throughput che catturino queste caratteristiche”. La Figura 6.24 mostra le prestazioni: numeri maggiori sono associati a prestazioni migliori. Il modello roofline aiuta a spiegare la differenza nelle prestazioni per questo caso di studio.

Dato che le specifiche associate alle prestazioni della GTX 280 sono superiori da 2,5 volte (frequenza di clock) a 7,5 volte (numero di core per chip) al Core i7, mentre le prestazioni variano da 2,0 volte inferiori (Solv) a 15,2 superiori (GJK), i ricercatori Intel hanno deciso di indagare la ragione di queste differenze.

- *Larghezza di banda della memoria.* La GPU ha una larghezza di banda della memoria di 4,4x, che contribuisce a spiegare come mai LBM e SAXPY vengono eseguiti rispettivamente 5,0 e 5,3 volte più velocemente. Il loro insieme di lavoro è di centinaia di megabyte per cui non può essere contenuto interamente nella cache del Core i7. Inoltre, dato che il loro scopo è accedere intensivamente alla memoria, questi kernel non fanno uso dell'accesso a blocchi alla cache (*vedi* Cap. 5). Perciò, la pendenza della roofline spiega le prestazioni ottenute. Anche SpVM ha un grosso insieme di lavoro, ma viene eseguito solo 1,9 volte più velocemente sulla GPU, poiché il calcolo in virgola mobile sulla GTX 280 è solamente di 1,5 volte più veloce che sul Core i7.
- *Larghezza di banda del calcolo.* Cinque degli altri kernel sono limitati dal calcolo: SGEMM, Conv, FFT, MC e Bilat, sui quali la GTX risulta più veloce rispettivamente di 3,9, 2,8, 3,0, 1,8 e 5,7 volte. I primi tre kernel utilizzano

Kernel	Unità di misura	Core i7 960	GTX 280	GTX 280/i7 960
SGEMM	GFLOP/s	94	364	3,9
MC	Miliardi cammini/s	0,8	1,4	1,8
Conv	Milioni pixel/s	1,25	3,5	2,8
FFT	GFLOP/s	71,4	213	3,0
SAXPY	GByte/s	16,8	88,8	5,3
LBM	Milioni accessi/s	85	426	5,0
Solv	Frame/s	103	52	0,5
SpMV	GFLOP/s	4,9	9,1	1,9
GJK	Frame/s	67	1020	15,2
Sort	Milioni elementi/s	250	198	0,8
RC	Frame/s	5	8,1	1,6
Search	Milioni ricerche/s	50	90	1,8
Hist	Milioni pixel/s	1517	2583	1,7
Bilat	Milioni pixel/s	83	475	5,7

**Figura 6.24 Prestazioni grezze e relative misurate sulle due piattaforme.** In questo studio, SAXPY viene utilizzata esclusivamente per misurare la larghezza di banda della memoria, per cui viene misurata in GByte/s e non in GFLOPS (tabella basata sulla Tabella 3 di Lee *et al.*, 2010).

l'aritmetica in virgola mobile in singola precisione, che sulla GTX è più veloce di un fattore da 3 a 6; MC utilizza la doppia precisione, e ciò spiega come mai risulti più veloce solamente di un fattore 1,8, dato che le prestazioni in doppia precisione su GTX sono migliori solo di un fattore 1,5. Bilat utilizza le funzioni trascendenti, che vengono supportate direttamente sulla GTX, mentre il Core i7 spende due terzi del tempo nel calcolo di queste funzioni; il risultato è che l'esecuzione su GTX risulta di 5,7 volte più veloce. Questa osservazione ci aiuta a capire il valore del supporto hardware alle operazioni contenute nel vostro codice: in questo caso operazioni in doppia precisione e calcolo delle funzioni trascendenti.

- **Benefici della cache.** Il *ray casting* (RC) risulta più veloce di solo 1,6 volte sulla GTX perché l'utilizzo della cache a blocchi del Core i7 evita che questo kernel venga limitato dalla memoria (*vedi* parr. 5.4 e 5.14), come succede sulla GPU. La gestione a blocchi della cache è utile anche per il kernel *Search*: se gli alberi degli indici sono piccoli e possono essere contenuti nella cache, il Core i7 risulta 2 volte più veloce, altrimenti le prestazioni di questo kernel sono limitate dalla larghezza di banda della memoria. Complessivamente, questo kernel viene eseguito sulla GTX 280 1,8 volte più velocemente. La gestione a blocchi delle cache è utile anche per il kernel *Sort*: anche se molti programmatori non eseguiranno ordinamenti sui processori SIMD, questi potrebbero essere scritti mediante una primitiva di ordinamento a 1 bit, chiamata *split*. Tuttavia l'algoritmo *split* esegue molte più istruzioni della sua versione scalare e come risultato questo kernel viene eseguito su un Core i7 1,25 volte più velocemente che su GTX 280. Si noti che la cache è utile anche per gli altri kernel sul Core i7, dato che la gestione a blocchi della cache fa sì che i kernel SGEMM, FFT e SpVM siano limitati dal calcolo. Questa osservazione sottolinea ancora una volta l'importanza dell'ottimizzazione che mira a utilizzare la cache a blocchi (*vedi* Cap. 5).
- **Gather-scatter.** Le estensioni multimediali SIMD non sono di grande aiuto quando i dati sono sparsi nella memoria principale: le prestazioni ottimali si ottengono solamente quando i dati sono allineati di 16 byte in 16 byte. Perciò il kernel *GJK* non riceve grandi vantaggi dalle istruzioni SIMD del Core i7. Come si è detto in precedenza, le GPU offrono la modalità di indirizzamento gather-scatter come tutte le architetture vettoriali, che non viene invece supportata dalla maggior parte delle estensioni SIMD; inoltre, il controllore della memoria di una GPU è anche in grado di raggruppare gli accessi alla stessa pagina di DRAM (*vedi* par. 5.2). Tutto ciò consente a una GTX 280 di eseguire il kernel *GJK* a una velocità di ben 15,2 volte superiore al Core i7; questa differenza di velocità è la massima riscontrata tra i kernel di Figura 6.22. Questa osservazione rinforza l'importanza della funzionalità gather-scatter di accesso alla memoria contenuta nelle GPU e mancante nelle estensioni SIMD.
- **Sincronizzazione.** Le prestazioni della sincronizzazione sono limitate dagli aggiornamenti atomici, che sono responsabili del 28% del tempo totale di esecuzione su un Core i7, nonostante sia disponibile un'istruzione di fetch e incremento dell'indirizzo. Di conseguenza, il kernel *Hist* viene eseguito solo 1,7 volte più velocemente su una GTX 280. Il kernel *Solv* risolve un insieme di vincoli indipendenti con un gruppo limitato di istruzioni, seguite da una barriera di sincronizzazione. Il Core i7 beneficia delle istruzioni atomiche e di un sistema di controllo della consistenza della memoria che consente di ottenere il risultato corretto anche se non tutti gli accessi alla gerarchia delle memorie sono stati completati. Non possedendo questo sistema di consistenza, la GTX 280 deve lanciare alcuni gruppi di istruzioni dal processore

di sistema, e ciò porta la GTX 280 ad avere una velocità di esecuzione che è la metà del Core i7. Questa osservazione sottolinea come le prestazioni della sincronizzazione possano rivelarsi particolarmente importanti in alcuni problemi su dati paralleli.

È interessante notare come le debolezze dell'architettura Tesla GTX 280, identificate dai kernel selezionati dai ricercatori Intel, siano state risolte già nell'architettura NVIDIA successiva: l'architettura Fermi ha incrementato le prestazioni di calcolo in virgola mobile in doppia precisione, ha operazioni atomiche più veloci e ha introdotto le memorie cache. È anche interessante notare come la funzionalità di accesso alla memoria gather-scatter, supportata dalle architetture vettoriali decenni prima dell'avvento delle istruzioni SIMD, sarebbe molto utile per l'efficacia delle istruzioni SIMD, cosa che era stata pre detta prima che questa comparazione fosse disponibile: i ricercatori Intel hanno notato che in 6 dei 14 kernel le istruzioni SIMD del Core i7 sfruttavano meglio il parallelismo sui dati se avessero a disposizione la funzionalità gather-scatter di accesso alla memoria. Questo studio sicuramente sancisce anche l'importanza della gestione a blocchi della memoria cache.

Ora che abbiamo visto diversi risultati di benchmark su diversi multiprocessori, ritorniamo alla nostra funzione DGEMM per vedere nel dettaglio come possiamo modificare il codice C per sfruttare appieno le architetture multiprocessori.

## 6.12 | Come andare più veloce: processori multipli e moltiplicazione di matrici

Questo paragrafo contiene l'ultimo passo, il più importante, del nostro viaggio nell'incremento progressivo delle prestazioni adattando la funzione DGEMM all'hardware sottostante del Core i7 Intel (Sandy Bridge). Ciascun Core i7 contiene 8 core e il calcolatore che utilizziamo contiene 2 microprocessori Core i7; abbiamo perciò a disposizione 16 core sui quali eseguire DGEMM.

La Figura 6.25 mostra la versione OpenMP di DGEMM configurata per quest'architettura: si noti che la linea 30 è l'unica linea aggiunta rispetto alla Figura 5.48 per fare sì che questo codice venga eseguito su un multiprocessore: un comando pragma OpenMP che istruisce il calcolatore a spargere su tutti i thread il lavoro richiesto dal ciclo `for` più esterno.

La Figura 6.26 mostra un classico grafico dello speed-up di un multiprocessore; viene mostrato l'incremento delle prestazioni rispetto all'utilizzo di un singolo thread in funzione del numero di thread. Questo grafico rende evidenti le difficoltà che occorre affrontare per ottenere che un problema scal si in modo forte, rispetto alla scalatura debole. Quando si riesce a mettere tutto nella cache di primo livello, come nel caso delle matrici  $32 \times 32$ , aumentare il numero di thread in realtà diminuisce le prestazioni: la versione a 16 thread di DGEMM risulta veloce la metà rispetto a quella a singolo thread in questo caso. Viceversa, la versione a 16 thread di DGEMM ottiene uno speed-up vicino a  $14\times$  per le due matrici di dimensioni maggiori, corrispondenti alle due linee rette inclinate di Figura 6.26.

La Figura 6.27 mostra l'incremento assoluto delle prestazioni all'aumentare del numero di thread da 1 a 16. DGEMM dopo tutte le ottimizzazioni viene eseguito a 174 GFLOPS su matrici  $960 \times 960$ . In confronto, la versione in C non ottimizzata della stessa funzione, riportata in Figura 3.22, era in grado di eseguire il codice a soli 0,8 GFLOPS. Le ottimizzazioni apportate via via nei Capitoli dal 3 al 6, miranti ad adattare il codice all'hardware sottostante, hanno prodotto uno speed-up di più di 200 volte!

```

1 #include <x86intrin.h>
2 #define ESPANDI (4)
3 #define DIM_BLOCCO 32
4 void esegui_blocco (int n, int si, int sj, int sk,
5                      double *A, double *B, double *C)
6 {
7     for (int i = si; i < si+DIM_BLOCCO; i+=ESPANDI*4)
8         for (int j = sj; j < sj+DIM_BLOCCO; j++) {
9             __m256d c[4];
10            for (int x = 0; x < ESPANDI; x++)
11                c[x] = _mm256_load_pd(C+i+x*4+j*n);
12            /* c[x] = C[i][j] */
13            for(int k = sk; k < sk+DIM_BLOCCO; k++)
14            {
15                __m256d b = _mm256_broadcast_sd(B+k+j*n);
16                /* b = B[k][j] */
17                for (int x = 0; x < UNROLL; x++)
18                    c[x] = _mm256_add_pd(c[x], /* c[x]+=A[i][k]*b */
19                                           _mm256_mul_pd(_mm256_load_pd(A+n*k+x*4+i), b));
20            }
21
22            for (int x = 0; x < UNROLL; x++)
23                _mm256_store_pd(C+i+x*4+j*n, c[x]);
24            /* C[i][j] = c[x] */
25        }
26    }
27
28 void dgemm (int n, double* A, double* B, double* C)
29 {
30 #pragma omp parallel for
31     for (int sj = 0; sj < n; sj += DIM_BLOCCO)
32         for (int si = 0; si < n; si += DIM_BLOCCO)
33             for (int sk = 0; sk < n; sk += DIM_BLOCCO)
34                 esegui_blocco(n, si, sj, sk, A, B, C);
35 }

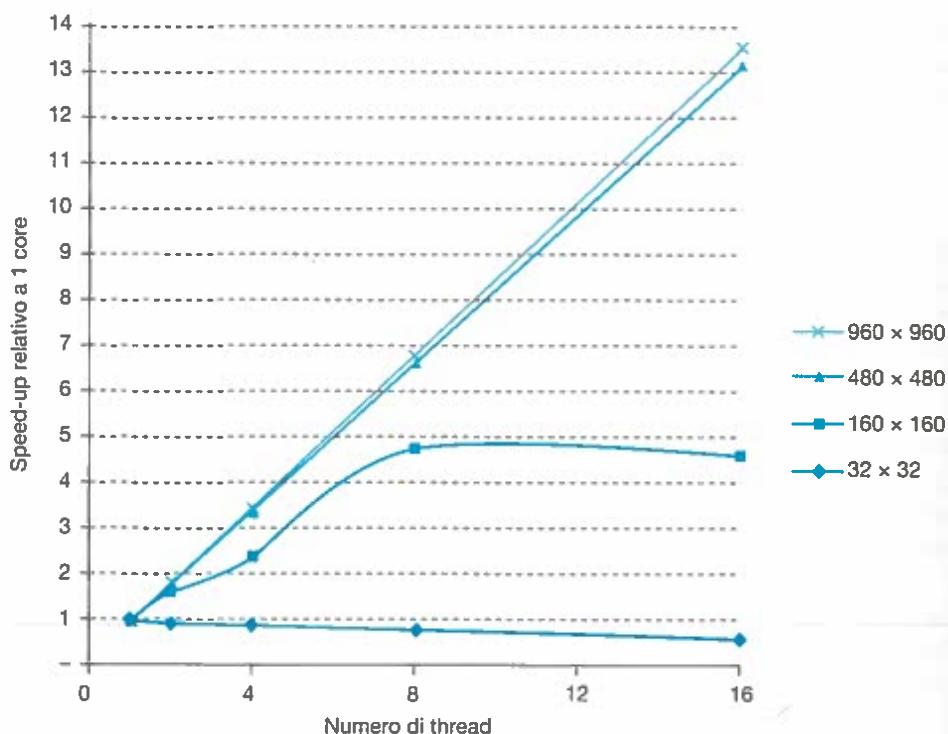
```

**Figura 6.25** Versione OpenMP di DGEMM di Figura 5.48. La linea 30 è l'unica linea di codice OpenMP e fa sì che il ciclo for più esterno venga eseguito in parallelo. Questa linea di codice è l'unica differenza rispetto alla Figura 5.48.

