

# Appunti di Architetture Avanzate dei Calcolatori

Matteo Gianello

3 luglio 2014

Quest'opera è stata rilasciata con licenza Creative Commons Attribuzione - Non commerciale - Condividi allo stesso modo 3.0 Unported. Per leggere una copia della licenza visita il sito web  
<http://creativecommons.org/licenses/by-nc-sa/3.0/deed.it> .

# Indice

<b>1 Pipelining</b>	<b>4</b>
1.1 Concetti base . . . . .	4
1.1.1 Reduced Instruction Set nei processori MIPS . . . . .	4
1.1.2 Esecuzione delle istruzioni . . . . .	7
1.1.3 Implementazione base di un MIPS . . . . .	8
1.2 Pipelining . . . . .	11
1.2.1 Implementazione di una Pipeline . . . . .	13
1.3 Il problema del "Hazard" . . . . .	16
1.3.1 Data Hazard . . . . .	16
1.3.2 Altri tipi di Data Hazard . . . . .	18
1.4 Analisi delle performance . . . . .	19
<b>2 Tecniche di predizione dei salti</b>	<b>22</b>
2.1 Il problema del Control Hazard . . . . .	22
2.2 Tecniche di predizione dei salti . . . . .	23
2.2.1 Tecniche di predizione statiche . . . . .	25
2.2.2 Tecniche di predizione dinamiche . . . . .	27
2.3 Speculazione . . . . .	31
<b>3 Instruction Level Parallelism</b>	<b>34</b>
3.1 Tipi di Hazards sui dati . . . . .	34
3.2 Parallelismo a livello di istruzione . . . . .	35
3.2.1 ILP in pratica . . . . .	37
3.2.2 Esecuzione super-scalare e VLIW . . . . .	38
3.3 Scoreboard . . . . .	39
3.4 Algoritmo di Tomasulo . . . . .	42
3.4.1 Gli stadi dell'algoritmo di Tomasulo . . . . .	43
3.4.2 Alcuni dettagli . . . . .	44
3.4.3 Tomasulo in pratica . . . . .	44
3.4.4 Tomasulo vs Scoreboard . . . . .	44
3.5 Register Renaming . . . . .	45
3.5.1 Renaming implicito . . . . .	45
3.5.2 Renaming esplicito . . . . .	46
<b>4 Static Multiple-Issue Processor: Approccio VLIM</b>	<b>49</b>
4.1 Processori VLIW . . . . .	49
4.1.1 Alcuni esempi . . . . .	51
4.2 Code scheduling VLIW . . . . .	53
4.3 List-based scheduling . . . . .	54
4.4 Global e Local scheduling . . . . .	54
<b>5 Reorder Buffer</b>	<b>60</b>
5.1 Struttura del reorder buffer . . . . .	61
<b>6 Multithreading</b>	<b>63</b>
6.1 Processori embedded . . . . .	64
6.2 Multithreading e Multiprocessing . . . . .	65

<b>7 Gerarchie di memoria</b>	<b>68</b>
7.1 Caches . . . . .	68
7.1.1 Struttura e funzionamento di una cache . . . . .	69
7.1.2 Analisi delle performance . . . . .	76
7.2 Incremento delle performance . . . . .	77
7.2.1 Ridurre il miss rate . . . . .	78
7.2.2 Ridurre il miss penalty . . . . .	80
7.2.3 Riduzione dello hit time . . . . .	82
7.3 Virtual Memory . . . . .	83
7.3.1 Memory management unit . . . . .	84
<b>8 Mutua esclusione e sincronizzazione</b>	<b>86</b>
8.1 Mutua esclusione . . . . .	86
8.2 Dalle primitive di sincronizzazione ai metodi di sincronizzazione . . . . .	87
<b>9 Introduzione ai multiprocessori</b>	<b>91</b>
9.1 Gestione della memoria . . . . .	93
9.2 Il problema della coerenza della cache . . . . .	96
<b>10 Analisi delle performance</b>	<b>101</b>
10.1 Legge di Amdahl . . . . .	101
10.2 Analisi delle performance in un processore pipelined . . . . .	101
10.3 Analisi delle performance nelle gerarchie di memoria . . . . .	103

# Introduzione

Il corso di Architetture Avanzate dei Calcolatori si prefigge lo scopo di fornire una vista sulle più recenti architetture avanzate dei calcolatori, introducendo i meccanismi base delle microarchitetture che si possono ritrovare nei moderni microprocessori, ed infine fornire le ragioni dietro alle tecniche adottate nelle architetture dei computer.

## 1 Pipelining

### 1.1 Concetti base

Definiamo prima di tutto quali sono le principali caratteristiche dell'architettura MIPS partendo dalla definizione delle istruzioni utilizzate da questi calcolatori calcolatori. Esistono due tipi di istruzioni le CISC e le RISC; le istruzioni di tipo **CISC** (*Complex Instruction Set Computer*) sono un set di istruzioni esteso che permettono ai processori di eseguire operazioni molto complesse come somme tra operandi caricati direttamente dalla memoria centrale, le istruzioni **RISC** (*Reduced Instruction Set Computer*), invece, sono istruzioni semplici che possono essere eseguite in un unico ciclo di clock e ottimizzate per le performance sulle CPU CISC. Le architetture RISC sono anche dette architetture di tipo *LOAD/STORE*, in quanto le istruzioni non accedono direttamente ai dati in memoria ma accedono ai dati contenuti in registri del processore, solo due istruzioni permettono l'accesso alla memoria principale, queste due istruzioni sono:

- **load** che carica i dati dalla memoria ai registri.
- **store** che sposta i dati dai registri alla memoria.

Un'altra caratteristica fondamentale per le architetture RISC è l'utilizzo della *Pipeline* una tecnica di ottimizzazione basata sull'esecuzione sovrapposta di molteplici istruzioni sequenziali.

#### 1.1.1 Reduced Instruction Set nei processori MIPS

Vediamo ora quali sono le diverse istruzioni di tipo RISC e come sono rappresentate nel calcolatore

**Istruzioni ALU** Vediamo innanzitutto le istruzioni di somma ovvero quelle eseguite dalla ALU. Queste possono essere di due tipi, Una istruzione di tipo *R-Format* è un'istruzione che prende in considerazione due registri, tale tipo di operazione è applicabile solo alle istruzioni **add** di tipo registro-registro. Esistono poi le **addi** che viene chiamata anche somma *immediata* in quanto avviene tra un registro ed un valore costante. Tale tipo di istruzione è del tipo *I-Format*. Lo pseudo codice assembly delle due istruzioni è mostrato qui di seguito, inoltre possiamo anche vedere come vengono svolte le operazioni.

```
add $s1, $s2, $s3      # $s1 <- $s2 + $s3
addi $s1, $s2, 4        # $s1 <- $s2 + 4
```

In Figura 1 vediamo come è suddivisa un'istruzione di tipo *R-Format* I diversi campi indicano rispettivamente:

**op** identifica il tipo di istruzione ALU da eseguire

**rs** indica il registro nel quale è contenuto il primo operando

op	rs	rt	rd	shamt	funct
6 bit	5 bit	5 bit	5 bit	5 bit	6 bit

Figura 1: Esempio di istruzione ALU di tipo R-Format

**rt** indica il registro nel quale è contenuto il secondo operando

**rd** indica il registro di destinazione

**shamt** sta ad indicare i bit di shift amount

**funct** identifica i diversi tipi di istruzione

op	rs	rt	immediate
6 bit	5 bit	5 bit	16 bit

Figura 2: Divisione delle informazioni in un registro di un istruzione ALU di tipo diretto

Nella Figura 2 vediamo invece la suddivisione di un registro nel caso di un'operazione di ALU immediata. La suddivisione dei diversi campi è la seguente:

**op** identifica l'istruzione di tipo immediato

**rs** indica il registro nel quale è posizionato il primo operando

**rt** indica il registro di destinazione del risultato

**immediate** contiene il valore per l'operazione immediata nel range  $-2^{15}$  e  $+2^{15} - 1$

**Istruzioni LOAD/STORE** Le istruzioni di *load* e di *store* sono quelle che permettono di caricare e scaricare i valori dai registri della CPU alla memoria centrale e viceversa. Un esempio di codice assembly per le istruzioni load e store.

```
lw $s1, offset($s2) # $s1 <- M[$s2 + offset]
sw $s1, offset($s2) # M[$s2 + offset] <- $s1
```

In Figura 3 vediamo come sono strutturate le istruzioni di *load* e di *store*; come possiamo vedere anche queste istruzioni sono nel formato *I-Format*

La suddivisione del registro è la seguente:

op	rs	rt	offset
6 bit	5 bit	5 bit	16 bit

Figura 3: Struttura di un'istruzione tipo load/store

**op** identifica l'istruzione di tipo load o store

**rs** identifica il registro base

**rt** identifica il registro sorgente o destinazione per i dati delle operazioni di store o di load da o per la memoria

**offset** da sommare all'indirizzo contenuto in **rs** per calcolare l'indirizzo di memoria

**Istruzioni di salto** Per quanto riguarda le istruzioni di salto possiamo suddividerle in due categorie, i salti *condizionati*, che richiedono la verifica di una determinata condizione per decidere se effettuare un salto; oppure istruzioni di salto *incondizionato* che effettuano il salto sempre. Lo pseudo assembly di un'istruzione di salto condizionato è:

```
beq $s1, $s2, L1 #go to L1 if ($s1 == $s2)
bne $s1, $s2, L1 #go to L1 if ($s1 != $s2)
```

Le istruzioni di salto condizionato sono nel formato *I-Format*. La suddivisione del registro per questo tipo di istruzione è mostrata in Figura 4

op	rs	rt	address
6 bit	5 bit	5 bit	16 bit

Figura 4: Struttura di un'istruzione tipo branch condizionato

**op** identifica l'istruzione di tipo branch condizionale

**rs** identifica il primo registro da comparare

**rd** identifica il secondo registro da comparare

**address** identifica l'offset rispetto al PC che corrisponde all'indirizzo dell'etichetta

Per quanto riguarda il salto incondizionato il funzionamento è molto più semplice, quando si raggiunge l'istruzione di salto il PC punta direttamente all'istruzione indicata dall'etichetta. Tale semplicità si rispecchia nella struttura dell'istruzione che in questo caso è di tipo *J-Format*. Lo pseudo-codice assembly dell'istruzione di salto è:

```
j L1 #go to L1
jr $s1 #go to add. contenuto in $1
```

Un esempio di struttura di istruzione di salto è rappresentato in Figura 5 dove i valori indicano

op	address
6 bit	26 bit

Figura 5: Divisione delle informazioni in un registro di un istruzione tipo branch incondizionato rispettivamente

**op** identifica il tipo di istruzione

**address** identifica l'indirizzo della prossima istruzione da eseguire

Ricapitolando possiamo suddividere le istruzioni in tre categorie in base alla loro struttura e a come vengono eseguite:

- Tipo R (*Registro*)
  - Istruzione ALU
- Tipo I (*Immediate*)
  - ALU immediate
  - Istruzioni Load/Store
  - Istruzioni di salto condizionato
- Tipo J (*Jump*)
  - Istruzioni di salto incondizionato

Il perché di questa divisione lo si capisce molto facilmente dallo schema in Figura 6 nel quale vengono confrontati le diverse suddivisioni degli Instruction Register.

	31	26	25	21	20	16	15	11	10	6	5	0
R	op	rs		rt		rd		shamt		funct		
I	op	rs		rt						offset/immediate		
J	op					address						

Figura 6: Divisione dei registri nei diversi casi di operazione

### 1.1.2 Esecuzione delle istruzioni

Vediamo ora come possono essere implementate le diverse istruzioni in ambiente MIPS. Tutte le istruzioni possono essere implementate suddividendo l'esecuzione in cinque fasi distinte:

1. **Instruction Fetch Cycle:** durante questo ciclo viene inviato il contenuto del *Program Counter* all'*Instruction Memory* e viene prelevata l'istruzione corrispondente. Successivamente viene aggiornato il PC perché punti alla prossima istruzione aggiungendo 4 al valore attuale (le istruzioni sono di 4 bytes)
2. **Instruction Decode and Register Read Cycle:** in questo ciclo si decodifica l'istruzione corrente e si leggono dal *Register File* i registri necessari corrispondenti ai registri specificati nei campi dell'istruzione. Si fa inoltre l'estensione del segno nel caso sia necessario.
3. **Execution Cycle:** In questo ciclo la ALU effettua le operazioni sugli elementi che sono stati preparati nel ciclo precedente. In base alle istruzioni la ALU esegue le seguenti operazioni
  - Istruzione ALU registro-registro: la ALU esegue le operazioni sugli operandi che ha letto dal *Register File*

- Istruzioni ALU immediate: la ALU esegue l'operazione specificate sul primo operando letto dal *Register File* e sul operando immediato al quale è stata applicata l'estensione di segno.
- Istruzioni Load/Store la ALU aggiunge all'indirizzo base l'offset per calcolare l'indirizzo effettivo.
- Istruzioni di salto condizionato: la ALU compara i due registri e calcola l'indirizzo target del salto da aggiungere al PC

4. **Memory Access (ME):** durante questo ciclo le istruzioni di *Load* effettuano la lettura dalla memoria usando l'indirizzo effettivo calcolato al ciclo precedente, le istruzioni di *Store* scrivono nella memoria i dati provenienti dal registro, infine, le istruzioni di *Branch* aggiornano il valore del Program Counter con l'indirizzo target calcolato al passo precedente, nel caso in cui la condizione sia verificata.
5. **Write-Back Cycle (WB):** in questo ciclo le istruzioni di *Load* scrivono i dati letti dalla memoria nel registro di destinazione mentre le istruzioni di *ALU* scrivono il risultato delle operazioni nei registri di destinazione.

Come possiamo notare le diverse operazioni non usano sempre tutti i cicli appena descritti ma solitamente (tranne nel caso della *load*) attraversano solo alcune fasi come possiamo vedere dallo schema in Figura 7.

Mentre nella tabella di Figura 8 possiamo vedere le latenze di ogni operazione nel caso di tempo di ciclo uguale ad 1 ns.

#### **ALU Instructions: op \$x, \$y, \$z**

Instr. Fetch & PC Increm.	Read of Source Regs. \$y and \$z	ALU OP (\$y op \$z)	Write Back of Destinat. Reg. \$x
---------------------------	----------------------------------	---------------------	----------------------------------

#### **Load Instructions: lw \$x, offset(\$y)**

Instr. Fetch & PC Increm.	Read of Base Reg. \$y	ALU Op. (\$y+offset)	Read Mem. M(\$y+offset)	Write Back of Destinat. Reg. \$x
---------------------------	-----------------------	----------------------	-------------------------	----------------------------------

#### **Store Instructions: sw \$x, offset(\$y)**

Instr. Fetch & PC Increm.	Read of Base Reg. \$y & Source \$x	ALU Op. (\$y+offset)	Write Mem. M(\$y+offset)
---------------------------	------------------------------------	----------------------	--------------------------

#### **Conditional Branch: beq \$x, \$y, offset**

Instr. Fetch & PC Increm.	Read of Source Regs. \$x and \$y	ALU Op. (\$x-\$y) & (PC+4+offset)	Write PC
---------------------------	----------------------------------	-----------------------------------	----------

Figura 7: Cicli eseguiti da ogni operazione

### 1.1.3 Implementazione base di un MIPS

Vediamo ora come potrebbe essere una semplice implementazione di un *Data Path* in un MIPS. Come notiamo dalla Figura 9 abbiamo che la parte di memoria dedicata alle istruzioni (*Instruction Memory*) è composta da 16 righe e 32 colonne. Ogni riga rappresenta un'istruzione di 32 bit. La colonna 0 rappresenta il bit di segno, le colonne 1-5 rappresentano l'indirizzo immagine, le colonne 6-11 rappresentano i dati immagine, le colonne 12-15 rappresentano i dati di controllo, le colonne 16-23 rappresentano i dati di controllo, le colonne 24-27 rappresentano i dati di controllo, le colonne 28-31 rappresentano i dati di controllo.

Instruction Type	Instruct. Mem.	Register Read	ALU Op.	Data Memory	Write Back	Total Latency
ALU Instr.	2	1	2	0	1	6 ns
Load	2	1	2	2	1	8 ns
Store	2	1	2	2	0	7 ns
Cond. Branch	2	1	2	0	0	5 ns
Jump	2	0	0	0	0	2 ns

Figura 8: Latenza delle diverse operazioni

*tion Memory*) è di sola lettura ed è separata dalla memoria dedicata ai dati (*Data Memory*). Inoltre abbiamo 32 registri organizzati in un *Register File* (RF) con 2 porte di lettura (le due frecce che escono sulla destra) e una porta in scrittura (la freccia in ingresso che punta al campo *Data*). La fase di *Instruction Fetch* invece richiede un adder il quale in uscita si connette al

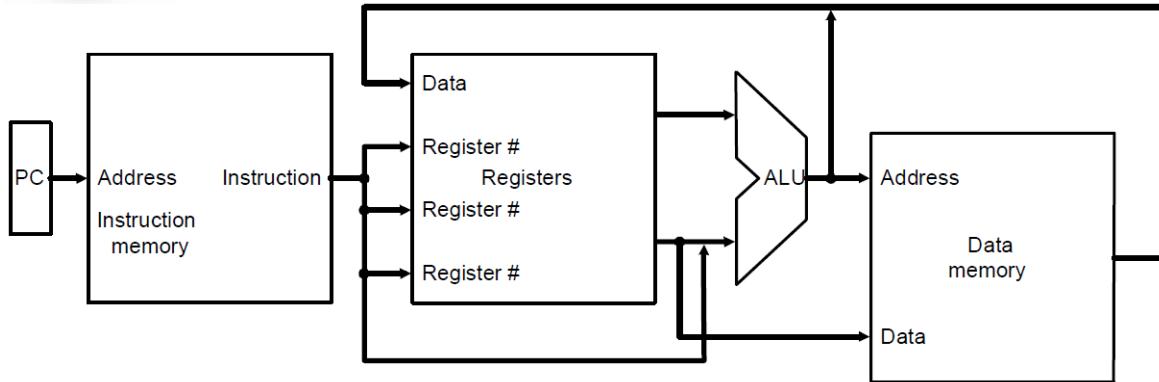


Figura 9: Esempio di implementazione di MIPS

PC mentre come ingressi riceve un valore costante  $4$  mentre all’altro ingresso riceve il valore corrente di del PC come possiamo vedere in Figura 10. Analizziamo ora in breve quale hardware

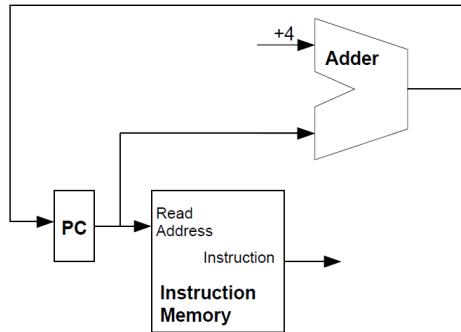


Figura 10: Hardware necessario per realizzare l’Instruction Fetch

è necessario per implementare le diverse operazioni che possono essere eseguite da un MIPS. Partiamo con l’analizzare un istruzione di tipo ALU come vediamo in Figura 11. Dal *Register File* escono due porte che sono connesse ad un’unità ALU la quale ha un’uscita *Result* che si

connette alla porta di scrittura del *RF* e un'uscita *Zero* per indicare eventuali anomalie. Inoltre la *ALU* ha un ingresso *OP* che serve a selezionare il tipo di operazione. Per quanto riguarda

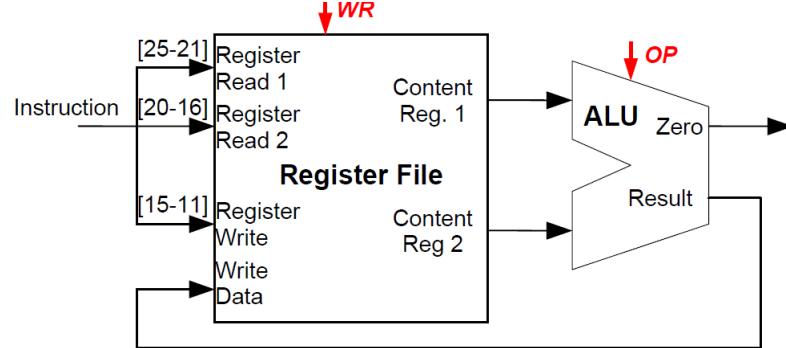


Figura 11: Hardware che implementa un istruzione di tipo aritmetico

l'istruzione di *load* e quella di *store* sono molto simili come si può vedere da Figura 12 e da Figura 13. Nel caso della *load* la *alu* calcola l'indirizzo di memoria da leggere il quale viene inviato alla memoria e il risultato della lettura è registrato tramite la porta write data del *RF*. Nel caso della *store* invece la *ALU* calcola l'indirizzo di destinazione della scrittura e tramite la porta write del *Data Memory* viene copiato il valore del registro. In entrambi i casi un'unità

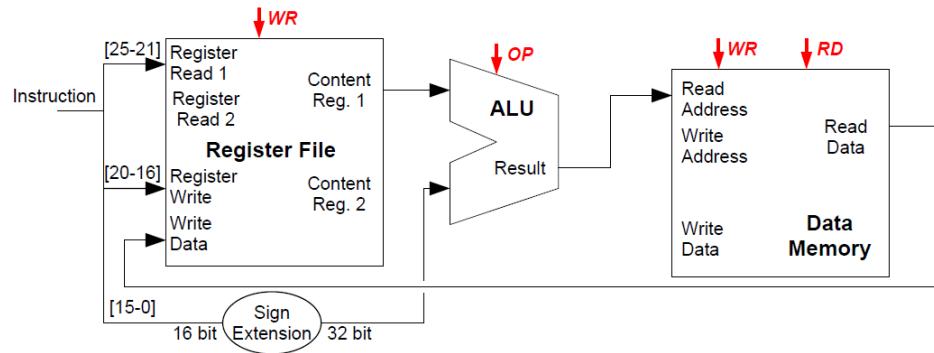


Figura 12: Hardware che implementa un istruzione di tipo load

particolare si occupa di eseguire l'estensione del segno in caso di bisogno.

Stabiliamo ora come viene implementato il clock del circuito; possiamo avere due possibilità, la prima è avere un unico ciclo di clock lungo quanto il percorso critico necessario per eseguire l'istruzione di *load* (la più lunga), la seconda è avere un ciclo di clock lungo quanto un singolo passaggio in uno dei componenti prima analizzati.

Analizziamo innanzitutto il caso di singolo ciclo, in questo caso il ciclo dovrà avere una durata pari al tempo necessario per eseguire un'istruzione di *load* che come abbiamo visto è pari a  $T = 8ns$  ( $f = 125 MHz$ ). Assumiamo quindi che ogni istruzione verrà eseguita in un singolo ciclo di clock, ogni modulo verrà utilizzato una sola volta per clock e quei moduli che dovrebbero essere utilizzati più di una volta dovranno essere duplicati. Inoltre dobbiamo tener conto anche delle differenze tra i diversi tipi di istruzioni, infatti, all'ingresso di scrittura dell'*RF* possiamo avere dati provenienti da una *ALU* e quindi di lunghezza [15-11] bit oppure dati provenienti da una *load/store* con una lunghezza di [20-16] bit questo richiede un *Multiplexer* all'ingresso dei registri nel *RF*. In secondo luogo al secondo ingresso della *ALU* possiamo avere avere il dato proveniente da un registro nel caso di operazioni *ALU* oppure l'offset per le istruzioni di

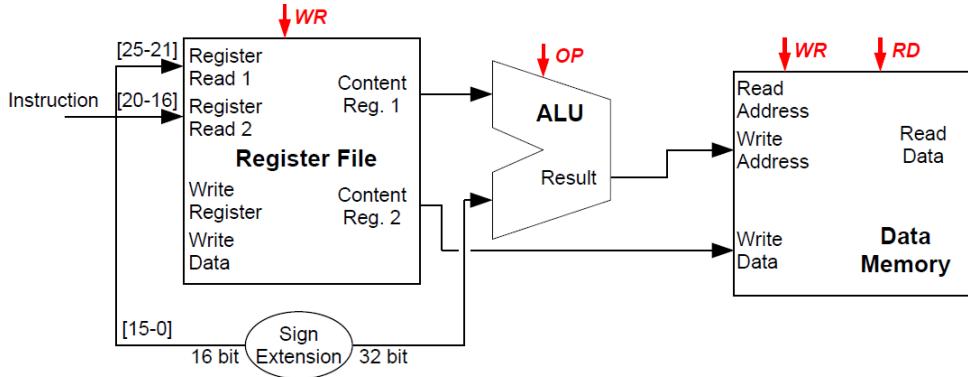


Figura 13: Hardware che implementa un'istruzione di tipo store

load/store, questo richiede un *MUX* al secondo ingresso della ALU. Infine i dati all'output del Destination Register possono arrivare sia dal risultato della ALU oppure dal Data Memory nel caso di load questo comporta l'utilizzo di un *MUX* all'ingresso in scrittura dei dati su RF. In Figura 14 vediamo l'implementazione completa di un MIPS a ciclo singolo con l'introduzione, oltre che dei MUX precedentemente specificati, anche di una ALU (parte alta della figura) che permette l'implementazione dei branch, e della logica di controllo (in rosso nella figura)

Veniamo ora al caso in cui il ciclo di clock sia di lunghezza pari al tempo necessario per un singolo modulo  $T = 2\text{ns}$  questo significa che per eseguire un'istruzione di load sono necessari 5 cicli di clock per un totale di  $10\text{ns}$ . Ogni fase dell'istruzione richiede un ciclo di clock ma questo permette la condivisione dei moduli tra diverse istruzioni in differenti cicli di clock, anche se questo richiede l'inserimento di registri tra un'unità e l'altra.

## 1.2 Pipelining

Il pipelining è una tecnica di ottimizzazione basata sull'esecuzione multipla sovrapposta di istruzioni sequenziali. L'idea fondamentale è quella di sfruttare il parallelismo intrinseco delle istruzioni sequenziali in quanto l'esecuzione di una istruzione è suddivisa in fasi differenti (*pipelines stages*) che richiedono soltanto una piccola frazione di tempo per essere completate. I diversi stadi sono connessi in sequenza nella pipeline, un'istruzione entra da una parte procede attraverso i diversi stadi e esce all'altro capo come in una catena di montaggio.

I vantaggi di questa tecnica è che è completamente trasparente al programmatore, inoltre come in una catena di montaggio, il tempo necessario per eseguire un'istruzione è uguale al caso in cui l'istruzione sia eseguita senza pipeline; quello che la pipeline fa è incrementare il numero di istruzioni eseguite contemporaneamente e perciò aumentare la frequenza di completamento come vediamo in Figura 15

Il tempo necessario per far avanzare un'istruzione di una fase corrisponde ad un ciclo di clock, le diverse fasi perciò devono essere *sincronizzate*, il periodo di clock deve essere uguale al tempo di esecuzione della fase più lenta (nel nostro esempio 2ns). L'obiettivo è quello di bilanciare la lunghezza di ogni fase della pipeline in modo da avere uno *speedup ideale* uguale al numero di fasi della pipeline.

Nel caso ideale vediamo come la pipeline sia più efficiente sia dell'architettura a singolo ciclo che a quella multi-ciclo viste in precedenza. Nel caso di una CPU1 non pipeline con un unico ciclo di clock della durata di 8ns contro una CPU2 con una pipeline a 5 stadi e ciclo di 2ns abbiamo che:

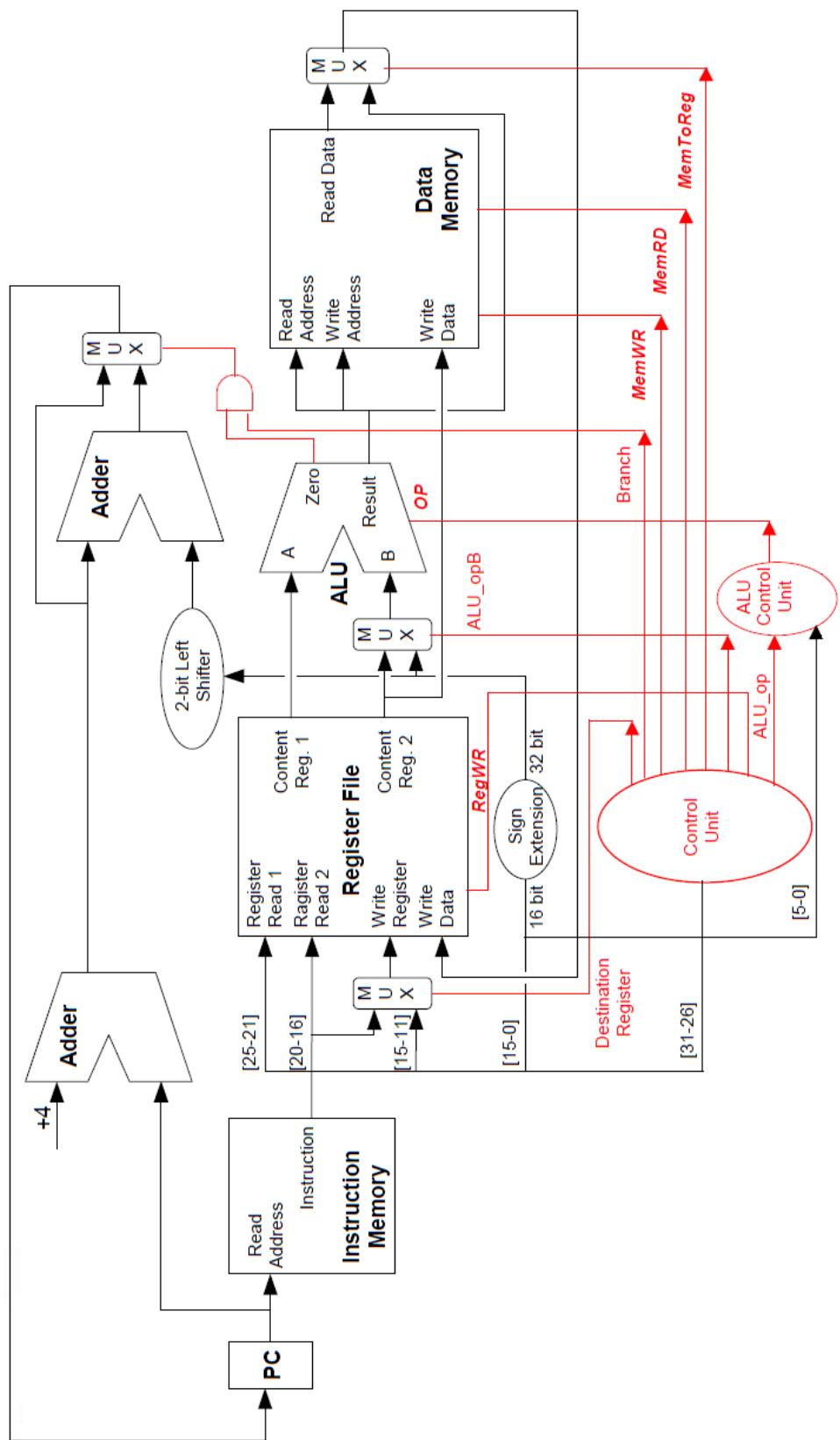


Figura 14: MIPS a singolo ciclo con logica di controllo

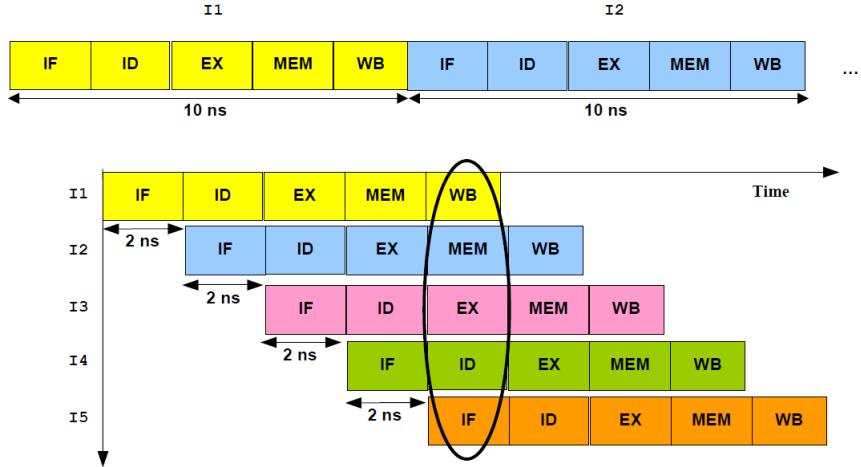


Figura 15: Confronto tra esecuzione sequenziale e pipelined

- la *latenza* ovvero il tempo necessario per eseguire una istruzione è peggiore nel caso di CPU2: 8ns vs 10ns
- il *throughput* tuttavia è notevolmente migliorato: 1 *istruzione*/8ns vs 1 *istruzione*/2ns

Nel caso di CPU3 multi-ciclo senza pipeline contro un'architettura CPU2 come quella descritta in precedenza abbiamo:

- la *latenza* resta invariata: 10 ns
- il *throughput* cresce di ben 5 volte: 1 *istruzione*/10ns vs 1 *istruzione*/2ns

### 1.2.1 Implementazione di una Pipeline

Innanzitutto vediamo quali fasi devono attraversare ciascuna operazione in quanto non tutte le fasi sono necessarie per tutte le operazioni, uno schema riassuntivo è specificato in Figura 16.

La divisione dell'esecuzione di una istruzione in 5 fasi implica che in ogni ciclo di clock cinque

IF Instruction Fetch	ID Instruction Decode	EX Execution	ME Memory Access	WB Write Back
<b>ALU Instructions: op \$x,\$y,\$z</b>				
Instr. Fetch & PC Increm.	Read of Source Regs. \$y and \$z	ALU Op. (\$y op \$z)		Write Back Destinat. Reg. \$x
<b>Load Instructions: lw \$x,offset(\$y)</b>				
Instr. Fetch & PC Increm.	Read of Base Reg. \$y	ALU Op. (\$y+offset)	Read Mem. M(\$y+offset)	Write Back Destinat. Reg. \$x
<b>Store Instructions: sw \$x,offset(\$y)</b>				
Instr. Fetch & PC Increm.	Read of Base Reg. \$y & Source \$x	ALU Op. (\$y+offset)	Write Mem. M(\$y+offset)	
<b>Conditional Branches: beq \$x,\$y,offset</b>				
Instr. Fetch & PC Increm.	Read of Source Regs. \$x and \$y	ALU Op. (\$x-\$y) & (PC+4+offset)	Write PC	

Figura 16: Fasi della pipeline necessarie ad ogni istruzione

istruzione sono in esecuzione questo comporta la necessità di inserire dei *registri* tra una fase e l'altra della pipeline per separare i diversi stage.

In Figura 17 vediamo una possibile implementazione di un'architettura MIPS pipelined con l'introduzione dei registri tra le fasi (in verde).

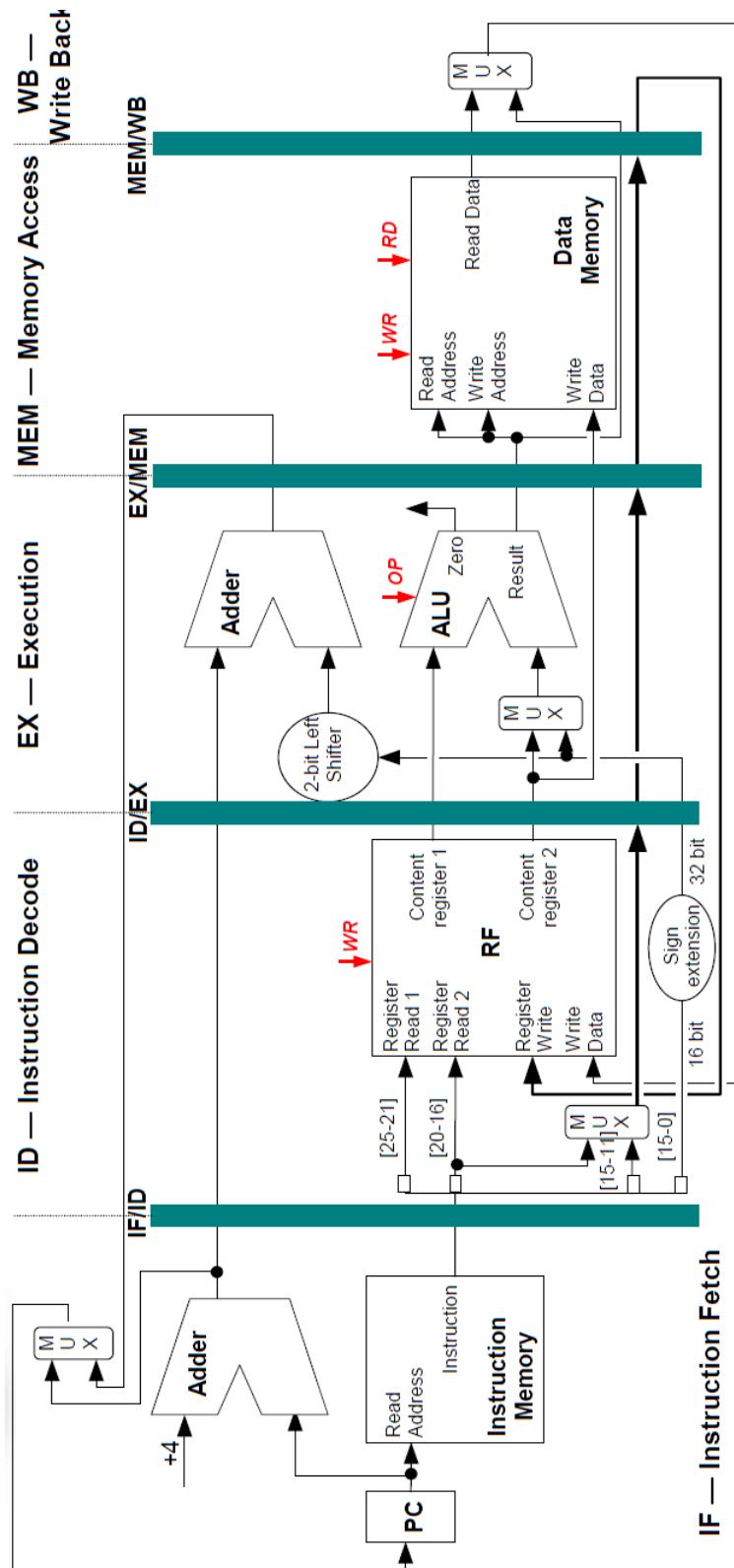


Figura 17: Schema di un MIPS con pipeline

### 1.3 Il problema del "Hazard"

Si ha un *hazard* quando vi è una dipendenza tra istruzioni diverse e la sovrapposizione dovuta al pipeline cambia l'ordine di dipendenza sugli operandi. Hazard previene l'esecuzione della prossima istruzione nel ciclo di clock designato ma così facendo riduce le performance allontanandole dallo speedup ideale.

Possiamo distinguere tre classi di *hazard*:

- **Structural Hazards:** si ha quando diverse istruzioni cercano di utilizzare la stessa risorsa simultaneamente (stessa memoria per istruzioni e dati)
- **Data Hazards:** si ha quando si cerca di utilizzare un risultato prima che questo sia pronto (istruzione dipendente dalla precedente che è nella pipeline)
- **Control Hazards:** si ha quando si deve prendere una decisione sulla esecuzione della prossima istruzione prima della valutazione di una condizione (branch condizionali)

Tra questi tre tipi di hazard il primo non può presentarsi nelle architetture MIPS in quanto lo spazio di memoria dedicato alle istruzioni e quello dedicato ai dati sono fisicamente separati.

#### 1.3.1 Data Hazard

Per quanto riguarda il data hazard si verifica quando sono in esecuzione nella pipeline due o più istruzioni *dipendenti*.

```
sub $2, $1, $3
and $12, $2, $5    #1° operando dipende dalla sub
or  $13, $6, $2    #2° operando dipende dalla sub
add $14, $2, $2    #1° & 2° operando dipendono dalla sub
sw   $15, 100($2)  #Il registro base dipende dalla sub
```

Come vediamo dall'esempio qui sopra e dalla sua esecuzione in Figura 18 abbiamo che le istruzioni successive alla **sub** debbono aspettare che la prima istruzione arrivi nella fase di *write-back* prima di poter utilizzare il dato come avviene per l'ultima istruzione evidenziata da una freccia verde. Esistono diversi meccanismi per far sì che queste dipendenze vengano soddisfatte, le

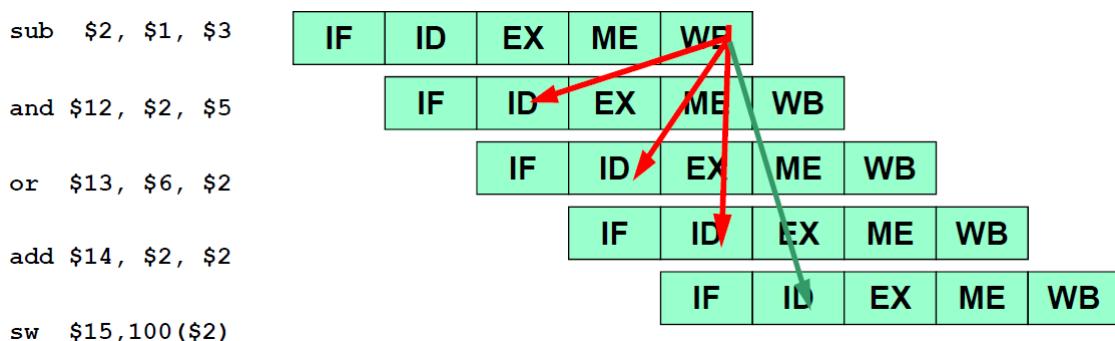


Figura 18: Esempio di data hazard

principali si possono suddividere in due categorie:

- **Tecniche di compilazione:** in questa categoria rientra il re-scheduling delle operazioni che consiste nell'inserire istruzioni indipendenti tra le istruzioni correlate in modo da permettere il calcolo dei valori necessari; nel caso non sia possibile inserire altre operazioni il compilatore inserisce delle **nop** ovvero delle operazioni che non fanno nulla *no operation*

- **Tecniche hardware:** in questa categoria rientrano la possibilità di inserire delle *bubbles* o degli stalli oppure le tecniche di *Data Forwarding* e di *Bypassing*

Vediamo innanzi tutto un esempio di inserimento di `nop` in Figura 19 Come vediamo l'inserimento

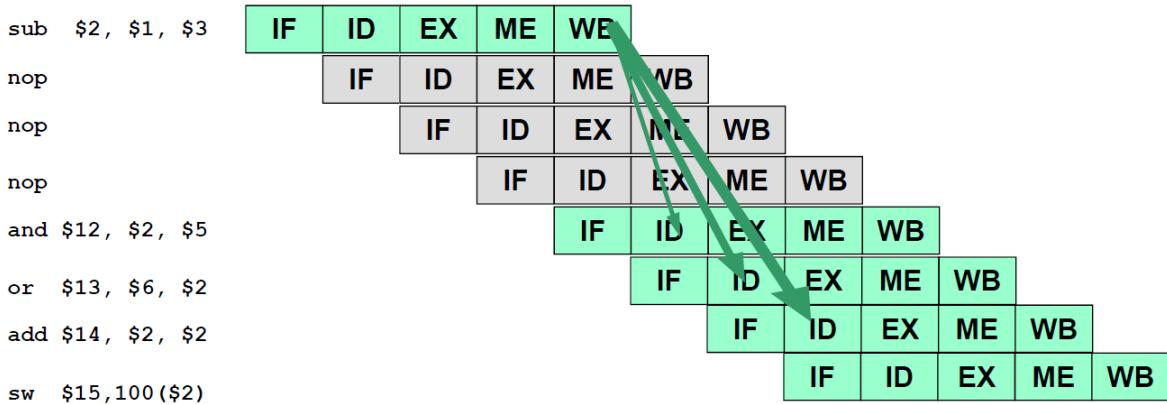


Figura 19: Esempio di uso `nop`

di `nop` peggiora lo speedup ideale, cosa che invece non succede se si applicano le tecniche di rescheduling in quanto non vengono inserite istruzioni inutili nell'esecuzione delle istruzioni ma viene modificato semplicemente l'ordine nel quale vengono eseguite.

Il caso di inserimento di stalli è molto simile a quello di inserimento delle `nop` la differenza sta nel fatto che si ferma l'esecuzione dell'istruzione dipendente il tempo necessario affinché l'istruzione in esecuzione renda disponibile il dato come vediamo in Figura 20, anche in questo caso abbiamo un peggioramento dello speedup ideale.

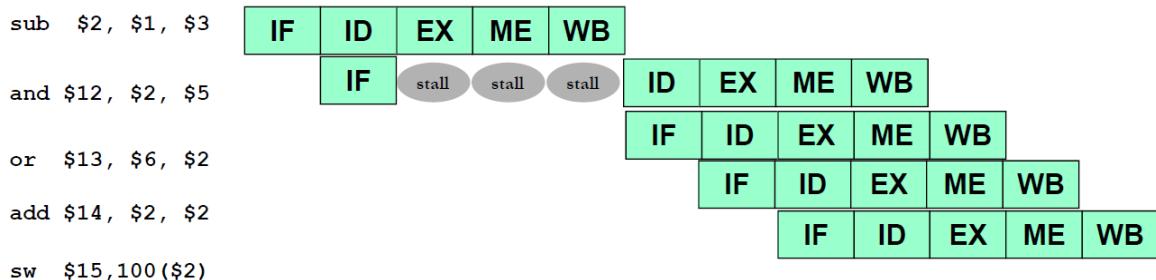


Figura 20: Esempio di uso degli stalli

**Data Forwarding** Il *data forwarding* è una tecnica hardware che comporta l'utilizzo dei risultati temporanei immagazzinati nei registri della pipeline, per fare ciò abbiamo bisogno di aggiungere dei *multiplexer* all'ingresso della ALU per selezionare la provenienza dei dati. In Figura 21 vediamo quali sono i collegamenti necessari per l'esecuzione di istruzioni aritmetiche dipendenti; questi path sono tre:

- **EX/EX path:** in figura da ALU ad ALU che risolve il problema di due istruzioni consecutive nella quale la seconda necessita del risultato della prima. Nell'esempio precedente la dipendenza `sub`→`and`
- **MEM/EX path:** in questo caso si risolve la dipendenza tra la prima e la terza istruzione (`sub`→`or`)

- **MEM/ID path:** risolve la dipendenza tra la prima e la quarta istruzione (**sub**→**add**)

Con l'introduzione di questi path si riesce a risolvere tutti i problemi di dipendenza per quanto riguarda le istruzioni di tipo aritmetico. In Figura 22 vediamo una possibile implementazione hardware di una pipeline con sistema di forwarding path. Esiste ancora una situazione però

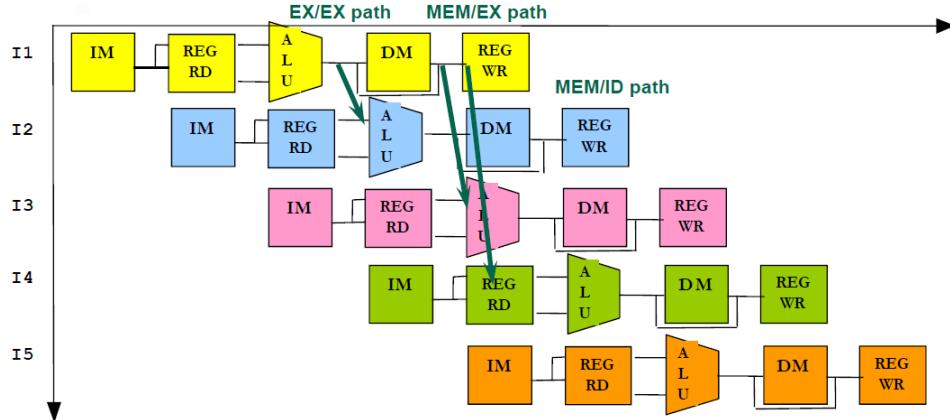


Figura 21: Esempio di forwarding path

in cui è necessario inserire degli stalli, questa situazione è dovuta alla sequenza di istruzioni seguenti:

```
L1: lw $s0, 4($t1)      #$s0<-M[4 + $t1]
L2: add $s5, $s0, $s1    #1° operand depends from L1
```

Questa dipendenza si crea in quanto il dato viene scritto in **\$s0** solo nella fase di WB mentre viene letto durante la fase di ID. In questo caso non si può fare molto se non sfruttare il forwarding path MEM/EX già analizzato prima anche se comunque è necessario introdurre uno stallo per risolvere la dipendenza. Nel caso invece in cui la dipendenza sia tra un'istruzione di *load* e una di *store* si può risolvere la dipendenza aggiungendo un *forwarding path* tra le due fasi di MEM. Questo path permette di risolvere la dipendenza senza dover introdurre altri stalli come si può vedere in Figura 23. Con l'architettura attuale, come abbiamo visto, nel caso di dipendenza tra una *load* ed un'istruzione aritmetica che legge un registro è necessario introdurre uno stallo; in quanto l'accesso in lettura avviene durante la fase di ID mentre la scrittura avviene durante la fase di WB. Nel caso, invece, di *pipeline ottimizzata* possiamo assumere che la fase di lettura avviene nelle seconda metà del ciclo di clock mentre la fase di scrittura nella prima metà; in questo modo nel caso in cui lettura e scrittura facciano riferimento allo stesso registro nello stesso ciclo di clock non è più necessario inserire degli stalli, e si può inoltre eliminare il forwarding path tra MEM e ID.

### 1.3.2 Altri tipi di Data Hazard

Fino ad ora abbiamo analizzato solo un tipo di dipendenza sui dati, questo tipo di dipendenza è chiamato **RAW (Read After Write)**, e si ha quando l'istruzione  $n+1$  cerca di leggere un registro prima che l'istruzione  $n$  abbia finito di scrivere tale registro.

Esistono tuttavia altri due tipi di *data hazard* che sono:

- Write After Write (WAW)
- Write After Read (WAR)

La dipendenza di tipo WAW si ha quando un istruzione  $n+1$  di tenta di scrivere un registro il quale non è ancora stato scritto dall'istruzione  $n$ . Tale tipo di dipendenza avviene solo nel caso in cui la nostra pipeline preveda la possibilità di fasi di memorizzazione o di esecuzione multi-ciclo come mostrato in Figura 24 e 25 le quali portano alla terminazione delle istruzioni fuori ordine. Per quanto riguarda le dipendenze di tipo WAR si hanno quando l'istruzione  $n+1$  tenta di scrivere un registro prima che questo sia stato letto da un'istruzione  $n$ , nel caso di architettura MIPS però tale tipo di dipendenza non può mai verificarsi in quanto la lettura avviene nella fase ID mentre la scrittura nella fase WB.

## 1.4 Analisi delle performance

L'utilizzo della pipeline aumenta il throughput della CPU ma non riduce il tempo di esecuzione della singola istruzione, anzi solitamente aumenta la latenza di ogni istruzione bisogna quindi bilanciare il numero di fasi con l'overhead dovuto alla pipeline.

Definito  $IC = \text{Instruction Count}$  ovvero il numero di istruzioni eseguite, possiamo determinare il numero di cicli di clock necessari per completare queste operazioni operazione. Tale valore è uguale a:

$$\#Clock\ Cycle = IC + \#Stall\ Cycles + 4$$

Dividendo tale valore per il numero di operazioni otteniamo:

$$\begin{aligned} CPI &= \text{Clock Per Instruction} = \#Clock\ Cycle/IC = \\ &= (IC + \#Stall\ Cycles + 4)/IC \\ MIPS &= f_{clock}/(CPI * 10^6) \end{aligned}$$

Come visto fino ad ora la CPI ideale per la pipeline è 1 ma gli stalli degradano le performance. Abbiamo così che la CPI media è data da:

$$\begin{aligned} \text{Ave. } CPI &= \text{Ideal } CPI + \#Stall\ per\ Instruction \\ &= 1 + \#Stall\ per\ Instruction \end{aligned}$$

Possiamo misurare il miglioramento delle performance dato dall'introduzione della pipeline come

$$\begin{aligned} \text{Pipeline SpeedUp} &= \frac{\text{Ave. Exec. Time Unpipelined}}{\text{Ave. Exec. Time Pipelined}} = \\ &= \frac{\text{Ave. CPI Unp.}}{\text{Ave. CPI Pipe}} \times \frac{\text{Clock Cycle Unp.}}{\text{Clock Cycle Pipe}} \end{aligned}$$

Se ignoriamo l'overhead sul tempo di clock e assumiamo che i diversi stage siano perfettamente bilanciati possiamo ridefinire lo speedup come

$$\text{SpeedUp}_{\text{pipeline}} = \frac{\text{Ave. CPI Unp.}}{1 + \#Stall\ per\ Instruction}$$

Nel caso ideale nel quale tutte le istruzioni richiedano lo stesso numero di cicli questi sono uguali al numero di fasi della pipeline e possiamo riscrivere la precedente come:

$$\text{SpeedUp}_{\text{pipeline}} = \frac{\text{Pipeline Depth}}{1 + \#Stall\ per\ Instruction}$$

Nel caso ideale in cui non ci siano stalli vediamo come le performance migliori tanto è più profonda (maggior numero di fasi) la pipeline.

Nel caso in cui si abbiano dei salti condizionati le performance peggiorano in base alla penalità del branch, infatti:

$$\text{SpeedUp}_{\text{pipeline}} = \frac{\text{Pipeline Depth}}{1 + \text{Branch Frequency} \times \text{Branch Penalty}}$$

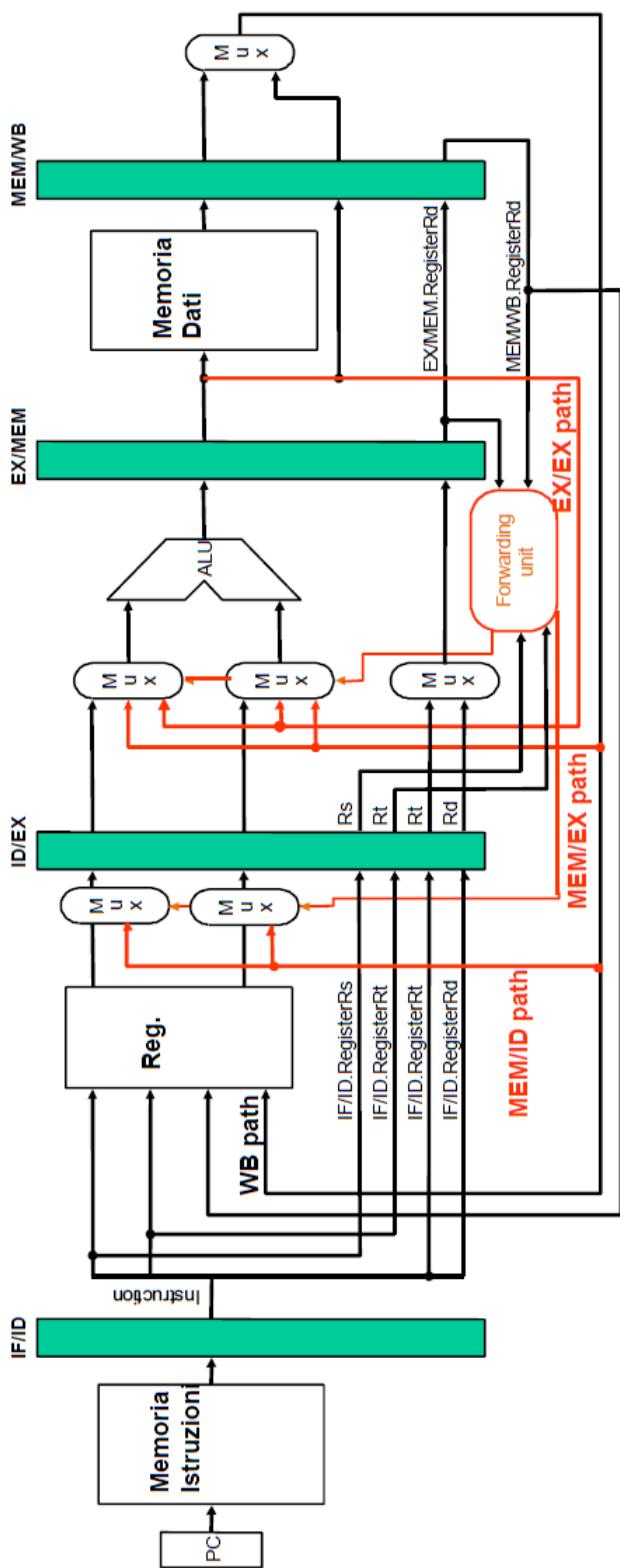


Figura 22: Schema MIPS con forwarding

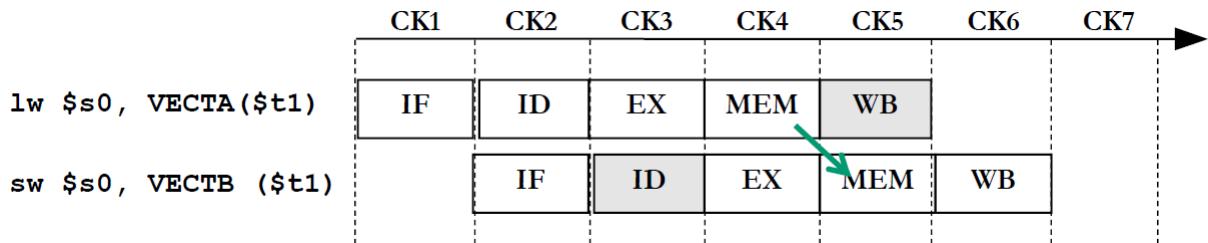


Figura 23: Forwarding path tra due fasi di memorizzazione successive.

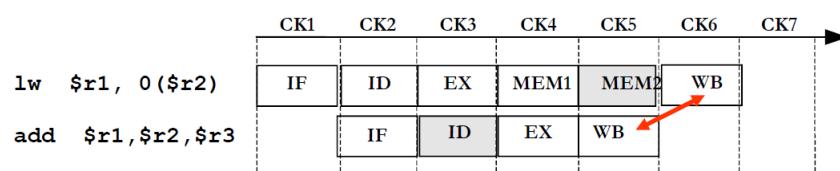


Figura 24: Fase di memorizzazione multiciclo che porta alla creazione di dipendenze di tipo WAW

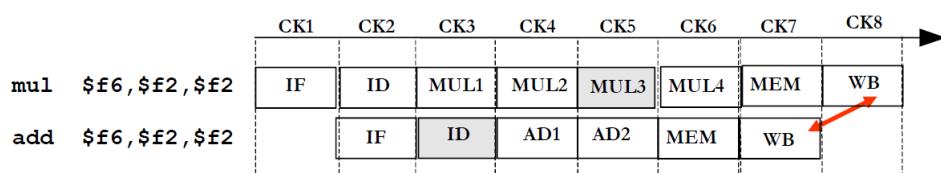


Figura 25: Fase di esecuzione multiciclo che porta alla creazione di dipendenze di tipo WAW

op	rs	rt	address
6 bit	5 bit	5 bit	16 bit

Figura 26: Esempio di istruzione di branch

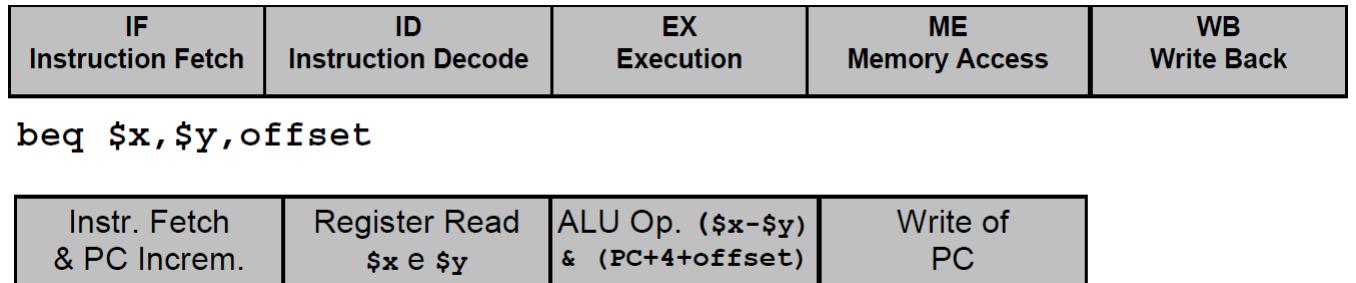


Figura 27: Suddivisione dell'esecuzione di un'istruzione di salto nelle varie fasi di una pipe

## 2 Tecniche di predizione dei salti

Come abbiamo visto nel Capitolo 1 i branch condizionati sono istruzioni di salto che vengono eseguite soltanto se viene soddisfatta la condizione. L'indirizzo di destinazione del branch viene sostituito nel program counter al posto dell'indirizzo dell'istruzione sequenziale successiva. Nel caso di ambiente MIPS possiamo distinguere due tipi di branch:

- **beq:** *branch on equal* che richiede che i valori nei registri da confrontare siano uguali.
- **bne:** *branch on not equal* che richiede che i valori nei registri da confrontare siano diversi.

Come abbiamo visto nel capitolo precedente le istruzioni di salto condizionato sono nel formato *I-Format* e un'istruzione di questo tipo è suddivisa nei diversi campi come mostrato in Figura 26. Dove il campo **address** indica l'indirizzo relativo rispetto al program counter che punta all'etichetta di salto.

Questo tipo di istruzione sfrutta solo quattro dei cinque stadi della pipeline come mostrato in Figura 27; durante l'instruction fetch si recupera l'istruzione da eseguire e si aggiorna il program counter all'istruzione sequenziale successiva, successivamente durante la fase di instruction decode si leggono i due registri da confrontare. Durante la fase di execution la ALU compara i due registri e calcola il valore di destinazione del salto. Durante la fase *Memory Access* si decide in base al valore della comparazione effettuata dalla ALU se aggiornare il PC con il valore del salto.

### 2.1 Il problema del Control Hazard

Il *control hazard* è il problema di decidere quale istruzione eseguire prima che la condizione di salto sia valutata. I problemi di *control hazard* nascono ogni qualvolta nella pipeline sia necessario modificare il valore del PC. Tali problemi riducono perciò la velocità della pipeline riducendo lo speedup ideale a causa di introduzioni di stalli nella pipeline.

Per alimentare la pipeline è necessario prelevare una nuova istruzione ad ogni ciclo di clock ma la decisione se effettuare o non effettuare un salto avviene solo durante lo stage *MEM*. Questo ritardo nel determinare l'istruzione successiva corretta è chiamato *Control Hazard* o *Conditional*

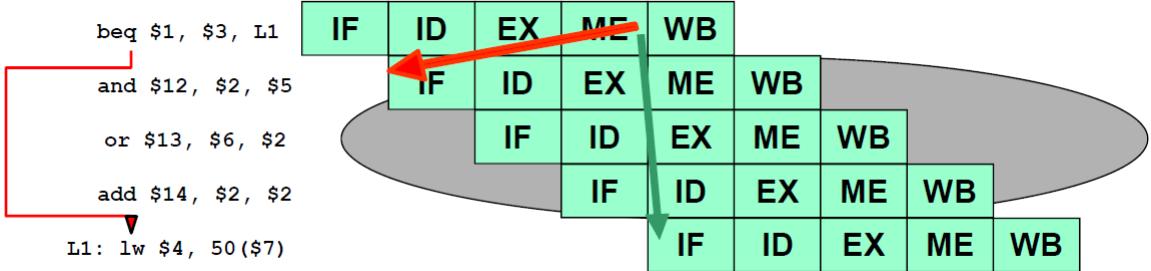


Figura 28: Esempio di esecuzione di un istruzione di salto

#### *Branch Hazard.*

Analizziamo ora l'esempio di Figura 28 in questo esempio la prima istruzione è un salto condizionato che viene valutato solo nella fase si MEM; durante l'esecuzione di tale istruzione vengono prelevate anche le tre istruzioni successive per continuare ad alimentare la pipeline. Se il salto non viene eseguito l'esecuzione è corretta e può proseguire, nel caso in cui, invece, il salto venga eseguito allora diventa necessario effettuare il *flush* delle tre istruzioni prelevate durante l'esecuzione del branch. Una possibile soluzione al problema è quella di attendere la decisione del salto prima di effettuare qualsiasi altra operazione; questo comporta l'inserimento di stalli nella pipeline; più precisamente sono necessari:

- tre stalli senza forwarding
- due stalli con forwarding

Nel caso in cui il salto non venga eseguito, tuttavia, la penalità di tre cicli di stallo non è giustificata. Un'altra soluzione è quella di assumere che il salto non sia mai eseguito e quindi scartare le tre istruzioni nel caso in cui il salto venga preso.