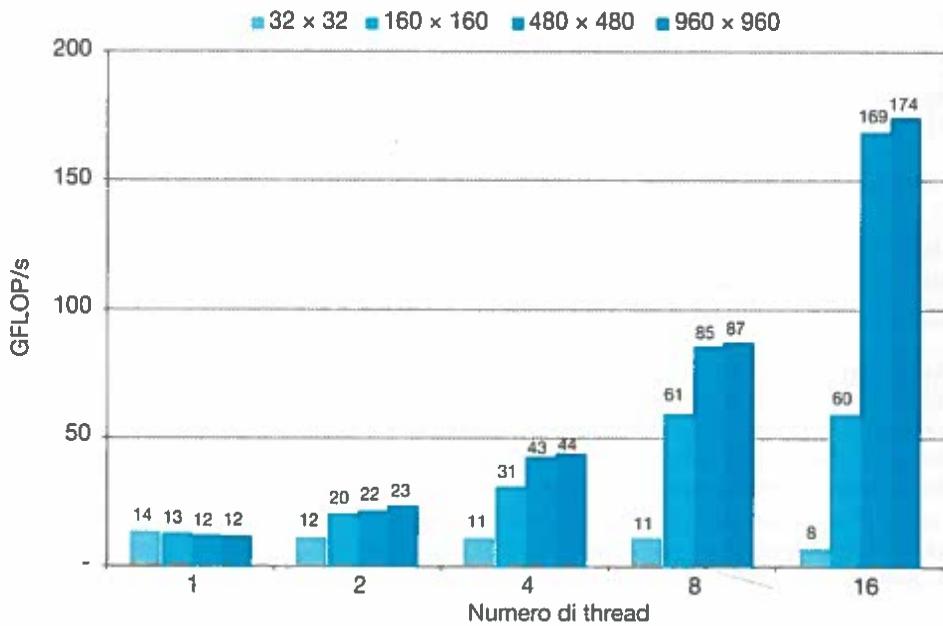
Il prossimo paragrafo illustrerà gli errori e i trabocchetti associati all'elaborazione sui multiprocessori. Il campo delle architetture dei calcolatori è pieno di progetti di calcolo parallelo che li hanno ignorati.

**Approfondimento.** Questi risultati sono stati ottenuti con la modalità Turbo disattivata. Utilizzando un chip contenente due core, in modo non sorprendente, possiamo ottenere uno speed-up di 1,27 (cioè 3,3/2,6) sia quando utilizziamo 1 thread (un core solo dei due) oppure 2 thread (1 core per ogni chip). All'aumentare del numero di thread e quindi del numero di core attivi, i benefici della modalità Turbo diminuiscono, dato che diminuisce la potenza aggiuntiva che può essere utilizzata da ogni singolo thread. Per 4 thread, lo speed-up medio ottenuto con la modalità Turbo è di 1,23, per 8 thread è di 1,13 e per 16 thread è di 1,11.

**Approfondimento.** Sebbene il Sandy Bridge supporti due thread hardware per ogni core, non riusciamo a ottenere un incremento di prestazioni utilizzando 32 thread. Il motivo è che un singolo componente hardware AVX viene condiviso da due thread multiplexati sul singolo core; perciò assegnare due thread a ogni core, di fatto, riduce le prestazioni perché introduce il sovraccarico del multiplexer.



**Figura 6.26** Aumento delle prestazioni relative a un singolo thread all'aumentare del numero di thread. Il modo più corretto per presentare questi grafici è riportare le prestazioni relative alla versione migliore del programma per un processore singolo. Questo grafico è relativo alle prestazioni del codice di Figura 5.48 senza includere i comandi pragma di OpenMP.



**Figura 6.27** Prestazioni di DGEMM all'aumentare del numero di thread per matrici di dimensioni diverse. L'aumento delle prestazioni rispetto al codice non ottimizzato di Figura 3.22 per matrici 960 × 960 quando vengono utilizzati 16 thread è di ben 212 volte!

## 6.13 | Errori e trabocchetti

L'enorme interesse destato dall'elaborazione parallela ha fatto nascere molti errori e trabocchetti. Di seguito ne elenchiamo quattro.

**Errore:** la Legge di Amdahl non si applica ai calcolatori paralleli.

Nel 1987 il responsabile di un'organizzazione di ricerca sostenne che la Legge di Amdahl era stata contraddetta da un sistema multiprocessore. Per cercare di capire perché i mezzi di comunicazione riportarono questa notizia, consideriamo la formulazione della Legge di Amdahl [1967]:

*Una conclusione abbastanza ovvia che può essere tratta a questo punto è che lo sforzo necessario per raggiungere elevate velocità di elaborazione parallela è sprecato, a meno che non sia accompagnato da miglioramenti di un ordine di grandezza molto simile nella velocità di elaborazione sequenziale.*

Quest'affermazione è ancora valida: la parte di un programma non interessata dall'architettura parallela limiterà comunque le prestazioni. Un'interpretazione di questa legge conduce al seguente lemma: esistono diverse porzioni di un programma che vengono comunque eseguite in modo sequenziale, cosicché deve esserci un limite superiore di natura economica al numero di processori, per esempio 100. Se si osservasse un incremento di velocità lineare con 1000 processori, si dimostrerebbe che questo lemma è falso e quindi si contraddirrebbe la Legge di Amdahl.

L'approccio seguito dai ricercatori era basato sulla scalabilità debole: anziché cercare di andare 1000 volte più veloce sullo stesso insieme di dati, essi cercarono di eseguire una quantità di lavoro 1000 volte superiore in un tempo confrontabile. Per il loro algoritmo, la parte sequenziale del programma era costante e indipendente dalla dimensione dei dati in ingresso, mentre la parte restante era completamente parallela e, di conseguenza, riuscirono a ottenere un incremento di velocità lineare con 1000 processori, ma non confutarono per questo la Legge di Amdahl.

Naturalmente la Legge di Amdahl si applica anche ai processori paralleli. Ciò che questa ricerca ha sottolineato è che l'utilizzo principale di calcolatori più veloci è quello di lavorare su problemi di dimensioni maggiori. Occorre però fare attenzione che i problemi considerati siano veramente importanti e non siano soltanto dei problemi che impegnano molti processori utilizzati per giustificare l'acquisto di calcolatori più costosi.

**Errore:** le prestazioni di picco ricalcano il comportamento delle prestazioni osservate.

Le aziende che producono supercomputer hanno utilizzato questa metrica per pubblicizzare i loro prodotti, e questo errore si è rivelato ancora più grave per le macchine parallele. Non solo i vendori utilizzano le prestazioni di picco (praticamente mai raggiungibili) dei singoli nodi monoprocessoressi, ma moltiplicano queste prestazioni per il numero totale di processori, assumendo un incremento di velocità perfettamente lineare con il numero dei processori! La Legge di Amdahl suggerisce come sia difficile raggiungere qualsiasi prestazione di picco. Il modello roofline aiuta a considerare le prestazioni di picco nella giusta prospettiva.

**Trabocchetto:** non sviluppare il software per sfruttare l'architettura multiprocessore o non ottimizzarlo per questo tipo di architettura.

C'è una lunga storia di applicazioni software che rimangono indietro rispetto ai processori paralleli, forse perché i problemi software sono più difficili da

*Per oltre un decennio i soliti profeti hanno continuato a sostenere che il modello strutturale basato su un singolo calcolatore aveva raggiunto i propri limiti e che qualunque ulteriore progresso poteva essere realizzato solo interconnettendo molti calcolatori che lavorassero in maniera coordinata per risolvere i problemi ... È stato dimostrato che l'approccio basato su un solo processore continua a essere valido...*

Gene Amdahl, "Validity of the single processor approach to achieving large scale computing capabilities", Spring Joint Computer Conference, 1967

risolvere. Forniamo ora un esempio di come siano sottili questi problemi (ma gli esempi potrebbero essere molti altri).

Un problema frequente si verifica quando il software progettato per un monoprocessore deve essere adattato a un ambiente multiprocessore. Per esempio, il sistema operativo SGI<sup>1</sup> inizialmente proteggeva la tabella delle pagine con un unico lucchetto, ipotizzando che l'allocazione delle pagine fosse un'operazione poco frequente. Nei monoprocessori questo non rappresenta un problema per le prestazioni, mentre nei multiprocessori può diventare uno dei colli di bottiglia principali per alcuni programmi. Si consideri un programma che utilizza un grande numero di pagine inizializzate all'avvio (come fa UNIX per le molte pagine allocate staticamente). Si supponga ora che il programma venga parallelizzato, in modo che più processi debbano allocare le pagine. Poiché l'allocazione delle pagine richiede l'utilizzo della tabella delle pagine, che è bloccata quando è utilizzata da un processo, anche le procedure del kernel di un SO che utilizza thread multipli per le procedure del SO stesso verrebbero eseguite in modo sequenziale se i diversi processi cercassero di allocare le loro pagine contemporaneamente, cosa che ci si aspetta che accada proprio in fase di inizializzazione.

Questa serializzazione forzata degli accessi alla tabella delle pagine elimina il parallelismo nell'inizializzazione e ha un impatto rilevante sulle prestazioni complessive. Questo collo di bottiglia riguarda anche il parallelismo a livello di task. Per esempio, si supponga di suddividere un programma a esecuzione parallela in diversi task separati e di eseguire ciascun task su un processore diverso; in questo modo non c'è condivisione tra i vari task. Questo è esattamente ciò che ha fatto un utente, il quale riteneva che il problema delle prestazioni fosse dovuto a una condivisione non voluta dei task, oppure a interferenze sulla memoria. Invece, era il blocco della tabella delle pagine a costringere in sequenza i diversi task, quindi, anche se i task erano indipendenti, le prestazioni risultavano scarse.

Questo trabocchetto illustra quanto possano essere sottili i problemi di prestazioni che si possono verificare nell'eseguire i programmi sui multiprocessori. Come molti altri componenti software cruciali, gli algoritmi e le strutture dati del SO devono essere ripensati per il contesto dei multiprocessori. Poder bloccare sottoinsiemi più piccoli della tabella delle pagine, per esempio, consentirebbe di eliminare questo problema.

*Errore: si possono ottenere buone prestazioni con il calcolo vettoriale senza bisogno di una grande larghezza di banda della memoria.*

Come abbiamo visto con il modello roofline, la larghezza di banda della memoria è molto importante in tutte le architetture. DAXPY richiede 1,5 accessi alla memoria per ogni operazione in virgola mobile, e questo rapporto è tipico di molti programmi scientifici. Anche se il costo delle operazioni in virgola mobile è trascurabile su un Cray-1, le prestazioni di DAXPY non potrebbero essere aumentate sui vettori e matrici delle dimensioni considerate, perché il limite al Cray-1 proviene dalla memoria. Le prestazioni del Cray-1 sul benchmark Linpack aumentano significativamente quando il compilatore utilizza l'approccio a blocchi per modificare la strategia di calcolo in modo tale che gli elementi siano mantenuti nei registri vettoriali. Questo approccio consente di diminuire il numero di accessi alla memoria per ogni operazione in virgola

<sup>1</sup> SGI sta per Silicon Graphics, una società che produceva workstation grafiche ad alte prestazioni negli anni Novanta [N.d.T.].

mobile e di aumentare le prestazioni di quasi un fattore due! E la larghezza di banda della memoria di un Cray-1 diventa sufficiente per eseguire un ciclo che prima richiedeva più larghezza di banda, esattamente come predice il modello roofline.

## 6.14 Note conclusive

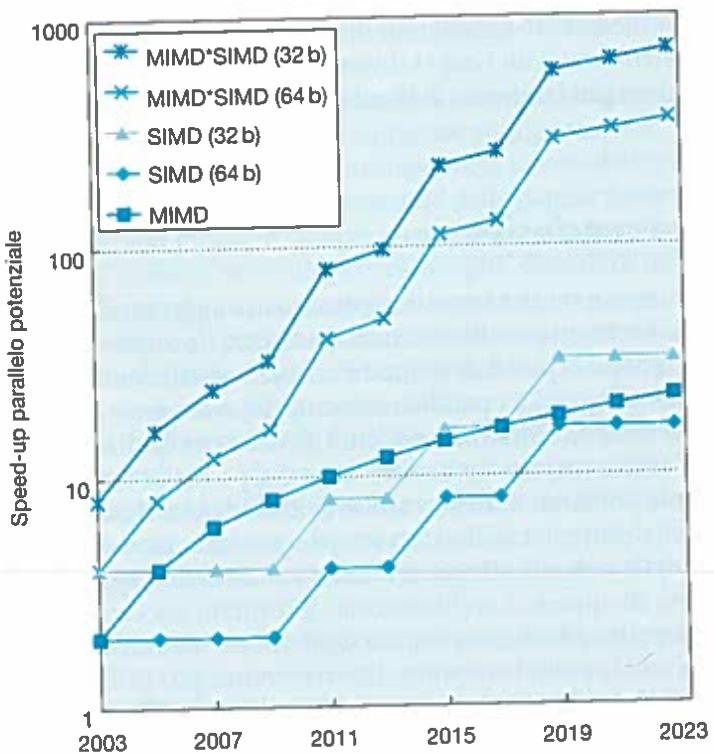
Il sogno di costruire i calcolatori semplicemente aggregando vari processori è nato con le architetture stesse; ciononostante, i progressi nel costruire e utilizzare i processori paralleli in modo efficace ed efficiente sono stati lenti: i grossi problemi legati alla parallelizzazione del software e il lungo processo di evoluzione delle architetture dei multiprocessori (al fine di aumentarne l'usabilità e l'efficienza) hanno limitato lo sviluppo dei processori paralleli. In questo capitolo abbiamo discusso molte delle sfide che riguardano il software, compresa la difficoltà nello scrivere programmi che ottengano un buon incremento della velocità di esecuzione, a causa della Legge di Amdahl. La grande varietà di approcci architetturali, il limitato successo e la breve vita di molte architetture realizzate hanno contribuito ad aumentare le difficoltà legate allo sviluppo del software. Descriveremo più in dettaglio la storia dello sviluppo dei multiprocessori nel paragrafo 6.15 . Per approfondire gli argomenti trattati in questo capitolo, potete consultare il Capitolo 4 della quinta edizione del testo *Computer Architecture: A Quantitative Approach*, per quanto riguarda il confronto tra GPU e CPU, e il Capitolo 6 per i grandi centri di elaborazione dati.

Come abbiamo accennato nel Capitolo 1, nonostante questo lungo passato costellato anche di insuccessi, l'industria della tecnologia dell'informazione ha ora legato il suo futuro al calcolo parallelo. E anche se potrebbe ancora succedere che questo tentativo fallisca, come altri tentativi sono falliti in passato, ci sono alcune ragioni per ritenere che non sarà così:

- Il *software come servizio* (SaaS) sta chiaramente crescendo di importanza e i cluster hanno avuto molto successo nel fornire questi servizi. Garantendo la ridondanza ad alto livello, anche attraverso centri di elaborazione dati distribuiti su aree geografiche molto ampie, questi servizi sono disponibili 24 ore al giorno, tutti i giorni dell'anno, per utenti che si trovano in ogni parte del mondo.
- Crediamo che i calcolatori dei grandi centri di calcolo stiano cambiando gli obiettivi e i principi su cui sono basati i loro servizi, proprio come i clienti dei dispositivi mobili stanno cambiando gli obiettivi e i principi della progettazione dei microprocessori. Entrambi questi settori stanno rivoluzionando anche l'industria del software. Sia le prestazioni per euro che le prestazioni per joule guidano lo sviluppo dei calcolatori dei grandi centri di calcolo e dei dispositivi mobili, e il parallelismo è la chiave per ottenere il loro miglioramento.
- Le operazioni SIMD e vettoriali sono particolarmente adatte alle applicazioni multimediali, che giocano un ruolo principale nell'era post-PC; esse condizionano il vantaggio della maggior facilità di utilizzo da parte dei programmati rispetto alla classica programmazione MIMD e sono più efficienti dal punto di vista energetico delle operazioni MIMD. Per mettere nella giusta prospettiva il confronto tra SIMD e MIMD, la Figura 6.28 confronta il numero di core di processori MIMD con il numero di operazioni a 32 bit o 64 bit per ciclo di clock nella modalità SIMD dei calcolatori x86 nei diversi anni. Nei calcolatori x86, ci aspettiamo di vedere due nuovi core per chip circa ogni due anni e la larghezza delle istruzioni SIMD raddoppiata ogni

*Stiamo dedicando lo sviluppo di tutti i nostri prodotti futuri alle architetture multicore. Crediamo che questo sia un punto di svolta per l'industria. [...] Questa non è una competizione. Questa è una svolta radicale nell'approccio al calcolo...*

Paul Otellini, presidente Intel,  
Intel Developers Forum, 2004



**Figura 6.28** Speed-up potenziale ottenibile attraverso il parallelismo nelle architetture Intel x86, MIMD, SIMD, e miste SIMD e MIMD, degli ultimi anni. Abbiamo ipotizzato che due core per ogni chip vengano aggiunti ogni due anni alle architetture MIMD e che il numero di operazioni nelle architetture SIMD raddoppi ogni quattro anni.

quattro anni. In base a queste ipotesi, nella prossima decade lo speed-up ottenuto dal parallelismo SIMD sarà il doppio del parallelismo MIMD. Data l'efficacia delle istruzioni SIMD per il multimediale e l'aumento della sua importanza nell'era post-PC, questa enfasi sulle architetture SIMD può essere considerata appropriata. Perciò, è tanto importante capire il parallelismo SIMD quanto il parallelismo MIMD, anche se il secondo finora ha ricevuto un'attenzione maggiore.