Una terza soluzione è quella di aggiungere delle risorse hardware per permettere di:

- comparare i registri
- calcolare l'indirizzo di destinazione del branch
- aggiornare il valore del PC

il prima possibile nella catena della pipeline. Nei processori MIPS tutto questo avviene durante lo stage ID come mostrato in Figura 29 Utilizzando queste tecniche si riesce a ridurre al minimo il costo per recuperare la corretta esecuzione di un branch nel caso di scelta sbagliata, riducendo a uno il numero degli stalli da introdurre.

Tuttavia queste tecniche comportano tuttavia una riduzione delle performance quantificabile tra il 10% e il 30% in base alla frequenza dei salti. Tali perdite di performance possono essere ridotte ulteriormente tramite alcune tecniche.

## 2.2 Tecniche di predizione dei salti

In generale il problema dei salti diventa importante quando si tratta di processori con delle pipeline profonde, dove il costo di una predizione errata è molto alto. Lo scopo principale delle tecniche di predizione dei salti è quello di predire il prima possibile il risultato di un istruzione di salto.

Le performance di una tecnica di predizione si possono misurare in:

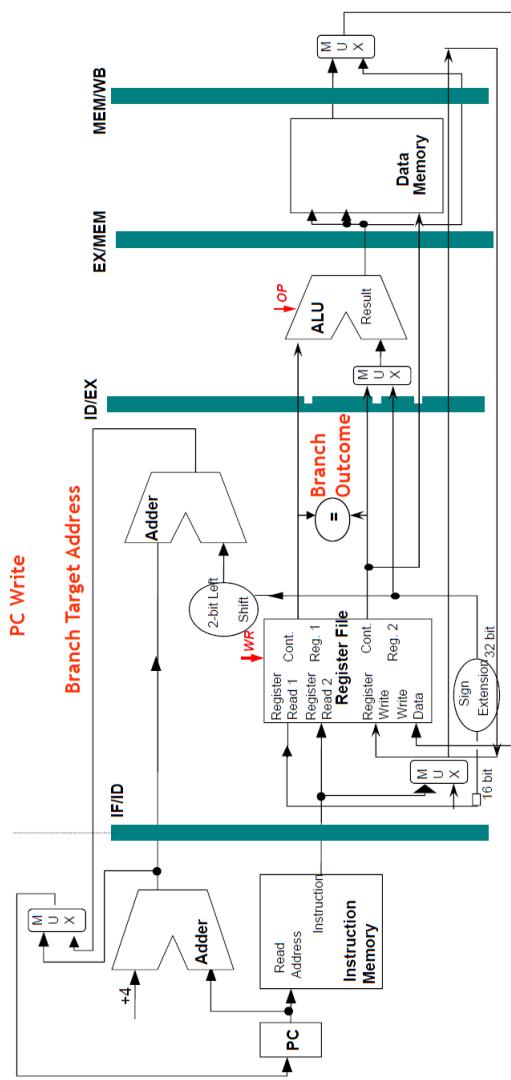


Figura 29: Hardware aggiuntivo per risolvere i problemi di controllo

- **Accuratezza:** misurata in termini di percentuale di predizioni sbagliate.
- **Costo:** misurato come tempo perso nel caso di predizione sbagliata.

Le tecniche di predizione dei salti possono essere suddivise in due categorie:

- *Tecniche di predizione statiche:* le azioni intraprese per il branch sono prefissate e uguali per tutta l'esecuzione e determinate a tempo di compilazione.
- *Tecniche di predizione dinamiche:* in questo caso le decisioni variano a seconda dell'esecuzione.

### 2.2.1 Tecniche di predizione statiche

Le tecniche di predizione statiche sono utilizzate soprattutto in quei processi dove ci si aspetta che i salti siano altamente predicibili. Alcune tecniche statiche di predizione dei salti sono:

- Branch Always Not Taken (Predicted-Not-Taken)
- Branch Always Taken (Predicted-Taken)
- Backward Taken Forward Not Taken (BTFNT)
- Profile-Driven Prediction
- Delayed Branch

**Branch Always Not Taken** In questa particolare tecnica assumiamo che il salto non venga mai intrapreso e le istruzioni vengono prelevate sequenzialmente e il flusso prosegue come se il salto non venga intrapreso. Se la condizione nello stage ID non viene soddisfatta la predizione è corretta e perciò non abbiamo perdita di performance.

Se la condizione nello stage ID risulta soddisfatta allora la predizione è errata e il salto viene effettuato. A questo punto dobbiamo effettuare il flush delle successive istruzioni che sono già state messe in esecuzione sostituendole con delle `nop` e riprendere l'esecuzione inserendo nella pipeline la prima istruzione del salto. Tutto questo porta ad una penalità di un ciclo di clock.

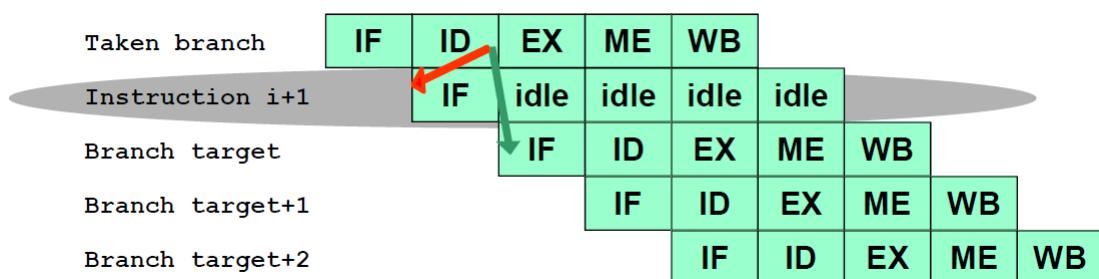


Figura 30: Esempio di penalità dovuto ad una predizione sbagliata

**Branch Always Taken** In alternativa al tipo di predizione precedente si può considerare che ogni salto sia sempre eseguito. Ogni qualvolta che un salto è decodificato e il suo indirizzo di destinazione è calcolato allora si assume che il salto sia eseguito e si introducono nella pipeline le istruzioni puntate dall'indirizzo di destinazione. Questo tipo di previsione ha senso in quelle pipeline dove l'indirizzo target è calcolato prima della comparazione dei registri.

Nelle pipeline di tipo MIPS noi non conosciamo l'indirizzo di destinazione prima della valutazione delle condizioni di salto così non vi è alcun vantaggio dall'utilizzo di questa tecnica.

**Backward Taken Forward Not Taken (BTFTNT)** La predizione di questa tecnica si basa sulla direzione del salto, ovvero se i salti sono all'indietro sono previsti come eseguiti (come ad esempio nei cicli) salti in avanti sono considerati come non eseguiti.

**Profile-driven Prediction** Questo tipo di predizione si basa su dati raccolti da precedenti esecuzioni del programma utilizzando alcune funzioni del compilatore.

**Delayed Branch Technique** In questo tipo di tecnica il compilatore schedula una particolare istruzione indipendente dal salto in un campo chiamato **branch delay slot**. L'istruzione in questo slot viene eseguita ogni qualvolta che il salto viene eseguito oppure no. Se assumiamo il ritardo dovuto ad un salto pari ad un ciclo di clock allora gli slot necessari per le istruzioni sono uno. Un esempio di questa tecnica è mostrato in Figura 31 dove nel *branch delay slot* viene inserita un'istruzione di add indipendente dal ciclo. Sia nel caso che il salto sia eseguito sia nel

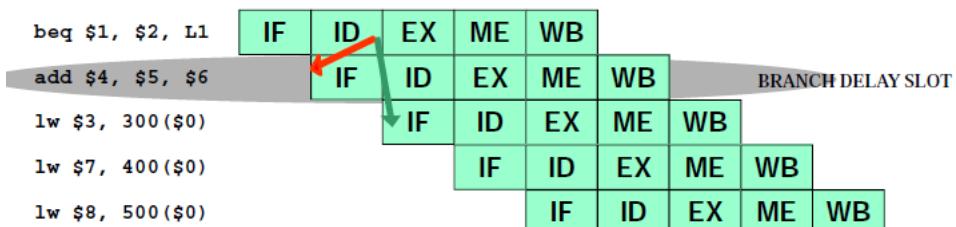


Figura 31: Esempio di utilizzo del branch delay slot

caso non sia eseguito l'istruzione dopo quella di salto è sempre quella del *delay slot* tuttavia nel caso il salto sia eseguito dopo l'istruzione del *delay slot* l'esecuzione prosegue con le istruzioni del salto viceversa nel caso non sia eseguito allora l'esecuzione prosegue con le istruzioni successive. Il compilatore deve essere in grado di selezionare l'istruzione da inserire nel *delay slot* in modo che essa sia valida ed utile. Ci sono quattro modi per selezionare tale istruzione:

- From Before
- From Target
- From Fall-Through
- From After

La tecnica *From Before* prevede di selezionare un'istruzione indipendente selezionata tra quelle che precedono il salto in tale modo l'istruzione viene sempre eseguita. Il metodo *From Target* prevede la copia dell'istruzione puntata dal salto; questa tecnica è preferibile quando il salto ha un'alta probabilità di essere eseguito. Un esempio di utilizzo di questa tecnica è mostrato

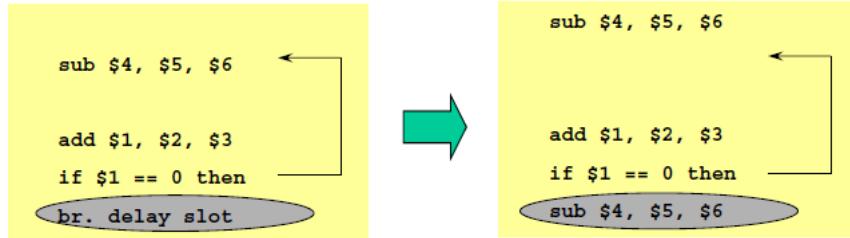


Figura 32: Esempio di selezione dell’istruzione *From Target*

in Figura 32. La tecnica *From Fall-Through* è contrapposta alla tecnica *From Target* infatti questa prevede che l’istruzione selezionata per il delay slot sia la prima istruzione che verrebbe prelevata nel caso di salto non eseguito come si vede dalla Figura 33. Questo tipo di tecnica è

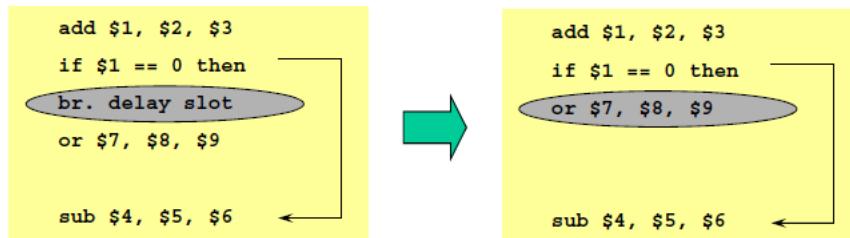


Figura 33: Esempio di uso della tecnica *From Fall-Through*

preferibile quando il salto ha un alta probabilità di essere non eseguito. Perché le ultime due tecniche risultino corrette è necessario che le istruzioni eseguite quando il salto non prende la direzione prevista risultino ininfluenti rispetto alla normale esecuzione del programma. Questo è possibile ad esempio se l’istruzione nel delay slot agisce su dei registri che sono inutilizzati nel caso di risultato inaspettato del salto.

In generale il compilatore è in grado di assegnare il 50% dei *delay slot* con istruzioni valide e utili, all’altro 50% viene riempito con istruzioni di tipo *nop*. Nel caso di pipeline più profonda il tempo necessario per valutare il salto è maggiore di un ciclo di clock e di conseguenza aumenta il numero di *delay slot* da riempire; diventa sempre più difficile perciò riempire tali slot con istruzioni valide e utili. Le principali limitazioni che nascono sono dovute alle restrizioni sulle istruzioni che possono essere schedulate e sull’abilità del compilatore di predire staticamente il risultato del salto.

Per migliorare l’abilità del compilatore di riempire i *delay slot* molti processori hanno introdotto il **cancelling or nullifying branch** ovvero salti nei quali è anche compresa la predizione della direzione del salto. Quando la predizione è verificata allora il contenuto del *delay slot* è eseguito in caso contrario viene eseguita una *nop*.

### 2.2.2 Tecniche di predizione dinamiche

L’idea di base in questo tipo di tecnica è quella di utilizzare il risultato di esecuzioni di salti passate per predire i salti futuri. Possiamo utilizzare dell’hardware per predire dinamicamente il risultato del salto; la predizione dipende dal comportamento dei salti durante l’esecuzione. Il meccanismo di predizione dinamica si basa su due meccanismi che interagiscono tra loro:

- **Branch Outcome Predictor:** che tenta di predire la direzione del branch.

- **Branch Target Predictor:** che calcola l'indirizzo di destinazione del salto nel caso in cui la condizione del salto abbia un risultato positivo.

Questi due moduli sono usate dall'unità di Instruction Fetch per predire la prossima istruzione da prelevare dall'I-Cache. Nel caso in cui il salto non venga eseguito il PC viene semplicemente incrementato, nel caso in cui, invece il branch venga preso il branch target predictor fornisce l'indirizzo di destinazione. Esiste inoltre una tabella chiamata **Branch History Table** la quale contiene un bit per ogni predizione passata che indica se il salto è stato preso oppure no. L'indice della tabella è basato su di una piccola porzione dell'indirizzo dell'istruzione di salto. Se la previsione è corretta e si prosegue nella direzione della previsione. Se la previsione risulta sbagliata il bit di previsione viene invertito e riportato nella tabella. Ogni accesso in tabella è un *hit* anche se tuttavia il bit di predizione può essere stato modificato da un altro salto con la stessa porzione di indirizzo. Una predizione errata avviene quando una predizione risulta

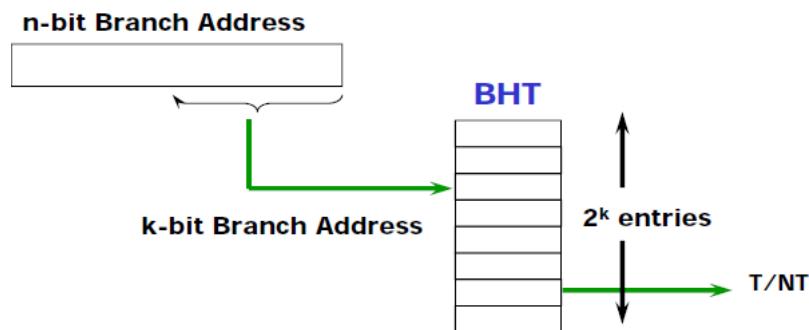


Figura 34: Esempio di *Branch History Table*

sbagliata oppure quando un indice è puntato da due differenti salti e la previsione si riferisce all'altro salto, per risolvere questo problema è sufficiente incrementare il numero di righe della BHT o utilizzare una funzione di hash.

Nel caso di una BHT ad un solo bit e considerando per esempio un salto di un loop che solitamente è eseguito la BHT sbaglierà previsione due volte, nel caso non sia eseguito e nel caso in cui venga eseguito il salto subito dopo non essere stato eseguito. Per meglio specificare i due casi abbiamo:

- All'ultima iterazione quando il salto non viene eseguito anche se la predizione indicava il contrario
- Quando rientriamo nel loop alla fine della prima interazione la nostra predizione indica che dovremmo uscire (ultima interazione del loop precedente) mentre in realtà effettueremo il salto.

Per ovviare a questo problema sfruttiamo una BHT a due bit nella quale sono necessarie due predizioni sbagliate consecutive per cambiare la nostra predizione. Per ogni indice della tabella i due bit sono utilizzati per indicare un dei quattro stati della macchina a stati finita in Figura 35 Tale tecnica si può generalizzare fino ad utilizzare una tabella con  $n$  bit per ogni record. Il valore che può assumere questo record va da 0 a  $2^n - 1$ ; quando il valore del record diventa uguale ad almeno la metà del suo valore massimo ( $2^n - 1$ ) la predizione del salto indica che esso deve essere eseguito, altrimenti la predizione sarà che non deve essere eseguito. Nello schema precedente il contatore veniva incrementato quando il salto veniva intrapreso e decrementato quando non veniva intrapreso.

Tuttavia anche se una generalizzazione è possibile gli studi hanno dimostrato che una una tabella

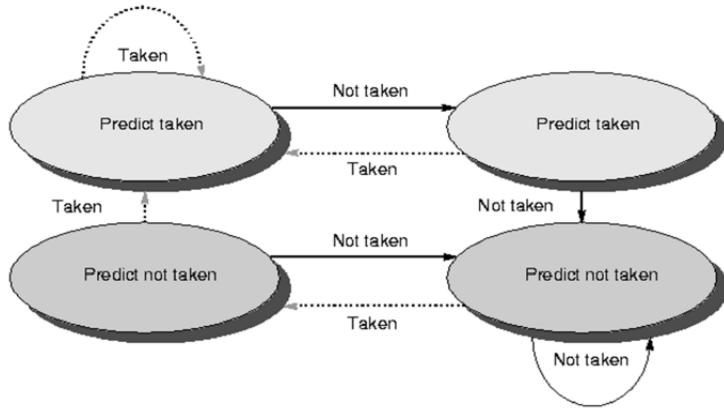


Figura 35: Macchina a stati finiti per una BHT a 2 bit

a 2 bit fornisce dati più che soddisfacenti. Ad esempio per una architettura IBM *SPEC89* con una BHT con 4K record di 2 bit l'accuratezza nella predizione varia dal 99% all'82%.

La BHT a 2 bit utilizza solo i risultati delle esecuzioni precedenti di un singolo salto per predirne il risultato di quel salto. L'idea di base però è che il comportamento dei salti recenti è che essi sono correlati con il comportamento di altri salti e perciò la predizione può essere influenzata da tali comportamenti. Un esempio è mostrato in Figura 36 Un predittore che utilizza il comportamento

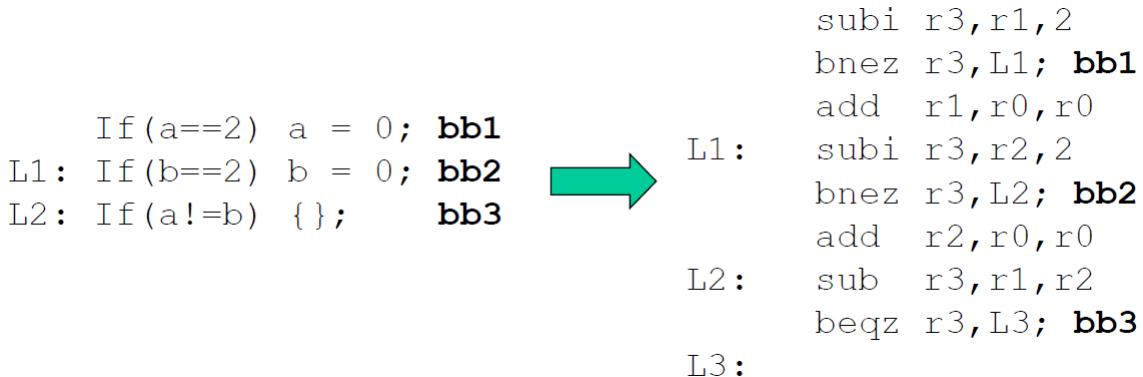


Figura 36: Esempio di salti correlati

di altri salti per effettuare una predizione è chiamato **Correlating Predictors** o anche **2-Level Predictors**. Un esempio di un *(1,1) Correlating Predictors* è un predittore ad un bit con un bit di correlazione, ovvero il comportamento dell'ultimo salto è utilizzato per scegliere una coppia di predittori ad un bit come mostrato in Figura 37. Si registrano le esecuzioni degli ultimi  $k$  salti; la predizione si basa sull'esecuzione del salto precedente selezionando la BHT ad un bit appropriata:

- Una predizione è usata nel caso in cui l'ultimo salto è stato eseguito.
- L'altra predizione è utilizzata se l'ultimo salto non è stato intrapreso.

In generale l'esecuzione dell'ultimo salto non riguarda la stessa istruzione sulla quale si cerca di fare una predizione come normalmente accade nei loop semplici.

In generale possiamo costruire un predittore correlato con  $(m,n)$  dove  $m$  indica gli ultimi  $m$

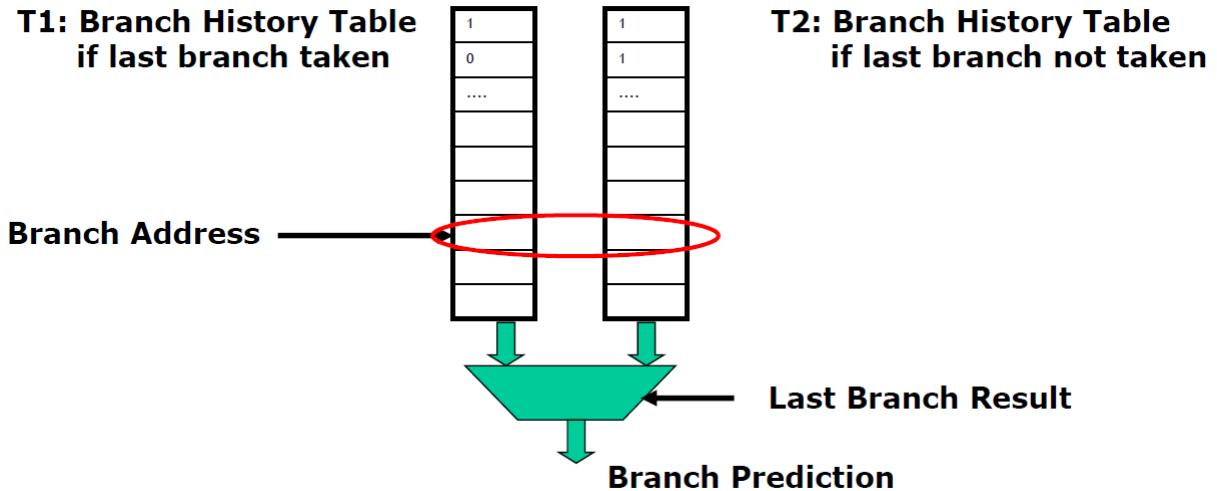


Figura 37: Esempio di predittore correlato di tipo (1,1)

salti da analizzare selezionando  $2^m$  BHT ognuna delle quali è un predittore a  $n$  bit. Il branch prediction buffer può essere indicizzato concatenando la parte finale del branch address con gli  $m$  bit della *global history*. Un esempio di un predittore correlato è un predittore (2,2) dove si hanno quattro BHT a 2 bit tra i quali scegliere e si usano 2 bit dalla global history per selezionare quale utilizzare come mostrato in Figura 38. Un predittore a due bit non correlato

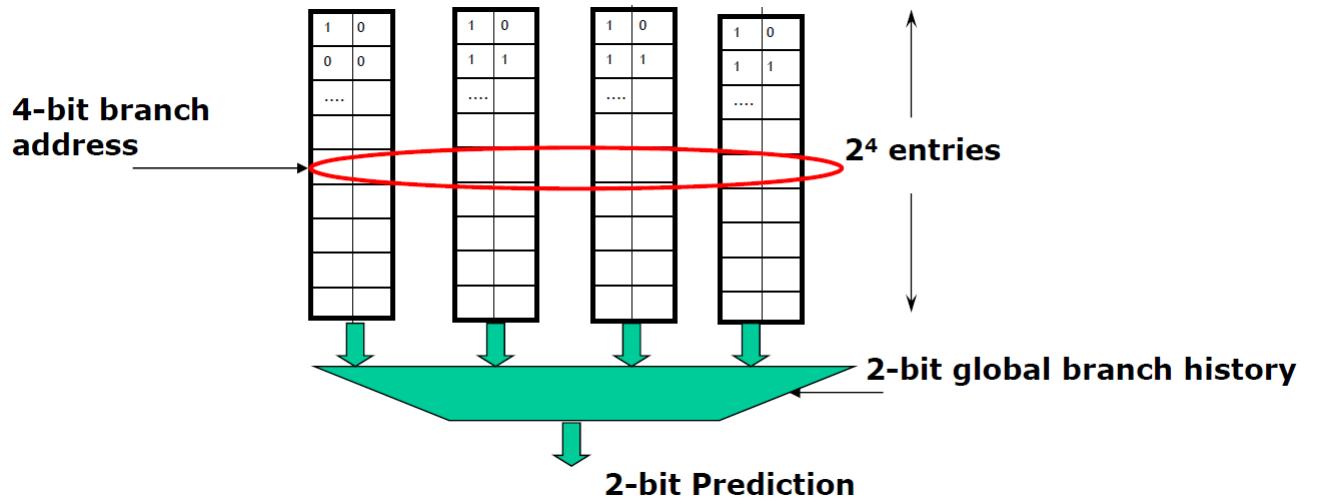


Figura 38: Esempio di predittore correlato (2,2)

non è altro che un predittore correlato con i valori (0,2); a questo punto possiamo confrontare le performance nel caso di un predittore semplice a 2 bit con una tabella di 4K entità e un predittore correlato (2,2) con tabelle di 1K entità. Come vediamo dal grafico in Figura 39 il predittore correlato è, in molti casi più efficace di un predittore semplice mentre nei casi peggiori ne egualga le performance. Un'altra tecnica di predizione dinamica è quella del **Predittore di salto adattativo a due livelli** (*Two-Level Adaptative Branch Predictors*) nel quale un primo livello di storia viene memorizzato in uno shift register a  $k$  bit chiamato **Branch History Register (BHR)** il quale memorizza i risultati degli ultimi  $k$  salti. Il secondo livello di storicizzazione è

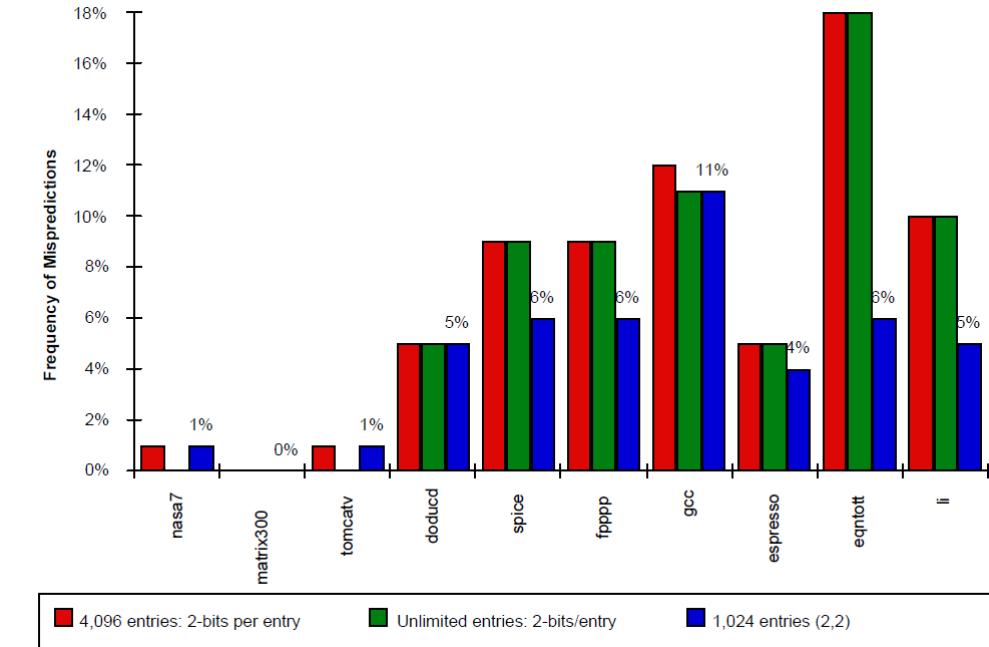


Figura 39: Comparazione delle performance per predittore non correlato e predittore corellato

memorizzato in una tabella con record di due bit chiamata **Pattern History Table (PHT)** per indicare la predizione. La BHR è utilizzata per indicizzare la PHT e selezionare i due bit da utilizzare; per selezionare quale dei due bit utilizzare si utilizza lo stesso principio utilizzato per i predittori a due bit semplici. Un'evoluzione di questo predittore è il **predittore GShare** dove le informazioni dell'indirizzo locali vengono correlate con quelle globali tramite un'operazione di XOR. Un ultimo elemento importante nella predizione dinamica è il **Branch Target Buffer** la quale è una cache nella quale vengono memorizzate l'indirizzo di destinazione del salto per le istruzioni dopo il salto. Accediamo al BTB nello stage IF utilizzando l'indirizzo dell'istruzione da prelevare per indicizzare la cache. Un possibile esempio di record è mostrato in Figura 42. L'indirizzo di destinazione del salto è espresso sempre in modo relativo al PC. La struttura del BTB è mostrata in Figura 43; in tale buffer dobbiamo memorizzare solo gli indirizzi per quei salti che vengono eseguiti.

### 2.3 Speculazione

Senza tecniche di predizione di salto il parallelismo risulta molto limitato e si riduce all'analisi dei **basic block** ovvero a pezzi di codice nei quali non entrano o escono dei salti. Tuttavia possiamo azzardare alcune supposizioni di parallelismo tra diversi blocchi base effettuando delle speculazioni. Tramite le speculazioni possiamo recuperare ed eseguire le istruzioni come se le nostre predizioni siano corrette gestendo in seguito il caso in cui non siano corrette. Tale speculazione può essere supportata sia dal compilatore che dall'hardware.

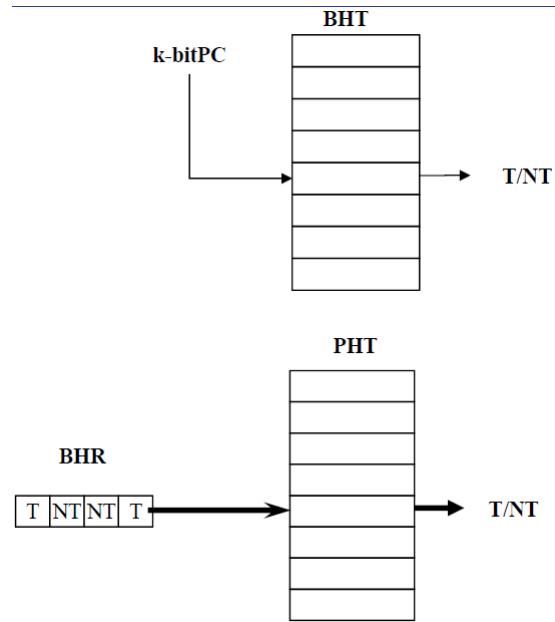


Figura 40: Esempio di predittore adattativo

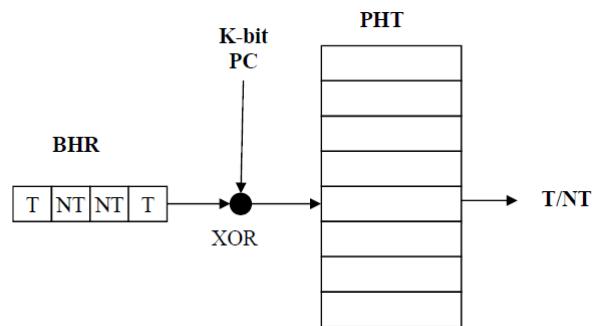


Figura 41: Esempio di predittore GShare

Exact Address of a Branch	Predicted target address

Figura 42: Esempio di record di un Branch Target Buffer

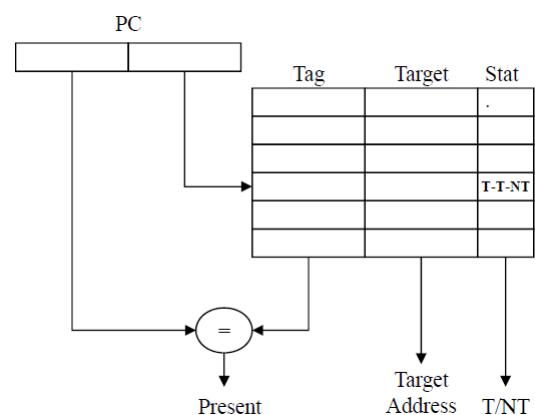


Figura 43: Struttura di un Branch Target Buffer

### 3 Instruction Level Parallelism

Come abbiamo visto nei capitoli precedenti in una macchina fornita di pipeline possiamo dimostrare come il numero di clock necessari per eseguire un'istruzione sia:

$$CPI_{pipeline} = CPI_{ideale} + \text{Stalli Strutturali} + \text{Stalli Data Hazard} + \\ \text{Stalli di Controllo} + \text{Stalli di Memoria}$$

La riduzione di uno qualsiasi dei termini sulla destra da in modo che  $CPI_{pipeline}$  si avvicini sempre più al  $CPI_{ideale}$ . Il caso migliore si ha quando il throughput è massimo ed è uguale a 1

$$IPC_{ideale} = 1; CPI_{ideale} = 1$$

Tuttavia si hanno dei limiti alle performance dovuti ai diversi tipi di *rischi (Hazards)*, questi possono essere di diversa natura:

- **Strutturali:** si possono risolvere tramite l'introduzione di nuovo hardware.
- **Dati:** necessitano di *forwarding* e di una schedulazione del codice a livello di compilazione.
- **Controllo:** Early evaluation, Branch Delay Slot, Predizione dei salti statica e dinamica.

Inoltre per le pipeline più profonde (superpipelining) questi problemi sono accentuati.  
In questo capitolo ci concentreremo su come incrementare il  $CPI$  oltre il valore ideale; per fare ciò però dobbiamo prima capire quali sono gli *hazard* sui dati che possiamo incontrare.

#### 3.1 Tipi di Hazards sui dati

Gli hazard innanzitutto sono quei conflitti che avvengono a tempo di esecuzione e sono generati da dipendenze a livello di istruzione. Consideriamo l'esecuzione di un'istruzione generale di questo tipo:

$$r_k \leftarrow (r_i) \text{ op } (r_j)$$

Possiamo avere tre tipi di dipendenza a livello di istruzione come mostrato in Figura 44

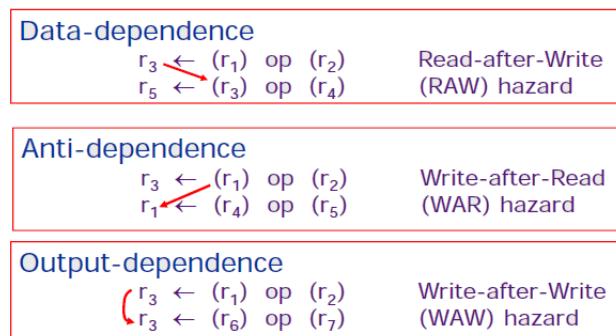


Figura 44: Esempi di dipendenza dei dati

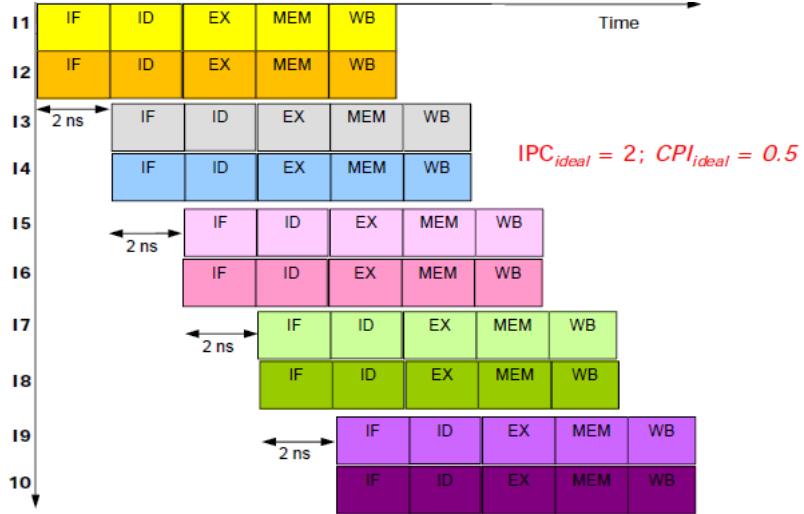


Figura 45: Esecuzione di istruzione in una pipeline dual-issue

### 3.2 Parallelismo a livello di istruzione

Per raggiungere livelli di performance maggiori è necessario estrarre dai programmi maggiore parallelismo, questo si traduce in pipeline con più uscite (*multiple-issue*). Per fare ciò è necessario individuare e risolvere le dipendenze, inoltre è utile riordinare (*schedule*) le istruzioni in modo da ottenere un maggiore parallelismo a tempo di esecuzione compatibilmente con le risorse a disposizione.

Per dare una definizione formale del parallelismo a livello di istruzione (*ILP*) possiamo dire che

$$ILP = Sfruttare\ la\ potenziale\ esecuzione\ sovrapposta\ di\ istruzione\ non\ correlate$$

Tale sovrapposizione è possibile tutte le volte che:

- Non abbiamo degli *Hazard* di tipo strutturale
- Non abbiamo *Hazard* di tipo RAW, WAR oppure WAW
- Non abbiamo *Hazard* di controllo

Un esempio di sistema *dual-issue* è mostrato in Figura 45 e Figura 46. In pipeline di tipo multiple-issue il  $CPI_{ideale}$  risulta essere  $< 1$  ad esempio considerando il caso ottimale di un processore *2-issue* abbiamo che ad ogni ciclo di clock vengono completate 2 istruzioni questo significa che:

$$IPC_{ideale} = 2; CPI_{ideale} = 0.5$$

Uno degli aspetti più critici nel caso di ILP è la determinazione delle dipendenze tra le istruzioni; infatti, se due istruzioni sono dipendenti tra loro esse non possono essere eseguite in parallelo ma dovranno essere eseguite in sequenza o al più in parziale sovrapposizione. Possiamo distinguere tre tipi di dipendenza:

- Dipendenza dei dati (Vera dipendenza)
- Dipendenza dei nomi
- Dipendenza di controllo

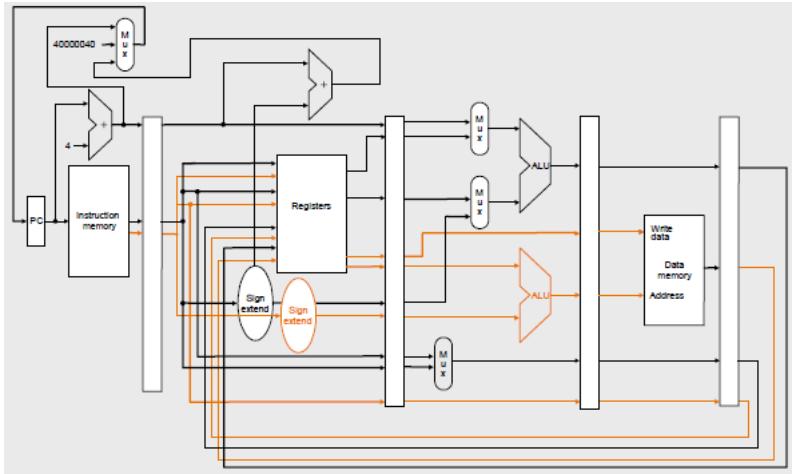


Figura 46: Schema hardware per una pipeline dual-issue con una unità ALU/BR e una unità load/store

**Dipendenza dei nomi** Tale dipendenza avviene quando due istruzioni usano lo stesso registro o la stessa area di memoria (chiamata *nome*) ma non esiste alcun flusso di dati tra queste due istruzioni. Possiamo individuare due tipi di dipendenza dai nomi tra due istruzioni **i** e **j** nelle quali **i** precede **j**

- **Antidipendenza:** quando l'istruzione *j* scrive un registro che l'istruzione **i** legge; l'ordine originale delle istruzioni deve essere preservato per essere sicuri che **i** legga il valore corretto (WAR).
- **Output dipendenza:** quando **i** e **j** scrivono lo stesso registro o la stessa area di memoria; l'ordine delle istruzioni deve essere rispettato per essere sicuri che il valore finale sia quello scritto da **j**

In realtà la dipendenza dai nomi non è una vera e propria dipendenza in quanto non vi è alcuno scambio di valori tra le istruzioni; se il *nome* usato nell'istruzione può essere cambiato non ci sono conflitti. Tuttavia per quanto riguarda le aree di memoria è molto più difficile localizzare questo tipo di conflitto infatti due indirizzi possono essere diversi ma puntare alla stessa area (*memory disambiguation*) mentre una rinominazione dei registri risulta molto più semplice. La rinominazione può essere effettuata sia staticamente dal compilatore sia in modo dinamico dall'hardware.

**Dipendenze dei dati** Le dipendenze dei dati possono potenzialmente generare dei *Data Hazard* ma l'impatto che questi hazard hanno in termini di stalli e tecniche di eliminazione degli stalli sono caratteristiche specifiche della pipeline e non dipendono dalla dipendenza. Le *RAW* sono le uniche vere dipendenze sui dati che abbiamo. Le dipendenze sono una caratteristica del programma mentre gli *hazard* sono specifiche della pipeline.

**Dipendenze di controllo** Le dipendenze di controllo sono determinate dall'ordinamento delle istruzioni e sono preservate da due proprietà:

- Le istruzioni devono essere eseguite nell'ordine del programma per assicurare che un'istruzione che si trova prima di un salto venga eseguita prima del salto.

- Individuazione degli *hazard di controllo* per assicurare che un’istruzione che dipende da un salto non sia eseguita prima di conoscere la direzione del salto.

Sebbene preservare il controllo delle dipendenze sia un modo semplice per preservare l’ordine del programma esso non è così essenziale da dover essere preservato.

### 3.2.1 ILP in pratica

Dalla trattazione appena fatta possiamo ricavare due proprietà importanti per verificare la correttezza di un programma (e che in realtà preservano sia le dipendenze dei dati che quelle di controllo):

- **Data flow:** il flusso dei valori dei dati attraverso le istruzioni deve produrre il risultato corretto.
- **Exception behavior:** preservare il comportamento delle eccezioni che significa che qualsiasi cambiamento nell’ordine di esecuzione delle istruzioni non deve cambiare come le eccezioni sono sollevate dal programma.

Esistono due tecniche fondamentali per supportare e implementare l’*ILP*, lo **scheduling dinamico** che dipende dall’hardware per localizzare il parallelismo e lo **scheduling statico** che fa affidamento sul software per individuare possibili parallelismi. La prima soluzione è quella più utilizzata in ambito desktop e server.

Consideriamo ora un processore di tipo *single-issue* lo stage IF precede quello EXE e le istruzioni possono essere prelevate sia dall’*Instruction Register* sia da una coda di istruzioni pendenti. La fase di esecuzione può richiedere più cicli di clock in base al tipo di operazione.

**Scheduling dinamico** Il problema principale è quello che non si può nascondere una dipendenza dai dati senza causare uno stallo nell’esecuzione della pipeline. Una soluzione è quella di permettere alle istruzioni situate dopo lo stallo di procedere; l’hardware deve riordinare dinamicamente l’esecuzione delle istruzioni per dar in modo di ridurre gli stalli. Per fare ciò è necessario permettere un’esecuzione fuori ordine e una fase di commit finale.

L’hardware riordina l’esecuzione delle istruzioni per ridurre il numero degli stalli della pipeline mentre mantiene il *data flow* e *exception behavior*. I vantaggi principali di questa tecnica sono il fatto di permettere una gestione di alcuni casi in cui esistono delle dipendenze non note al tempo della compilazione, inoltre permette di semplificare la complessità del compilatore ed infine permette di compilare il codice affinché esso venga eseguito efficientemente anche su pipeline diverse. Questi vantaggi tuttavia hanno un costo che comporta un significativo incremento della complessità dell’hardware, un incremento dei consumi e può generare delle eccezioni imprecise. Tale soluzione è utilizzata soprattutto nei *Processori Super-scalari*

**Scheduling statico** In questo caso il compilatore utilizza dei sofisticati algoritmi per individuare e organizzare il codice in modo da sfruttare l’*ILP*, per fare ciò analizza i **basic block** e ricerca il parallelismo in questi seppur esso sia minimo (15% - 25%), ed inoltre sfrutta il parallelismo tra diversi *basic block*. Un limite a questa tecnica è però la dipendenza dai dati presente nei vari blocchi base. Tipicamente questa tecnica è sfruttata dai processori **VLIW** (*Very Long Instruction Word*) i quali si aspettano del codice privo di dipendenze dal compilatore.

Il limite principale di questa tecnica è dato dall’impredicibilità dei salti, dalla latenza della memoria, dalla dimensione del codice e dalla complessità del compilatore.

### 3.2.2 Esecuzione super-scalare e VLIW

Passando al caso *multi-issue* la questione si complica anche se le prestazioni migliorano come vediamo in Figura 47, infatti si vogliono eseguire più istruzioni per ogni ciclo di clock; per fare ciò è necessario prelevare più istruzioni per ciclo dall'IR e questo dipende dalla banda a disposizione. La cosa più difficile però risulta essere l'analisi delle dipendenze dei dati e di controllo per le istruzioni da eseguire. In Figura 48 vediamo come è strutturato un processore super-scalare nel

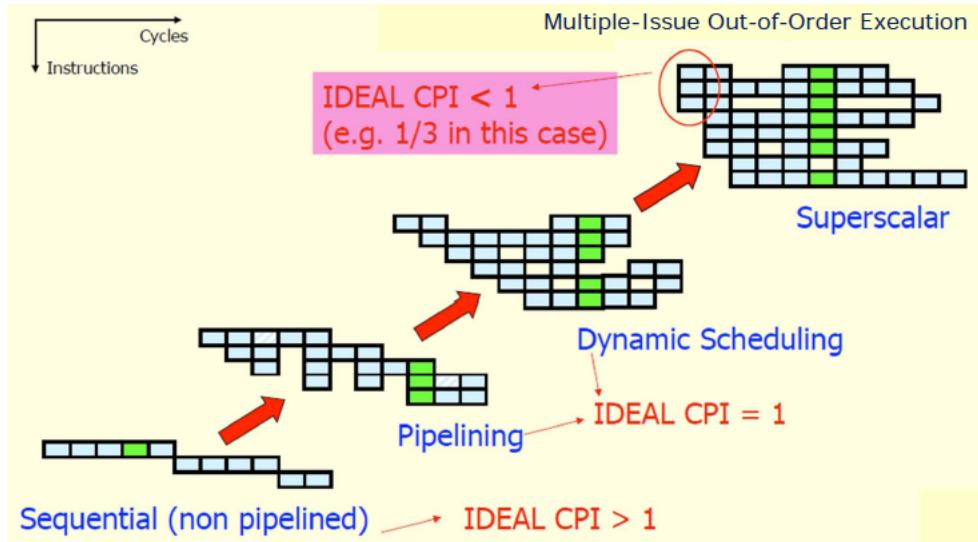


Figura 47: Confronto di prestazioni tra architetture

quale possiamo notare le diverse unità per l'esecuzione di istruzioni parallele come le due ALU o l'unità per le istruzioni di load/store, e l'unità per il riordino delle istruzioni.

Per decidere quali istruzioni mandare in esecuzione si utilizza lo *Scheduler dinamico* il quale per

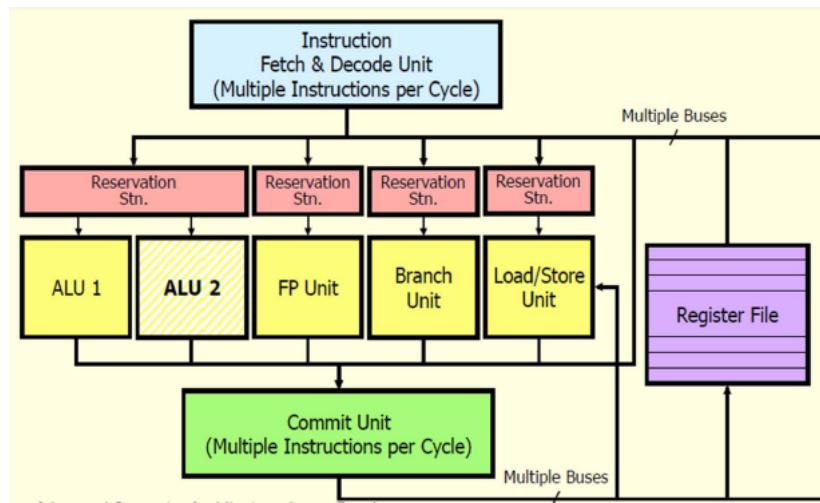


Figura 48: Struttura di un processore superscalare

ogni ciclo di clock analizza quali sono le dipendenze e per fare ciò la sua complessità è molto alta, nell'ordine del quadrato delle possibili istruzioni come mostrato in Figura 49. Esiste un limite al numero di istruzioni che possono essere analizzate durante un singolo ciclo di clock, infatti

In-flight Instructions ( $kR$ )												
Fetched Instructions To Execute ( $R$ )												
✗	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓
✗	✓	✓	✓	✗	✓	✓	✗	✓	✓	✗	✓	✗
✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓
✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Figura 49: Tabella delle dipendenze in uno scheduler dinamico

i limiti principali dei processori super-scalari riguardano tutti lo scheduler dinamico in quanto esso è molto costoso in termini di area in quanto la sua logica è molto complessa, il tempo di clock dipende dal tempo di analisi delle istruzioni ed infine la verifica del design dello scheduler è molto complessa. Queste limitazioni portano a delle limitazioni in termini di istruzioni che possono essere eseguite simultaneamente. Attualmente in realtà esistono dei processori a 6-issue anche se molto rari, più normalmente sono di tipo 4-issue o minori in quanto è troppo difficile trovare 8 o addirittura 16 istruzioni indipendenti in un singolo ciclo.

La famiglia di processori multi-issue che sfruttano lo scheduler statico sono, invece, i processori VLIW (*Very Long Instruction Word*) nei quali è il compilatore che decide cosa far fare ad ogni istruzione per ogni ciclo di clock come si vede in Figura 50. I processori di tipo super-scalare

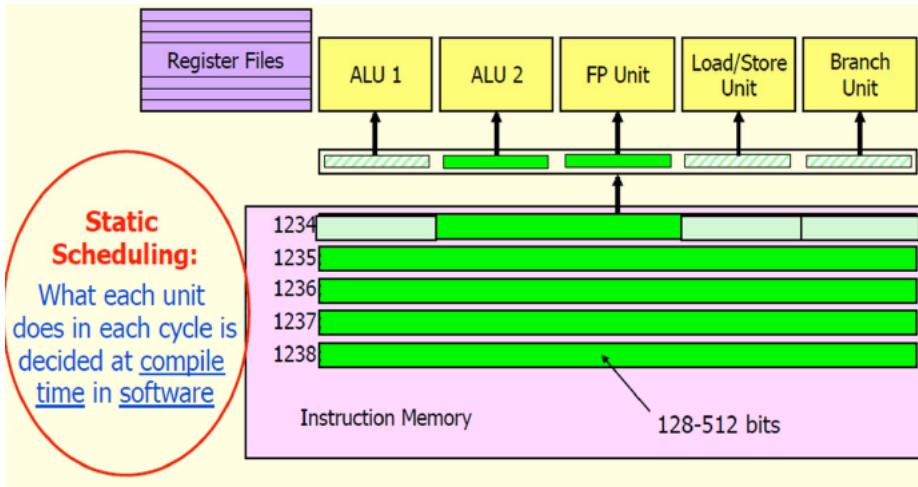


Figura 50: Scheduler statico nel caso di processori VLIW

sono utilizzati soprattutto in ambito desktop e server, mentre i processori VLIW sono utilizzati soprattutto in ambito embedded in quanto la decisione sull'esecuzione è presa a compile-time e non è necessario aggiungere dell'area per lo scheduler dinamico.

### 3.3 Scoreboard

Come abbiamo visto fino ad ora lo scheduling dinamico è il meccanismo più utilizzato nei sistemi general purpose per sfruttare il parallelismo tra le istruzioni. Tale meccanismo però ha diverse tecniche di implementazione, una di queste è lo **Scoreboard**. Lo *Scoreboard* divide il normale stage ID in due stage per permettere l'esecuzione delle istruzioni nell'ordine più efficiente. Questi due stage sono:

- Lo stage *issue* che si occupa di decodificare le istruzioni e di verificare eventuali hazard strutturali.
- Lo stage *read operands* (RR) che si occupa di risolvere i data hazard ritardando eventualmente la lettura dei registri.

Grazie a questo meccanismo le istruzioni vengono eseguite quando non hanno alcuna dipendenza o non esistono degli hazard strutturali non verifica però una eventuale priorità delle istruzioni. Per spiegare la struttura dello *Scoreboard* dobbiamo innanzitutto definire quando un'istruzione si dice in **esecuzione**. Possiamo distinguere quando un'istruzione inizia la sua esecuzione e quando essa la termina durante questi due istanti l'istruzione si dice in *esecuzione*. In una pipeline possono esistere molte istruzioni in esecuzione nello stesso momento, questo richiede che esistano più unità funzionali. Quello che noi analizzeremo è la struttura di un **CDC 6600** nel quale le istruzioni vengono immesse in ordine ma vengono eseguite e completate in un ordine casuale. L'architettura di questo sistema è mostrata in Figura 51 Nella pipeline di uno scoreboard

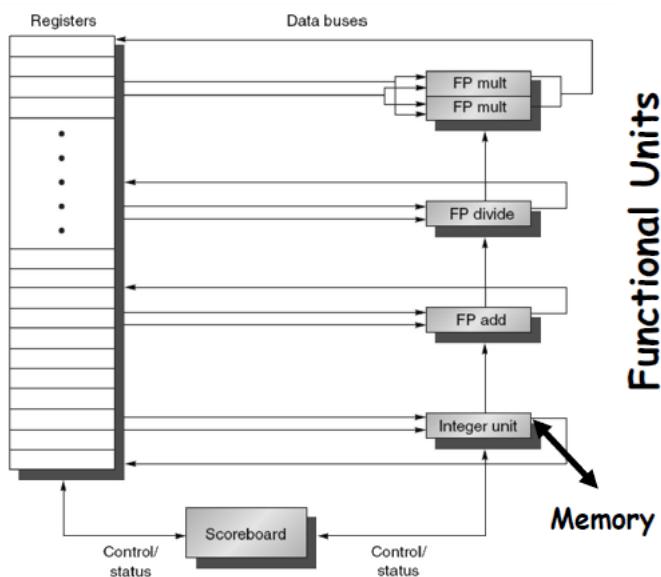


Figura 51: Architettura di un sistema *Scoreboard*

le fasi ID, EXE e WB sono sostituite da quattro stage. Lo stage ID è diviso in due parti, la prima chiamata **Issue** la quale decodifica l'istruzione e controlla eventuali hazard strutturali, la seconda chiamata **Read Operands** che attende fino alla risoluzione degli hazard sui dati. Lo scoreboard fa in modo che le istruzioni eseguite siano prive di dipendenze. Come abbiamo visto le istruzioni vengono prelevate in ordine dallo stage *Issue* ma da quell'istante esse possono essere ordinate in qualsiasi modo, infatti durante lo stage *read operands* esse possono essere bloccate o bypassate, inoltre, la latenza di esecuzione può variare tra le diverse unità funzionali.

Un completamento fuori ordine può portare però a hazard di tipo WAR o WAW che tuttavia possono essere facilmente risolti infatti per i WAR è sufficiente:

- Bloccare il *write back* finché i registri non vengono letti
- Effettuare la lettura dei registri soltanto durante la fase di *Read Operands*

Mentre per risolvere i WAW è sufficiente individuare l'hazard e bloccare l'esecuzione delle istruzioni dipendenti successive fino a quando esse non vengono concluse.

L'individuazione e la risoluzione degli hazard è centralizzata nello *scoreboard*; ogni istruzione attraversa lo scoreboard dove viene aggiornata una tabella delle dipendenze, a questo punto lo scoreboard determina quando l'istruzione può leggere i registri ed iniziare la sua esecuzione. Se un'istruzione non può cominciare immediatamente la sua esecuzione lo scoreboard monitora ogni cambiamento e decide quando l'istruzione può andare in esecuzione. Infine, lo scoreboard, controlla quando l'istruzione scrive il risultato dentro i registri di destinazione.

Un problema che si viene a creare quando si accetta il completamento delle istruzioni fuori ordine è quello della preservazione del comportamento delle eccezioni; una soluzione è quella di assicurarsi che nessuna istruzione possa generare una eccezione finché il processore non conosce esattamente l'istruzione che ha sollevato l'eccezione. Un'eccezione si dice *imprecisa* se lo stato del processo nell'istante in cui viene sollevata una eccezione non è uguale a quello del caso in cui l'istruzione sia eseguita in ordine; un'eccezione imprecisa si può verificare quando la pipeline ha già completato delle istruzioni che sequenzialmente si trovano dopo l'istruzione che ha sollevato l'eccezione, oppure se non ha ancora completato istruzioni che sequenzialmente la precedono.

**I quattro stadi dello Scoreboard** Analizziamo ora in dettaglio i quattro stadi dello Scoreboard. Il primo stadio è quello di *Issue* nel quale le istruzioni vengono decodificate e si verificano gli eventuali hazard strutturali e quelli di WAW. Le istruzioni lasciano questo stage in ordine, se l'unità funzionale necessaria per eseguire l'istruzione è libera e non esistono altre istruzioni che hanno lo stesso registro di destinazione (no WAW) lo stage issue inoltra l'istruzione all'unità funzionale e aggiorna la sua struttura interna. In caso contrario lo stage ferma l'istruzione fino a quando tutti gli hazard sono stati risolti. Ogni qualvolta lo stage di issue ferma un'istruzione il buffer tra IF e Issue si riempie. Nel caso il buffer sia singolo IF si blocca e attende l'Issue nel caso invece il buffer sia a coda l'IF si blocca solo nel caso in cui la coda sia piena.

Lo stage *Read Operands* attende la risoluzione dei data hazard e legge i registri. Un registro risulta disponibile quando non ci sono istruzioni attive precedenti che scrivono su di esso oppure se un'unità funzionale scrive il suo risultato in tale registro. Quando i registri sono disponibili lo scoreboard comunica all'unità funzionale di leggere i registri. Gli RAW hazard vengono risolti dinamicamente in questa fase.

Lo stage *execution* è lo stage composto dalle unità funzionali; durante questo stage le unità funzionali svolgono le diverse operazioni, quando il risultato è pronto viene notificato allo scoreboard. Le unità funzionali sono caratterizzate da una latenza variabile, il tempo per iniziare l'esecuzione è variabile in quanto è il tempo necessario per leggere i diversi registri, i tempi per effettuare una load/store variano in base ai cache HIT/MISS.

L'ultimo stage è denominato *Write Result*, in questa fase si controllano eventuali hazard di tipo WAR e si termina l'esecuzione scrivendo il risultato nel registro di destinazione.

**Lo Scoreboard in pratica** Lo scoreboard è formato da tre unità fondamentali:

- Instruction status
- Functional Unit status: il quale indica lo stato delle unità fondamentali tramite diversi indici

**Busy:** indica se la FU è occupata oppure no

**Op:** indica l'operazione da eseguire

**$F_i$ :** registro di destinazione

**$F_j, F_k$ :** registri sorgenti

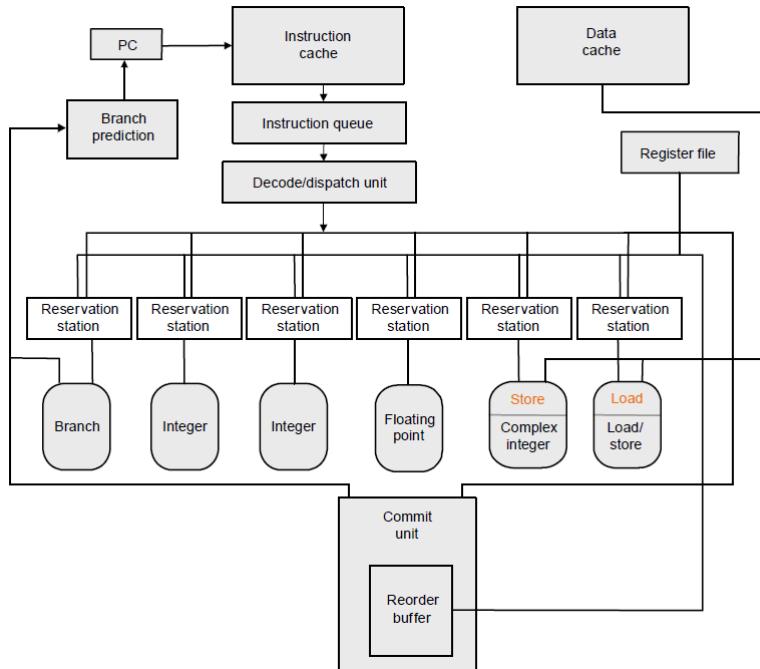


Figura 52: Architettura per l'algoritmo di Tomasulo

$Q_j, Q_k$ : indicano quale FU sta tenendo occupato il registro corrispondente

$R_j, R_k$ : sono dei flag che indicano lo stato dei registri sorgenti.

- Register Result status: indica quale FU dovrà scrivere sui registri di destinazione.

Per un esempio completo si rimanda alle slide dell'insegnante.

### 3.4 Algoritmo di Tomasulo

Un altro tipo di algoritmo da sfruttare per lo scheduling dinamico è l'algoritmo di *Tomasulo* introdotto da IBM tre anni dopo il CDC 6600 lo scopo di tale algoritmo è sempre quello di ottenere prestazioni elevate senza l'utilizzo di speciali compilatori.

A differenza dello Scoreboard tuttavia il meccanismo di controllo e di buffer è distribuito sulle diverse unità funzionali. I buffer associati alle unità funzionali sono chiamate *Reservation Station*. I registri nelle istruzioni sono sostituiti da valori o puntatori alle reservation station è possibile così il renaming dei registri. In questo modo si evitano anche eventuali hazard di tipi WAR e WAW; inoltre esistono più *RS* che registri e questo permette delle ottimizzazioni che il compilatore non può fare. Infine il risultato delle FU non transita dai registri bensì viene diffuso a tutte le altre FU tramite un *Common Data Bus*. Le operazioni di Load e Store sono trattate come normali FU . L'architettura necessaria per implementare l'algoritmo di Tomasulo è mostrata in Figura 52 mentre la struttura di una unità funzionale è mostrata in Figura 53 I vari campi che compongono una reservation station sono:

**Tag:** identifica quale RS è coinvolta.

**Busy:** identifica se la RS è occupata.

**OP:** Identifica il tipo di operazione eseguita dal componente.

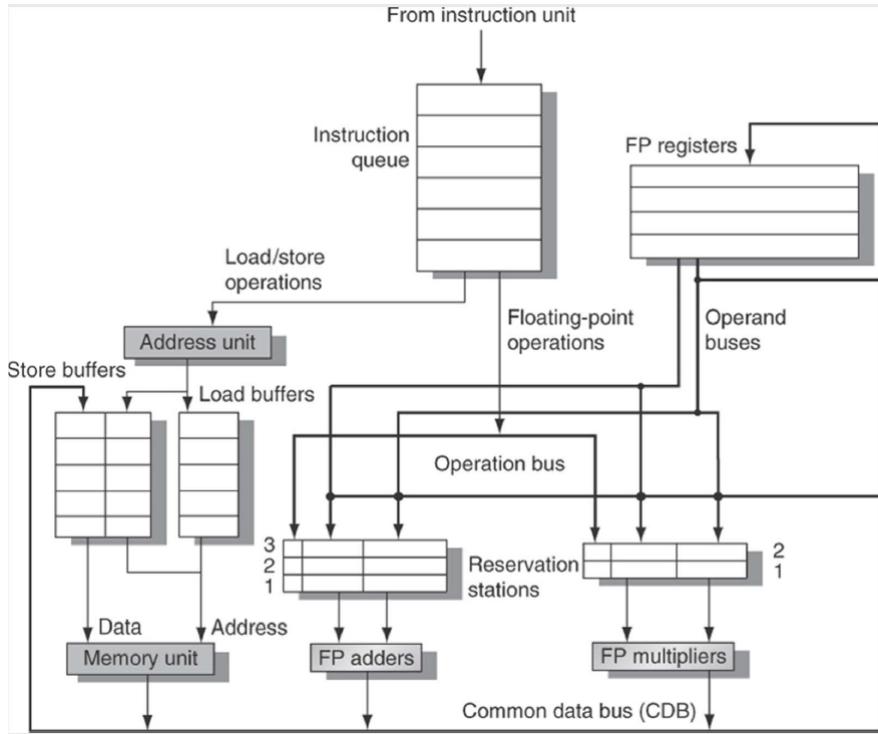


Figura 53: Architettura di un'unità funzionale

$V_j, V_k$ : Valori contenuti nei registri; per la *load*  $V_j$  contiene il valore dell'offset.

$Q_j, Q_k$ : Puntatori alle RS che producono i valori  $V_j, V_k$  se il valore è uguale a 0 l'operando è già disponibile.

Si noti che solo uno dei due campi  $V$  e  $Q$  è disponibile per ogni operando in un determinato istante.

Il register file e lo Store buffer hanno un campo *Value* ( $V$ ) e uno *Puntatore* ( $Q$ ) il quale punta al numero della *reservation station* che produce il risultato da immagazzinare, se questo puntatore è uguale a zero allora non ci sono istruzioni attive e il valore contenuto nel RF/Buffer è quello corretto. Nel caso di Load buffer abbiamo un campo *Address* ( $A$ ) e un campo *Busy*; il campo  $A$  mantiene le informazioni riguardanti gli indirizzi calcolati nelle operazioni di Load/Store, all'inizio contiene le informazioni sull'offset poi, una volta calcolato, contiene l'indirizzo effettivo.