- L'utilizzo del calcolo parallelo in ambito scientifico e ingegneristico è molto diffuso. Questi settori hanno bisogno di una potenza di calcolo pressoché illimitata; inoltre, molte applicazioni utilizzate in questi campi presentano un grado elevato di concorrenza intrinseca nei calcoli. Attualmente, i cluster dominano in queste aree applicative: per esempio, analizzando la graduatoria dei 500 calcolatori più veloci (Top 500) del 2012, l'80% dei calcolatori più veloci sui benchmark Linpack era costituito da cluster.
- Tutti i produttori di microprocessori desktop e server stanno costruendo multiprocessori con prestazioni sempre più elevate; rispetto al passato, però, non esiste una strada già tracciata da seguire per aumentare le prestazioni delle applicazioni sequenziali. Perciò, i programmatore *devono* parallelizzare il loro codice o scrivere nuovi programmi che consentano l'elaborazione parallela.
- In passato si usavano diverse metriche per valutare le prestazioni di microprocessori e multiprocessori. Quando aumentavano le prestazioni dei monoproessori, i progettisti si ritenevano soddisfatti se le prestazioni delle applicazioni a singolo thread crescevano proporzionalmente alla radice quadrata dell'area di silicio del chip; era quindi accettabile un incremento sublineare in funzione delle risorse. Il successo di un multiprocessore veniva

decretato quando la sua velocità risultava cresciuta *linearmente* con il numero dei processori, supponendo implicitamente che il costo dell'acquisto e dell'amministrazione di  $n$  processori fosse  $n$  volte quello del singolo processore. Ora che il parallelismo viene realizzato sullo stesso chip (grazie ai processori multicore), possiamo utilizzare la stessa metrica dei processori tradizionali (andamento sublineare) per valutare i multiprocessori.

- Il successo della compilazione Just-In-Time durante l'esecuzione e dell'autoadattamento rende plausibile l'idea di un software che si adatti autonomamente, per sfruttare l'aumento del numero di core del chip; ciò fornirebbe una flessibilità che non è disponibile quando ci si limita ai compilatori statici.
- Oggi il movimento dell'*open source* è diventato una parte significativa dell'industria del software. Questo movimento ha molti meriti, anche perché le soluzioni meglio ingegnerizzate possono essere concepite quando la condivisione delle idee degli sviluppatori non deve sottostare a certi vincoli legali. Inoltre, esso favorisce l'innovazione, dato che invita a modificare il vecchio software e ad accogliere positivamente nuovi linguaggi e nuovi prodotti. Questo approccio "open" (aperto) potrebbe rivelarsi estremamente utile in questo periodo di rapidi cambiamenti.

Per motivare il lettore a seguire questa rivoluzione, abbiamo dimostrato le potenzialità del parallelismo concretamente su una funzione che moltiplica matrici su un Core i7 Intel (Sandy Bridge) nei paragrafi 3.8, 4.12, 5.14, 6.12.

- Il parallelismo a livello di dati descritto nel Capitolo 3 ha consentito di aumentare le prestazioni di un fattore 3,85 eseguendo in parallelo quattro operazioni in virgola mobile a 64 bit utilizzando gli operandi a 256 bit delle istruzioni AVX, dimostrando così la validità dell'approccio SIMD.
- Il parallelismo a livello di istruzioni descritto nel Capitolo 4 ha consentito di aumentare le prestazioni di un ulteriore fattore 2,3 espandendo i cicli 4 volte in modo da fornire allo scheduler hardware, che invia a esecuzione le istruzioni fuori ordine, più istruzioni da inviare a ogni ciclo di clock.
- Le ottimizzazioni della cache del Capitolo 5 ci hanno consentito di aumentare ancora le prestazioni di un fattore tra 2,0 e 2,5 utilizzando la cache a blocchi per ridurre le miss della cache.
- Il parallelismo a livello di thread di questo capitolo ci ha consentito di aumentare le prestazioni quando le matrici non possono essere contenute interamente nella cache L1 di un altro fattore che varia da 4 a 14 utilizzando tutti i 16 core del nostro microprocessore multicore, dimostrando la validità dell'approccio MIMD. Abbiamo ottenuto ciò aggiungendo una singola linea che contiene un comando *pragma OpenMP*.

Utilizzare le idee di questo libro e adattare il software al calcolatore considerato ha comportato l'aggiunta di 24 linee di codice alla funzione DGEMM originaria. Per matrici di dimensione  $32 \times 32$ ,  $160 \times 160$ ,  $480 \times 480$  e  $960 \times 960$ , lo speed-up complessivo delle prestazioni ottenuto mettendo in pratica le idee descritte aggiungendo progressivamente fino a 24 linee di codice è stato rispettivamente pari a un fattore 8, 39, 129 e 212!

Questa rivoluzione parallela dell'interfaccia hardware/software è forse la sfida più importante degli ultimi 60 anni nel campo dell'informatica. Fornirà molte opportunità di ricerca e di business all'interno e al di fuori del settore della tecnologia dell'informazione, e le società che domineranno nell'era dei multicore non è detto che siano le stesse che hanno dominato nell'era delle architetture monoprocessoressi. Dopo avere compreso il trend di sviluppo dell'hardware e avere appreso come adattare il software all'hardware sottostante, potrete essere tra gli innovatori che coglieranno le opportunità che sicuramente si presenteranno nei prossimi anni. E noi saremo ben lieti di beneficiare delle vostre invenzioni!

## 6.15 Inquadramento storico e approfondimenti

Questo paragrafo, disponibile online , descrive la ricca e spesso disastrosa storia dei multiprocessori degli ultimi 50 anni.

### Bibliografia

B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, R. Sears. Benchmarking cloud serving systems with YCSB, In: *Proceedings of the 1st ACM Symposium on Cloud computing*, June 10-11, 2010, Indianapolis, Indiana, USA, doi:10.1145/1807128.1807152.

G. Regnier, S. Makineni, R. Illikkal, R. Iyer, D. Minturn, R. Huggahalli, D. Newell, L. Cline, and A. Foong. TCP onloading for data center servers. *IEEE Computer*, 37(11):48–58, 2004.

## 6.16 | Esercizi

**6.1** Stilare un elenco delle attività quotidiane che si svolgono in un tipico giorno feriale. Per esempio, alzarsi, fare la doccia, vestirsi, fare colazione, asciugarsi i capelli, lavarsi i denti ecc. L'elenco deve contenere almeno 10 attività.

**6.1.1** [5] <6.2> Dire quali attività sfruttano già una forma di parallelismo: per esempio, lavare tutti i denti contemporaneamente invece di lavare un dente alla volta, o portare un libro alla volta in università invece di metterli nello zaino e portarli tutti assieme. Spiegare perché alcune attività non sono eseguite in parallelo.

**6.1.2** [5] <6.2> Dire quali attività potrebbero essere svolte in modo concorrente (per esempio, fare colazione e ascoltare il giornale radio). Per ciascuna attività, definire quali altre attività possono essere svolte assieme ad essa.

**6.1.3** [5] <6.2> Considerando la risposta all'Esercizio 6.1.2, come si potrebbero modificare i sistemi utilizzati (per esempio, le docce, i vestiti, le TV, le automobili), per riuscire a eseguire più attività in parallelo?

**6.1.4** [5] <6.2> Stimare quanto tempo si risparmierebbe eseguendo il maggior numero di attività in parallelo.

**6.2** Si supponga di dover cucinare 3 torte ai mirtilli. Gli ingredienti sono i seguenti:

1 tazza di burro ammorbidente

1 tazza di zucchero

4 uova grandi

1 cucchiai da tè di estratto di vaniglia

1/2 cucchiai da tè di sale

1/4 cucchiai da tè di noce moscata

1 1/2 tazza di farina

1 tazza di mirtilli

La ricetta per una singola torta è la seguente:

Passo 1: scaldare il forno a 160 °C. Ungere e cospargere di farina la tortiera.

Passo 2: in un recipiente grande mescolare con un frullino burro e zucchero a velocità intermedia fino a quando l'impasto non diventa spumoso e leggero. Aggiungere uova, vaniglia, sale e noce moscata. Frullare l'impasto fino a quando non diventa omogeneo e amalgamato. Ridurre la velocità del frullino e aggiungere 1/2 tazza alla volta di farina, mescolando fino a quando l'impasto non diventa uniforme.

Passo 3: disporre delicatamente i mirtilli sull'impasto spargendoli uniformemente. Cuocere in forno per 60 minuti.

**6.2.1** [5] <6.2> Il compito consiste nel cuocere 3 torte nella maniera più efficiente possibile. Si supponga di avere a disposizione un forno (grande abbastanza da contenere una sola torta), un recipiente, una tortiera grande e un frullino. Proporre la sequenza di attività che consenta di cucinare le 3 torte nel minor tempo possibile, identificando i colli di bottiglia nell'eseguire questo compito.

**6.2.2** [5] <6.2> Si supponga ora di avere a disposizione 3 recipienti, 3 tortiere e 3 frullini; quanto più velocemente si riesce a eseguire questo compito, avendo a disposizione più risorse?

**6.2.3** [5] <6.2> Si supponga ora di avere due amici che aiutano a cucinare e di avere a disposizione un

forno più grande che possa contenere 3 torte. Come si può modificare la sequenza delle attività proposta nell'Esercizio 6.2.1?

**6.2.4 [5] <6.2>** Confrontare la preparazione delle torte con il completamento di 3 iterazioni di un ciclo su un calcolatore parallelo. Identificare il parallelismo a livello dei dati e il parallelismo a livello dei task nella preparazione delle torte.

**6.3** Molte applicazioni eseguite sui calcolatori richiedono la ricerca e l'ordinamento di insiemi di dati. È stato sviluppato un gran numero di algoritmi di ricerca e ordinamento con l'obiettivo di ridurre il tempo di esecuzione di queste attività di routine. In questo esercizio verrà spiegato come rendere il più possibile paralleli questi algoritmi.

**6.3.1 [10] <6.2>** Si consideri il seguente algoritmo di ricerca binaria; si tratta di un algoritmo classico del tipo "divide et impera" che cerca un certo numero  $X$  all'interno di un vettore di  $N$  elementi,  $A$ , e restituisce l'indice dell'elemento cercato.

```
RicercaBinaria(A[0..N-1], X) {
    inferiore = 0
    superiore = N-1
    while (inferiore <= superiore) {
        metà = (superiore + inferiore) / 2
        if (A[metà] > X)
            superiore = metà - 1
        else if (A[metà] < X)
            inferiore = metà + 1
        else
            return metà // elemento trovato
    }
    return -1 // elemento non trovato
}
```

Si supponga di avere a disposizione  $Y$  core su un processore multicore per eseguire la procedura Ricerca-Binaria. Si supponga che  $Y$  sia molto più piccolo di  $N$ . Esprimere lo speed-up atteso per diversi valori di  $Y$  e  $N$  e riportare questi valori su un grafico.

**6.3.2 [5] <6.2>** Si supponga quindi che  $Y$  sia uguale a  $N$ . Che effetto ha questa ipotesi sulle conclusioni riportate nel problema precedente? Se l'obiettivo fosse ottenere il maggior speed-up possibile, cioè disporre di un algoritmo che scalì in modo forte, spiegare come si potrebbe modificare questo frammento di codice per raggiungerlo.

**6.4** Si consideri il seguente frammento di codice C:

```
for (j=2; j<=1000; j++)
    D[j] = D[j-1] + D[j-2];
```

Il codice RISC-V corrispondente a questo frammento di codice è il seguente:

```
li      x5, 8000
add   x12, x10, x5
addi  x11, x10, 16
CICLO: fld   f0, -16(x11)
        fld   f1, -8(x11)
        fadd.d f2, f0, f1
        fsd   f2, 0(x11)
        addi  x11, x11, 8
        ble   x11, x12, CICLO
```

La latenza di un'istruzione è il numero di cicli di clock tra l'istruzione e un'istruzione che utilizza il suo risultato. Si supponga che le istruzioni in virgola mobile abbiano le seguenti latenze (in numero di cicli di clock):

fadd.d	fld	fsd
4	6	1

**6.4.1 [10] <6.2>** Quanti cicli di clock occorrono per eseguire questo codice?

**6.4.2 [10] <6.2>** Riordinare il codice per ridurre gli stalli. Ora, quanti cicli occorrono per eseguire il nuovo codice? (Suggerimento: potete rimuovere degli stalli modificando l'offset dell'istruzione fsd).

**6.4.3 [10] <6.2>** Quando un'istruzione di un'iterazione del ciclo dipende dal contenuto di una variabile scritta in un'iterazione precedente dello stesso ciclo, si dice che c'è una *dipendenza introdotta dal ciclo* tra due iterazioni. Identificare le dipendenze introdotte dal ciclo contenute nel frammento di codice precedente e le variabili dipendenti. Identificare i registri coinvolti analizzando il codice assembler. Si può ignorare la variabile di ciclo  $j$ .

**6.4.4 [15] <6.2>** Riscrivere il codice utilizzando i registri per trasportare i dati tra un'iterazione e l'altra dei cicli (invece di scrivere e ricaricare i dati dalla memoria principale). Mostrare dove questo codice va in stallo e calcolare il numero di cicli di clock necessari per eseguire questo codice. Notare che per questo problema si dovrà utilizzare la pseudo-istruzione assembler "fmv.d rd, rs1", che scrive il contenuto del registro in virgola mobile  $rs1$  nel registro in virgola mobile  $rd$ . Si supponga che  $fmv.d$  venga eseguita in un solo ciclo di clock.

**6.4.5 [10] <6.2>** L'espansione dei cicli è stata descritta nel Capitolo 4. Applicare tale tecnica a questo ciclo e ottimizzare il ciclo in modo tale che il ciclo ottenuto gestisca tre iterazioni del ciclo originale. Mostrare dove questo codice va in stallo e calcolare il numero di cicli di clock necessari per eseguire questo codice.

**6.4.6 [10] <6.2>** Lo srotolamento del ciclo dell'Esercizio 6.4.5 funziona bene perché volevamo un multiplo di tre iterazioni. Cosa succede se il numero di iterazioni non è noto al momento della compilazione? Con che efficienza si può gestire un numero di iterazioni che non sia un multiplo del numero di iterazioni del ciclo srotolato?

**6.4.7 [15] <6.2>** Si supponga di eseguire questo codice su un sistema a due nodi con memoria distribuita a scambio di messaggi; si supponga anche di utilizzare la procedura di scambio di messaggi descritta nel paragrafo 6.7. In questo paragrafo è stata introdotta l'operazione `send(x, y)`, che invia al nodo `x` il valore `y`, e l'operazione `receive()`, che attende il valore inviato al nodo. Si supponga che l'operazione di `send` richieda un ciclo di `clock` per essere eseguita e quindi che le istruzioni successive (in coda sul nodo) possano passare al ciclo successivo; si supponga anche che occorrono diversi cicli di `clock` perché il dato inviato arrivi al nodo di destinazione. Le istruzioni `receive` mettono in stallo l'esecuzione sul nodo in cui vengono eseguite fino a che il messaggio non è arrivato a quel nodo. Si può utilizzare un tale sistema per accelerare l'esecuzione del codice di questo esercizio? Se così fosse, quale sarebbe la massima latenza che si potrebbe tollerare per ricevere un'informazione? Se così non fosse, spiegare perché.

**6.5** Si consideri il seguente algoritmo ricorsivo di ordinamento e unione (Merge Sort); si tratta di un altro classico algoritmo di tipo “*divide et impera*”. L'algoritmo Merge Sort è stato descritto per la prima volta da John von Neuman nel 1945; l'idea di base è di dividere una lista non ordinata `x` di `m` elementi in due sottoliste lunghe circa la metà della lista originaria. Si ripete questa operazione per ciascuna sottolista fino a quando non si ottengono liste di dimensione unitaria. Quindi, le sottoliste di dimensione 1 vengono unite in un'unica lista ordinata.

```
Mergesort(m)
var list sinistra, destra, risultato
if length(m) ≤ 1
    return m
else
    var metà = length(m) / 2
    for ciascun x in m fino a metà
        add x a sinistra
    for ciascun x in m dopo metà
        add x a destra
    sinistra = Mergesort(sinistra)
    destra = Mergesort(destra)
    risultato = Merge(sinistra, destra)
    return risultato
```

Il passo “merge” viene eseguito dal frammento seguente di codice:

```
Merge(sinistra, destra)
var list risultato
while length(sinistra) > 0 and length(destra) > 0
    if primo(sinistra) ≤ primo(destra)
        append primo(sinistra) a risultato
        sinistra = resto(sinistra)
    else
        append primo(destra) a risultato
        destra = resto(destra)
    if length(sinistra) > 0
        append resto(sinistra) a risultato
    if length(destra) > 0
        append resto(destra) a risultato
return risultato
```

**6.5.1 [10] <6.2>** Si supponga di avere a disposizione `Y` core su un processore multicore per eseguire l'algoritmo MergeSort e si supponga che `Y` sia molto più piccolo della lunghezza di `m`, `length(m)`. Esprimere la percentuale di incremento della velocità che ci si può attendere per diversi valori di `Y` e di `length(m)`. Riportare questi valori su un grafico.

**6.5.2 [10] <6.2>** Si supponga quindi che `Y` sia uguale a `length(m)`. Come cambierebbe la conclusione del problema precedente? Se l'obiettivo fosse ottenere il maggior incremento di velocità possibile, cioè di fare in modo che l'algoritmo scali in modo forte, come si potrebbe modificare questo frammento di codice?

**6.6** La moltiplicazione di matrici gioca un ruolo importante in un gran numero di applicazioni. Due matrici possono essere moltiplicate tra loro solamente se il numero di colonne della prima matrice è uguale al numero di righe della seconda.

Si supponga di avere una matrice `A` di dimensioni `m × n` che deve essere moltiplicata per una matrice `B` di dimensioni `n × p`. Il prodotto sarà una matrice `C` di dimensioni `m × p` e si scriverà  $C = AB$  oppure  $C = A \cdot B$ , dove l'elemento nella posizione  $(i, j)$  di `C` sarà indicato con  $c_{ij}$  e sarà dato dalla seguente somma di prodotti:

$$c_{i,j} = \sum_{k=1}^n a_{i,k} \times b_{k,j} = a_{i,r} b_{r,j} = a_{i,1} b_{1,j} + a_{i,2} b_{2,j} + \dots + a_{i,n} b_{n,j}$$

per ogni elemento nella posizione  $(i, j)$ , con  $1 \leq i \leq m$  e  $1 \leq j \leq p$ . Vogliamo ora vedere se possiamo rendere parallelo il calcolo della matrice `C`. Si supponga che le matrici siano scritte contiguamente in memoria per colonne:  $a_{1,1}, a_{2,1}, a_{3,1}, a_{4,1}, \dots$  ecc.

**6.6.1 [10] <6.5>** Si supponga di calcolare  $C$  su una macchina a core singolo con memoria condivisa e su una macchina con 4 core con memoria condivisa. Calcolare l'incremento di velocità atteso sulla macchina con 4 core trascurando le problematiche legate alla memoria.

**6.6.2 [10] <6.5>** Risolvere l'Esercizio 6.6.1 supponendo che gli aggiornamenti degli elementi della matrice  $C$  provochino una miss della cache causata da una falsa condivisione dei blocchi quando vengono scritti due elementi consecutivi sulla stessa riga, cioè con lo stesso indice di riga  $i$ .

**6.6.3 [10] <6.5>** Come si potrebbe risolvere il problema delle false condivisioni identificate nel problema precedente?

**6.7** Si considerino le seguenti quattro istruzioni che vengono eseguite contemporaneamente sui quattro core di un *processore simmetrico multicore* (SMP). Si supponga che, prima di eseguire queste istruzioni, le variabili  $x$  e  $y$  contengano entrambe 0.

Core 1:  $x = 2;$   
 Core 2:  $y = 2;$   
 Core 3:  $w = x + y + 1;$   
 Core 4:  $z = x + y;$

**6.7.1 [10] <6.5>** Quali sono i valori che possono assumere  $w$ ,  $x$ ,  $y$  e  $z$ ? Spiegare come si arriva a ciascun insieme di valori. Occorre esaminare tutti i possibili intrecci tra le istruzioni.

**6.7.2 [5] <6.5>** Come si potrebbe rendere l'esecuzione deterministica ed essere certi, quindi, di ottenere solamente un insieme di valori?

**6.8** Il problema della cena dei filosofi è un classico problema di sincronizzazione e di concorrenza, in genere definito nel modo seguente. I filosofi sono seduti attorno a un tavolo e svolgono una delle seguenti due attività: mangiare e pensare; quando mangiano non pensano e quando pensano non mangiano. Sul tavolo c'è un'insalatiera con degli spaghetti, e tra un filosofo e l'altro c'è una forchetta. Le loro abitudini impongono l'uso di due forchette per mangiare, in particolare quella che hanno alla loro destra e alla loro sinistra. I filosofi tra loro non parlano.

**6.8.1 [10] <6.7>** Descrivere lo scenario in cui nessun filosofo mangia e muoiono tutti di fame. Quale potrebbe essere la sequenza di azioni che porta a questa situazione?

**6.8.2 [10] <6.7>** Descrivere come si possa risolvere questo problema introducendo il concetto di priorità.

Possiamo assicurare che tutti i filosofi saranno trattati allo stesso modo? Motivare la risposta.

Si supponga ora che ci sia un cameriere incaricato di assegnare le forchette a ciascun filosofo: nessun filosofo può impugnare la forchetta fino a quando il cameriere non gli dà il permesso. Il cameriere conosce la situazione di tutte le forchette; inoltre, se imponiamo ai filosofi di prendere sempre la forchetta alla loro sinistra prima di quella che sta alla loro destra, si possono anche evitare situazioni di stallo.

**6.8.3 [10] <6.7>** Le richieste ai camerieri da parte dei filosofi possono essere implementate sia come code di richieste sia come tentativi di richiesta periodici. Se si utilizzano le code, le richieste sono gestite nell'ordine con cui vengono ricevute; il problema delle code risiede nel fatto che potrebbe succedere di non riuscire a servire sempre il filosofo la cui richiesta si trova in cima alla coda per mancanza di risorse (non ci sono forchette disponibili). Descrivere lo scenario con cinque filosofi in cui si utilizza una coda di richieste e in cui il servizio non viene assegnato (nessun filosofo mangia) anche quando ci sono due forchette disponibili per uno dei filosofi la cui richiesta non è in cima alla coda.

**6.8.4 [10] <6.7>** Se si implementano le richieste al cameriere con chiamate periodiche fino a quando le risorse non diventano disponibili, si riesce a risolvere il problema descritto nell'Esercizio 6.8.3? Motivare la risposta.

**6.9** Si considerino le seguenti tre organizzazioni della CPU:

CPU SS: un microprocessore superscalare con 2 core che fornisce la possibilità di avviare all'esecuzione le istruzioni fuori ordine su due unità funzionali (UF). Un solo thread alla volta può essere eseguito su ogni core.

CPU MT: un processore multithread a grana fine che consente di eseguire le istruzioni su due thread in modo concorrente su due unità funzionali, anche se solamente le istruzioni di un singolo thread possono essere avviate all'esecuzione a ogni ciclo di clock.

CPU SMT: un processore SMT (a "multithreading simultaneo") che consente di eseguire in modo concorrente le istruzioni di due thread; cioè, il processore contiene due unità funzionali e le istruzioni di uno o di entrambi i thread possono essere avviate all'esecuzione a ogni ciclo di clock.

Si supponga che ci siano due thread, X e Y, che devono essere eseguiti su queste CPU e che i thread comprendano le seguenti operazioni:

Thread X	Thread Y
A1 – richiede 3 cicli per essere eseguita	B1 – richiede 2 cicli per essere eseguita
A2 – non ci sono dipendenze	B2 – è in conflitto con B1 su un'unità funzionale
A3 – è in conflitto con A1 su un'unità funzionale	B3 – dipende dal risultato di B2
A4 – dipende dal risultato di A3	B4 – non ci sono dipendenze e richiede 2 cicli per essere eseguita

Si supponga che tutte le istruzioni richiedano un singolo ciclo di clock per essere eseguite a meno che non sia specificato altrimenti o che si verifichi un hazard.

**6.9.1 [10] <6.4>** Si supponga di avere a disposizione 1 CPU SS. Quanti cicli di clock sono richiesti per eseguire questi due thread? Quanti pacchetti di istruzioni non vengono riempiti a causa degli hazard?

**6.9.2 [10] <6.4>** Si supponga ora di avere 2 CPU MT. Quanti cicli sono richiesti per eseguire questi due thread? Quanti pacchetti di istruzioni non vengono riempiti a causa degli hazard?

**6.9.3 [10] <6.4>** Si supponga ora di avere 1 CPU MT. Quanti cicli sono richiesti per eseguire questi due thread? Quanti pacchetti di istruzioni non vengono riempiti a causa degli hazard?

**6.9.4 [10] <6.4>** Si supponga ora di avere 1 CPU SMT. Quanti cicli sono richiesti per eseguire questi due thread? Quanti pacchetti di istruzioni non vengono riempiti a causa degli hazard?

**6.10** Il software di virtualizzazione viene utilizzato in modo massiccio per ridurre i costi della gestione dei moderni server ad alte prestazioni. Società come VMWare, Microsoft e IBM hanno tutte sviluppato un insieme di prodotti per la virtualizzazione. Il concetto generale, descritto nel Capitolo 5, è che si può introdurre un livello di supervisione (*hypervisor*) tra l'hardware e il sistema operativo per consentire a più sistemi operativi di condividere lo stesso hardware fisico. Il livello di supervisione è quindi responsabile dell'allocazione delle risorse della CPU e della memoria, oltre che della gestione dei servizi tipicamente gestiti dal sistema operativo (quale l'I/O).

La virtualizzazione fornisce una visione astratta dell'hardware sottostante al sistema operativo ospitato e al software applicativo; questo richiede di ripensare la progettazione dei sistemi multicore e multiprocessore per supportare la condivisione della CPU e della memoria da parte di più sistemi operativi concorrenti.

**6.10.1 [30] <6.4>** Selezionare due supervisori disponibili sul mercato, comparare ed evidenziare le differenze nel modo di virtualizzare e gestire l'hardware sottostante, cioè la CPU e la memoria.

**6.10.2 [15] <6.4>** Discutere quali modifiche potrebbero essere apportate alle piattaforme multicore del futuro per soddisfare meglio le richieste di risorse fatte a queste architetture. Per esempio, il multithreading può giocare un ruolo determinante per diminuire la competizione per le risorse di calcolo?

**6.11** Si vuole eseguire questo ciclo nella maniera più efficiente possibile. Sono a disposizione due macchine diverse: una MIMD e una SIMD.

```
for (i=0; i < 2000; i++)
    for (j=0; j < 3000; j++)
        X_vettore[i][j] = Y_vettore[i][j] + 200;
```

**6.11.1 [10] <6.3>** Mostrare la sequenza di istruzioni RISC-V che verrebbero eseguite su ciascuna CPU di una macchina MIMD dotata di quattro CPU. Quale sarebbe l'incremento di velocità per questa macchina MIMD?

**6.11.2 [20] <6.3>** Scrivere un programma assembler utilizzando un'estensione SIMD del RISC-V di vostra ideazione per eseguire il ciclo su una macchina SIMD larga 8, cioè dotata di 8 unità funzionali parallele. Confrontare il numero di istruzioni eseguite sulla macchina SIMD e sulla MIMD.

**6.12** Un array sistolico è un esempio di macchina MISD: si tratta, infatti, di una rete di pipeline o di un "fronte d'onda" di elementi di elaborazione. Ciascuno di questi elementi non richiede un program counter, poiché l'esecuzione viene avviata non appena arrivano i dati. Gli array sistolici sincronizzati effettuano i calcoli attraverso passi prestabiliti, in ciascuno dei quali ogni processore alterna una fase di calcolo e una di comunicazione.

**6.12.1 [10] <6.3>** Esaminare le implementazioni proposte di array sistolici, le cui descrizioni possono essere trovate su Internet o in letteratura. Provare a programmare il ciclo descritto nell'Esercizio 6.11 utilizzando il modello MISD e discutere le difficoltà incontrate.

**6.12.2 [10] <6.3>** Elenicare le differenze che sussistono tra le macchine MISD e le macchine SIMD, con particolare riferimento al parallelismo dei dati.

**6.13** Si supponga di volere eseguire sulla GPU NVIDIA 8800 GTX descritta in questo capitolo il ciclo DAXPY implementato nell'assembler del RISC-V e riportato a pagina 444. Per questo esercizio si supponga che tutte le operazioni matematiche siano eseguite su numeri in virgola mobile in singola precisione; per questo motivo chiameremo il ciclo SAXPY. Si supponga, inoltre, che le diverse istru-

zioni richiedano il seguente numero di cicli di clock per essere eseguite:

Load	Store	Add.S	Mult.S
5	2	3	4

**6.13.1** [20] <6.6> Descrivere come si possano costruire dei warp per il ciclo SAXPY in modo da sfruttare gli otto core forniti da un singolo multiprocessore.

**6.14** Scaricare l'SDK e il Toolkit CUDA dalla pagina web <https://developer.nvidia.com/cuda-toolkit>. Assicurarsi di scaricare la versione "emurelease" (modalità emulazione) del codice, in modo da non dover avere a disposizione fisicamente una scheda NVIDIA per svolgere questo esercizio. Compilare i programmi di esempio forniti con l'SDK e verificare che vengano eseguiti sull'emulatore.

**6.14.1** [90] <6.6> Utilizzando l'esempio "template" dell'SDK come punto di partenza, scrivere un programma CUDA che esegua le seguenti operazioni vettoriali:

$$1) \mathbf{a} - \mathbf{b} \text{ (sottrazione vettore-vettore)}$$

$$2) \mathbf{a} \cdot \mathbf{b} \text{ (prodotto scalare tra due vettori)}$$

Il prodotto scalare fra i due vettori  $\mathbf{a} = [a_1, a_2, \dots, a_n]$  e  $\mathbf{b} = [b_1, b_2, \dots, b_n]$  è definito come:

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

Scrivere un programma per ogni operazione, eseguirlo e verificare la correttezza del risultato.

**6.14.2** [90] <6.6> Se si ha a disposizione una GPU, effettuare un'analisi delle prestazioni dei programmi scritti esaminando il tempo di calcolo della versione per GPU e della versione per CPU dei due programmi, utilizzando vettori di dimensioni molto diverse. Spiegare i risultati ottenuti.