### 3.4.1 Gli stadi dell'algoritmo di Tomasulo

**Il primo stadio: ISSUE** Durante questo stadio si preleva un'istruzione  $I$  dalla testa della coda delle istruzioni (**in-order issue**). Si controlla se la RS associata a quell'istruzione è vuota altrimenti si attende (controllo sugli structural hazard). Nel caso gli operandi non siano ancora disponibili si tiene traccia della FU che produce tali operandi (puntatore  $Q$ ). Durante questa fase si effettua anche il *Rename* dei registri in modo tale da evitare eventuali *WAR*, infatti, se un'istruzione  $I$  scrive un registro  $R_x$  e un'istruzione  $K$  già prelevata legge tale registro  $K$  conosce già il valore di  $R_x$  e lo ha immagazzinato nella sua RS oppure conosce quale operazione produce tale valore. Inoltre si evitano anche eventuali *WAW* in quanto le istruzioni sono prelevate in ordine.

**Secondo stadio: Esecuzione** Quando entrambi gli operandi sono disponibili si esegue l'operazione; nel caso in cui non siano pronti invece si controlla il *Common Data Bus* in attesa del risultato; ritardando l'esecuzione si evitano eventuali RAW. Si noti come più istruzioni possono diventare eseguibili allo stesso istante per la stessa FU a questo punto bisogna verificare la disponibilità dell'unità funzionale. Le RAW hazard sono molto meno incisive in quanto sono gestite a livello di RS e non è necessario attendere il *write back* sul register file.

Per quanto riguarda le istruzioni di load e di store l'esecuzione avviene in due passi, nel primo passo si calcola l'effettivo indirizzo di destinazione e si memorizza nel buffer, nel secondo passo in caso di load si esegue l'operazione non appena l'unità è disponibile nel caso della store si attende invece che il dato da immagazzinare sia disponibile.

Per preservare il comportamento dell'esecuzione nessuna istruzione può cominciare l'esecuzione fino a quando i branch precedenti non sono stati eseguiti. Se si usano tecniche di predizione la CPU deve conoscere se la predizione è corretta prima di procedere.

**Terzo stadio: Write result** Quando un risultato è disponibile esso viene scritto sul *Common Data Bus* e da qui copiato sia sul RF sia su tutte le RS che attendono questo risultato. Il *Common Data Bus* è un bus di tipo data+source composto da 64 bit di dati e 4 bit per la sorgente in questo modo le FU possono effettuare un lookup associativo.

#### 3.4.2 Alcuni dettagli

Le operazioni di load e di store attraversano un'unità funzionale per il calcolo dell'indirizzo effettivo prima di procedere con le vere e proprie operazioni di load e di store. La load necessita di una seconda fase per accedere alla memoria mentre invia il risultato al RF e alle RS durante lo stage *Write Result*. La store, invece completa la sua esecuzione durante lo stage *Write Result* nel quale scrive i dati sul buffer. Le operazioni di load e di store possono essere eseguite in ordine differente purché esse accedano a differenti aree di memoria, in caso contrario possono presentarsi problemi di WAR (scambio tra load e store), di RAW (scambio tra store e load) oppure di WAW (scambio tra due store) invece le load possono essere scambiate liberamente. Per identificare questo tipo di anomalie il calcolo degli indirizzi di tutte le operazioni deve essere calcolato dalla CPU e quindi secondo l'ordine del programma.

Prendiamo il caso di una load eseguita fuori ordine con una store precedente ed assumiamo che il calcolo dell'indirizzo sia eseguito con l'ordine del programma. Quando l'indirizzo della load è calcolato esso viene comparato con i campi *A* dello *Store Buffer* nel caso vi sia un match la load non viene inviata allo Load Buffer fino a quando il conflitto non è risolto. Le operazioni di store invece verificano la presenza di conflitti sia nello Store Buffer che nel Load Buffer.

#### 3.4.3 Tomasulo in pratica

Per un esempio completo si rimanda alle slide dell'insegnante per il funzionamento dell'algoritmo.

#### 3.4.4 Tomasulo vs Scoreboard

Al contrario dello scoreboard l'algoritmo di Tomasulo ha una finestra di prelevamento delle istruzioni minore (5 vs 12) in entrambi i casi non si hanno hazard di tipo strutturale nel prelevamento delle istruzioni, nel caso di Tomasulo questi sono bloccati a livello di RS mentre nel caso dello Scoreboard a livello di FU. Tomasulo è più efficiente per quanto riguarda la risoluzione di WAW e di WAR che vengono risolti tramite renaming mentre per lo Scoreboard sono necessari alcuni stalli; inoltre, in Tomasulo il risultato di una FU è distribuito a tutte le altre tramite

il Common Data Bus mentre nello scoreboard è necessario attendere che il risultato sia scritto nei registri di destinazione. Il controllo in Tomasulo è distribuito ed è possibile effettuare un loop unrolling al contrario dello Scoreboard. Tuttavia i limiti di Tomasulo risiedono nella sua complessità e alla limitazione delle prestazioni del CCommon Data Bus; inoltre il parallelismo è ridotto a causa dei salti.

## 3.5 Register Renaming

### 3.5.1 Renaming implicito

Analizziamo ora un piccolo codice di esempio che rappresenta un ciclo

```
Loop: LD F0 0 R1
      MULTD F4 F0 F2
      SD F4 0 R1
      SUBI R1 R1 #8
      BNEZ R1 Loop
```

ed assumiamo che la moltiplicazione abbia una latenza di 4 cicli, che la prima volta che viene effettuata la load si abbia un overhead di 8 cicli (cache miss) mentre nelle successive la latenza sia di 1 ciclo (cache hit), infine, assumiamo che la branch prediction predica che il salto sia effettuato.

Come abbiamo visto nel paragrafo precedente l'algoritmo di Tomasulo fornisce il *register renaming* in modo implicito tramite le *reservation station* le quali bufferizzano gli operandi delle istruzioni per evitare eventuali problemi di WAW e di WAR. Questo permette a iterazioni diverse di utilizzare registri fisici diversi (**dynamic loop unrolling**), inoltre, permette di sostituire i registri statici con dei puntatori dinamici che fanno in modo di incrementare praticamente la dimensione del register file. Questo permette alle istruzioni di procedere e, tramite l'uso della branch prediction di prelevare più istruzioni di iterazioni diverse.

Per un esempio di questo meccanismo si rimanda alle slide della professoressa riportiamo di seguito solo il passaggio dell'algoritmo al ciclo di clock 14 (Figura 54) nella quale si nota come sia stato possibile effettuare un loop unrolling e come tale loop venga gestito in modo implicito infatti le operazioni del primo loop hanno tutte come registro base  $R1 = 80$  mentre il secondo loop abbia come registro base  $R1 = 72$ .

Il problema di questa tecnica è che necessita di prelevare le istruzioni *in ordine* in quanto un prelevamento fuori ordine ci può portare ad avere WAR e RAW che in realtà non esistono. Tuttavia il meccanismo funziona bene nel caso di prelevamento di un'unica istruzione. La situazione cambia completamente nel caso di prelevamento di istruzioni multiple in un singolo ciclo di clock, infatti, è necessario disporre di porte multiple per la *rename table* e dobbiamo essere in grado di effettuare il rename su diverse istruzioni contemporaneamente. Dobbiamo inoltre fornire istruzioni a più *reservation station* nello stesso ciclo di clock e questo comporta l'utilizzo di 2x porte in lettura e 1x porta in scrittura. Il prelevamento delle istruzioni in sequenza è il vero collo di bottiglia nel caso di istruzioni multiple per singolo ciclo di clock.

Il completamento fuori ordine riduce notevolmente la nostra possibilità di avere delle eccezioni *precise* in quanto il register file può contenere risultati di istruzioni successive e magari non contenere risultati di istruzioni precedenti e non ancora completate. In questo contesto sarebbe necessario effettuare un *rollback* del register file in modo da avere un'eccezione *precisa* ovvero nella quale tutte le istruzioni precedenti a quella che ha generato l'eccezione hanno committato il loro risultato e nessuna istruzione successiva ha committato il risultato

Instruction status:				Exec Write			Busy Addr	Fu
ITER	Instruction	j	k	Issue	CompResult			
1	LD F0	0	R1	1	9	10	Load1	No
1	MULTD F4	F0	F2	2	14		Load2	No
1	SD F4	0	R1	3			Load3	Yes 64
2	LD F0	0	R1	6	10	11	Store1	Yes 80 Mult1
2	MULTD F4	F0	F2	7			Store2	Yes 72 Mult2
2	SD F4	0	R1	8			Store3	No

Reservation Stations:				S1	S2	RS	Code:
Time	Name	Busy	Op	Vj	Vk	Qj	Qk
	Add1	No					LD F0 0 R1
	Add2	No					MULTD F4 F0 F2
	Add3	No					SD F4 0 R1
0	Mult1	Yes	Multd M[80] R(F2)				SUBI R1 R1 #8
1	Mult2	Yes	Multd M[72] R(F2)				BNEZ R1 Loop

Clock	R1	F0	F2	F4	F6	F8	F10	F12	...	F30
14	64	Fu	Load3	Mult2						

Figura 54: Algoritmo di Tomasulo all’istante 14

### 3.5.2 Renaming esplicito

Come abbiamo appena visto tomasulo fornisce il renaming in modo implicito tramite l’utilizzo delle reservation station. Quello che vogliamo fare ora è capire come implementare un renaming dei registri in modo esplicito. Per fare questo innanzitutto dobbiamo tener conto che dobbiamo utilizzare un maggior numero di registri fisici rispetto a quelli specificati dall’ISA. Il principio chiave è quello di allocare una nuova destinazione fisica per ogni istruzione che scrive un risultato. Questa tecnica è molto simile alla trasformazione effettuata dal compilatore chiamata *Static Single Assignment (SSA)* ma in questo caso viene effettuata dallo hardware. Con questa tecnica si rimuovono tutte le possibilità di avere WAR o WAW. In Figura 55 vediamo come il *Register File Fisico* sia molto più grande di quello standard, inoltre, notiamo la presenza di una tabella denominata *Freelist* nella quale sono memorizzati i registri fisici che non sono utilizzati e che sono quindi *liberi*. Per implementare questo meccanismo è sufficiente tenere traccia dell’associazione dei registri tramite una *tabella di traduzione* come mostrato in Figura 56. Quando un’istruzione deve scrivere un risultato in un registro esso viene sostituito da un nuovo registro preso dalla *freelist*. Un registro torna a far parte della *freelist* quando non è più usato da nessuna istruzione. I vantaggi di questa tecnica sono il disaccoppiamento del concetto di renaming da quello di scheduling, la pipeline è esattamente uguale a quella dello MIPS, con la possibilità di implementare scheduling dinamici (Tomasulo o Scoreboard) la possibilità di prelevare più istruzioni per singolo ciclo di clock; inoltre, è possibile utilizzare meccanismi di forwarding e bypassing. Un altro vantaggio è il fatto che, in caso di eccezione, è immediato ricostruire l’esatto stato al tempo del breakpoint in quanto basta solo effettuare la sostituzione dei valori della *rename table*.

**Renaming in pratica** Quando un’istruzione viene prelevata si rinominano tutti i registri relativi agli operandi, gli operandi vengono letti dal RF (reale o esteso) oppure via CDB. Alla fine dell’esecuzione un *reorder buffer* forza un commit ordinato delle istruzioni senza però rinominare i risultati.

Per effettuare queste operazioni sono necessarie alcune caratteristiche, prima fra tutti la disponibilità di una tabella di traduzione, un register file fisico molto più grande di quello ISA nel

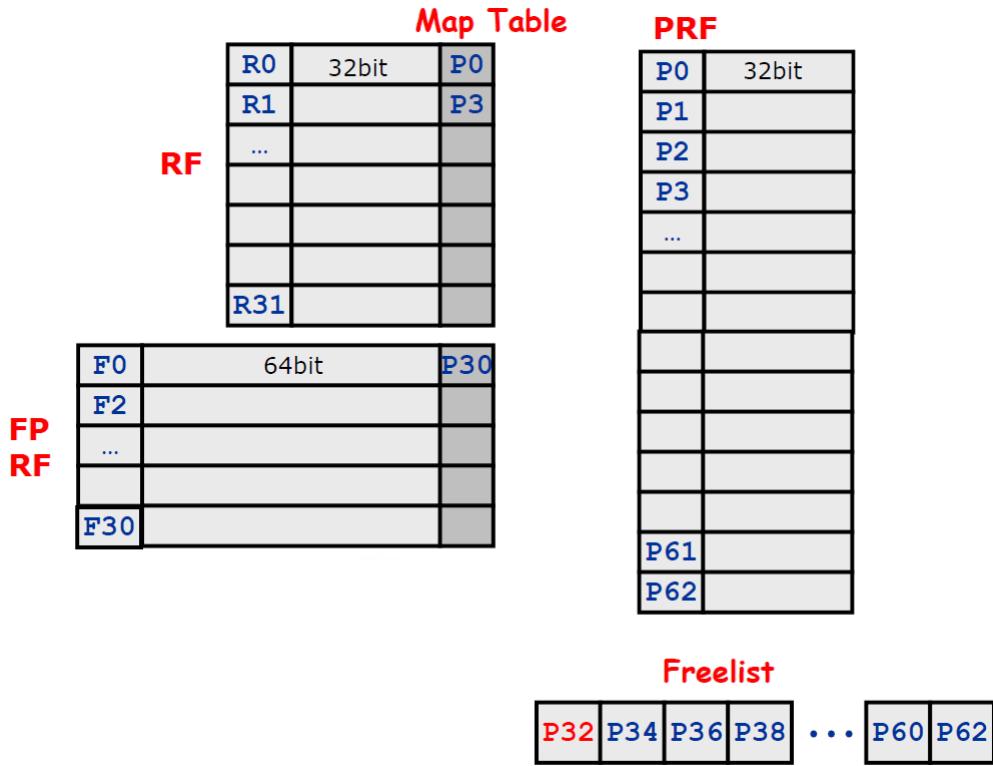


Figura 55: Associazione tra RF e RF fisico

quale sia possibile capire quali sono i registri *liberi*.

Utilizzando il *Register Renaming* possiamo semplificare lo scheduler Scoreboard i quattro stage che lo compongono diventano:

**Issue:** preleva e decodifica le istruzioni, controlla eventuali hazard di tipo strutturale, alloca nuovi registri fisici per il risultato:

- Le istruzioni vengono prelevate nell'ordine del programma per verificare eventuali conflitti
- Non si prelevano ulteriori istruzioni nel caso non vi siano registri fisici liberi.
- Si inserisce uno stallo fino a quando i conflitti strutturali non sono stati risolti.

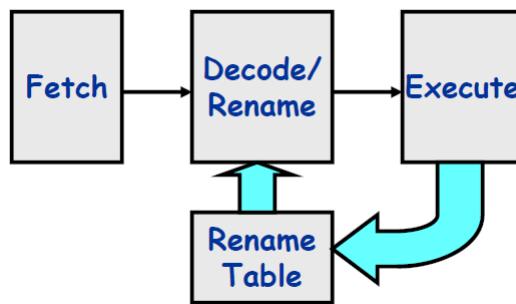


Figura 56: Meccanismo di register renaming

**Read Operands:** Si attende fino a quando sono risolti eventuali conflitti RAW dopo di che si prosegue con la lettura degli operandi.

**Execution:** Si eseguono le operazioni nelle unità funzionali specificate.

**Write result:** I risultati vengono scritti nei registri.

Come si nota non si effettuano controlli per eventuali conflitti di tipo WAR o WAW in quanto risolti automaticamente dal register renaming.

Per un esempio di esecuzione di uno Scoreboard con l'utilizzo del renaming esplicito si rimanda alle slide del corso.

## 4 Static Multiple-Issue Processor: Approccio VLIM

Fino ad ora abbiamo analizzato tecniche che permettono di estrapolare il parallelismo dalle istruzioni solo nel caso in cui queste siano prelevate dalla memoria in modo sequenziale una alla volta; con questo meccanismo abbiamo che il CPI massimo che possiamo ottenere è 1 ovvero possiamo eseguire, nel migliore dei casi, al massimo una istruzione per ciclo.

Introduciamo ora una categoria di processori che, invece, sono in grado di prelevare più di una istruzione per ogni ciclo di clock. Questi processori possono essere a *scheduler dinamico* ovvero possono prelevare un numero diverso di istruzioni ad ogni ciclo di clock, oppure a scheduler statico che preleva un numero prefissato di istruzioni ad ogni ciclo.

Il numero di istruzioni che si possono prelevare ad ogni ciclo può variare da un minimo di 1 ad un massimo di 8 il CPI in questo caso diventa  $CPI = 1/\#istruzioni\ prelevate$ . Questo tipo di processore viene definito *processore superscalare*. Lo scheduler può essere implementato esclusivamente tramite lo hardware anche se il compilatore può migliorare notevolmente la qualità dello scheduler. Lo hardware risistema le istruzioni in esecuzione per ridurre il numero degli stalli mentre mantiene il flusso dei dati e il comportamento delle eccezioni, i vantaggi principali sono la possibilità di gestire casi di dipendenze sconosciute al tempo della compilazione, l'utilizzo di un compilatore semplificato ed infine la possibilità per il codice compilato di essere eseguito su pipeline diverse; questi vantaggi sono ottenuti al costo di una maggiore complessità dello hardware e di un maggiore consumo energetico.

Lo scheduler statico utilizza compilatori con algoritmi sofisticati per estrarre ILP da codice sorgente e individuando quando due istruzioni possono essere eseguite in parallelo. Tuttavia il problema principale è che questa analisi può essere effettuata solo tra *basic block*, ovvero tra piccoli segmenti di codice sequenziale privi di salti ad eccezione del punto iniziale e nessun salto in uscita se non alla fine. Tipicamente questi blocchi hanno una lunghezza compresa tra le 4 e le 7 istruzioni. Un altro fattore che limita la quantità di ILP che si può estrarre dal codice è la dipendenza dei dati; il compilatore tuttavia può, in certa misura, eliminare alcune false dipendenze in modo da aumentare il parallelismo. Per aumentare notevolmente le performance dobbiamo estrarre il parallelismo tra diversi basic block.

Come primo passo bisogna determinare le dipendenze tra le istruzioni in quanto queste dipendenze determinano il livello di parallelismo del programma. Come abbiamo visto esistono tre tipi di dipendenza:

- Dipendenza dei dati.
- Dipendenza dei nomi (WAR e WAW)
- Dipendenze di controllo

### 4.1 Processori VLIW

Come abbiamo visto la ricerca delle dipendenze tramite hardware e lo scheduling dinamico richiedono un grande consumo di area e di energia. L'idea generale è quella di ridurre questi due fattori spostando sul compilatore la decisione di quali operazioni possono essere eseguite in parallelo. Queste operazioni parallele sono raggruppate dal compilatore in un unico pacchetto chiamato *bundle* così che l'hardware non debba controllare eventuali dipendenze.

Il compilatore deve essere certo che non vi siano dipendenze tra le istruzioni inserite nel bundle, tuttavia può indicare quando una dipendenza può presentarsi.

Il vantaggio di questo approccio è che si semplifica notevolmente l'hardware si ha un notevole risparmio sul consumo di energia e si ottengono buone performance grazie a ottimizzazioni del

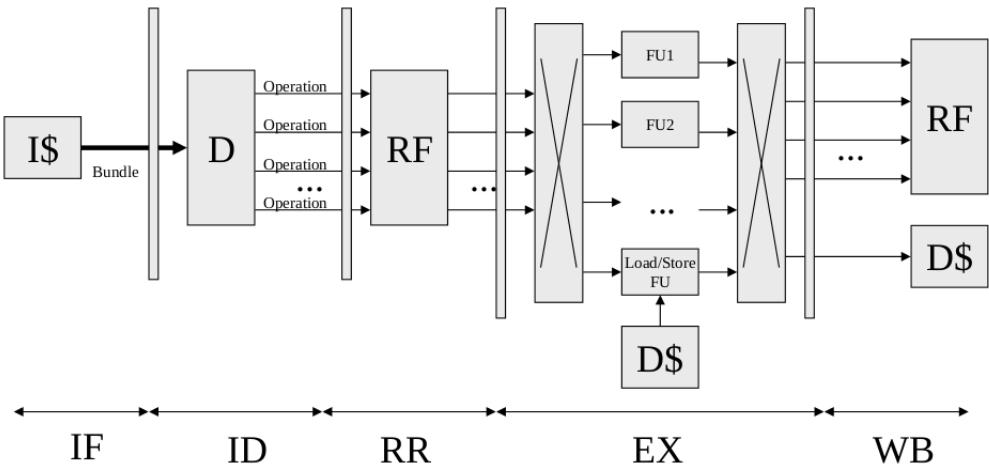


Figura 57: Architettura di una pipeline per VLIW

compilatore.

Un singolo pacchetto è in realtà un'istruzione molto grande (64, 128 o più bits) nella quale sono inserite più operazioni. In principio i processori VLIW erano molto rigidi sul formato delle istruzioni e richiedevano la ricompilazione del programma nel caso di utilizzo su diversi tipi di hardware.

L'istruzione lunga è composta da una serie di campi chiamati *slot* corrispondenti ognuno ad un'unità funzionale; ad esempio una *5-issue VLIW* è un'istruzione lunga che contiene 5 operazioni ad esempio una operazione intera o di salto, due operazioni con la virgola e due operazioni di load/store. In questo modo la fase di decodifica si riduce alla decodifica di ogni istruzione come si può vedere nella Figura 57. Sfortunatamente le operazioni VLIW non sono di tipo atomico, ovvero non avvengono in un solo ciclo di clock perciò la latenza delle operazioni deve essere nota al compilatore ad esempio:

```
I [C=A*B, ...];
I+1 [nop, ...];
I+2 [X=C*F, ...];
```

In questo caso l'operazione ha una latenza di 2 perciò il compilatore inserisce una *nop* al secondo ciclo in quanto *C* non è ancora pronta; se il compilatore schedulasse la seconda moltiplicazione nel secondo ciclo si comprometterebbe la corretta esecuzione del programma. Le dipendenze di tipo WAW e WAR sono risolte dal compilatore e non dallo hardware tenendo conto della latenza delle diverse FU. Per quanto riguarda, invece, le dipendenze di tipo RAW i processori superscalari inseriscono delle *nop* oppure eseguono se possibile, delle istruzioni successive. Nei processori VLIW sono inserite dal compilatore durante la fase di scheduling, idealmente sono usate, se possibile, solo istruzioni non coinvolte in dipendenze, altrimenti sono generate delle *nop*. Tutte le dipendenze vengono risolte dal compilatore, il quale fornisce anche informazioni riguardo alle predizioni dei salti; l'unico tipo di dipendenza che rimane da risolvere sono quelle di controllo che viene evitata dallo hardware abortendo l'esecuzione delle operazioni in caso di predizione incorretta.

Un esempio di meccanismo VLIW è presentato in Figura 58. I vantaggi dell'utilizzo delle VLIW sono il fatto che il compilatore può analizzare il codice ad un livello più alto rispetto a quello che fa lo hardware, in questo modo, tramite sofisticati algoritmi, si può estrapolare un maggiore

SLOT1: LD/ST Ops	SLOT2: LD/ST Ops	SLOT3: Integer Ops	SLOT4: Integer+Branch Ops
lw\$2, BASEA(\$4)	lw \$3, BASEB(\$4)	NOP	NOP
NOP	NOP	NOP	NOP
NOP	NOP	addi \$2,\$2,INC1	addi \$3,\$3,INC2
NOP	NOP	addi \$4, \$4, 4	NOP
NOP	NOP	NOP	NOP
NOP	NOP	NOP	bne \$4,\$7, L1

Figura 58: Esempio di utilizzo di Very Long Instruction Word con queste FU: 2 LD/ST, 1 Int e 1 Int/Branch

parallelismo tra le diverse istruzioni e quindi aumentare le performance. Inoltre, le istruzioni hanno dei campi prefissati è quindi più facile decodificarle. Infine, si riduce notevolmente la complessità dello hardware tramite una superficie minore e quindi un minor consumo di energia e la possibilità di introdurre più FU. Le difficoltà all'applicazione su larga scala di questo meccanismo sono la necessità di avere una tecnologia di compilazione che individui ed estrapoli il parallelismo anche oltre i singoli *basic block*. Inoltre la dimensione del codice aumenta di molto a causa delle numerose `nop` che vengono introdotte; infine, la complessità del Register File e della circuiteria di trasporto verso le FU è incrementata di molto. L'aspetto più importante che però limita l'utilizzo della tecnologia VLIW è l'incompatibilità binaria, ad esempio architettura con lo stesso ISA ma diversi bundle VLIW sono incompatibili, ma anche architetture con lo stesso ISA e con lo stesso bundle ma latenze differenti sono incompatibili. L'unica soluzione a questo problema è la *Just In Time Compilation* ma è molto costosa. Perciò, in molti casi, la VLIW è utilizzata solo nei sistemi embedded dove la compatibilità binaria non è un fattore rilevante.

#### 4.1.1 Alcuni esempi

Analizziamo ora alcuni esempi dei più diffusi processori VLIW presenti in commercio.

**STMicroelectronics ST200** Processore embedded pensato per l'utilizzo e la gestione dei media; esso è un processore VLIW di tipo cluster con quattro cluster eseguiti con un singolo PC. Ogni cluster ha 4 slot e può eseguire perciò 4 istruzioni contemporaneamente, 1 salto, 1 load/store, due moltiplicazioni. Ogni cluster ha a disposizione un banco di 64 registri.

**NXP Trimedia** È un processore per i media con cinque unità di esecuzione tali unità richiedono 15 porte in lettura e 5 in scrittura, ogni unità necessita di tre porte in lettura per leggere i due operandi e un *guard operand* il quale condiziona l'esecuzione di ogni operazione.

**VLIW vs EPIC: Intel IA-64** Un evoluzione del VLIW è stato EPIC (*Explicitly Parallel Instruction Computer*) che al contrario del VLIW permette una flessibilità del formato delle istruzioni inoltre indica quando un istruzione non può essere eseguita in parallelo alle successive. Una prima implementazione di questa tecnologia è stato *Intel Itanium* il quale permetteva un alto parallelismo e una pipeline profonda con una frequenza di clock molto bassa 800MHz. Inoltre, aveva 128 registri a 64-bit per gli interi e 128 registri a 82 bit per i numeri con la virgola. I registri di tipo integer sono configurati per aiutare ed accelerare le chiamate a procedura usando lo stack dei registri, questo meccanismo è simile a quello utilizzato nei RISC e nelle architetture

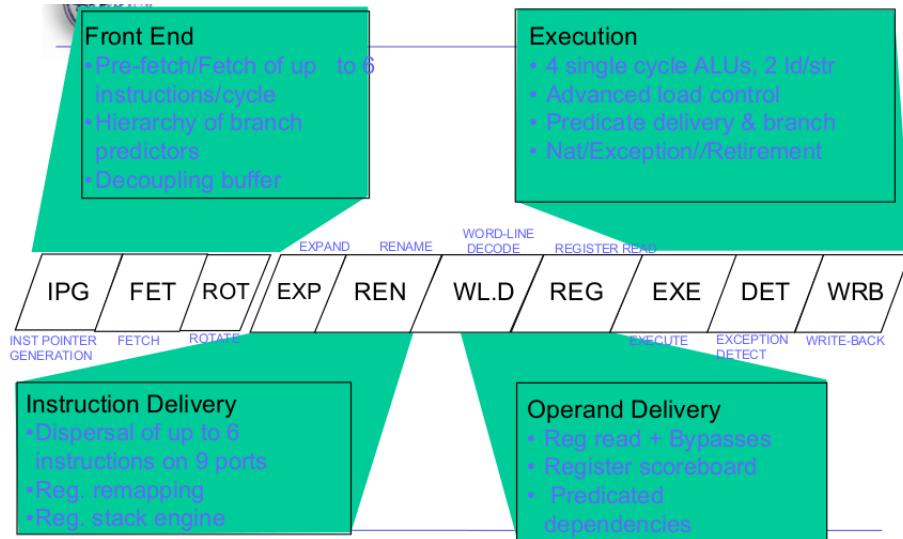


Figura 59: Stage della pipeline di un processore Itanium

SPARC. I registri dallo 0 al 31 sono sempre accessibili tramite il loro indirizzo reale, i registri 32-128 sono allocati come register stack ed ogni procedura ne può allocare alcuni (da 0 a 96) rinominando i registri fisici.

Un *instructions group* è una sequenza di istruzioni consecutive senza dipendenze dei dati, queste istruzioni perciò possono essere eseguite in parallelo se sono disponibili le risorse hardware. La lunghezza dei gruppi è arbitraria ma il compilatore deve separare esplicitamente i due gruppi inserendo uno *stop* tra due istruzioni che appartengono a due gruppi diversi. Nel IA-64 le istruzioni sono codificate in bundle di lunghezza 128 bit, ogni bundle è composto da un campo *template* di 5 bit e da 3 istruzioni di 41 bits.

Il processore Itanium ha una pipeline composta da 10 stadi come possiamo vedere in Figura 59. Come vediamo dalla Figura 59 possiamo raggruppare i diversi stage in quattro gruppi:

**Front-end:** composto dagli stage IPG, Fetch e Rotate, precaricano 32 byte per clock (2 bundle) in un buffer di precarico, il quale mantiene 24 istruzioni, in questa fase si possono effettuare delle predizioni tramite un predittore adattativo multi livello.

**Instruction delivery:** formato dagli stage EXP e REN i quali distribuiscono le istruzioni alle 9 unità funzionali ed effettuano il renaming dei registri.

**Operand delivered:** formato dagli stage WLD e REG, accede ai registri e verifica eventuali dipendenze.

**Execution:** composto dagli ultimi tre stage (EXE, DET e WRB) si occupa di eseguire le istruzioni individuando eventuali eccezioni ed inoltre effettua il write-back dei risultati.

**Processore Crusoe** Il processore Crusoe è un processore VLIW con esecuzione sequenziale formato da 64 registri di tipo integer e da 32 registri per i *floating point*. Esso è formato da una pipeline a 6 stadi per gli integer: 2 fetch, 1 di decodifica, 1 register read, 1 execution e 1 write back; e da una pipeline a 10 stadi per le operazioni in virgola mobile, che comprendono 4 stadi supplementari per la fase di esecuzione.

I processori Crusoe hanno a disposizione 5 unità funzionali:

- ALU;
- Compute che è un unità che può occuparsi di due ALU contemporaneamente o una floating point o un'operazione sui media.
- Memory: che implementa le operazioni di load e di store.
- Branch: per eseguire un istruzione di salto
- Immediate: una istruzione a 32 bit utilizzata immediatamente da un'altra operazione

## 4.2 Code scheduling VLIW

L'obiettivo principale dello scheduling nei processori VLIW è quello di stabilire staticamente l'ordine di esecuzione delle istruzioni nel codice oggetto in modo che queste vengano eseguite in modo corretto ed efficiente.

Il meccanismo base per schedulare le istruzioni in modo efficiente è quello di dividere il codice in blocchi base a questo punto per ogni blocco base si costruisce un grafico delle dipendenze come quello mostrato in Figura 60 Un grafico delle dipendenze cattura tutti i tipi di dipendenze

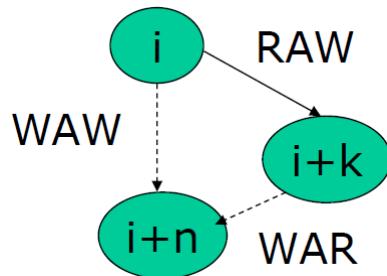


Figura 60: Esempio di grafico delle dipendenze tra le istruzioni di un blocco base

(WAR, WAW e RAW) le dipendenze WAR e WAW però sono dipendenze solo di nome in quanto sono dovute al riuso dei registri o delle variabili. Dal grafico delle dipendenze si può individuare il cosiddetto *critical path* ovvero il percorso più lungo in un grafico delle dipendenze; questo percorso determina il minimo tempo di esecuzione come possiamo vedere dalla Figura 61. Dal

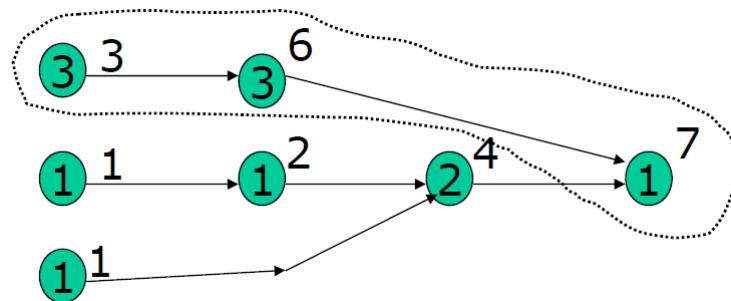


Figura 61: Esempio di *critical path*

quale si possono calcolare la lunghezza di ogni percorso tramite la formula

$$LP(i) = \text{Max}(LP(\text{Pred}(i))) + \text{Latency}(i)$$

mentre il *critical path* è dato dal valore massimo tra tutte le lunghezze

$$LCP = \text{Max}(LP(i))$$

Lo scheduler in teoria dovrebbe programmare l'esecuzione di ogni operazione del *critical path* in modo da minimizzare il tempo di esecuzione e parallelamente schedulare le altre istruzioni, questo però è possibile soltanto nel caso di processori con risorse infinite. Con risorse finite invece il tempo di esecuzione dipende anche da come sono schedulate le rimanenti operazioni. Inoltre, uno scheduler ottimo analizza in modo esaustivo lo spazio degli schedule disponibili per minimizzare il tempo di esecuzione, ma tale analisi è un problema NP-completo siamo perciò costretti ad usare dei meccanismi euristici.

### 4.3 List-based scheduling

In questo tipo di scheduling per ogni ciclo di clock viene selezionata un'istruzione da un *ready set* che può essere inserita in uno degli slot. Un'istruzione si definisce *ready* quando tutti i suoi predecessori sono già stati schedulati e tutti gli operandi necessari sono disponibili. All'inizio dello scheduling tutte le operazioni all'inizio del grafico sono inserite nel *ready set* e ad ogni ciclo si cerca di schedulare tutte le istruzioni nel set nel caso più istruzioni siano presenti si schedula quella con la priorità maggiore.

Per implementare questo meccanismo è necessario però tenere traccia di quali risorse sono occupate per fare questo si utilizza una *Resource Reservation Table* ovvero una tabella che indica quali risorse sono occupate in un determinato istante di tempo.

Un esempio di tale sistema è mostrato in Figura 62

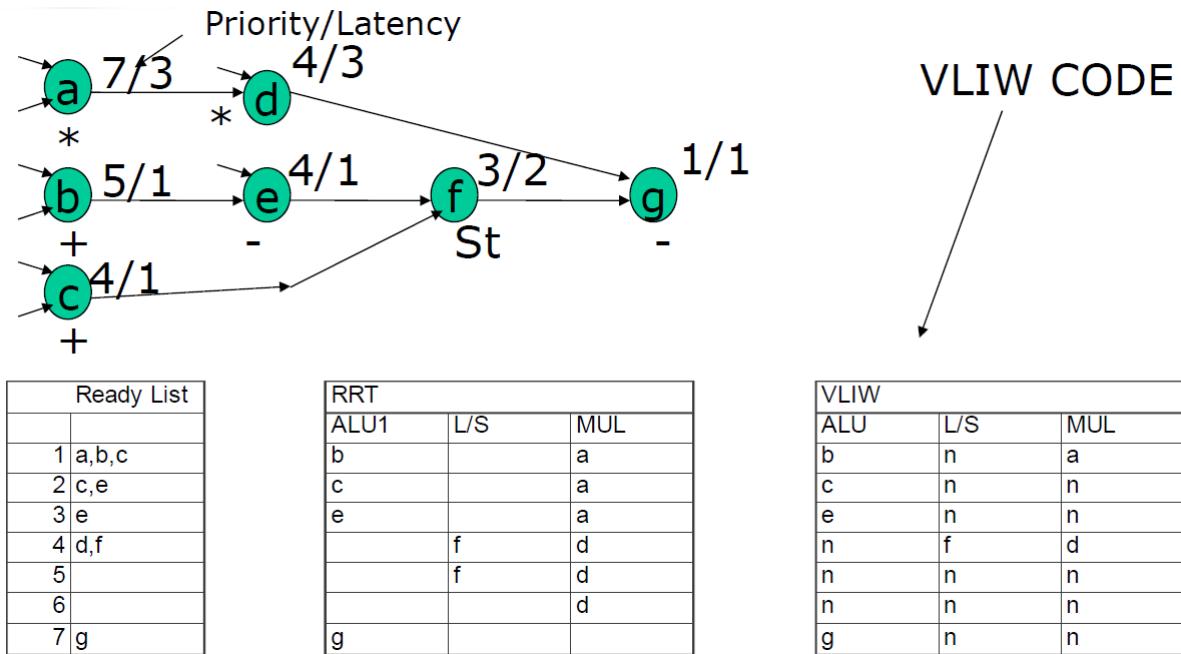


Figura 62: Esempio di scheduling List-Based con utilizzo di Reservation Table

### 4.4 Global e Local scheduling

Per esplorare tutto il possibile parallelismo presente è necessario che il compilatore espanda la dimensione del basic block o che scheduli parallelamente istruzioni che appartengono a basic

block differenti. Le tecniche di *local scheduling* operano su di un singolo basic block e possono essere tecniche di *Loop Unrolling* o di *Software Pipelining*. Le tecniche di *global scheduling* espandono la ricerca del parallelismo al di fuori del singolo basic block e alcune tecniche sono *Trace Scheduling* e il *Superblock Scheduling*.

**Loop unrolling** Il *loop unrolling* è una tecnica di scheduling locale che permette al compilatore di incrementare la quantità di parallelismo disponibile in un basic block. Per fare ciò il compilatore replica il corpo di un loop diverse volte (dipende dal fattore di *unrolling*) sistemando poi il codice di terminazione del ciclo. Per effettuare questa operazione però il compilatore deve prima testare l'indipendenza tra le diverse iterazioni.

Il *loop unrolling* permette di incrementare il numero di operazioni effettuate ad ogni ciclo minimizzando però il numero di salti da effettuare; inoltre, aumentando il numero di operazioni si aumenta la lunghezza del blocco base permettendo al compilatore di effettuare uno scheduling più efficiente. Gli svantaggi di questa tecnica sono però l'aumento della dimensione del codice e il numero di registri richiesto per l'esecuzione.

Un esempio di loop unrolling è mostrato nelle Figura 63(a) e Figura 63(b) dove si mostra un unrolling di fattore di srotolamento di 4.

<pre> Loop: LD    F0, 0 (R1)       NOP       ADD   F4, F0, F2       NOP       NOP       SD    F4, 0 (R1)       SUBI R1, R1, #8       NOP       BNE   R1, R2, LOOP       NOP   </pre>	<pre> Loop: LD    F0, 0 (R1)       ADD   F4, F0, F2       SD    F4, 0 (R1)       LD    F0, -8 (R1)       ADD   F4, F0, F2       SD    F4, -8 (R1)       LD    F0, -16 (R1)       ADD   F4, F0, F2       SD    F4, -16 (R1)       LD    F0, -24 (R1)       ADD   F4, F0, F2       SD    F4, -24 (R1)       SUBI R1, R1, #32       BNE   R1, R2, LOOP   </pre>
(a) loop base	(b) loop unrolled

Figura 63: Esempio di unrolling con fattore 4

**Loop-carried dependences** Il loop-carried è una tecnica incentrata sull'analisi delle dipendenze presenti tra gli operandi all'interno delle diverse interazioni di un loop. Si hanno delle dipendenze tra due interazioni di un loop quando un valore in una interazione dipende da un secondo valore prodotto in una interazione precedente. Un esempio di questa dipendenza è dato dal seguente codice:

```

for (i = 6; i < 100; i++)
{
    Y[i] = Y[i-5] + Y[i];
}
  
```

In questo caso ogni iterazione dipende dalla 5 iterazione precedente e quindi le iterazioni *i*, *i+1*, *i+2*, *i+3*, *i+4* sono indipendenti (*i+5* dipende dall'iterazione *i*). Trovando il fattore di indipendenza si può trovare effettuare un *unrolling* del ciclo come segue:

```
for (i = 6; i < 100; i=i+5)
{
    Y[i] = Y[i-5] + Y[i];
    Y[i+1] = Y[i-4] + Y[i+1];
    Y[i+2] = Y[i-3] + Y[i+2];
    Y[i+3] = Y[i-2] + Y[i+3];
    Y[i+4] = Y[i-1] + Y[i+4];
}
```

Si riesce così ad estendere il blocco base ma il fattore di *unrolling* non può essere maggiore di 5.

**Loop peeling e fusion** La tecnica di *peeling & fusion* è una tecnica che consiste nell'eliminare (sbucciare) da un ciclo le interazioni superflue in modo da poterlo fondere con un secondo ciclo; ad esempio prendiamo in considerazione due cicli:

```
for (i = 0; i < 102; i++) b[i] = b[i-2] + c;
for (j = 0; j < 100; j++) a[j] = a[j] * 2;
```

In questo caso sbucciamo il primo ciclo delle ultime due iterazioni e lo fondiamo poi con il secondo ciclo, otteniamo così:

```
for (i = 0; i < 100; i++)
{
    b[i] = b[i-2] + c;
    a[i] = a[i] * 2;
}
b[100] = b[98] + c;
b[101] = b[99] + c;
```

dove abbiamo un ciclo che fonde i due precedenti e due operazioni aggiuntive che servono a compensare le due iterazioni mancanti del primo.

**Software pipeline** Supponiamo di avere un loop nel quale ad ogni iterazione possiamo identificare delle istruzioni indipendenti differenti come quello mostrato in Figura 64 A questo punto possiamo riorganizzare queste istruzioni in un nuovo loop nel quale ad ogni iterazione si eseguono delle istruzioni provenienti da diverse iterazioni. Questa tecnica può essere considerata un po come un *symbolic loop unrolling*.

Prendiamo in esempio il seguente loop dove sono presenti delle dipendenze interne al loop.

```
for(i = 0; i < 100; i++)
{
    A[i] = B[i];
    A[i] = A[i]+1;
    C[i] = A[i];
}
```

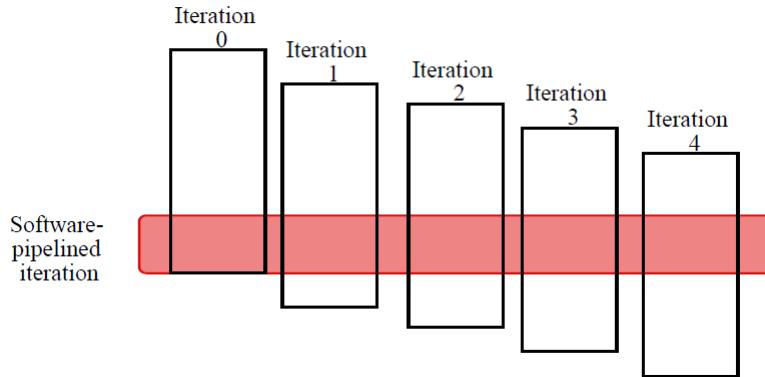


Figura 64: Loop con istruzioni indipendenti ad ogni iterazione

In Figura 65 vediamo lo svolgimento delle diverse iterazioni; i riquadri nella figura indicano il nuovo ciclo che si viene a creare che viene mostrato in Figura 66. Il vantaggio di questa tecnica è che consuma meno spazio del loop unrolling, in quanto non vi è la necessità di duplicare il codice. Si riempie e si svuota la pipe una volta sola per ogni loop al contrario del caso del loop unrolling che si riempie e si svuota ad ogni interazione. Infine tale tecnica può essere associata anche al loop unrolling per incrementarne le prestazioni.

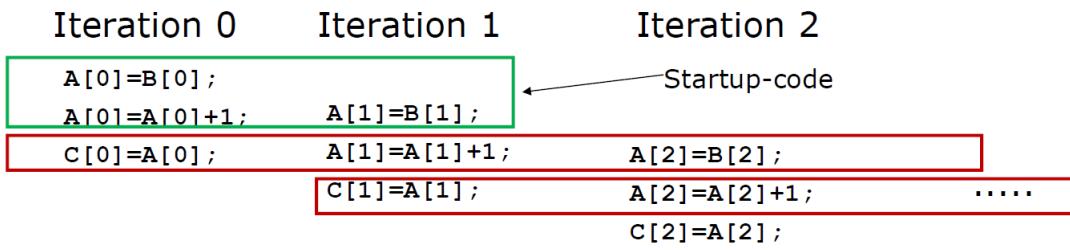


Figura 65: Svolgimento delle diverse iterazioni

**Trace scheduling** I sistemi di local scheduling funzionano solo quando i loop contengono un singolo basic block, quando i corpi dei loop contengono dei flussi di controllo allora è necessario espandere lo scheduling attraverso i diversi basic block tramite tecniche di *global scheduling*. Il *trace scheduling* è la prima tecnica che analizzeremo e consiste nel tentare di trovare del parallelismo attraverso i salti condizionati. Si compone di due passi, il primo consiste nel trovare una sequenza di blocchi base composta dalla più lunga sequenza di istruzioni, la seconda fase consiste nel compattare questa sequenza in poche istruzioni VILW inserendo delle istruzioni di compensazione in caso di predizione errata. Questa tecnica è una forma di speculazione del compilatore.

**Superblock scheduling** Questa è una tecnica di ottimizzazione del *trace scheduling* che consiste nel creare un *super blocco* formato da diversi blocchi base con un unico punto d'entrata e molteplici flussi di controllo in uscita come mostrato in Figura 67.

**Hardware support** Tutte le tecniche viste fino ad ora si applicano quando il comportamento dei salti è molto predicibile altrimenti il controllo delle dipendenze limita di molto il parallelismo. Per aggirare questo problema possiamo estendere il set di istruzioni per includere delle

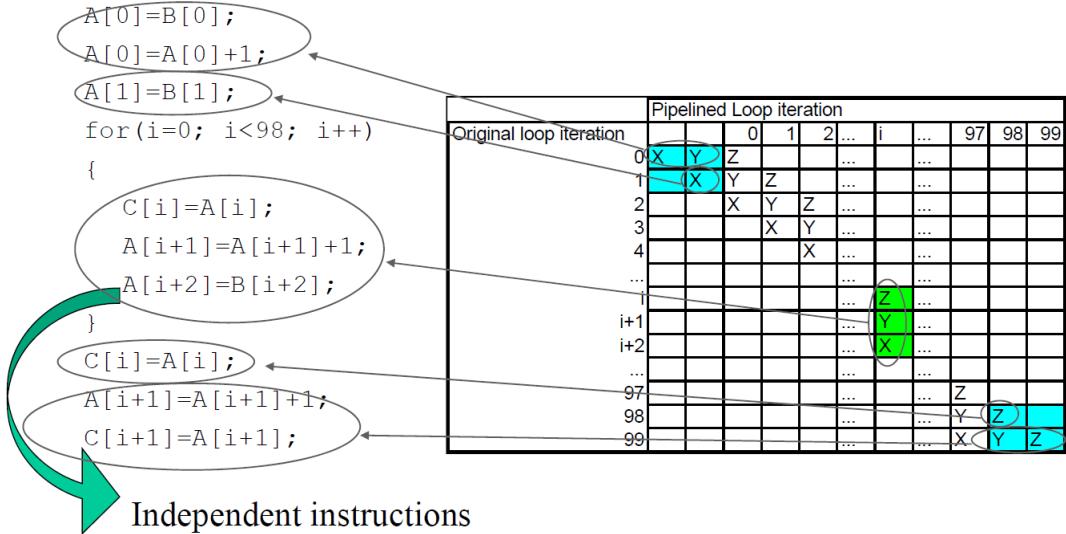


Figura 66: Nuovo codice che implementa il loop

istruzioni condizionali; utilizzare la speculazione del compilatore con il supporto dell'hardware per permettere di creare codice speculativo mantenendo il comportamento delle eccezioni. Un esempio è l'*esecuzione condizionale* che utilizza un particolare tipo di istruzione come quella seguente:

(p) op Rd, R1, R2

dove p è un predicato booleano che condiziona l'operazione op, infatti, op viene *committata* soltanto se p risulta *vera*.

Avendo effettuato queste modifiche possiamo a questo punto modificare il codice tramite delle *if-conversion* che trasformano i salti in sequenze di istruzioni condizionali e le dipendenze di controllo si trasformano in dipendenze sui dati eliminando così i salti. I vantaggi sono l'eliminazione dei problemi di *miss-prediction* e l'aumento delle dimensione dei basic block. Questa soluzione diventa efficiente se le predizioni errate e quindi la loro penalità è considerevole e se i salti sono sbilanciati e la parte eseguita più frequentemente è quella più lunga.

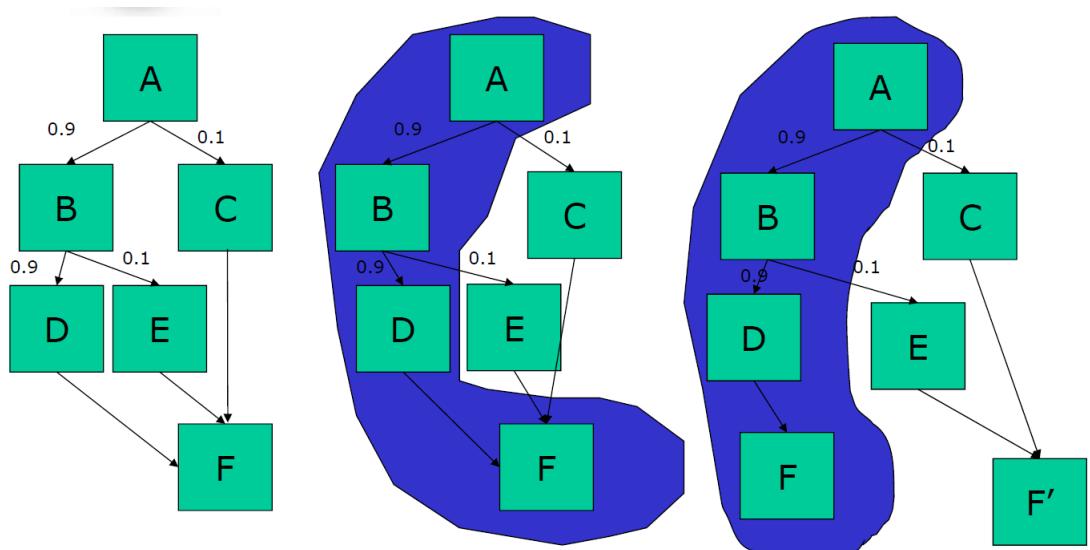


Figura 67: Realizzazione di un superblocco

## 5 Reorder Buffer

Nei capitoli precedenti abbiamo visto come, sia Tomasulo che lo Scoreboard, prevedano il prelievo delle istruzioni in ordine ma che poi l'esecuzione e il completamento di tali istruzioni avvenga fuori ordine. Ovvero esiste un disaccoppiamento tra il prelevamento e l'esecuzione delle istruzioni. Un altro punto che abbiamo analizzato superficialmente è la risoluzione dei salti, infatti, noi facevamo affidamento sul fatto che il risultato del branch fosse controllato da un'operazione tra interi e che quindi fosse un'operazione *veloce*. Nel caso in cui, invece, il ciclo dipenda da un'operazione più lenta come una moltiplicazione perdiamo tutti i nostri vantaggi come nel caso del codice seguente:

```
Loop:    LD      F0  0   R1
          MULTD  F4  F0  F2
          SD      F4  0   R1
          SUBI   R1  R1  #8
          BNEZ   R1  Loop
```

Il primo problema è nella predizione del salto, infatti, una corretta previsione diventa fondamentale per mantenere delle buone prestazioni. Oltre alla predizione sui salti l'architettura deve prevedere qualsiasi altro tipo di dipendenza come quella sui dati. Tutte queste previsioni sono effettuate dallo hardware; l'idea di base è quella di prelevare ed eseguire delle istruzioni dipendenti da un salto prima che il risultato di questo salto sia conosciuto, ovvero permettere alle operazioni di essere eseguite fuori ordine ma è necessario che esse siano completate *in ordine* tutto questo per prevenire che un'operazione venga *committata* prima che tutte le sue precedenti non siano concluse. Questo significa che un'operazione deve essere committata solo quando essa non è più *speculativa*; il meccanismo che permette questo tipo di controllo è il *ReOrder Buffer (ROB)* che mantiene il risultato delle istruzioni che hanno completato la loro esecuzione ma che non possono essere ancora committate.

Il risultato di un salto è predetto e il programma viene eseguito come se la previsione fosse corretta (senza speculazione non si ha la fase di esecuzione). Per fare ciò però sono necessari dei meccanismi per manipolare i casi in cui la previsione è sbagliata. La speculazione hardware permette estende lo scheduling dinamico al di fuori dei blocchi base.

La *speculazione hardware* combina tre idee:

**Dynamic Branch Prediction:** che permette di selezionare quale ramo del salto dovrà essere eseguito prima che il risultato del salto sia conosciuto.

**Speculazione:** che permette di eseguire delle istruzioni prima che le dipendenze di controllo siano eseguite.

**Scheduling dinamico:** che supporta l'esecuzione fuori ordine ma il completamento in ordine.

Essenzialmente il modello basato sulla speculazione hardware è un modello basato sul *data flow* ovvero, l'esecuzione di un'istruzione inizia quando i suoi operandi sono disponibili.

La speculazione hardware è stata introdotta per estendere e supportare l'algoritmo di Tomasulo, in particolare per separare la fase di commit da quella di esecuzione è stato introdotto il *Reorder Buffer*. Il meccanismo del *Reorder Buffer* è abbastanza semplice, le istruzioni vengono mantenute in un ordine di tipo FIFO esattamente come vengono prelevate, per ogni record del ROB si mantengono il valore del PC, del registro di destinazione, del risultato e l'eventuale stato dell'eccezione. Quando un'istruzione completa la sua esecuzione il risultato viene inserito nel corrispettivo campo del ROB. Una volta completata l'esecuzione si forniscono i risultati alle

altre istruzioni ma si utilizzano i valori dei tag del ROB invece di utilizzare le reservation station. Un'istruzione effettua il commit solo quando è pronta ed è in cima al ROB, solo a quel punto i valori vengono copiati nei registri.

Oltre a questo il Reorder Buffer è comodo per effettuare delle speculazioni, infatti, esso permette di eseguire delle istruzioni senza conseguenze nel caso in cui il branch non sia chiuso; questo meccanismo è chiamato *boosting*. È l'insieme dell'utilizzo di dynamic scheduling e branch prediction che permette ad un'istruzione di essere eseguita prima che sia conosciuto il valore di un salto.

Il fatto di eseguire le istruzioni fuori ordine ma di completarle in ordine permette di prevenire azioni dannose come l'aggiornamento di stati non consistenti o eccezioni.

## 5.1 Struttura del reorder buffer

Come abbiamo detto il reorder buffer mantiene lo stato delle istruzioni che hanno completato la loro esecuzione ma che non sono ancora state committate, inoltre permette di scambiare il risultato di un'istruzione tra le diverse istruzioni in esecuzione. Tutto questo permette un'esecuzione delle istruzioni fuori ordine ma un commit di tali istruzioni in ordine. Un esempio di utilizzo del ROB lo abbiamo nell'algoritmo di Tomasulo come mostrato in Figura 68 Il *Reorder*

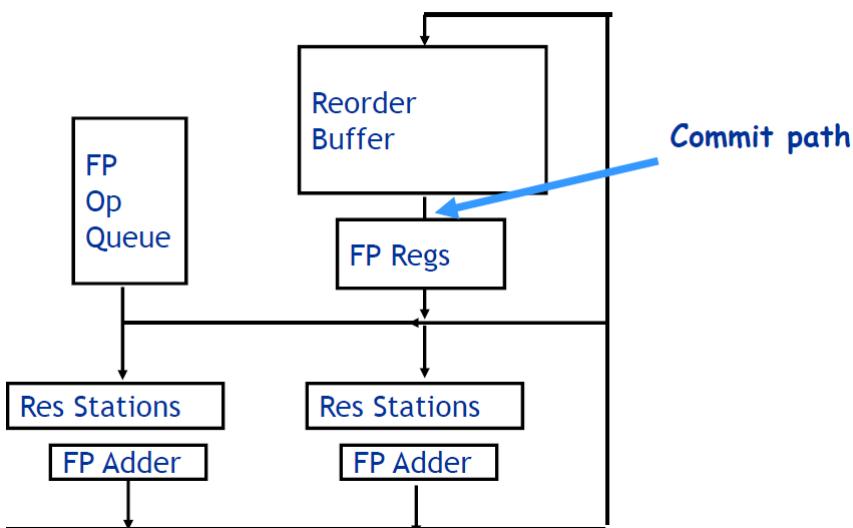


Figura 68: Struttura dell'algoritmo di Tomasulo con Reorder Buffer

*Buffer* rimpiazza completamente lo *Store Buffer*, la funzione di renaming delle *Reservation Station* adesso è effettuata direttamente dal ROB, le RS ora sono utilizzate solamente per immagazzinare istruzioni e operandi da passare alle FU per diminuire gli hazard strutturali. I puntatori adesso puntano direttamente alle entità del ROB.

Ogni entità contenuta nel *Reorder Buffer* contiene a sua volta 4 campi:

**Instruction Type:** Identifica il tipo di istruzione da eseguire.

**Destination:** Contiene l'indirizzo del registro di destinazione (ALU e load) o l'indirizzo di memoria (store).

**Value:** Mantiene il valore del risultato fino a quando l'istruzione non è committata.

**Ready:** Indica che l'istruzione ha completato la sua esecuzione.

LE istruzioni nel Reorder Buffer sono inserite nell'ordine del programma, quando un'istruzione viene prelevata essa è allocata in sequenza, essa può essere in tre stati:

- **i issued**
- **x in esecuzione**
- **f completata**

Un'istruzione viene committata solo quando tutte le precedenti istruzioni sono state committate e tale istruzione si trova nello stato **f**. Il reorder buffer ha una struttura circolare con due puntatori, uno che indica la coda delle istruzioni (punto in cui la successiva istruzione prelevata sarà memorizzata) e uno che ne indica la testa (che tiene traccia della prossima istruzione da committare). Un esempio di tale struttura è mostrato in Figura 69. Nel caso di algoritmo di

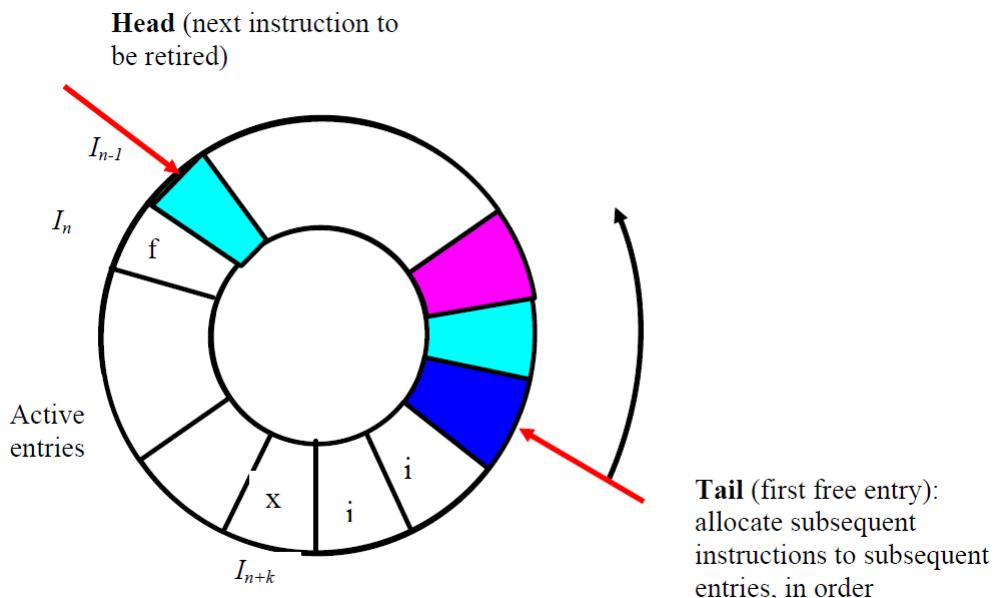


Figura 69: Struttura di un Reorder Buffer

Tomasulo speculativo con Reorder Buffer i passi dell'algoritmo variano leggermente:

**Issue:** durante questa fase viene prelevata l'istruzione dalla coda delle istruzioni; se il ROB contiene uno slot libero allora si preleva l'istruzione e si copia nel reorder buffer e nelle reservation station.

**Execution:** si eseguono le operazioni sugli operandi, quando sono entrambi disponibili si esegue l'operazione altrimenti si controlla il *Common Data Bus* in attesa dei risultati.

**Write Result:** durante questa fase i risultati delle esecuzioni vengono scritti sul CDB e nel ROB inoltre le reservation station vengono marcate come disponibili.

**Commit:** vengono aggiornati i diversi registri con il valore contenuto nel ROB, quando l'istruzione puntata dal puntatore di testa del ROB ha completato la sua esecuzione ed è marcata come conclusa allora i dati vengono copiati nei registri e l'istruzione viene rimossa dal ROB; nel caso di predizione errata le istruzioni vengono cancellate.

Per la fase di commit esistono tre possibili sequenze:

1. **Normal commit:** l'istruzione raggiunge la testa del ROB il risultato è presente nel buffer allora esso viene copiato nei registri e l'istruzione viene rimossa dalla coda.
2. **Store commit:** come la precedente solo che il risultato viene scritto in memoria anziché nei registri.
3. **Incorrect prediction:** l'istruzione è un salto la cui predizione è però errata, questo fa sì che il ROB sia svuotato e che la predizione ricominci dall'istruzione successiva corretta.

Per quanto riguarda il comportamento delle eccezioni esse sono riconosciute solo quando l'istruzione che genera l'eccezione è committata, questo mantiene il comportamento dell'eccezione. Per un esempio completo del funzionamento dell'algoritmo di Tomasulo con Reorder Buffer si rimanda alle slide del corso.

## 6 Multithreading

Fino ad ora abbiamo visto il parallelismo intrinseco che esiste tra le istruzioni, ovvero quello che viene comunemente chiamato *ILP*. Il problema dell' ILP è che comporta meccanismi talvolta molto costosi come il *dynamic scheduling* che richiede una grande quantità di logica e limita inoltre il clock.

In realtà una macchina ideale dovrebbe avere le seguenti caratteristiche:

**Register renaming:** dovrebbe avere un numero di registri infinito più un meccanismo di buffering degli operandi per evitare tutti i problemi di WAW e WAR.

**Branch prediction (*perfetta*):** ovvero una predizione dei salti che non commette errori e una lista illimitata di istruzioni da poter eseguire.

**Memory-address alias analysis:** gli indirizzi sono conosciuti e le *store* possono effettuare le loro operazioni prima che le *load* provino che gli indirizzi non sono uguali.

**Unlimited issue:** la CPU può prelevare un numero di istruzioni arbitrario per ogni ciclo di clock, ricercando tali istruzioni in qualsiasi punto del codice.

**One cycle latency for all instruction:** tutte le istruzioni hanno un solo ciclo di latenza questo comporta che ogni istruzione dipendente può essere schedulata nel ciclo successivo.

**Perfect cache:** tutte le *load* e le *store* vengono eseguite in un singolo ciclo di clock e non avvengono mai *cache miss*.

Oltre a queste caratteristiche fisiche una macchina ideale deve anche disporre di uno *scheduling dinamico perfetto* il che comporta il fatto di poter prelevare le istruzioni in modo arbitrariamente lontano, la possibilità di rinominare tutti i registri, la capacità di determinare quali dipendenze dati esistono nel set di istruzioni prelevato ed infine fornire un numero adeguato di unità funzionali replicate per poter eseguire tutte le istruzioni pronte.

In realtà nelle CPU attualmente in commercio non si possono avere più di due riferimenti a memoria per ciclo, inoltre vi sono alcune limitazioni sul numero di bus e sul numero di porte del *register file*; tutte queste limitazioni definiscono un limite al numero di operazioni che può essere prelevato durante un singolo ciclo di clock.

Appena introdotto i processori super scalari erano di tipo 2-issue e si velocemente in processori

4-issue. Attualmente possiamo trovare molto raramente dei processori 6-issue ma non superiori in quanto risulta molto complicato decidere 8 o più istruzioni indipendenti da eseguire ad ogni ciclo, il calcolo è molto complesso e la frequenza del processore dovrebbe essere diminuita. Gli svantaggi dei processori super scalari sono la grande quantità di logica necessaria per decidere l'indipendenza delle istruzioni, ed inoltre, non è scalare infatti per aumentare la finestra di prelevamento è necessario ridurre la frequenza di clock.