**6.15** AMD ha recentemente annunciato che ha integrato un'unità di elaborazione grafica nei suoi core x86 in un unico chip, anche se i due diversi core utilizzeranno un clock a frequenza diversa. Si tratta di un esempio di sistema multiprocessore eterogeneo. Uno degli elementi chiave del progetto è l'elevata velocità di trasmissione dei dati tra la CPU e la GPU. Prima dell'architettura Fusion di AMD, erano richiesti dei canali di comunicazione appositi per fare comunicare il chip della CPU e quello della GPU. Attualmente, il progetto prevede di utilizzare più canali PCI express (almeno 16) per facilitare la comunicazione.

**6.15.1** [25] <6.6> Confrontare la larghezza di banda e la latenza associate a queste due tecnologie di interconnessione.

**6.16** La Figura 6.14b mostra una topologia di interconnessione a forma di  $n$ -cubo di ordine 3 che collega otto nodi. Una caratteristica interessante della topologia di interconnessione di rete  $n$ -cubo è la sua capacità di fornire la connettività anche in presenza di collegamenti interrotti.

**6.16.1** [10] <6.8> Sviluppare un'equazione che calcoli quanti collegamenti di un  $n$ -cubo (dove  $n$  è l'ordine del cubo) possono essere interrotti prima di non riuscire più a garantire il collegamento a ogni nodo.

**6.16.2** [10] <6.8> Confrontare la resistenza ai guasti di un  $n$ -cubo rispetto a una rete completamente connessa con lo stesso numero di nodi. Per le due topologie, disegnare un grafico che mostri l'affidabilità in funzione del numero di collegamenti che si possono interrompere.

**6.17** I benchmark rientrano in un campo di studio che comprende l'identificazione dei carichi di lavoro rappresentativi da eseguire su specifiche piattaforme di calcolo, con lo scopo di comparare in modo oggettivo le prestazioni di diversi sistemi. In questo esercizio confronteremo due tipi di benchmark: i benchmark CPU Whetstone e la suite di benchmark PARSEC. Tutti i programmi della PARSEC dovrebbero essere liberamente disponibili in Internet. Si supponga di eseguire più copie dei programmi Whetstone e di confrontarne l'esecuzione con quella del benchmark PARSEC su uno qualsiasi dei sistemi descritti nel paragrafo 6.11.

**6.17.1** [60] <6.10> Quali sono le differenze intrinseche tra i due tipi di carichi di lavoro quando sono eseguiti su questi sistemi multicore?

**6.17.2** [60] <6.10> Utilizzando il modello roofline, determinare quanto i risultati ottenuti dall'esecuzione di questi benchmark dipendano dal livello di condivisione dei dati e di sincronizzazione richiesta dai rispettivi carichi di lavoro.

**6.18** Quando si eseguono calcoli con matrici sparse, la latenza delle gerarchie delle memorie diventa molto importante. Alle matrici sparse manca la località spaziale del flusso continuo di dati che si riscontrano tipicamente nelle operazioni matriciali; per questo motivo, sono state proposte nuove forme di rappresentazione delle matrici.

Una delle prime forme di rappresentazione delle matrici sparse è il formato delle matrici sparse di Yale (*Yale Sparse Matrix Format*). In questo formato una matrice sparsa  $M$  di dimensioni  $m \times n$  viene memorizzata come se fosse una singola riga utilizzando tre vettori. Sia  $R$  il numero di elementi di  $M$  diversi da 0; si può costruire un vettore  $A$ , di lunghezza  $R$ , che

contiene tutti gli elementi diversi da zero. Costruiamo quindi un secondo vettore  $IA$  di lunghezza  $m + 1$ , cioè pari al numero di righe più 1:  $IA(i)$  conterrà l'indice di  $A$  corrispondente al primo elemento diverso da zero di ogni riga  $i$ . La riga  $i$ -esima della matrice originaria si estende quindi tra  $A(IA(i))$  e  $A(IA(i+1) - 1)$ . Il terzo vettore  $JA$  contiene gli indici di colonna di ciascun elemento di  $A$ , per cui la sua lunghezza è pari a  $R$ .

**6.18.1** [15] <6.10> Si consideri la matrice sparsa  $X$  qui riportata. Scrivere il codice C che memorizza questa matrice nel formato delle matrici sparse di Yale.

Riga 1	[1, 2, 0, 0, 0, 0]
Riga 2	[0, 0, 1, 1, 0, 0]
Riga 3	[0, 0, 0, 0, 9, 0]
Riga 4	[2, 0, 0, 0, 0, 2]
Riga 5	[0, 0, 3, 3, 0, 7]
Riga 6	[1, 3, 0, 0, 0, 1]

**6.18.2** [10] <6.10> Si supponga che ciascun elemento della matrice  $X$  sia un numero in virgola mobile in singola precisione. Calcolare la quantità di memoria necessaria per memorizzare la matrice precedente nel formato delle matrici sparse di Yale.

**6.18.3** [15] <6.10> Eseguire la moltiplicazione matriciale di  $X$  per il seguente vettore colonna  $Y$ :

$$[2, 4, 1, 99, 7, 2]$$

Inserire il calcolo in un ciclo e misurare il tempo di esecuzione. Assicurarsi di ripetere più volte la moltiplicazione, in modo da ottenere una misura affidabile dei tempi. Confrontare i tempi di esecuzione di una codifica semplice delle matrici con il formato delle matrici sparse di Yale.

**6.18.4** [15] <6.10> Sareste in grado di proporre una rappresentazione delle matrici sparse che sia più efficiente in termini di memoria richiesta e di incremento del tempo di calcolo?

**6.19** Ci aspettiamo che nei calcolatori delle prossime generazioni saranno inserite piattaforme di calcolo eterogenee (cioè costruite con CPU eterogenee). Abbiamo già iniziato a vedere alcune di queste piattaforme nel mondo dei processori embedded, per esempio DSP in virgola mobile e CPU che funzionano da microcontrollori sono già integrate in moduli contenenti più chip.

Si supponga di avere a disposizione tre tipi di CPU:

CPU A – Una CPU multicore con velocità moderata, dotata di un'unità in virgola mobile, che può eseguire più istruzioni per ciclo di clock.

CPU B – Una CPU a singolo core, veloce, intera (cioè senza unità in virgola mobile), che può eseguire un'istruzione per ciclo di clock.

CPU C – Una CPU vettoriale lenta con capacità di elaborazione in virgola mobile, che può eseguire più copie di una stessa istruzione per ciclo di clock.

Si supponga che i processori lavorino alle seguenti frequenze:

CPU A	CPU B	CPU C
1 GHz	3 GHz	250 MHz

La CPU A può eseguire 2 istruzioni per ciclo di clock, la CPU B può eseguire 1 istruzione per ciclo di clock e la CPU C può eseguirne 8 (purché si tratti della stessa istruzione eseguita otto volte). Si supponga, inoltre, che tutte le operazioni possano essere completate in un singolo ciclo di clock se non si verificano hazard.

Tutte e tre le CPU hanno la possibilità di eseguire operazioni aritmetiche sui numeri interi; la CPU B non può eseguire direttamente le operazioni sui numeri in virgola mobile. Le CPU A e B sono dotate di un insieme di istruzioni simile a quello del processore RISC-V. La CPU C può eseguire in virgola mobile solamente le operazioni di somma e sottrazione e le operazioni di load e store relative alla memoria. Si supponga che tutte le CPU abbiano accesso a una memoria condivisa e che la sincronizzazione non abbia alcun costo.

L'operazione con cui vogliamo confrontare le CPU consiste nel comparare due matrici  $X$  e  $Y$ , ciascuna contenente  $1024 \times 1024$  elementi in virgola mobile. Il risultato di questa operazione sarà il numero di elementi per i quali il valore contenuto nella matrice  $X$  è maggiore di quello contenuto nella matrice  $Y$ .

**6.19.1** [10] <6.11> Descrivere come si potrebbe partizionare il problema su queste tre diverse CPU per ottenere le prestazioni più elevate.

**6.19.2** [10] <6.11> Quale tipo di istruzione aggiungereste alla CPU vettoriale C per ottenere prestazioni migliori?

**6.20** Questo quesito analizza la lunghezza delle code che si verifica in un sistema data la massima frequenza di elaborazione delle transazioni e la latenza media osservata per ogni transazione. La latenza comprende sia il tempo per fornire il servizio (che viene calcolato a partire dalla frequenza massima) sia il tempo di attesa in coda.

Si supponga che un sistema di calcolo con quattro core debba elaborare le transazioni contenute in un database con una frequenza di richieste al secondo costante nel tempo. Si supponga, inoltre, che ciascuna transazione richieda, in media, un tempo fisso per essere elaborata. La tabella seguente mostra coppie

di valori: latenza delle transazioni e frequenza di elaborazione.

Latenza media delle transazioni	Massima frequenza di elaborazione delle transazioni
1 ms	5000/s
2 ms	5000/s
1 ms	10 000/s
2 ms	10 000/s

Per ciascuna coppia di valori riportati nella tabella precedente, rispondere alle domande che seguono.

**6.20.1** [10] <6.11> In media, quante richieste devono essere elaborate in ogni istante?

**6.20.2** [10] <6.11> Se si passasse a un sistema con 8 core, come cambierebbe idealmente il throughput (cioè, quante transazioni/secondo riuscirebbe a elaborare il calcolatore)?

**6.20.3** [10] <6.11> Spiegare perché è raro che si riesca a ottenere questo incremento di velocità semplicemente aumentando il numero di core.

### Risposte alle domande di autovalutazione

**Paragrafo 6.1, pagina 436** – Falso. Il parallelismo a livello di task può aiutare le applicazioni sequenziali e le applicazioni sequenziali possono essere modificate per essere eseguite su hardware parallelo, anche se è più difficile.

**Paragrafo 6.2, pagina 441** – Falso. La scalatura *debole* può compensare una parte seriale di un programma che altrimenti limiterebbe la scalabilità, ma non la scalatura *forte*.

**Paragrafo 6.3, pagina 448** – Vero. Tuttavia mancano di alcune caratteristiche vettoriali, quali la funzione gather-scatter, e i registri vettoriali che migliorano l'efficienza delle architetture vettoriali. Come menzionato nella sezione *Approfondimento*, le estensioni SIMD AVX2 offrono, su dati sparsi, operazioni di load (gather) ma *non* di store (scatter). La generazione Haswell dei microprocessori x86 è stata la prima a supportare le istruzioni AVX2.

**Paragrafo 6.4, pagina 451** – 1. Vero. 2. Vero.

**Paragrafo 6.5, pagina 455** – Falso. Dato che un indirizzo condiviso è un indirizzo *fisico*, più task, ciascuno nel suo spazio di indirizzamento *virtuale*, possono essere eseguiti bene su un multiprocessore con memoria condivisa.

**Paragrafo 6.6, pagina 461** – Falso. I chip di DRAM grafica sono apprezzati per la loro maggiore larghezza di banda.

**Paragrafo 6.7, pagina 466** – 1. Falso. Inviare e ricevere un messaggio è una sincronizzazione implicita, così come un modo per condividere i dati. 2. Vero.

**Paragrafo 6.8, pagina 469** – Vero.

**Paragrafo 6.10, pagina 479** – Vero. È verosimile che occorrono delle innovazioni a tutti i livelli dell'hardware e del software per vincere la scommessa dell'industria sul calcolo parallelo.



# Indice analitico

## A

Acronimi, 7  
*add immediate* (addi), 64  
Affidabilità, 10  
AFR (frequenza annua media di malfunzionamenti), 363  
Algoritmo Bubble Sort, 126  
Aliasing, 388  
AMAT (tempo medio di accesso alla memoria), 348  
  calcolo del, 348  
Amdhal, legge di, 44  
AND, 80  
  immediato (andi), 81  
Architettura  
  dell'insieme di istruzioni (ISA), 18  
  implementazione, 19  
  metamorfosi della, 37  
  registri, 300  
  RISC-V, 78, 166-167  
  scalare, 445  
  vettoriale, 443  
  x86, 131  
ASCII (*American Standard Code for Information Interchange*), 95  
  codifica dei caratteri, 96  
  confronto con codifica binaria, 95  
Assemblatore, 111  
Assembler, 11  
  linguaggio, 11  
Astrazione, 9  
Augmented reality, 2  
Autocalibrazione, 362

## B

Banda, 20  
Barriera dell'energia, 35  
Benchmark, 41  
  paralleli, 470, 471  
  per i multiprocessori, 470  
  SPEC per la CPU, 41  
  SPEC sull'assorbimento di potenza, 42  
Bit, 11  
  di indirizzamento, 379  
  di presenza (sticky), 190  
  di segno, 68  
  di una cache, 336  
  di utilizzo, 379  
  di validità, 333  
  in ultima posizione (ulp), 189  
  meno significativo, 66  
  più significativo, 66  
Bitmap, 14  
Blocco (linea), 322  
  come si individua nelle gerarchie delle memorie, 396  
  dove può essere posizionato nelle gerarchie delle memorie, 395  
  quale deve essere sostituito in caso di miss della cache, 397  
Blocco di base, 84  
*branch if equal* (beq), 81  
*branch if not equal* (bne), 81  
Bubble Sort, 126  
Buffer  
  degli indirizzi di salto, 277

- di predizione dei salti, 276
- di riordino, 293
- di scrittura, 340
- Bytecode Java, 118
  
- C**
- Cache
  - a mappatura diretta, 331, 332
  - accesso alla, 333
  - bit di una, 336
  - coerenza della, 405
  - come misurare e migliorare le prestazioni, 345
  - come trovare un blocco nella, 353
  - completamente associativa, 348
  - determinazione delle prestazioni, 346
  - gestione delle miss della, 339
  - indirizzata fisicamente, 388
  - indirizzata virtualmente, 388
  - miss della, 339
  - miss e associatività nelle, 350
  - multilivello, 355, 356, 357
  - set-associativa, 349
  - split, 343
- Calcolatori
  - architettura, 18
  - aritmetica, 151
    - divisione, 160
    - moltiplicazione, 155
    - somme e sottrazioni, 151
  - caratteristiche, 3
  - componenti, 13
  - comunicazione con altri calcolatori, 20
  - embedded, 3
  - introduzione, 2
  - linguaggio dei, 53
  - operandi dell'hardware, 58
  - operazioni svolte dall'hardware dei, 54
  - organizzazione con i 5 componenti classici, 15
  - per grandi centri di calcolo, 463
  - prestazione, 24
    - tipi di, 3
  - Cambio di contesto, 390
  - Cammino vettoriale, 447
  - Campo
    - indice, 335
    - tag, 335, 336
      - dimensione del campo, 355
  - Caricamento anticipato (*prefetching*), 422
  - Caricamento di una costante su 32 bit, 100
  - Carico di lavoro, 41
  - Chiamante, 87
  - Chiamata, 87
    - di sistema (*system call*), 389
  - Chip, 16, 23
    - processo produttivo, 23
  - Ciclo, 83
    - while, 83
  - Ciclo di clock, 28
    - numero di cicli per istruzione (CPI), 30, 240
  - Ciclo di stallo
    - della memoria, 345, 346
    - in lettura, 345
    - in scrittura, 345
  - Cifra binaria, 11, 65
  - Circuito integrato, 16, 21
    - a grandissima scala di integrazione (VLSI), 22
    - costo di un, 24
  - Clock
    - ciclo di, 28
    - ciclo per istruzione (CPI), 30
    - periodo di, 28
  - Cloud computing, 5
  - Cluster, 435, 463
  - CMOS (*Complementary Metal Oxid Semiconductors*), 36
  - Codice
    - di Hamming, 365
    - di rilevamento degli errori, 365
    - macchina, 72
    - operativo (cod op), 73, 223
  - Codifica
    - ASCII, 95-96
    - Unicode, 98, 99
  - Collegamento *lazy*, 117
  - Competizione sui dati, 107
  - Compilatore, 11, 57, 110
    - Compilatori Just In Time (JIT), 119
    - Complemento a 1, 71
    - Complemento a 2, 67
    - Completamento in ordine, 294
    - Concetto di programma memorizzato, 54, 78
  - Conversione
    - da binario a decimale, 68
    - da binario a esadecimale, 73
    - da esadecimale a binario, 73
  - Coprocessore, 188
  - Core dump, 105
  - Core i7 Intel, 41
    - 920, 299
    - 960, 480
    - CPI, 302
    - gerarchie delle memorie, 410
    - pipeline, 297, 301
    - prestazioni, 302
    - programmi benchmark su, 42
    - specifiche, 297
    - valutazione, 41
  - Cortex-A53 ARM, 297
    - CPI, 299
    - gerarchie delle memorie, 410
    - pipeline, 297, 298
    - specifiche del, 297
  - Costrutti
    - case/switch, 85
    - if, 81, 83
    - if-then-else, 85
  - CPU (*central processing unit*), 7, 16
    - prestazioni della, 29

**D**

- Deadlock*, 109  
**DGEMM** (moltiplicazione generale tra matrici, in doppia precisione), 185, 305, 484  
 prestazioni di tre versioni su matrici  $32 \times 32$ , 306  
**Difetto**, 23  
**DIMM** (*Dual Inline Memory Modules*), 327  
**Dipendenza effettiva**, 291  
**Dipendenze**  
 false (o nominali), 291  
 nella pipeline, 262, 264  
 rilevamento delle, 262  
**Dischi RAID**, 409  
**Disco magnetico** (hard disk), 19  
**Disco rigido**, 19  
**Dispositivi**  
 di input, 13  
 di output, 13  
**Dispositivo mobile** (PMD), 5  
**Distanza di Humming**, 365  
**Divisione**, 160  
 algoritmi, 161, 163  
 di numeri dotati di segno, 164  
 dividendo, 161  
 divisore, 161  
 nel RISC-V, 165  
 quoziente, 161  
 resto, 161  
 veloce, 165  
**DRAM**, 325

**E**

- Eccezione**, 170, 278  
 abilitazione, 391  
 disabilitazione, 391  
 dovuta a un malfunzionamento dell'hardware in un'istruzione add, 283  
 gestione nella pipeline, 281  
 gestione nelle architetture RISC-V, 280  
 imprecisa, 284  
 in un calcolatore dotato di pipeline, 282  
 precisa, 284  
 unità di elaborazione e segnali di controllo, 282  
**Efficienza energetica**, 296  
**Elemento**  
 combinatorio, 210  
 dell'unità di elaborazione, 213  
 di stato, 211  
**Esecuzione**  
 fuori ordine, 294  
 parallela, 285  
 slot di, 286  
**Espansione dei cicli**, 290  
 con pipeline a esecuzione parallela, 291  
**Esponente**, 169  
**Estensione**  
 avanzata dell'x86, 193  
 del segno, 216

**F**

- multimediale, 446  
**SIMD**, 193, 442  
**EX** (esecuzione dell'indirizzo), 245, 249, 251  
 hazard nello stadio, 266  
**F**  
 Falsa condivisione, 408  
**FastMATH Intrinsics**, 342  
 TLB del, 384  
**File eseguibile**, 113  
**File oggetto**, 112  
 eseguire il link di, 114  
 intestazione del, 114  
**Flush** (di istruzioni), 273  
**Formato**  
 di tipo I, 74  
 di tipo R, 74  
**Frame buffer**, 14  
**Frame**  
 della procedura, 92  
 pointer, 92  
**free**, 93  
**Frequenza**  
 delle hit (hit rate), 322  
 delle miss (miss rate), 323  
 di clock di processori Intel x86, 35  
**Frequenza annua media di malfunzionamenti** (AFR), 363  
**Frequenza di miss**  
 globale, 361  
 locale, 361  
**Funzione stato prossimo**, 403

**G**

- Gerarchia delle memorie**, 10, 320, 322  
 affidabilità, 363  
 come individuare un blocco, 396  
 complesso delle operazioni, 387  
 dove posizionare un blocco, 395  
 malfunzionamento, 363  
 schema comune, 395  
 struttura di base, 322  
**Global pointer**, 92  
**GPU** (unità di elaborazione grafica), 455  
 di NVIDIA, 467  
 guida ai termini utilizzati per le, 461  
 prospettiva delle, 460  
 strutture della memoria, 459  
**Grado di associatività**, 355

**H**

- Hamming**  
 codice di, 365  
 codice di correzione degli errori, 366  
 distanza di, 365  
**Hamming, Richard**, 365  
**Hazard**, 235  
 e stallo, 268

- nello stadio EX, 266  
 nello stadio MEM, 266  
 strutturali, 235  
 sui dati, 236
  - di una load, 238
  - propagazione o stallo, 261, 268
  - sul controllo (o sui salti condizionati), 239, 272
- Heap**, 93
- Hit**
  - frequenza delle, 322
  - tempo di, 323
- I**
- ID** (decodifica dell'istruzione), 244, 249
- IF** (fetch dell'istruzione), 244, 248
- Indirizzamento**
  - con interleaving, 327
  - immediato, 104
  - modalità del RISC-V, 104, 105
  - nei salti, 101
  - relativo al program counter, 102, 104
  - RISC-V di un campo immediato e di un indirizzo ampio, 100
  - spazio di, 372
  - tramite base e spiazzamento, 104
  - tramite registro, 104
- Indirizzo**, 60
  - base, 60
  - di destinazione del salto, 216
  - di ritorno, 87
  - fisico, 373
  - traduzione, 374
  - virtuale, 373
- Intel x86**
  - codifica delle istruzioni, 138
  - estensione avanzata del, 193
  - evoluzione, 131
  - formati tipici delle istruzioni, 138
  - istruzioni che creano problemi alla virtualizzazione, 422
  - istruzioni tipiche, 137
  - modalità di indirizzamento, 134
  - operazioni su numeri interi, 134
  - operazioni tipiche, 137
  - registri, 134
- Intensità aritmetica**, 473
- Interfaccia binaria delle applicazioni (ABI)**, 19
- Interrupt**, 170, 278
  - imprecisi, 284
  - precisi, 284
  - vettorizzati, 280
- iPad2**
  - componenti, 17
  - schedina del processore, 17
- ISA** (architettura dell'insieme di istruzioni), 18
- Istruzioni**, 11
  - componenti di base, 33
  - composizione delle, 34
  - dell'architettura x86, 131
  - di tipo R (aritmetico-logiche), 214
  - di trasferimento dati, 60
  - formato delle, 72
  - insieme delle, 53
  - jump and link, 87
  - latenza, 308
  - lh (*load half unsigned*), 98
  - linguaggio dei calcolatori, 53
  - lui (*load upper immediate*), 100
  - nop, 269
  - numero di, 31, 33b
  - pacchetto di, 287
  - rappresentazioni nel calcolatore, 71
  - riavviabili, 392
  - riordinamento dinamico delle, 292
  - RISC-V, 106, 145
  - frequenza negli SPEC CPU2006, 201
  - sh (*store half*), 98
- Istruzioni RISC-V**
  - codifica, 75
  - codop, 73
  - funz3, 73
  - funz7, 73
  - rd, 73
  - rs1, 73
  - rs2, 73
- J**
- Java Virtual Machine (JVM)**, 118
- Java**
  - caratteri e stringhe, 98
  - gerarchia della traduzione del codice, 118
- K**
- Kernel**, modalità, 389
- L**
- Larghezza**
  - di banda totale della rete, 467
  - di bisezione, 467
- Latenza**, 244
  - dell'istruzione, 308
  - di rotazione, 330
  - di utilizzo, 289
- LCD** (schermo a cristalli liquidi), 14
- Legge**
  - dei rendimenti decrescenti, 44
  - di Amdahl, 44
  - di Moore, 8
- Libreria a caricamento dinamico (DLL)**, 116-117
- Linguaggio**
  - assembler, 11, 110
  - dei calcolatori, 53
  - di programmazione ad alto livello, 12
  - macchina, 11, 72, 104
- Linker**, 113
- Linux**, 65

Livellamento dell'usura, 328

*load*, 60

*load byte* (1b), 69

*load byte unsigned* (1bu), 68

*load doubleword* (1d), 60

Loader, 116

Località

principio di, 321

spaziale, 321

temporale, 321

Lock, 452

## M

Macchina

di Mealy, 404

di Moore, 404

Macchine a stati finiti (FSM), 401, 402

controllori basati su, 403

per il controllo semplificato della cache, 404

Macchine virtuali (VM), 369

mancanza di supporto da parte dell'architettura

dell'insieme di istruzioni, 371

requisiti del monitor, 370

*malloc*, 93

Mantissa, 169

Mappatura dell'indirizzo, 374

Maschera, 80

Mascheratura, 22

Media geometrica, 42

MEM (accesso alla memoria dati), 245, 249, 250

hazard nello stadio, 266

Memoria, 16

a disco, 328

cache, 17, 330

di massa, 19

dinamica ad accesso casuale, 16

flash, 19, 328

gerarchia, 320, 322

indirizzabile per contenuto, 353

indirizzo di, 60

locale, 458

non volatile, 19

principale, 19

statica ad accesso casuale (SRAM), 17

tecnologie della, 325

virtuale, 372

per un insieme ampio di indirizzi virtuali, 380

volatile, 19

Metodo della ripartenza anticipata, 339

Metodologia di temporizzazione, 211

Mezza parola (halfword), 98

Microarchitettura, 300

Microprocessore quadcore, 38

Migrazione, 407

MIMD, 441

MIPS (*Million Instructions Per Second*), 45

Miss

della cache, 339

di capacità, 399

di conflitto (miss di collisione), 399

frequenza, 332

obbligate (da partenza a freddo), 399

penalità, 323

Miss rate (frequenza delle miss), 323

in funzione della dimensione del blocco della cache,  
338

Modalità di memorizzazione

automatica, 92

statica, 92

Modalità kernel, 389

Modello delle tre C, 399

Modulo e segno, rappresentazione, 67

Moltiplicazione, 155

algoritmo di, 158

decimale in virgola mobile, 180

di numeri dotati di segno, 158

e somma integrate, 190

in virgola mobile, 178

moltiplicando, 155

moltiplicatore, 155

nel RISC-V, 159

prodotto, 155

veloce, 159

Moore

legge di, 8

macchina di, 404

MTBF (tempo medio tra due guasti), 364

MTTF (tempo medio di funzionamento prima di un  
malfunzionamento, 363

MTTR (tempo medio di riparazione), 364

Multicore, 451

Multiplexer, 209

Multiprocessore, 434

a memoria condivisa (SMP), 435, 451

con memoria ad accesso non uniforme (NUMA), 452

con memoria ad accesso uniforme (UMA), 452

passaggio dai sistemi uniprocessore a, 37

Multithreading

a grana fine, 448

a grana grossa, 448

hardware, 448

simultaneo (SMT), 449

## N

*nop (not operation)*, 269

NOT, 80

Notazione

polarizzata, 71

scientifica, 168

NUMA (multiprocessore ad accesso non uniforme), 452

Numeri

con e senza segno, 65

esadecimali, 72

*integers*, 69

in virgola mobile, 168

normalizzati, 168

*unsigned integers*, 69

**N**umero

- del registro da leggere, 215
- del registro di lettura, 215
- del registro di scrittura, 215

**N**umero di istruzioni per ciclo di clock, 286**NVIDIA**

- confronto con un Core i7 960 Intel, 480
- GPU, 457
- strutture di memoria delle GPU, 458

**O****O**ffset, 61**O**penMP, 454, 472**O**perandi

- allocati in memoria, 60
- del calcolatore, 58
- immediati ampi, 100
- immediati o costanti, 63
- RISC-V, 55

**O**perazioni logiche, 79

- AND, 80
- NOT, 80
- OR, 80
- XOR, 80

**O**pteron, 475

- AMD, 477
- X2, 475
- X4, 475

**O**R, 80

- esclusivo immediato (`xori`), 81
- immediato (`ori`), 81

**O**verflow, 67, 169

- condizioni per somma e sottrazione, 153

**P****P**acchetto di istruzioni, 287**P**age fault, 373, 375, 377**P**arallelismo, 9

- a livello di attività, 435
- a livello di dati, 442
- a livello di istruzioni (ILP), 285
- a livello di parola, 192
- a livello di processo, 435

**P**arallelizzazione

- a due vie, 288
- dinamica dell'esecuzione, 286
- statica dell'esecuzione, 286, 287
- statica dell'ISA del RISC-V, 288

**P**arola (word), 58**P**arola doppia (doubleword), 58**P**enalità di miss, 323**P**eriodo di clock, 28**P**ersonal computer (PC), 3**P**iastrina, 23**P**ipeline, 230

- analoga con il bucato, 231
- avanzate, 296

**C**onfronto tra processori con e senza, 232**C**ore i7, 301**D**dipendenze, 262**G**estione delle eccezioni, 281**H**azard, 235**S**trutturali, 235**S**ui dati**C**ontrollo, 239**I**mpatto sulle istruzioni di salto condizionato, 272**I**ntroduzione, 230**I**struzioni per architetture dotate di, 235**P**rogettazione dell'insieme di istruzioni per architetture con, 235**R**appresentazione grafica, 253**D**iagrammi a più cicli di clock, 253, 255, 256**D**iagrammi a singolo ciclo di clock, 253, 256**R**iordinare il codice per evitare, 238**E**sempio, 240**R**iorganizzazione dinamica del codice, 292**S**alto condizionato nella, 274**S**tallo della, 238**U**nità di controllo della, 257**P**ixel, 14**P**op, 88**P**otenza relativa, 36**P**recisione**D**oppia, 170**S**ingola, 170**P**redittore**C**orrelato, 277**D**i salto a torneo, 277**P**redizione, 9, 241, 276**B**uffer di, 276**C**cili e, 276**D**ei salti, 241**D**inamica dei salti, 275**P**restazioni**C**omponenti di base, 33**D**efinizione, 25**D**ella CPU, 29**D**i sistema, 28**D**i un calcolatore, 24**D**i una cache, 346**E**quazione di misura delle, 31**M**isura delle prestazioni associate alle istruzioni, 30**M**isura delle, 28**M**odelli delle, 472**R**elative, 27**P**rincipio di località, 321**P**rocedura, 86**A**nnidata, 90**A**zzera, 127**C**onfronto tra le versioni, 129**V**ersione che utilizza i puntatori, 128**V**ersione che utilizza vettore e indice, 127**D**i invio, 462**D**i ricezione dei messaggi, 462**F**oglia, 90**F**rame della, 92

- ordina, 121
  - allocazione dei registri, 121
  - chiamata alla procedura contenuta nella, 123
  - codice del corpo, 121
  - completa, 124
  - passaggio dei parametri nella, 123
  - salvataggio dei registri nella, 124
- scambia, 120
- Processo, 448
- Processor cores*, 4
- Processore, 6, 16, 206
  - crescita nelle prestazioni, 39
  - dotato di parallelizzazione dinamica dell'esecuzione, 291
  - FastMATH Intrinsity, 342
  - multicore, 6, 451
  - parallelo, 434
  - vettoriale, 446
- Progetto logico, convenzioni, 210
- Program counter (contatore di programma), 87, 213
- Programma memorizzato, 78
- Programmi a esecuzione parallela, 435
  - difficoltà nel creare, 436
- Propagazione (bypassing), 236
  - nel caso di due istruzioni, 236
  - rappresentazione grafica, 237
- Protocollo di snooping, 407
- Pseudoistruzioni, 111
- Pthreads, 472
- Push, 88
  