## 6.1 Processori embedded

Al contrario dei processori super scalari i processori pensati per prodotti embedded devono essere economici e consumare poco energia. La maggior parte dei processori per i sistemi embedded è progettata da:

- ARM
- STMicroelectronics

**ARM Cortex-A8** Basato sull'architettura ARMv7 è il processore presente sul System-on-Chip A4 di Apple. È un processore di tipo dual-issue con esecuzione in ordine. Utilizzato nel SoC A4 con manifattura a 45nm e frequenza di 1GHz è stato utilizzato nel primo iPad, sull'iPhone4 e sull'iPod Touch.

**ARM Cortex-A9** Basato sull'architettura ARMv7 è un processore di tipo dual core presente sul System-on-Chip A5 di Apple. È un processore di tipo dual-issue con esecuzione in ordine. Utilizzato nel SoC A5 con manifattura a 45nm e successivamente 35nm e frequenza di 1GHz è stato utilizzato nell'iPad2, sull'iPhone 4S e sull'iPad Mini e sulle Apple TV di terza generazione.

**Apple A6 SoC** Il SoC A6 è stato introdotto nel settembre 2012 per l'arrivo dell' iPhone5 è basato sull'architettura ARMv7 personalizzata da Apple, comprende un processore dual core da 1.3GHz e un GPU triple-core PowerVR SGX, inoltre incorpora una RAM da 1GB LPDDR2-1066 che permette di incrementare la banda di memoria teorica fino ad un valore di 8.5 GB/s. Confrontando due processori Intel possiamo vedere come per quanto riguarda i sistemi embedded lo scopo principale sia quello di mantenere contenuto il consumo di energia, infatti:

- Intel i7 920: processore a 4 cores a 2.66 GHz; consumo medio di energia **130W**
- Intel Atom 230 (*embedded*): 1 core a 1.66 GHz; consumo medio di energia **4W**

Il problema per i sistemi embedded è trovare il giusto compromesso tra risparmio di energia e performance. Per fare questo utilizzano, a differenza dei sistemi *general purpose*, uno *scheduling statico* ovvero uno scheduling effettuato a *compile-time* e le istruzioni di tipo VLIW. Questo meccanismo permette di decidere quando e dove eseguire le istruzioni durante la compilazione del programma. Questo permette di ridurre il design dell'hardware.

Le difficoltà che si riscontrano sono tuttavia molteplici e diverse, ad esempio, il compilatore deve essere molto evoluto per individuare il parallelismo tra le istruzioni e schedularle sulle diverse unità funzionali. Inoltre, esiste una *incompatibilità binaria* a causa delle ottimizzazioni architettoniche che effettua il compilatore.

## 6.2 Multithreading e Multiprocessing

Un’alternativa per incrementare le performance è quella di sfruttare il parallelismo dei *thread* al posto del semplice ILP. Un *thread* è parte di un processo con istruzioni e dati propri, esso può essere parte di un programma con molti processi oppure può essere un programma indipendente; ogni thread ha un suo *stato* necessario per la sua esecuzione. Un esempio di come un thread può esistere è mostrato in Figura 70. I thread vengono creati dai programmatori o dal sistema

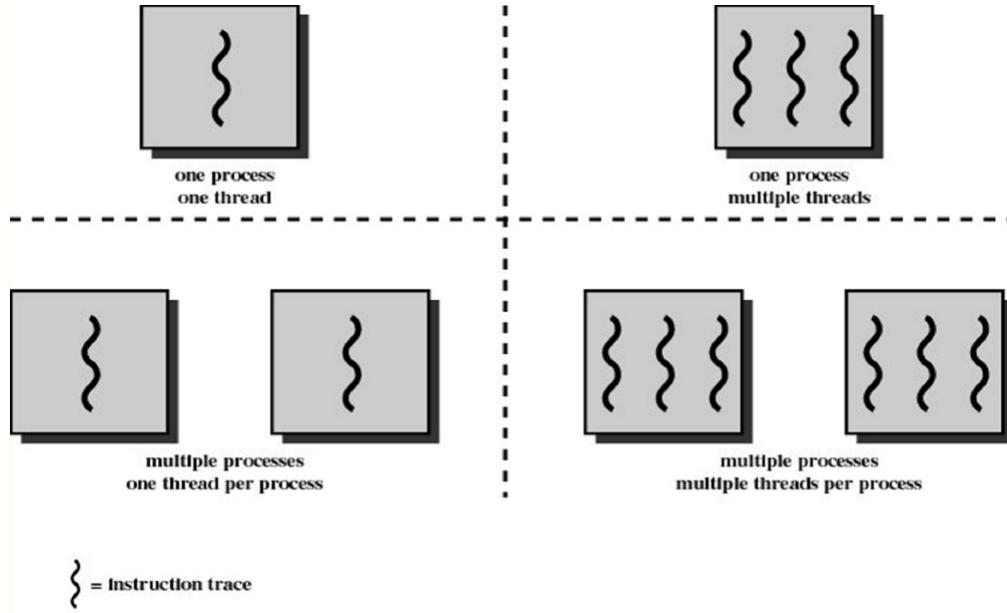


Figura 70: Esempio di esistenza dei thread

operativo, ad ogni thread viene associata una specifica porzione di computazione che può essere formata da poche istruzioni oppure da molte righe di codice. Un processo scambia tra i diversi thread, quando uno va in stallo un altro va in esecuzione, lo stato di ogni thread deve essere salvato mentre un altro thread è in esecuzione; sono così necessari register file e PC multipli. Il multithreading permette a più thread di condividere le unità funzionali di un singolo processore, lo spazio degli indirizzi di memoria è condiviso attraverso meccanismi di memoria virtuale, l’HW supporta l’abilità di cambiare tra i diversi threads in modo veloce e più efficientemente di quanto si possa fare in uno scambio di contesto tra processi.

Esistono diversi tipi di multithreading in ambienti super scalari.

**Coarse-grained:** meccanismo che prevede che quando un thread si blocca, ad esempio per una lettura del disco, un altro thread va in esecuzione.

**Fine-grained:** il switching tra i diversi thread è effettuato ad ogni istruzione

**Simultaneous:** più thread utilizzano slot di prelevamento multipli in un singolo ciclo di clock.

Analizziamo ora i diversi tipi di multithreading partendo dal caso di un processore super scalare senza multithreading; l’utilizzo degli *issue slot* è limitato dalla mancanza di ILP, inoltre nel nostro esempio abbiamo anche uno stallo dovuto ad una cache miss che lascia il programma in sospeso. Il nostro esempio è mostrato in Figura 71(a) Nel caso di multithreading super scalare di tipo coarse-grained come quello in Figura 71(b) i lunghi stalli vengono nascosti mandando in esecuzione un nuovo thread che occupi le risorse del processore, questo meccanismo riduce il

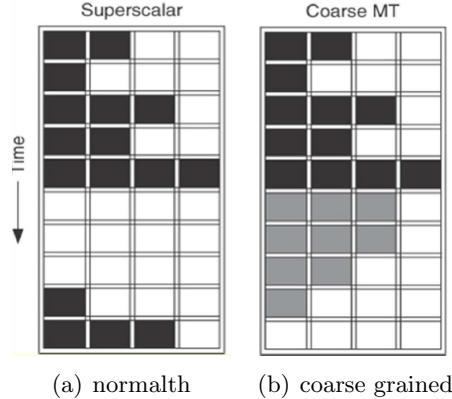


Figura 71: Esempio di processore super scalare 47 e di processore super scalare con multithreading di tip coarse-grained 71(b)

numero di cicli in cui il processore è inattivo; tuttavia le limitazioni dovute alle ILP continuano a far sì che alcuni slot siano vuoti, quando c'è uno stallo è necessario svuotare la pipeline prima di iniziare con il nuovo thread, inoltre il nuovo thread necessita di un periodo di start-up nel quale non si completano operazioni e si riduce perciò il throughput del sistema. Date queste limitazioni il coarse-grained MT è applicabile solo quando il tempo per il riempimento della pipeline è molto minore dei tempi di stallo.

Un'alternativa al coarse-grained è il *fine-grained multithreading* in questo caso il MT preleva un istruzione da un thread diverso ad ogni ciclo, ovvero l'esecuzione di thread multipli viene intervallata in un circolo *round-robin* e si saltano quei thread che sono bloccati. Il processore deve essere in grado di cambiare thread ad ogni ciclo ma questo permette di nascondere gli stalli mandando in esecuzione altri thread e non considerando quello bloccato. Tuttavia il tempo di esecuzione di un thread è rallentato perché un thread pronto deve comunque aspettare l'esecuzione di altri thread. Anche in questo caso inoltre abbiamo degli issue slot vuoti dovuti alla limitazione dell'ILP. Un esempio di *fine-grained MT* è mostrato in Figura 72(a) Un esempio del

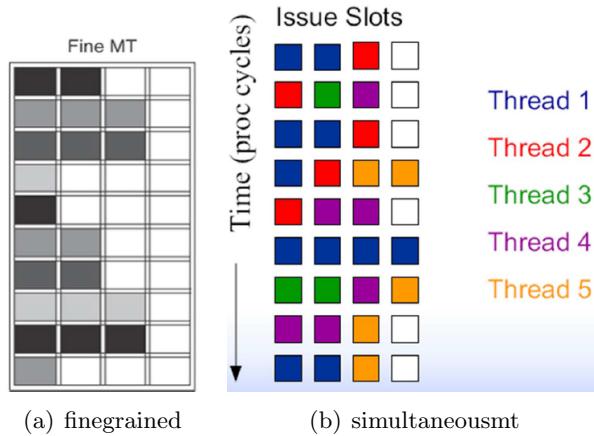


Figura 72: Esempio di fine-grained MT 72(a) e di simultaneous MT 72(b)

*fine-grained multithread* è il **Sun Niagara T1** un processore ad 8 core nel quale un singolo core può gestire fino ad un massimo di 4 thread per un totale di 32 thread gestiti. Ad ogni ciclo di clock ogni core esegue un thread diverso tra i quattro disponibili. Quando un thread è bloccato

il core sul quale è eseguito manda in esecuzione un altro thread solo quando tutti e quattro i thread di un core sono bloccati allora il core risulta in stallo.

In pratica il multithreading permette di nascondere eventi di latenze molto lunghe e mantenere occupate le unità funzionali.

A questo punto però ci chiediamo perché non sfruttare sia ILP che il TLP contemporaneamente. Per fare questo utilizziamo il *simultaneous multithreading* come mostrato in Figura 72(b). In questo meccanismo più thread utilizzano i diversi issue slot nello stesso ciclo di clock, la risoluzione delle dipendenze avviene tramite scheduling dinamico, il register renaming permette di utilizzare identificativi univoci per mischiare istruzioni provenienti da thread diversi. L'utilizzo degli issue slot è limitato solo dalla loro occupazione. L'unico difetto di questa tecnica è che inevitabilmente si compromette il tempo di esecuzione del singolo thread. Il fattore principale che ha spinto l'introduzione del SMT è il fatto che le CPU moderne hanno più unità funzionali di quanto un singolo thread può sfruttare, tale implementazione è la più comune nei processori Intel Core i7.

Per implementare tale meccanismo bisogna però migliorare anche l'unità di *Fetch* in quanto essa deve essere in grado di prelevare le istruzioni da più thread ed inoltre decidere da quali thread prelevare. I vantaggi di questa tecnica sono i ridotti salti non risolti, i minori problemi di load miss ed infine la minore quantità di istruzioni in coda.

## 7 Gerarchie di memoria

Per incrementare le performance di un computer bisogna considerare anche il suo sistema di memoria, bisogna infatti dare l'illusione di avere un sistema di memoria che sia simultaneamente ampio e veloce, inoltre bisogna fornire dati al processore ad alta frequenza.

Per fornire queste due caratteristiche dobbiamo sfruttare alcune caratteristiche, prima caratteristica da sfruttare è quella della **località**. La località è una caratteristica dei dati e può essere di due tipi:

**Località temporale:** quando vi è un riferimento ad un elemento di memoria, la tendenza è quella di riferirsi allo stesso elemento nei momenti successivi (ad esempio in un loop).

**Località spaziale:** quando vi è un riferimento ad un elemento di memoria la tendenza è quella di accedere ai dati vicini nel tempo seguente, come nel caso di istruzioni sequenziali.

Un'altra soluzione è quella di sfruttare delle gerarchie di memoria come quelle mostrate in Figura 73, ogni livello di questa gerarchia ha dimensioni e velocità diverse implementate attraverso diverse tecnologie. Lo scopo è quello di fornire all'utente una grande quantità di memoria ad un costo contenuto ma fornendo comunque un tempo di accesso ridotto grazie alle tecnologie più veloci.

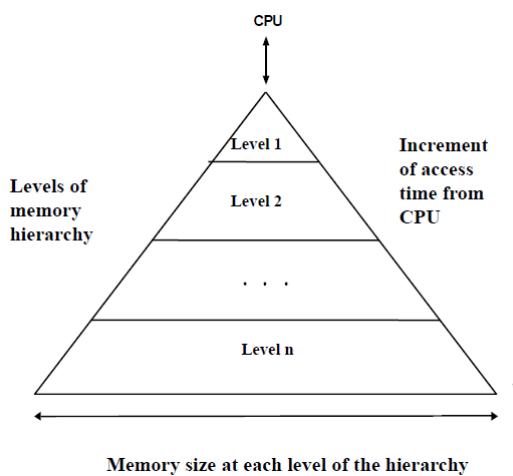


Figura 73: Esempio di gerarchia di memoria

### 7.1 Caches

La cache è un meccanismo che implementa i due obiettivi prefissati, prima di tutto sfrutta entrambi i tipi di località, sfrutta la *località temporale* mantenendo i dati acceduti più recentemente e sfrutta la *località spaziale* prelevando dal livello superiore un blocco di dati nell'intorno del dato acceduto. In Figura 74 vediamo i livelli di caches e le gerarchie di memoria usati in due sistemi differenti, un sistema server ed un device mobile. Il design della gerarchia di memoria è diventato importante con i recenti processori multi-core, infatti la banda necessaria cresce con l'aumentare del numero di core, ad esempio un processore *Core i7* alla frequenza di 3.2GHz genera due riferimenti a memoria per ogni core per clock che moltiplicati per i quattro core sono 25.6 miliardi di riferimenti a 64-bit al secondo più 12.8 miliardi di istruzioni a 128 bit per un totale di 409.6 GB/s. La banda di una memoria DRAM è di solo 25 GB/s ovvero il 6%

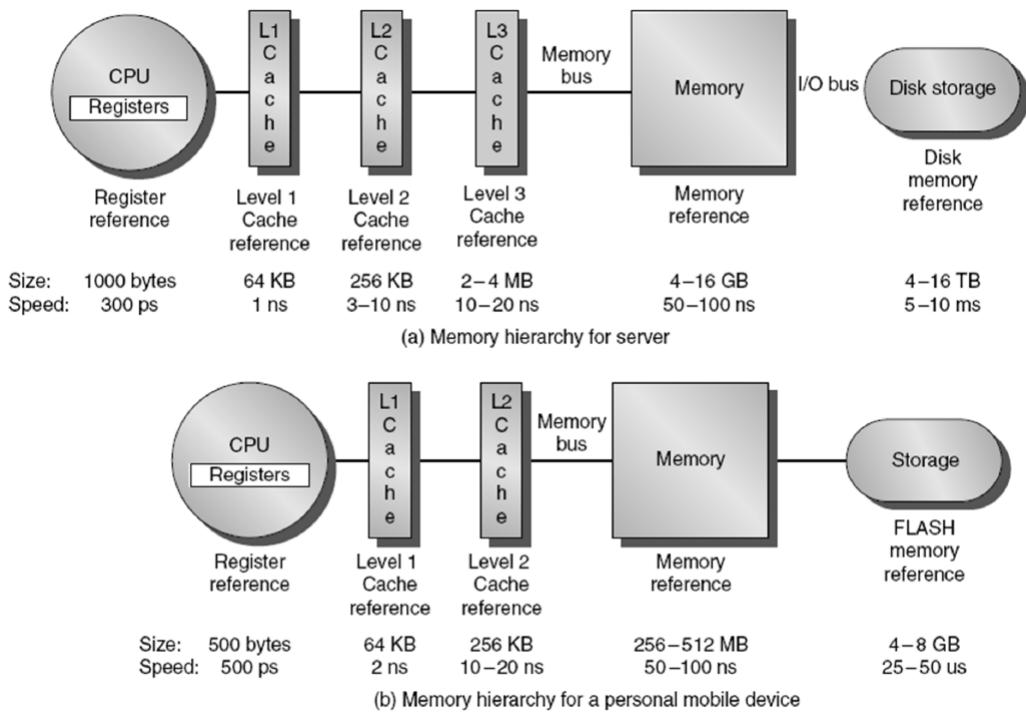


Figura 74: Sistema di memoria in un sistema server e in un device mobile.

per implementare un sistema efficiente è perciò necessario avere una cache multiport, due livelli di cache per singolo core e un terzo livello di cache condiviso. Un esempio di tale sistema è mostrato in Figura 75.

### 7.1.1 Struttura e funzionamento di una cache

Come abbiamo appena visto le gerarchie di memoria sono composte da diversi livelli i dati tuttavia sono copiati soltanto tra due livelli adiacenti. Per semplificare la nostra trattazione consideriamo soltanto due livelli, la cache e la memoria principale. La cache (livello **superiore**) è piccola e veloce ma molto costosa. La memoria centrale (livello **inferiore**) è di grandi dimensioni poco costosa e più lenta rispetto alla cache. La minima porzioni di dati che può essere copiata nella cache è chiamata **blocco** o **cache line**. Per sfruttare la *località spaziale* è necessario che la dimensione di un blocco sia un multiplo della dimensione di una parola. Per determinare il numero di blocchi che possono essere caricati in cache basta applicare la formula:

$$\# \text{ blocchi cache} = \text{cache size}/\text{block size}$$

Di seguito alcune definizioni che ci aiuteranno nella nostra esposizione.

**Cache Hit** Si ha un *cache hit* se ad una richiesta di dati si trova il dato già presente in cache.

**Cache Miss** Quando si richiede un dato ma tale dato non è presente in nessun blocco di cache allora si ha un *cache miss*. Per trovare tale blocco è necessario accedere al livello che sta più in basso nella gerarchia di memoria. In caso di miss è necessario bloccare la CPU e richiedere il blocco alla memoria centrale, copiare il blocco in cache e poi ripetere l'accesso alla cache per effettuare uno *hit*.

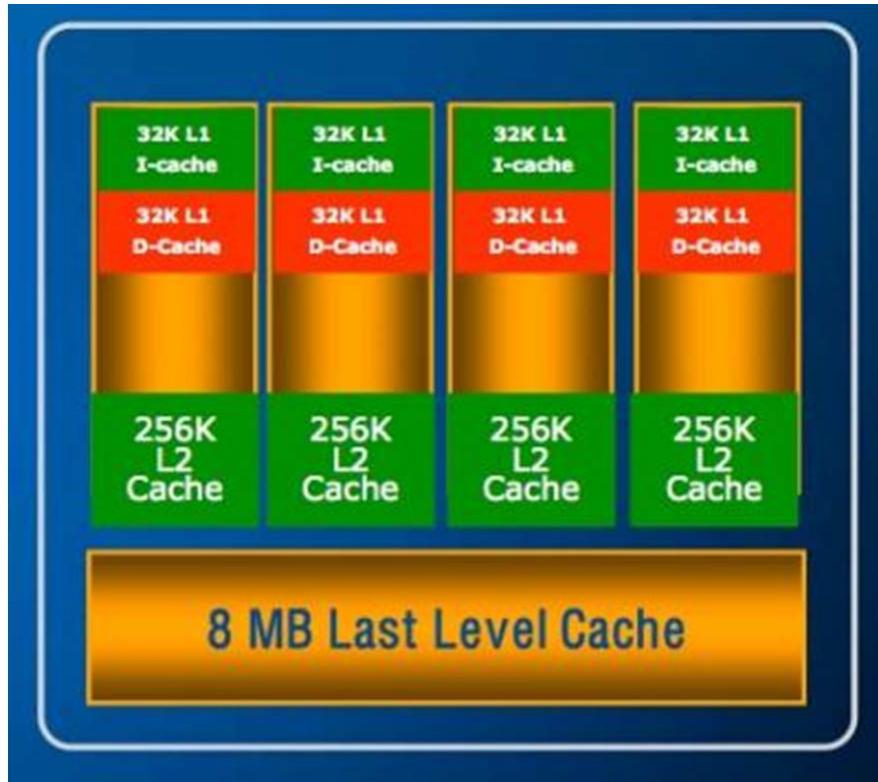


Figura 75: Architettura cache di un processore Intel Core i7

**Hit Rate** L'*hit rate* è il numero di accessi in memoria che trovano il dato in cache rispetto al numero di accessi totali.

$$Hit\ Rate = \frac{\# Hit}{\# Memory\ Accesses}$$

**Hit Time** L'*hit time* è il tempo necessario per accedere al dato quando esso è presente nel livello più alto della cache, questo tempo include il tempo per decidere se il dato è presente.

**Miss Rate** Il *miss rate* è il numero di accessi a memoria che non trovano il dato nel livello più alto e che quindi necessitano di un accesso ad un livello più basso di memoria rispetto al numero di accessi totali. Data questa definizione abbiamo che:

$$Miss\ Rate + Hit\ Rate = 1$$

**Miss Penalty** Il *miss penalty* è il tempo necessario per accedere al livello inferiore della memoria e rimpiazzare il blocco contenente il dato al livello superiore.

**Miss Time** Tempo totale per accedere ad un dato che non si trova nel livello più alto della cache è dato da:

$$Miss\ Time = Hit\ Time + Miss\ Penalty$$

**Tempo medio di accesso** Il *tempo medio di accesso* alla memoria (AMAT) è dato dalla percentuale di hit per il tempo di accesso più la percentuale di miss per il relativo tempo di



Figura 76: Struttura di un record di cache

accesso.

$$AMAT = Hit\ Rate * Hit\ Time + Miss\ Rate * Miss\ Time$$

Dato che

$$Miss\ Time = Hit\ Time + Miss\ Penalty$$

e

$$Miss\ Rate + Hit\ Rate = 1$$

allora possiamo scrivere il tempo medio di accesso come:

$$AMAT = Hit\ Time + Miss\ Rate * Miss\ Penalty$$

Analizziamo ora come sono strutturate le cache e le loro tecnologie di implementazione. Innanzitutto ogni record di una cache contiene:

**Bit di validità:** per indicare se la posizione corrente contiene dei dati validi oppure no, durante il *bootstrap* del sistema tutti i record vengono marchiati come *invalidi*.

**Tag:** Questo campo contiene un valore che identifica inequivocabilmente l'indirizzo di memoria corrispondente ai dati immagazzinati.

**Data:** Contiene una copia dei dati della memoria.

Tale struttura è rappresentata in Figura 76. Esistono tre metodi per implementare un sistema di cache, essi sono rispettivamente la cache *direct map*, la cache *fully associative* e la cache *associativa a n vie*. Questa distinzione di tipologie è dovuta al problema del **piazzamento dei blocchi** ovvero come assegnare la posizione in memoria cache ad un blocco che viene caricato in memoria centrale.

**Cache direct map** Nelle cache di tipo *direct map* ogni indirizzo di memoria corrisponde ad uno ed un solo blocco in cache tale blocco di cache è determinato da:

$$(Block\ Address)_{cache} = (Block\ Address)_{mem} \ mod \ (\# \ cache \ blocks)$$

Un esempio di tale meccanismo è mostrato in Figura 77 L'indirizzamento avviene nel modo descritto dalla Figura 78 dove dato un indirizzo di  $N$  bit viene suddiviso in 4 campi:

**Byte offset:** che serve ad identificare il byte dato all'interno della parola, nel caso in cui la memoria non sia indirizzabile al byte allora  $B = 0$

**Byte offset:** che serve ad identificare la parola all'interno del blocco , nel caso in cui il blocco abbia dimensione di una parola allora  $K = 0$

**Index:** identifica il blocco tramite  $M$  bit, dove

$$M = \log_2 \# \ Block$$

**Tag:** necessario per comparare il blocco selezionato tramite l'index, la dimensione del tag è data da:

$$N - (M + K + B)$$

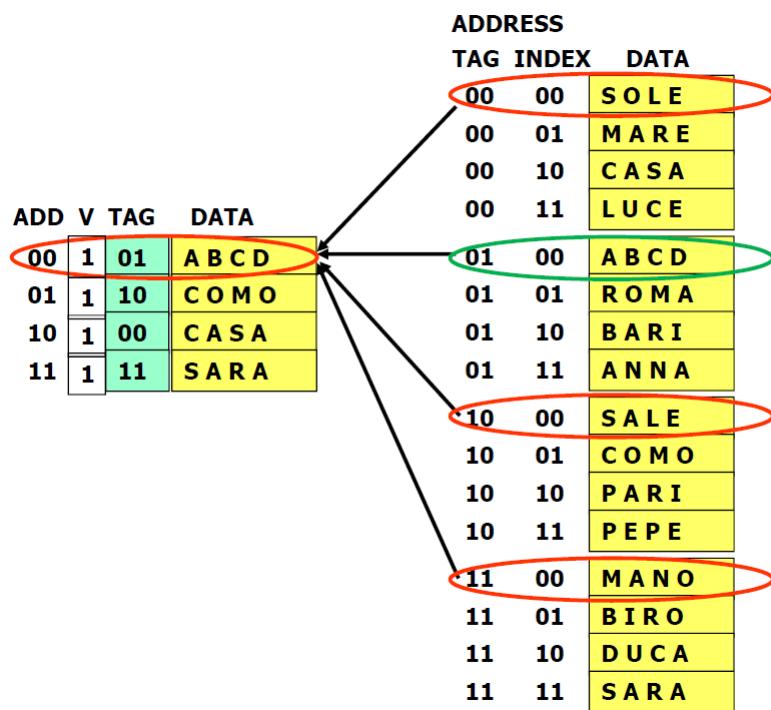


Figura 77: Esempio di piazzamento di blocchi nel caso di cache di tipo direct map

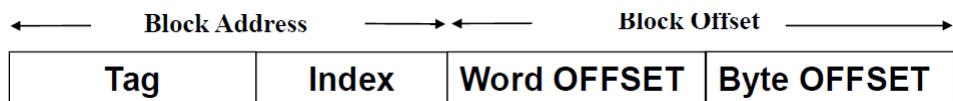


Figura 78: Indirizzamento in un sistema di cache direct map

**Cache completamente associativa** In una cache di tipo *fully associative* un blocco di memoria può essere posizionato in un qualsiasi blocco cache, durante la ricerca di un blocco tutti i blocchi cache vengono controllati e il loro campo *tag* viene comparato con quello del blocco da cercare. Il campo *index* non esiste nella cache completamente associativa. I campi che compongono un record in una cache completamente associativa sono mostrati in Figura 79 e sono rispettivamente:

**Byte offset:** che serve ad identificare il byte dato all'interno della parola, nel caso in cui la memoria non sia indirizzabile al byte allora  $B = 0$

**Byte offset:** che serve ad identificare la parola all'interno del blocco , nel caso in cui il blocco abbia dimensione di una parola allora  $K = 0$

**Tag:** Serve ad identificare il blocco in cache ed è uguale a

$$N - (B + K)$$

Tag (28 bit)	WO (2 bit)	BO (2 bit)
--------------	------------	------------

Figura 79: Campi di una cache completamente associativa

Un esempio di cache completamente associativa è mostrato in Figura 80

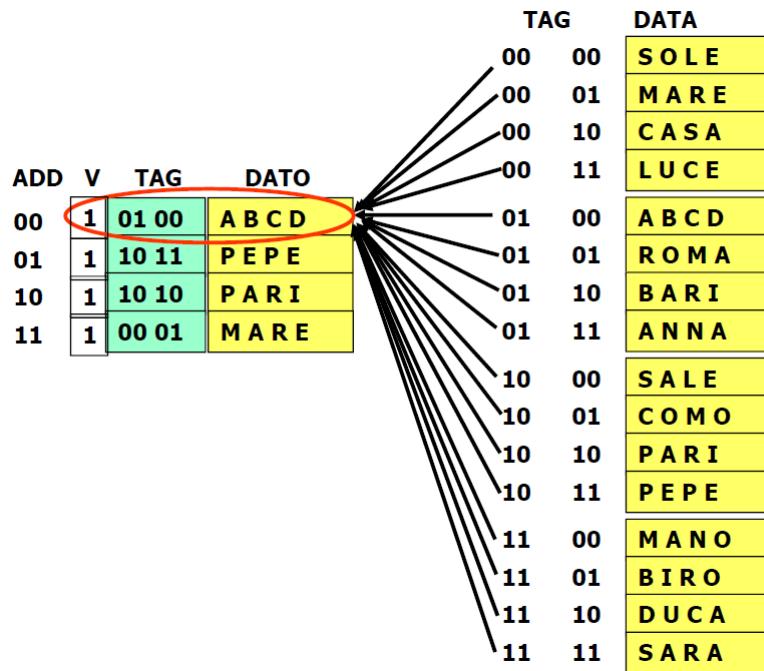


Figura 80: Esempio di una cache completamente associativa

**Cache associativa ad  $n$  vie** In questo caso la cache è composta da dei *set* ed ogni set è composto da un determinato numero di blocchi.

$$\# \text{ Blocchi} = \text{Cache size}/\text{Block size}$$

$$\# Sets = Cache size / (Block size * n)$$

Un blocco di memoria può essere posto in uno qualsiasi dei blocchi di un set e la ricerca deve avvenire in tutti i blocchi di quel set. Per calcolare in quale set un blocco dovrà essere posizionato bisogna utilizzare la formula:

$$(Set)_{cache} = (Block address)_{mem} \bmod (\text{Num sets in cache})$$

Un esempio di utilizzo di una cache set associativa è mostrato in Figura 81 L'indirizzamento in

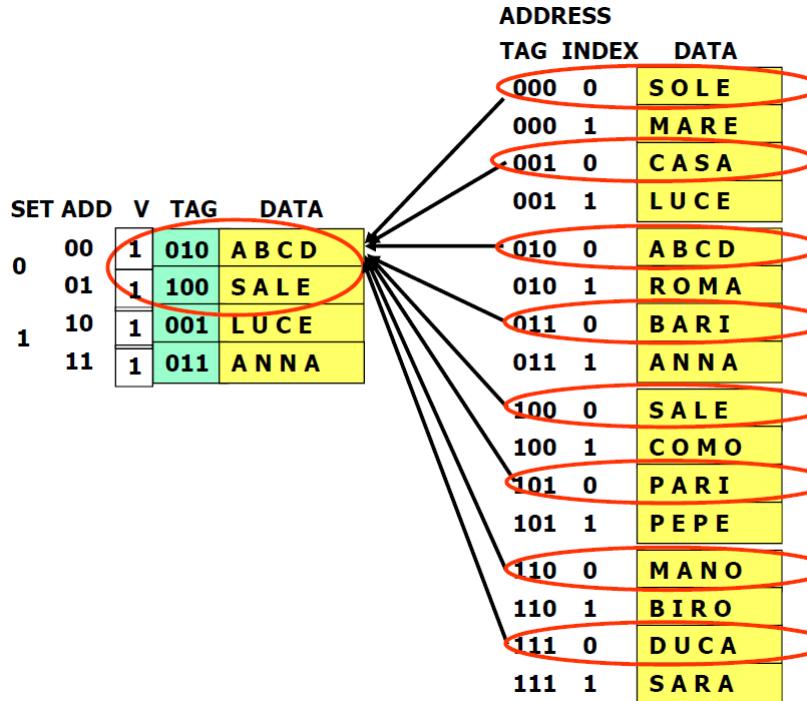


Figura 81: Esempio di una cache associativa a due vie

una cache associativa ad  $n$  vie avviene come nel caso di cache di tipo direct map ma questa volta il campo index serve ad identificare il set nel quale effettuare la ricerca del blocco come mostrato in Figura 82 Il problema principale che si ha quando si utilizza un meccanismo di cache

Tag	Index	Word OFFSET	Byte OFFSET
-----	-------	-------------	-------------

Figura 82: Indirizzamento in una cache associativa a due vie

è quello di decidere dove posizionare i blocchi di memoria quando questi sono copiati in cache in quanto la memoria è molto più grande di quella cache. I tre meccanismi di cache che abbiamo appena visto risolvono questo problema in modi diversi come possiamo vedere dalla Figura 83. Come vediamo nel caso di completamente associativa il blocco può essere posizionato in qualsiasi blocco cache, questo significa che per verificare che un blocco sia già presente in cache bisogna controllare tutti i blocchi della cache aumentando così il tempo di servizio. Nel caso di cache *direct map* il blocco può essere posizionato in un solo blocco di cache questo fa sì che la ricerca sia molto veloce ma può portare ad altri problemi (cosa succede se il blocco 4 contiene delle

istruzioni che sono contenute nel blocco 12?). Nel caso di cache associativa a  $n$  vie abbiamo una via di mezzo tra le precedenti infatti il blocco può essere posizionato in qualsiasi blocco all'interno del set ma la ricerca avviene soltanto all'interno del set e non su tutta la cache. La



Figura 83: Problema del *block placement* nei tre meccanismi di cache

soluzione al problema del piazzamento dei blocchi sembrerebbe allora quello dell'incremento del livello di associatività. Tuttavia pur riducendo il *miss rate* aumentano in modo esponenziale i costi di implementazione ed inoltre incrementa anche il tempo di *hit*. La decisione di quale tecnologia utilizzare dipende dal tradeoff di costi e riduzione di *miss rate*.