- Q**
- Quad word (parola quadrupla), 135
- Quicksort, 358
  
- R**
- Radix Sort, 358
- RAM (*Random Access Memory*), 7
- Raster refresh buffer, 14
- Record di attivazione, 92
- Register file, 214
- Register spilling (versamento dei registri), 63
- Regioni di mutua esclusione, 108
- Registri
  - base, 61
  - dell'architettura, 300
  - dell'x86, 134
  - della tabella delle pagine, 377
  - di pipeline, 263
    - dipendenze tra, 264
  - di utilizzo generale (GPR), 132
  - indice, 64
  - ridenominazione, 291
  - SCAUSA (registro causa di supervisione delle eccezioni), 280
  - SPEC, 391
  - utilizzo di più, 88
- Registro di lettura, numero del, 215
- Registro di scrittura, numero del, 215
- Rendimento, 23
- Replicazione, 407
- Rete
  - Ethernet, 20
  - geografica (WAN), 20
  - locale (LAN), 20
- Reti dei cluster, 469
- Reti di calcolatori, 466
  - completamente connesse, 467
  - crossbar, 468
  - implementazione delle topologie di, 468
  - larghezza di banda totale della rete, 467
  - larghezza di bisezione, 467
  - multistadio, 468
  - topologia ad anello, 467
- Ridondanza, 10
- Riduzione, 453
- Riordinamento
  - di un semplice frammento di codice per la parallelizzazione statica dell'esecuzione, 290
  - dinamico delle istruzioni, 292
- Riorganizzazione dinamica del codice in una pipeline, 292
- RISC-V
  - allocazione della memoria per programmi e dati, 94
  - architettura del, 78, 166-167
  - campi delle istruzioni, 73
  - caricamento di una costante su 32 bit, 100
  - codifica delle istruzioni, 75, 106
  - convenzioni di utilizzo dei registri, 94
  - divisione nel, 165
  - formati delle istruzioni, 107
  - gestione delle eccezioni, 280
  - implementazione di base del, 207, 208
  - implementazione hardware, 220
  - indirizzamento di un campo immediato e di un indirizzo ampio, 100
  - indirizzi di memoria effettivi nel, 62
  - istruzioni assembler in virgola mobile, 183
  - istruzioni in virgola mobile nel, 181, 182
  - linguaggio assembler, 55-56, 166-167
  - moltiplicazione nel, 159
  - operandi, 55
  - operandi in virgola mobile nel, 183
  - parallelizzazione statica dell'ISA del, 288
  - traduzione di assegnamenti in C in, 57
  - traduzione in linguaggio macchina delle istruzioni assembler, 75
- Ritardo di rotazione, 330
- Roofline, modello, 474, 481
  - Core i7 960 Intel
  - GPU NVIDIA, 480
  
- S**
- SaaS (software come servizio), 5, 464
- Salto
  - a indirizzi lontani, 104
  - buffer di predizione, 275

- condizionato, 81  
 hazard sul, 239  
 ipotizzare che non sia eseguito, 273  
 nella pipeline, 274  
 ridurre i ritardi associati al, 273  
 condizionato eseguito, 217  
 condizionato non eseguito, 217  
 predizione, 241, 275, 276  
**Scalabilità**  
 debole, 439  
 forte, 439  
**Scambio**  
 atomico, 108  
 di messaggi, 462  
**Scheduling**, 38  
**Schermo a cristalli liquidi (LCD)**, 14  
 a matrice attiva, 14  
**Scorrimento a destra aritmetico (srai)**, 80  
**Scrittura**  
 buffer di, 340  
 gestione della, 340  
**Seek**, 329  
**Seek time (tempo di ricerca)**, 329  
**Segmentazione**, 376  
**Segmento dei dati statici**, 93  
**Segmento di testo**, 93  
**Segnale**  
 asserito, 212  
 di controllo, 211  
 non asserito, 212  
**Semiconduttore**, 22  
**Server**, 3  
**Settori**, 328  
**Sfida**  
 dell'aumento di velocità, 438, 440  
 di velocità, 437  
**Shift**, 79  
 a destra logico immediato, 79  
 a sinistra logico immediato, 79  
**Silicio**, 22  
 lingotto cristallino di, 22  
**SIMD**, 441, 442  
 cammini di elaborazione, 457  
 negli x86, 443  
 schedulare di thread, 457  
 scheduler di blocchi, 457  
 thread di istruzioni, 457  
**Sincronizzazione**, 107, 452  
**SISD**, 441  
**Sistema operativo**, 10  
**Slot di esecuzione**, 286  
**SMT (multithreading simultaneo)**, 449  
**Software come servizio (SaaS)**, 5, 464  
**Software di sistema**, 10  
 ottimizzazione mediante elaborazione a blocchi, 357  
**Sommatori a salvataggio di riporto**, 159  
**Somme**, 151  
 binarie, 152  
 condizioni di overflow, 153  
**Sottrazioni**, 151  
**Spazio**  
 di swap, 378  
 virtuale, 373  
**Spazio di indirizzamento**, 372  
 protezione, 373  
**SPEC (Cooperativa per la Valutazione delle Prestazioni dei Sistemi)**, 41  
**SPECratio**, 41  
**Speculazione**, 286  
**Split cache**, 343  
**SPMD**, 441  
**Spostamento condizionato**, 278  
**SRAM**, 17, 325  
**Stack**, 88  
 allocazione, 93  
**Stack pointer**, 88  
 contenuto, 89  
**Stallo della pipeline**, 238  
 esempio, 240  
 modalità di inserimento, 270  
**Stazioni di prenotazione**, 293  
**Sticky bit (bit di presenza)**, 190  
**store doubleword (sd)**, 61  
**Supercomputer**, 3  
**Superscalare**, 291  
**Supporto hardware alle procedure**, 86

## T

- Tabella**  
 degli indirizzi di salto (di salto), 85  
 dei simboli, 112  
 della storia dei salti, 276  
 della verità, 222  
 delle pagine, 376, 378  
 delle pagine ombra, 392  
**Tag**, 332  
 dimensione del campo, 355  
**Tecnica iterativa di Newton**, 188  
**Tecnologia**  
 CMOS, 36  
 delle memorie, 325  
 DRAM, 325, 326, 327  
 SRAM, 325  
**Tempo**  
 assoluto, 28  
 di CPU, 345  
 di esecuzione, 25  
 di hit, 323  
 di ricerca, 329  
 di risposta, 25  
 di trasferimento, 330  
 medio di accesso alla memoria (AMAT), 348  
**Tempo di esecuzione della CPU**, 28  
 tempo di CPU sistema, 28  
 tempo di CPU utente, 28  
**Tempo medio di funzionamento prima di un malfunzionamento (MTTF)**, 363  
**Tempo medio di riparazione (MTTR)**, 364  
**Tempo medio tra due guasti (MTBF)**, 364

Temporizzazione  
metodologia di, 211  
sensibile ai fronti, 211, 212

Terabyte, 3

Termine indifferente, 222

Thread, 448

Throughput (larghezza di banda), 25  
e tempo di esecuzione, 26

Touchscreen, 15

Tracce, 328

Traduzione dell'indirizzo, 374

Transistor, 21

Translation lookaside buffer (TLB), 382, 383  
del processore FastMATH Intrinsity, 384  
gestione delle miss e dei page fault, 390

Turbo Mode, 35

## U

UMA (multiprocessore ad accesso uniforme), 452

Underflow, 169

Unicode, 98, 99

alfabeti rappresentati con, 99

Unità aritmetico-logica (ALU), 153  
unità di controllo della, 221

Unità di consegna, 292

Unità di controllo, 16, 209

associata a unità di elaborazione con pipeline, 244  
completamento della, 228  
della pipeline, 257

Unità di elaborazione dati (datapath), 16, 206

Unità di elaborazione, 218, 219

con pipeline, 244

durante l'esecuzione di un'istruzione branch equal,  
229

durante l'esecuzione di un'istruzione di load, 228

durante l'esecuzione di un'istruzione tipo R, 227

elementi della, 213

funzionamento, 226

progettazione, 218

realizzazione, 213

schema, 225

Unità di elaborazione grafica (GPU), 455

Unità di misura

binaria, 4  
decimale, 4

Utilizzo meno recente (LRU), 354

## V

Very Long Instruction Word (VLIW), 287

Vettore delle trappole del supervisore (STVEC), 285

Vincolo di allineamento, 62

Virgola mobile, 168

addizione binaria in, 177  
addizione in, 174  
istruzioni nel RISC-V, 181  
moltiplicazione in, 178  
numeri in, 168  
operandi nel RISC-V, 183  
rappresentazione in, 169  
standard IEEE 754, 170-171

VLIW (*Very Long Instruction Word*), 287

## W

Wafer, 22

Warehouse Scale Computers (WSC), 5

x86

codifica delle istruzioni, 138  
evoluzione, 131  
formati tipici delle istruzioni, 138  
istruzioni dell'architettura, 131  
istruzioni dell'architettura, 137  
istruzioni in virgola mobile delle estensioni  
SSE/SSE2, 193

operazioni su numeri interi, 134

operazioni tipiche, 137

registri e modalità di indirizzamento, 134

WB (write-back), 245, 249, 250, 341

Windows Microsoft, 65

Wireless, connessione, 5

Write-through, 340

## X

XOR, 80

# Manuale di riferimento RISC-V

## ISTRUZIONI INTERE DI BASE RV64I – In ordine alfabetico

Mnemonico	FMT Nome	Descrizione (in Verilog)
add, addw	R Somma (parola)	$R[rd] = R[rs1] + R[rs2]$
addi, addiw	I Somma immediata (parola)	$R[rd] = R[rs1] + cost$
and	R AND	$R[rd] = R[rs1] \& R[rs2]$
andi	I AND immediato	$R[rd] = R[rs1] \& cost$
auipc	I Somma la parte più significativa al PC	$R[rd] = PC + [cost, 12b:0]$
beq	SB Salta se uguale	$se(R[rs1] == R[rs2]) PC = PC + [cost, 1b:0]$
bge	SB Salta se maggiore uguale	$se(R[rs1] \geq R[rs2]) PC = PC + [cost, 1b:0]$
bgeu	SB Salta se maggiore uguale senza segno	$se(R[rs1] \geq R[rs2]) PC = PC + [cost, 1b:0]$
blt	SB Salta se minore	$se(R[rs1] < R[rs2]) PC = PC + [cost, 1b:0]$
bltu	SB Salta se minore senza segno	$se(R[rs1] < R[rs2]) PC = PC + [cost, 1b:0]$
bne	SB Salta se diverso	$se(R[rs1] != R[rs2]) PC = PC + [cost, 1b:0]$
csrrc	I Leggi&Cancella Reg Cont/Stat	$Rd = CSR_CSR = CSR & -R[rs1]$
csrrci	I Leggi&Cancella Reg Cont/Stat cost	$Rd = CSR_CSR = CSR & -cost$
csrrs	I Leggi&Imposta Reg Cont/Stat	$Rd = CSR_CSR = CSR! -R[rs1]$
csrrsi	I Leggi&Imposta Reg Cont/Stat cost	$Rd = CSR_CSR = CSR! -cost$
csrrw	I Leggi&Scrivi Reg Cont/Stat	$Rd = CSR_CSR = R[rs1]$
csrrwi	I Leggi&Scrivi Reg Cont/Stat cost	$Rd = CSR_CSR = cost$
ebreak	I Interrompi l'ambiente	Trasferisce controllo al debugger
ecall	I Chiamata all'ambiente	Trasferisce controllo al sistema operativo
fence	I Sincronizza thread	Sincronizza i thread
fence.i	I Sincronizza istruz e dati	Sincronizza le scritture di un flusso di istruz
jal	UI Salta e unisci ("jump and link")	$R[rd] = PC + 4:PC = PC + [cost, 1b:0]$
jalr	I Salta e unisci a registro	$R[rd] = PC + 4:PC = R[rs1] + cost$
lb	I Carica un byte	$R[rd] = [56'bM][7], M[R[rs1]+cost][7:0]$
lbu	I Carica un byte senza segno	$R[rd] = [56'b0M][7], M[R[rs1]+cost][7:0]$
ld	I Carica una parola doppia	$R[rd] = [56'bM], M[R[rs1]+cost][63:0]$
lh	I Carica una mezza parola	$R[rd] = [48'bM][15], M[R[rs1]+cost][15:0]$
lhu	I Carica una mezza parola senza segno	$R[rd] = [48'b0M], M[R[rs1]+cost][15:0]$
lui	U Carica la mezza parola superiore	$R[rd] = [32'b1'cost<31>, cost12'b0]$
lw	I Carica una parola	$R[rd] = [32'bM][7], M[R[rs1]+cost][15:0]$
lwu	I Carica una parola senza segno	$R[rd] = [32'b0M], M[R[rs1]+imm][31:0]$
or	R OR	$R[rd] = R[rs1]   R[rt]$
ori	I OR immediato	$R[rt] = R[rs1]   Cost$
sb	S Memorizza byte	$M[R[rs1]+cost][7:0] = R[rs2][7:0]$
sd	S Memorizza parola doppia	$M[R[rs1]+cost][63:0] = R[rs2][63:0]$
sh	S Memorizza mezza parola	$M[R[rs1]+cost][15:0] = R[rs2][15:0]$
sll, sllw	R Scorrimento a sinistra (parola)	$R[rd] = R[rs1] << R[rs2]$
slli, slliw	I Scorrimento a sinistra immediato (parola)	$R[rd] = R[rs1] << cost$
slt	R Imposta se minore di	$R[rd] = (R[rs1] < R[rs2]) ? 1 : 0$
slti	I Imposta se minore di (costante)	$R[rd] = (R[rs1] < cost) ? 1 : 0$
sltiu	I Imposta se minore di (costante, senza segno)	$R[rd] = (R[rs1] < cost) ? 1 : 0$
sltu	R Imposta se minore di (senza segno)	$R[rd] = (R[rs1] < R[rs2]) ? 1 : 0$
sra, sraw	R Scorrimento arith a destra (parola)	$R[rd] = R[rs1] >> R[rs2]$
srai, srawi	I Scorrimento arith a destra immediato (parola)	$R[rd] = R[rs1] >> cost$
srl, srliw	R Scorrimento a destra (parola)	$R[rd] = R[rs1] >> R[rs2]$
srls, srliw	I Scorrimento a destra immediato (parola)	$R[rd] = R[rs1] >> cost$
sub, subw	R Sottrazione (parola)	$R[rd] = R[rs1] - R[rs2]$
sw	S Memorizza parola	$M[R[rs1]+cost] = R[rs2]$
xor	R XOR	$R[rd] = R[rs1] ^ R[rs2]$
xori	I XOR immediato	$R[rd] = R[rs1] ^ cost$

### Note:

- La versione parola opera solo sui 32 bit più a destra significativi di un registro a 64 bit
- L'operazione viene svolta con interi senza segno e non in complemento a 2
- Il bit meno significativo di una jal viene impostato a 0
- Con segno. L'istruzione estende il bit di segno fino a riempire tutti i 64 bit del registro
- Replica il bit di segno fino a riempire tutti i bit del risultato in uno scorrimento a destra
- Moltiplicazione di un operando senza segno con un operando dotato di segno
- La versione in singola precisione utilizza 32 bit più a destra di un registro F a 64 bit
- La classificazione scrive una maschera di 10 bit che mostra quali proprietà siano vere (per es. -Inf, -0, +0, +Inf, denorm)
- Operazione atomica su memoria; non si può sovrapporre nulla tra la lettura e la scrittura della cella di memoria
- Il bit di segno viene esteso nel RISC-V

## NUCLEO DELLE ISTRUZIONI DI MOLTIPLICAZIONE

### Estensione della moltiplicazione RV64M

Mnemonico	FMT Nome	Descrizione (in Verilog)	Note
mul, mulw	R Moltiplicazione (parola)	$R[rd] = R[rs1] * R[rs2](63:0)$	1)
mulh	R Moltiplicazione mezza parola superiore	$R[rd] = R[rs1] * R[rs2](127:64)$	
mulhsu	R Moltiplicaz mezza parola superiore con/senza segno	$R[rd] = R[rs1] * R[rs2](127:64)$	6)
mulhu	R Moltiplicaz mezza parola superiore senza segno	$R[rd] = R[rs1] * R[rs2](127:64)$	2)
div, divw	R Divisione (parola)	$R[rd] = R[rs1] / R[rs2]$	1)
divu	R Divisione senza segno	$R[rd] = R[rs1] / R[rs2]$	2)
rem, remw	R Resto (parola)	$R[rd] = R[rs1] \% R[rs2]$	1)
remu, remuw	R Resto senza segno (parola)	$R[rd] = R[rs1] \% R[rs2]$	1,2)

### Estensioni virgola mobile RV64F e RV64M

2)	Mnemonico	FMT Nome	Descrizione (in Verilog)	Note
2)	fld, fldw	I Carica (parola)	$F[rd] = M[R[rs1]+cost]$	1)
2)	fsd, fsw	S Memorizza (parola)	$M[R[rs1]+cost] = F[rd]$	1)
2)	fadd.s, fadd.d	R Somma	$F[rd] = F[rs1] + F[rs2]$	7)
2)	fsub.s, fsub.d	R Sottrazione	$F[rd] = F[rs1] - F[rs2]$	7)
2)	fmul.s, fmul.d	R Moltiplicazione	$F[rd] = F[rs1] * F[rs2]$	7)
2)	fdiv.s, fdiv.d	R Divisione	$F[rd] = F[rs1] / F[rs2]$	7)
2)	fsqrts, fsqrtd	R Radice quadrata	$F[rd] = sqrt(F[rs1])$	7)
2)	fmadds, fmadd.d	R Moltiplicazione-somma	$F[rd] = F[rs1] * F[rs2] + F[rs3]$	7)
2)	fmsubs, fmsub.d	R Moltiplicazione-sottrazione	$F[rd] = F[rs1] * F[rs2] - F[rs3]$	7)
2)	fmnsub.s, fmnsub.d	R Moltiplicazione negativa-sottrazione	$F[rd] = -(F[rs1] * F[rs2] - F[rs3])$	7)
3)	fsignj.s, fsignj.d	R Provenienza segno	$F[rd] = F[rs2]<63>, F[rs1]<62:0>$	7)
3)	fsignjn.s, fsignjn.d	R Provenienza segno negativa	$F[rd] = -F[rs2]<63>, F[rs1]<62:0>$	7)
3)	fsignjk.s, fsignjk.d	R Provenienza segno con xor	$F[rd] = (F[rs1]<63> * F[rs2]<63>) & F[rs1]<62:0>$	7)
4)	fmin.s, fmin.d	R Minimo	$F[rd] = (F[rs1] < F[rs2]) ? F[rs1] : F[rs2]$	7)
4)	fmax.s, fmax.d	R Massimo	$F[rd] = (F[rs1] < F[rs2]) ? F[rs2] : F[rs1]$	7)
4)	feq.s, feq.d	R Confronta se uguale	$R[rd] = (F[rs1] == F[rs2]) ? 1 : 0$	7)
4)	flt.s, flt.d	R Confronta se minore	$R[rd] = (F[rs1] < F[rs2]) ? 1 : 0$	7)
4)	fel.s, fel.d	R Confronta se minore o uguale	$R[rd] = (F[rs1] \leq F[rs2]) ? 1 : 0$	7)
4)	fclass.s, fclass.d	R Classifica il tipo	$R[rd] = class(F[rs1])$	7,8)
4)	fmv.s.x, fmv.d.x	R Sposta da un intero	$F[rd] = R[rs1]$	7)
4)	fmv.x.s, fmv.x.d	R Sposta in un intero	$R[rd] = F[rs1]$	7)
4)	fcvt.s.d	R Converti da DP a SP	$F[rd] = single(F[rs1])$	7)
4)	fcvt.d.s	R Converti da SP a DP	$F[rd] = double(F[rs1])$	7)
4)	fcvt.s.w, fcvt.d.w	R Converti da un intero su 32 bit	$F[rd] = float(F[rs1])[31:0]$	7)
4)	fcvt.s.l, fcvt.d.l	R Converti da un intero su 64 bit	$F[rd] = float(F[rs1])[63:0]$	7)
4)	fcvt.s.wu, fcvt.d.wu	R Converti da un intero su 32 bit senza segno	$F[rd] = float(F[rs1])[31:0]$	2,7)
4)	fcvt.s.lu, fcvt.d.lu	R Converti da un intero su 64 bit senza segno	$F[rd] = float(F[rs1])[63:0]$	2,7)
4)	fcvt.s.w, fcvt.d.w	R Converti in un intero su 32 bit	$R[rd] = integer(F[rs1])[31:0]$	7)
4)	fcvt.s.l, fcvt.d.l	R Converti in un intero su 64 bit	$R[rd] = integer(F[rs1])[63:0]$	7)
4)	fcvt.s.wu, fcvt.d.wu	R Converti in un intero su 32 bit senza segno	$R[rd] = integer(F[rs1])[31:0]$	2,7)
4)	fcvt.s.lu, fcvt.d.lu	R Converti in un intero su 64 bit senza segno	$R[rd] = integer(F[rs1])[63:0]$	2,7)
4)	fcvt.s.w, fcvt.d.w	R Converti in un intero su 32 bit	$R[rd] = integer(F[rs1])[31:0]$	7)
4)	fcvt.s.l, fcvt.d.l	R Converti in un intero su 64 bit	$R[rd] = integer(F[rs1])[63:0]$	7)
4)	fcvt.s.wu, fcvt.d.wu	R Converti in un intero su 32 bit senza segno	$R[rd] = integer(F[rs1])[31:0]$	2,7)
4)	fcvt.s.lu, fcvt.d.lu	R Converti in un intero su 64 bit senza segno	$R[rd] = integer(F[rs1])[63:0]$	2,7)

### Estensione atomica RV64A

Mnemonico	FMT Nome	Descrizione (in Verilog)	Note	
amoadd.w, amoadd.d	R Somma	$R[rd] = M[R[rs1]], M[R[rs1]] = M[R[rs1]] + R[rs2];$	9)	
amoand.w, amoand.d	R AND	$R[rd] = M[R[rs1]], M[R[rs1]] = M[R[rs1]] \& R[rs2];$	9)	
amamax.w, amamax.d	R Massimo	$R[rd] = M[R[rs1]], if R[rs2] > M[R[rs1]] M[R[rs1]] = R[rs2];$	9)	
amomin.w, amomin.d	R Massimo senza segno	$R[rd] = M[R[rs1]], if R[rs2] > M[R[rs1]] M[R[rs1]] = R[rs2];$	2,9)	
amomin.w, amomin.d	R Minimo	$R[rd] = M[R[rs1]], if R[rs2] < M[R[rs1]] M[R[rs1]] = R[rs2];$	9)	
amomin.w, amomin.d	R Minimo senza segno	$R[rd] = M[R[rs1]], if R[rs2] < M[R[rs1]] M[R[rs1]] = R[rs2];$	2,9)	
1)	amoor.w, amoor.d	R OR	$R[rd] = M[R[rs1]], M[R[rs1]] = M[R[rs1]]   R[rs2];$	9)
1)	amoswap.w, amoswambio.d	R Scambio	$R[rd] = M[R[rs1]], M[R[rs1]] = R[rs2];$	2,9)
1)	amoxor.w, amoxor.d	R XOR	$R[rd] = M[R[rs1]], M[R[rs1]] = M[R[rs1]] ^ R[rs2];$	9)
1)	lr.w, lr.d	R Carica riservata	$R[rd] = M[R[rs1]], riserva M[R[rs1]]$	9)
sc.w, sc.d	R Memorizza condizionato	se riservato, $M[rs1] = M[rs2], R[rd] = 0, altrimenti R[rd] = 1$	9)	

## FORMATO DELLE ISTRUZIONI CORE

31	27	26	25	24	20	19	15	14	12	11	7	6	0		
R					funz7		rs2		rs1		funz3		rd	codop	
I					cost[11:0]				rs1		funz3		rd	codop	
S					cost[11:5]		rs2		rs1		funz3		cost[4:0]	codop	
SB					cost[1210:5]		rs2		rs1		funz3		cost[4:1][11]	codop	
U												cost[31:12]		rd	codop
W												cost[2010:11][11:9:12]		rd	codop

### PSEUDOISTRUZIONI

Mnemonico	Nome	Descrizione (In Verilog)	Tradotta in
beqz	Salta se uguale a zero	if R[rs1] == 0, PC=PC+(cost,1b'0)	beq
bnez	Salta se non zero	if R[rs1] != 0, PC=PC+(cost,1b'0)	bne
fabs.s,fabs.d	Valore assoluto	F[rd] = [F(rs1)<0] ? -F(rs1) : F(rs1)	fsgnx
fmv.s,fmv.d	Spostamento VM	F[rd] = F(rs1)	fsgnj
fneg.s,fneg.d	Spostamento negativo VM	F[rd] = -F(rs1)	fsgnjn
j	Salto incondizionato	PC = (cost,1b'0)	jal
jr	Salto a registro	PC = R(rs1)	jalr
la	Carica indirizzo	R[rd] = indirizzo	auipc
li	Carica immediato	R[rd] = costante	addi
mv	Sposta	R[rd] = R(rs1)	addi
neg	Negazione	R[rd] = -R(rs1)	sub
nop	Non operazione	R[rd] = R[0]	addi
not	Negazione logica	R[rd] = ~R(rs1)	xori
ret	Ritorna	PC = R[1]	jalr
seq	Imposta se zero	R[rd] = [[rs1] == 0] ? 1 : 0	situ
sne	Imposta se non zero	R[rd] = [[rs1] != 0] ? 1 : 0	situ

(3)

### NOME DEI REGISTRI, UTILIZZO, CONVENZIONI DI CHIAMATA

Registro	Nome	Utilizzo	Chi lo salva
x0	zero	La costante 0	N.A.
x1	ra	Indirizzo di ritorno	Chiamante
x2	sp	Puntatore a stack	Chiamato
x3	gp	Puntatore globale	---
x4	tp	Puntatore a thread	---
x5-x7	t0-t2	Temporanei	Chiamante
x8	s0 / _fp	Salvato/puntatore a frame	Chiamato
x9	s1	Salvato	Chiamato
x10-x11	a0-a1	Argomenti di funzione/valori restituiti	Chiamante
x12-x17	a2-a7	Argomenti di funzione	Chiamante
x18-x27	s2-s11	Registri salvati	Chiamato
x28-x31	t3-t6	Temporanei	Chiamante
f0-f7	ft0-ft7	Temporanei VM	Chiamante
f8-f9	fs0-fs1	Salvati VM	Chiamato
f10-f11	fa0-fa1	Argomenti di funzione/valori restituiti VM	Chiamante
f12-f17	fa2-fa7	Argomenti di funzione VM	Chiamante
f18-f27	fs2-fs11	Registri salvati VM	Chiamato
f27-f31	ft8-ft11	R[rd] = R(rs1) + R(rs2)	Chiamante

(4)

### ISTRUZIONI ORDINATE NUMERICAMENTE PER CODICE OPERATIVO

Mnemonico	FMT	Codice operativo	Funz3	Funz7 o cost	Eсадicimale
lb	I	0000011	000		03/0
lh	I	0000011	001		03/1
lw	I	0000011	010		03/2
ld	I	0000011	011		03/3
lbu	I	0000011	100		03/4
lhu	I	0000011	101		03/5
lwu	I	0000011	110		03/6
fence	I	0001111	000		0F/0
fence.w	I	0001111	001		0F/1
add!	I	0010011	000		13/0
slli	I	0010011	001	0000000	13/1/00
slti	I	0010011	010		13/2
sltiu	I	0010011	011		13/3
xci	I	0010011	100		13/4
srl1	I	0010011	101	0000000	13/5/00
sra1	I	0010011	101	0100000	13/5/20
ori	I	0010011	110		13/6
andi	I	0010011	111		13/7
auipc	U	0010111			17
addiw	I	0011011	000		1B/0
slliw	I	0011011	001	0000000	1B/1/00
srliw	I	0011011	101	0000000	1B/5/00
straw	I	0011011	101	0100000	1B/5/20
sb	S	0100011	000		23/0
sh	S	0100011	001		23/1
sw	S	0100011	010		23/2
sd	S	0100011	011		23/3
add	R	0110011	000	0000000	33/0/00
sub	R	0110011	000	0100000	33/0/20
sli	R	0110011	001	0000000	33/1/00
slt	R	0110011	010	0000000	33/2/00
situ	R	0110011	011	0000000	33/3/00
xor	R	0110011	100	0000000	33/4/00
srl	R	0110011	101	0000000	33/5/00
sra	R	0110011	101	0100000	33/5/20
cr	R	0110011	110	0000000	33/6/00
and	R	0110011	111	0000000	33/7/00
lui	U	0110111			37
addw	R	0110111	000	0000000	3B/0/00
subw	R	0110111	000	0100000	3B/0/20
sllw	R	0110111	001	0000000	3B/1/00
srlw	R	0110111	101	0000000	3B/5/00
straw	R	0110111	101	0100000	3B/5/20
beq	SB	1100011	000		63/0
bne	SB	1100011	001		63/1
bit	SB	1100011	100		63/4
bge	SB	1100011	101		63/5
bitu	SB	1100011	110		63/6
bgeu	SB	1100011	111		63/7
jarl	I	1100111	000		67/0
jal	W	1101111			6F
ecall	I	1110011	000	0000000000000000	73/0/000
ebreak	I	1110011	000	0000000000000001	73/0/001
CSRRW	I	1110011	001		73/1
CSRWS	I	1110011	010		73/2
CSRRC	I	1110011	011		73/3
CSRRI	I	1110011	101		73/5
CSRRII	I	1110011	110		73/6
CSRRCI	I	1110011	111		73/7

### STANDARD VIRGOLA MOBILE IEEE-754

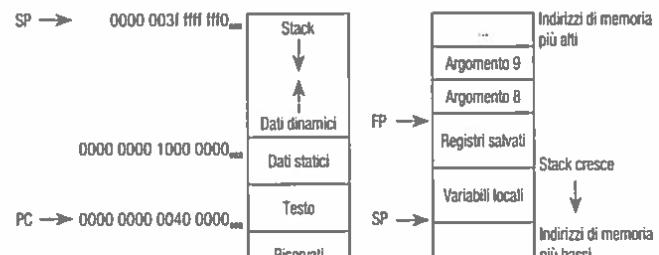
$$(-1)^{\text{S}} \times (1 + \text{Frazione}) \times 2^{(\text{Esponente} - \text{Polarizzazione})}$$

dove la polarizzazione per la mezza precisione è pari a 15, per la singola precisione è pari a 127, per la doppia precisione è pari a 1023, e per la precisione quadrupla è pari a 16383

### FORMATI RELATIVI A MEZZA, SINGOLA, DOPPIA E QUADRUPLA PRECISIONE

S	Esponente	Frazione	
15	14	10 9	0
S	Esponente	Frazione	
31	30	23 22	0
S	Esponente	Frazione	
63	62	52 51	0
S	Esponente	Frazione	
127	126	112 111	0

### ALLOCAZIONE MEMORIA



### PREFISSI DELLE DIMENSIONI E SIMBOLI

Dimensione	Prefisso	Simbolo	Dimensione	Prefisso	Simbolo
$10^3$	Chilo-	K	$2^{10}$	Kibi-	Ki
$10^6$	Mega-	M	$2^{20}$	Mebi-	Mi
$10^9$	Giga-	G	$2^{30}$	Gibi-	Gi
$10^{12}$	Tera-	T	$2^{40}$	Tebi-	Ti
$10^{15}$	Peta-	P	$2^{50}$	Pebi-	Pi
$10^{18}$	Esa-	E	$2^{60}$	Exbi-	Ei
$10^{21}$	Zetta-	Z	$2^{70}$	Zebi-	Zi
$10^{24}$	Yotta-	Y	$2^{80}$	Yobi-	Yi
$10^{-3}$	milli-	m	$10^{-11}$	femto-	f
$10^{-6}$	micro-	μ	$10^{-14}$	atto-	a
$10^{-9}$	nano-	n	$10^{-21}$	zepto-	z
$10^{-12}$	pico-	p	$10^{-34}$	yocto-	y