In caso di miss dobbiamo decidere quale blocco sostituire, in caso di completamente associativa i blocchi candidati sono tutti mentre nel caso di set associativa i blocchi candidati sono soltanto quelli del set corrispondente al blocco da inserire. Esistono diverse strategie per effettuare la selezione e sono:

- Random or pseudo random
  - LRU (Least Recently Used)
  - FIFO (First In First Out)

Un altro problema che si viene a creare quando dobbiamo sostituire un blocco è quello della *write policy* ovvero come ci dobbiamo comportare se il blocco che stiamo sostituendo è stato modificato. Esistono due tecniche principali:

**Write-through:** in questo caso ad ogni scrittura vengono aggiornati sia il blocco in cache che il blocco in memoria, questo meccanismo ci permette di non preoccuparci della sostituzione ed è di facile implementazione, tuttavia richiede la presenza di un *write buffer* come quello in Figura 84.

**Write-back:** questo meccanismo prevede che le modifiche ad un blocco siano scritte solamente in cache e tali modifiche siano scritte in memoria centrale solo quando il blocco in cache viene sostituito, per implementare questo meccanismo è necessario aggiungere ai campi della cache un *dirty bit* che indica se il blocco è stato modificato durante la sua permanenza in cache. Questo meccanismo permette scritture più veloci e anche scritture multiple su uno stesso blocco richiedono un'unica scrittura in memoria centrale.

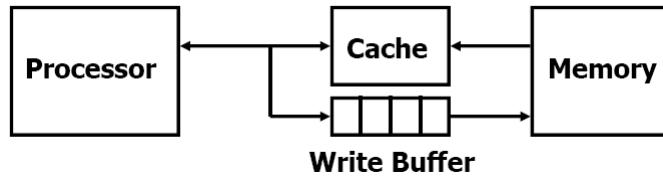


Figura 84: Utilizzo del write buffer

Per implementare un meccanismo di *write-through* efficace abbiamo visto come sia necessario un *write buffer*, ovvero un buffer di tipo FIFO che permette al processore di non attendere i tempi di scrittura della memoria. Il processore scrive i dati in cache e nel write buffer, un controllore scrive i dati del buffer nella memoria più bassa. Il problema della scrittura tuttavia non è ancora risolto in quanto il buffer si può saturare e a quel punto è necessario bloccare il processore in attesa della scrittura.

La scrittura di dati in memoria presenta un ulteriore problema, infatti può capitare che sia necessario scrivere su di un blocco che non è presente in cache, anche per questo caso esistono due soluzioni distinte:

**Write allocate:** in questo caso prima di scrivere il dato viene allocato in cache il blocco corrispondente allo spazio di memoria che dovrebbe essere scritto e la scrittura avviene sul blocco appena allocato.

**No write allocate:** in questo caso il dato da scrivere viene semplicemente inviato alla memoria di più basso livello senza allocare nuovi dati in cache.

Le tecniche di *write-back* o *write-through* e quelle di *write allocate* e di *no write allocate* possono teoricamente essere combinate in qualsiasi modo, ma in pratica si utilizzano solo le combinazioni *write-back* e *write allocate* e *write-through* con *no write allocate*.

### 7.1.2 Analisi delle performance

Come abbiamo visto nei capitoli precedenti abbiamo che

$$CPU_{time} = (CPU_{exec-cycles} + Memory\ stall\ cycle) \times T_{CLK}$$

dove:

$T_{CLK}$ : periodo di tempo del clock

$CPU_{exec-cycles}$ :  $IC \times CPI_{exec}$

$IC$ : Instruction count

$Memory\ stall\ cycle$ :  $IC \times Miss\ per\ instr \times Miss\ penalty$

A questo punto possiamo riscrivere la precedente equazione come:

$$IC \times (CPI_{exec} + Miss\ per\ instr \times Miss\ penalty) \times T_{CLK}$$

dove

$$Miss\ per\ instr = MemoryAccessPerInstruction \times Miss\ Rate$$

$$IC \times (CPI_{exec} + MAPI \times Miss\ Rate \times Miss\ penalty) \times T_{CLK}$$

Questa formula non tiene conto degli stalli dovuti ai conflitti della pipeline ma solo di quelli dovuti all'accesso in memoria ma sono questi quelli importanti per la nostra trattazione.  
A questo punto vogliamo ridurre il tempo di accesso ai dati ovvero ridurre *AMAT*

$$AMAT = Hit\ Time + Miss\ Rate * Miss\ Penalty$$

questo significa ridurre uno dei tre componenti della formula, lo *hit time*, il *miss rate* o il *miss penalty*. Per fare ciò introduciamo un secondo livello di cache, il primo livello (L1) abbastanza veloce da soddisfare il più veloce ciclo di clock, il secondo livello (L2) grande abbastanza da catturare la maggior parte degli accessi destinati alla memoria in modo da ridurre gli effetti del miss penalty. Il tempo medio di accesso per quanto riguarda il primo livello di cache è dato da:

$$AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1} \times Miss\ Penalty_{L1}$$

dove il *miss penalty* è dato da:

$$Miss\ Penalty_{L1} = Hit\ Time_{L2} + Miss\ Rate_{L2} \times Miss\ Penalty_{L2}$$

Otteniamo così che il tempo medio di accesso con due livelli di cache è dato da

$$AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1} \times (Hit\ Time_{L2} + Miss\ Rate_{L2} \times Miss\ Penalty_{L2})$$

Possiamo però fare una distinzione tra i miss rate possiamo definire un **local miss rate** ed un **global miss rate**. Per calcolare un *local miss rate* si calcolano il numero di miss nel livello di cache in esame diviso il numero di accessi in quel livello di cache otteniamo così un *miss rate* per L1 e un *miss rate* per L2. Nel caso di *global miss rate* invece si calcolano il numero di miss in quel livello di cache diviso il numero totale di accessi a memoria generati dalla CPU. Per il livello L1 il *miss rate* coincide con quello locale mentre per il livello L2 abbiamo che  $Miss\ Rate_{L1L2} = Miss\ Rate_{L1} \times Miss\ Rate_{L2}$ . Il miss rate di tipo globale è una misura più veritiera di quanti accessi a memoria effettuati dalla CPU arrivano realmente alla memoria centrale. Possiamo così riscrivere il tempo medio di accesso come:

$$AMAT = Hit\ Time_{L1} + Miss\ Rate_{L1} \times Hit\ Time_{L2} + Miss\ Rate_{L1L2} \times Miss\ Penalty_{L2}$$

A questo punto possiamo analizzare l'impatto dei memory miss sul tempo di esecuzione, questo impatto è dato da

$$CPU_{time} = IC \times (CPI_{exec} + MAPI \times MR_{L1} \times HT_{L2} + MAPI \times MR_{L1L2} \times MP_{L2}) \times T_{CLK}$$

## 7.2 Incremento delle performance

Fino ad ora abbiamo visto come sono costruite e come funzionano le diverse tipologie di cache, ora ci chiediamo come incrementarne le performance. Partendo dal tempo medio di accesso a memoria:

$$AMAT = HT + MR * MP$$

Da questa formula vediamo che per ridurre il tempo di accesso a memoria possiamo ottimizzare uno dei tre fattori ovvero:

- Ridurre il *miss rate*
- Ridurre il *miss penalty*
- Ridurre lo *hit time*

**Cache miss** Prima di iniziare la nostra trattazione vediamo quali tipi di cache miss possiamo incontrare, possiamo distinguerne quattro tipologie:

- Compulsory Misses
- Capacity Misses
- Conflict Misses
- Coherence Misses

La *compulsory misses* si ha all'avvio di un programma o del calcolatore quando tutti i blocchi di cache sono segnati come invalidi e per questo ogni nuovo caricamento di un blocco porta ad un miss. La *capacity misses* si ha quando la cache non è in grado di contenere tutti i blocchi necessari all'esecuzione tale tipo di miss decresce con l'aumentare della capacità della cache. La *conflict misses* avviene solo nel caso in cui la cache sia di tipo direct map o set associative ed avviene perché un blocco appena eliminato per far spazio ad un altro blocco sarà necessario nel breve periodo; tale miss è anche chiamata *collision misses*. L'ultima tipologia di miss è relativamente nuova ed è causata da problemi di coerenza tra cache in un ambiente multiprocessore.

### 7.2.1 Ridurre il miss rate

Per ridurre il miss rate in molti casi basta ridurre le capacity miss, il modo più semplice per farlo è aumentare la dimensione della cache, tuttavia questo incrementa il tempo di hit nonché l'area e il consumo di energia.

Un altro metodo è quello di incrementare la dimensione dei blocchi di cache, così facendo si sfrutta maggiormente la località spaziale, tuttavia aumentando troppo la dimensione del blocco si ottiene l'effetto opposto aumentando il miss rate inoltre l'aumento della dimensione del blocco di cache aumenta il miss penalty e riduce il numero di blocchi disponibili aumentando la possibilità di *conflict misses*.

Per ridurre la possibilità di *conflict misses* si può aumentare l'associatività della cache tuttavia questo comporta un aumento di logica con conseguente aumento dell'area di silicio e del consumo di energia, inoltre si ha anche un aumento del tempo di *hit*.

Un problema di difficile risoluzione è quello di ridurre i *conflict miss* mantenendo ridotti i tempi di hit come nel caso di cache *direct map*. Una soluzione è quella di introdurre un buffer nel quale immagazzinare i blocchi scartati dalla cache in modo da migliorare l'utilizzo della località temporale. Tale buffer viene chiamato *victim cache* ed è una cache di tipo completamente associativo di piccole dimensioni, infatti una victim cache di soli 4 registri associata ad una cache direct map di 4KB riduce fino al 95% dei conflitti. Il funzionamento della *victim cache* è molto semplice essa è situata tra la cache e il suo percorso di caricamento, ogni volta che avviene un miss si controlla se nella victim cache è presente il blocco mancante prima di inoltrare la richiesta al livello più basso. Nel caso in cui il blocco sia presente nella victim cache avviene una sostituzione tra un blocco in cache e il blocco necessario nella victim cache.

Analizziamo ora due metodi che ci permettono di ottenere un breve tempo di hit mantenendo comunque un livello di miss da conflitto dell'ordine di una cache a due vie, il primo meccanismo si applica alle cache di tipo *direct map* e viene denominato *divide cache*, quando avviene un miss si controlla l'altra metà della cache. Il secondo metodo denominato *way prediction* si applica alle cache di tipo associativo a due vie e consiste nell'utilizzare un bit extra per predire quale via sarà intrapresa durante il prossimo accesso in cache in modo da pre settare il multiplexer e ridurre i tempi di accesso, in caso la predizione sia corretta si ottiene un tempo di accesso pari allo *hit time* in caso di predizione sbagliata abbiamo un tempo di accesso più lento. Questa

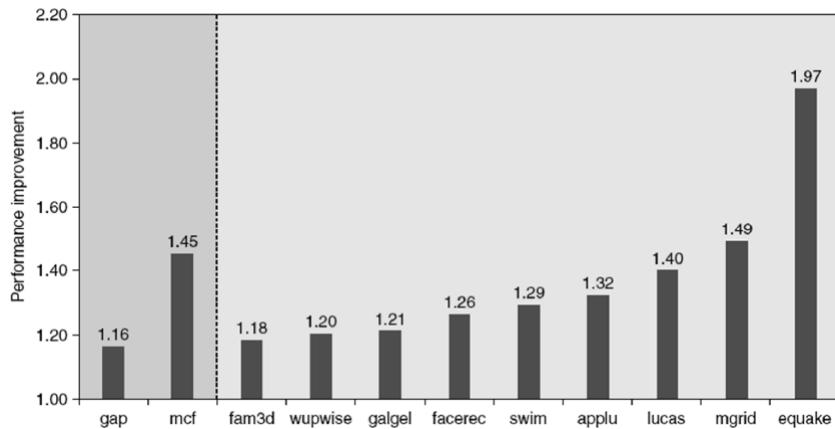


Figura 85: Incremento di performance dovuto al pre-fetching di due blocchi quando avviene un miss

tecnica oltre a ridurre il tempo di accesso riduce anche il consumo di energia. Tuttavia questo meccanismo è efficace per le cache che non sono a diretto contatto con il processore in quanto uno *slow hit* degrada notevolmente le performance, inoltre tale tecnica è efficace se si ha un buon predittore e quindi il numero di predizioni corrette è maggiore di quelle sbagliate altrimenti le performance degradano sensibilmente.

Un altro modo per ridurre il tempo di hit e il miss rate è quello di utilizzare delle tecniche hardware di *pre-fetching* ovvero la possibilità di prelevare più di un blocco nell'istante di prelevamento , questo meccanismo aumenta virtualmente la banda disponibile tuttavia nel caso vi sia un'interferenza tra *pre-fetching* e misses si ha un degrado delle performance. Un esempio dei vantaggi di questa tecnica è mostrato in Figura 85 nella quale si mostra l'incremento di performance grazie al prelevamento di due blocchi quando avviene un miss (includendo il blocco sequenzialmente successivo). Una tecnica simile è quella del *software pre-fetching* ovvero il compilatore inserisce delle istruzioni dette di pre-fetching per richiedere dei dati prima che essi siano necessari. Questo meccanismo può essere di due tipi:

**Register prefetch:** i dati vengono caricati nei registri

**Cache prefetch:** i dati vengono caricati in cache

Tuttavia questa tecnica è efficace solo se il costo in termini di tempo di esecuzione e aumento della complessità del codice è minore del vantaggio che si ha in termini di riduzioni di misses. Esiste, infine un ramo di ottimizzazioni che riguardano la riduzione dei misses tramite l'utilizzo di tecniche di ottimizzazione effettuate dal compilatore. Tali tecniche sono:

**Mearging arrays:** che permette di incrementare la località spaziale di una struttura di elementi rispetto a due array.

**Loop Interchange:** incrementa la località spaziale scambiando l'innestamento dei loop per permettere un accesso ai dati in modo più continuo.

**Loop Fusion:** incrementa la località spaziale combinando due cicli indipendenti che eseguono lo stesso loop su alcune variabili sovrapposte.

**Bloccking:** incrementa la località temporale tramite l'accesso a blocchi di dati ripetutamente al posto di accedere per righe o per colonne.

Un esempio di *measuring array* è mostrato in Figura 86 nel quale si riducono i conflitti tra le variabili `val` e `key` incrementando la località spaziale. Un esempio di *loop interchange* è mo-

```
/* Before: 2 sequential arrays */
int val[SIZE];
int key[SIZE];
/* After: 1 array of structures */
struct merge {
    int val;
    int key;
};
struct merge merged_array[SIZE];
```

Figura 86: Esempio di measuring array

strato in Figura 87 nel quale si ha un accesso a memoria con un avanzamento di 100 parole anziché 5000 aumentando così la località spaziale. Nell'esempio di Figura 88 viene mostrato

```
/* Before */
for (k = 0; k < 100; k = k+1)
    for (j = 0; j < 100; j = j+1)
        for (i = 0; i < 5000; i = i+1)
            x[i][j] = 2 * x[i][j];

/* After */
for (k = 0; k < 100; k = k+1)
    for (i = 0; i < 5000; i = i+1)
        for (j = 0; j < 100; j = j+1)
            x[i][j] = 2 * x[i][j];
```



Figura 87: Esempio di loop interchange

come incrementare la località spaziale tra due cicli che lavorano su alcune variabili comuni. Per quanto riguarda il *blocking* la Figura 89(a) e 89(b) mostrano come si effettua un blocking ovvero si modifica l'accesso ad una matrice non più per righe o colonne ma tramite dei blocchi, tale meccanismo richiede più memoria ma migliora la località spaziale, nell'esempio il fattore  $B$  è chiamato *fattore di blocco*.

### 7.2.2 Ridurre il miss penalty

Come per il *miss rate* anche per il *miss penalty* esistono diverse tecniche per la riduzione di questo fattore. La prima tecnica che analizziamo è denominata *read priority* e si basa sull'idea di dare una maggiore priorità ai miss che riguardano un dato in lettura rispetto a dei miss che riguardano una scrittura, per implementare questa tecnica è necessario dimensionare adeguatamente il *write buffer*. Questa tecnica complica leggermente l'accesso a memoria in quanto il write buffer deve immagazzinare gli aggiornamenti se un area di memoria è richiesta da una read miss. Tuttavia questa tecnica ha notevoli svantaggi, il *write-through* con un write buffer

```
/* Before */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        d[i][j] = a[i][j] + c[i][j];
/* After */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
    {
        a[i][j] = 1/b[i][j] * c[i][j];
        d[i][j] = a[i][j] + c[i][j];}
```

Figura 88: Esempio di loop fusion

```

/* Before */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        r = 0;
        for (k = 0; k < N; k = k+1){
            r = r + y[i][k]*z[k][j];};
        x[i][j] = r;
};

/* After */
for (jj = 0; jj < N; jj = jj+B)
    for (kk = 0; kk < N; kk = kk+B)
        for (i = 0; i < N; i = i+1)
            for (j = jj; j < min(jj+B-1,N); j = j+1)
                r = 0;
                for (k = kk; k < min(kk+B-1,N); k = k+1) {
                    r = r + y[i][k]*z[k][j];
                    x[i][j] = x[i][j] + r;
                };

```

Figura 89: Esempio di utilizzo del blocking

può portare a dei conflitti di tipo RAW.

Un'altra tecnica è quella del *sub-block placement* nel quale non è necessario caricare un intero blocco quando si verifica un miss, si hanno dei *bit di validità* per ogni sotto blocco.

Due tecniche molto utilizzate sfruttano il fatto che ogni qualvolta si verifica una miss la CPU necessita solamente di una parola del blocco e quindi non è necessario attendere il caricamento dell'intero blocco prima di ricominciare l'esecuzione; queste due tecniche sono **early restart** che preleva dalla memoria il blocco nel suo normale ordine ma non appena la parola necessaria alla CPU è disponibile allora essa viene inviata e l'esecuzione riprende mentre il blocco termina il caricamento. La seconda tecnica denominata **Critical Word First** prevede di recuperare innanzitutto la parola necessaria all'esecuzione e dopo completare il caricamento del blocco. Queste due tecniche sono molto utili soprattutto per blocchi cache di grosse dimensioni, tuttavia può creare alcuni problemi nel caso di località spaziale se la CPU richiede le parole sequenzialmente seguenti.

La tecnica di **non-blocking cache** prevede che la cache continui a fornire dati alla CPU durante un evento di miss, per fare ciò è necessario che la CPU sia in grado di effettuare una esecuzione fuori ordine. Il funzionamento base prevede che la CPU prosegua la sua esecuzione dopo aver effettuato una richiesta che ha sollevato una miss. Un'evoluzione di questa tecnica ha permesso di sovrapporre più eventi di miss ed è denominata *hit under multiple miss* o *miss under miss*.

tuttavia questa tecnica aumenta di molto la complessità della cache ed inoltre richiede diversi banchi di memoria.

Una tecnica che abbiamo visto nella sezione precedente è quella di aumentare i livelli di memoria per diminuire il miss penalty per la spiegazione si rimanda a 7.1.2.

Una tecnica che permette di ridurre gli stalli dovuti alla scrittura è quella del **merging write buffer** che consiste nell'aggiornare il blocco che è eventualmente già presente nel *write buffer* piuttosto di inserire ulteriori blocchi.

### 7.2.3 Riduzione dello hit time

Come ultima analisi vediamo come ridurre lo *hit time* per migliorare il tempo di accesso alla memoria. Come prima cosa è utile utilizzare un primo livello di cache di dimensioni ristrette con una associatività molto bassa, tali caratteristiche permettono di ottenere un breve tempo di accesso in quanto il tempo di hit è dato dal tempo necessario per effettuare le seguenti tre operazioni:

- indirizzare il tag in memoria
- comparare i tag
- selezionare il set corretto

In una cache di tipo direct map si possono sovrapporre la comparazione dei tag e la trasmissione dei dati riducendo così il tempo d'accesso inoltre accedendo ad un minor numero di linee si consuma un minor quantitativo di energia. La Figura 90 mostra i tempi di accesso in base alla associatività e alla dimensione dei blocchi. Un'altra tecnica per ridurre il tempo di hit è

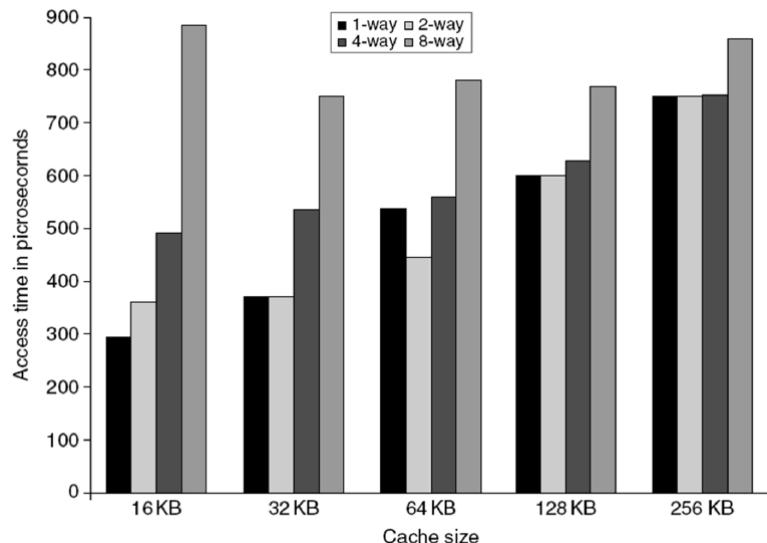


Figura 90: Tempi di hit per dimensione e livello di associatività di una cache.

quella di risolvere il problema della traduzione degli indirizzi. Ogni qualvolta la CPU effettua un cambio di contesto ovvero passa all'esecuzione di un altro processo è necessario svuotare la cache altrimenti si otterrebbero dei falsi hit. Una soluzione è quella di permettere che ogni blocco cache ha solo un indirizzo fisico che corrisponde agli n bits meno significativi di quelli virtuali, questa tecnica è chiamata *page coloring*. Per risolvere il problema del flush della cache è quello di aggiungere un tag identificativo del processo per verificare che il blocco al quale si

sta accedendo appartenga al processo in esecuzione.

Un'idea per incrementare la velocità di hit è quella di inserire una pipeline per effettuare il check del tag e aggiornare i dati in cache in due stage diversi così facendo si aumenta il throughput del sistema tuttavia aumenta anche la penalità in caso di predizione errata di un salto.

L'ultima tecnica si pone l'obiettivo di ridurre il tempo di write e si basa sull'idea che la maggior parte degli update avviene su una porzione molto piccola del blocco di cache, solitamente una *parola* per tale motivo è utile effettuare gli update solo di sotto-blocchi di dimensione di una parola.

Per un ricapitolazione delle diverse tecniche rimandiamo alla Figura 91 nella quale si mostrano i vantaggi e gli svantaggi nonché la complessità delle diverse soluzioni

	<i>Technique</i>	<i>MR</i>	<i>MP</i>	<i>HT</i>	<i>Complexity</i>
<b>miss rate</b>	Larger Block Size	+	-		0
	Higher Associativity	+		-	1
	Victim Caches	+			2
	Pseudo-Associative Caches	+			2
	HW Prefetching of Instr/Data	+			2
	Compiler Controlled Prefetching	+			3
	Compiler Reduce Misses	+			0
<b>miss penalty</b>	Priority to Read Misses		+		1
	Subblock Placement		+	+	1
	Early Restart & Critical Word 1st		+		2
	Non-Blocking Caches		+		3
	Second Level Caches		+		2
<b>hit time</b>	Small & Simple Caches	-		+	0
	Avoiding Address Translation			+	2
	Pipelining Writes			+	1

Figura 91: Principali tecniche di incremento delle performance

### 7.3 Virtual Memory

La memoria virtuale è un meccanismo che serve a tradurre gli indirizzi virtuali in indirizzi fisici, ovvero ad effettuare il *memory mapping*. Praticamente la memoria virtuale è un meccanismo che tratta la memoria come una cache per il disco in questo caso però i blocchi vengono chiamati *pagine*.

Il concetto di memoria virtuale è stato introdotto innanzitutto per separare lo spazio degli indirizzi del processore dalla dimensione reale della memoria fisica come mostrato in Figura 93. Il meccanismo di traduzione da indirizzo virtuale ad indirizzo reale è basato su una *Page Table* che viene associata ad ogni processo, ogni processo ha una sua page table e le page table di tutti i processi risiedono nella memoria fisica. Ogni page table mappa il *virtual page number (VPNs)* con il rispettivo *physical page number (PPNs)*, il VPN viene utilizzato come indice della *page table*. Un esempio di page table è mostrato in Figura 94 in tale esempio vediamo anche come sia presente per ogni record della tabella un bit di validità per distinguere quale di queste pagine sono presenti in memoria (bit di validità a 1) oppure immagazzinate sul disco (bit di validità a 0).

Technique	Hit time	Band-width	Miss penalty	Miss rate	Power consumption	Hardware cost/complexity	Comment
Small and simple caches	+		–	+		0	Trivial; widely used
Way-predicting caches	+			+		1	Used in Pentium 4
Pipelined cache access	–	+				1	Widely used
Nonblocking caches	+	+				3	Widely used
Banked caches	+			+		1	Used in L2 of both i7 and Cortex-A8
Critical word first and early restart			+			2	Widely used
Merging write buffer			+			1	Widely used with write through
Compiler techniques to reduce cache misses				+		0	Software is a challenge, but many compilers handle common linear algebra calculations
Hardware prefetching of instructions and data		+	+	–	2 instr., 3 data		Most provide prefetch instructions; modern high-end processors also automatically prefetch in hardware.
Compiler-controlled prefetching		+	+			3	Needs nonblocking cache; possible instruction overhead; in many CPUs

Figura 92: Principali tecniche di incremento delle performance

### 7.3.1 Memory management unit

Per incrementare le performance l'idea è quella di inserire una cache per velocizzare la traduzione degli indirizzi. Questa particolare cache prende il nome *Translation Look-Aside Buffer* ed è una piccola cache di circa 128-256 entità di tipo completamente associativo. Cosa succede in caso di miss nella TLB, nel caso di soluzione hardware lo hardware della MMU la controlla la page table corrente e nel caso di entità valida la trasferisce alla TLB in caso contrario comunica al kernel il *page fault*. Nel caso di soluzione software il processore riceve il TLB fault il kernel ritrova il record nella corrispettiva page table, nel caso di entità valida viene restituita l'entità in caso contrario il kernel richiama internamente la procedura di page fault.

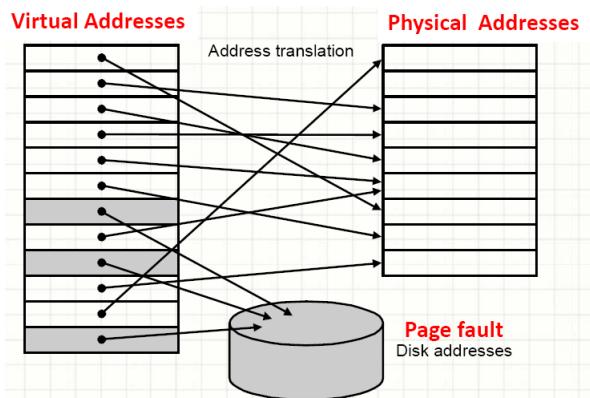


Figura 93: Traduzione da indirizzo virtuale ad uno fisico

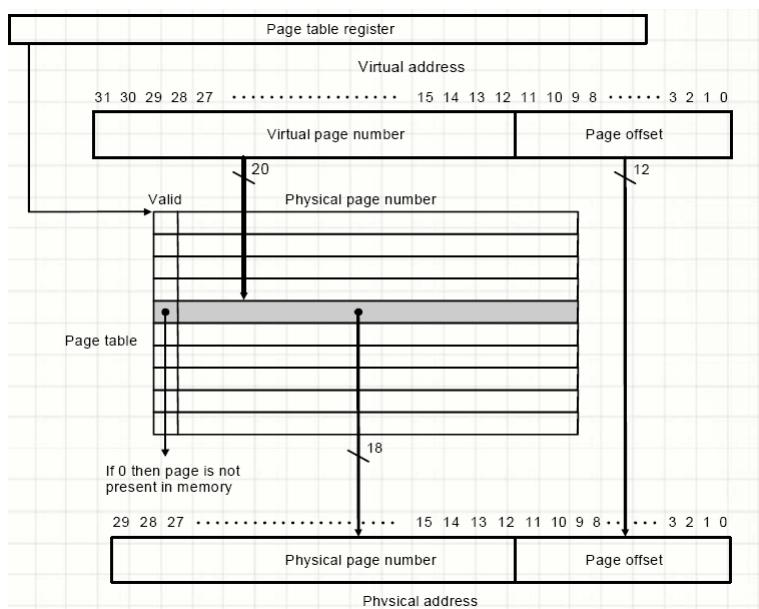


Figura 94: Esempio di page table

## 8 Mutua esclusione e sincronizzazione

Come prima cosa dobbiamo distinguere tra *mutua esclusione*, ovvero l'insieme di meccanismi che svolgono la funzione di proteggere dei dati condivisi da degli accessi concorrenti, e *sincronizzazione* che sono i meccanismi che permettono la coordinazione tra threads e processi che vengono eseguiti in parallelo. I meccanismi di sincronizzazione sono costruiti attorno a quelli di mutua esclusione.

Il problema principale quando parliamo di multithread o multi-processi è quello di ottenere un risultato deterministico, ovvero in esecuzioni differenti ottenere esattamente lo stesso risultato pur avendo inevitabilmente un accesso alle variabili in modo sempre diverso. Per ottenere questo risultato è necessario sfruttare delle tecniche di sincronizzazione che possono essere sia hardware che software.

Un aspetto fondamentale della sincronizzazione è che un processo o un thread deve essere in grado di eseguire determinate operazioni su delle determinate strutture dati senza la possibilità di essere interrotto; tale sequenza di operazioni è denominata *sezione critica*.

I meccanismi di sincronizzazione sono costruiti a livello utente tramite delle procedure software che sfruttano istruzioni di sincronizzazione fornite dal hardware. Queste istruzioni devono permettere la lettura e la modifica di aree di memoria in modo *atomico* ovvero senza essere interrotte. Esistono diverse soluzioni per implementare la mutua esclusione ad esempio l'uso della coppia *lock-unlock*, la sincronizzazione punto a punto tramite *flags* o la sincronizzazione globale tramite barriere.

Le principali caratteristiche che contraddistinguono un metodo di sincronizzazione da un altro sono principalmente tre:

**Metodo di acquisizione:** ovvero i modi in cui il processo tenta di acquisire i diritti per la sincronizzazione.

**Algoritmi di waiting:** i metodi utilizzati dal processo per attendere che la sincronizzazione sia disponibile.

**Metodi di rilascio:** i modi in cui un processo permette agli altri processi di procedere alla sincronizzazione.

Per quanto riguarda gli algoritmi si attesta esistono due alternative principali, la prima è quella di *busy-waiting* nel quale il processo attende tramite un ciclo che la variabile di sincronizzazione cambi il suo valore, nel secondo caso la soluzione è *bloccante* e in questo caso quando il processo entra in fase di attesa si blocca e lascia il processore ad un altro processo.

### 8.1 Mutua esclusione

La mutua esclusione è composta da due operandi *lock* e *unlock* i quali sono molto efficaci in caso di livelli di contesa delle risorse molto basso ma diventano molto inefficienti quando questi livelli di contesa si alzano. Esaminiamo due processi  $P_j$  e  $P_k$  i quali sono eseguiti su due nodi distinti e modificano entrambi la stessa struttura dati  $D$  le due modifiche non hanno precedenza ma devono avvenire in modo atomico. Lo schema di questo esempio è mostrato in Figura 95. Un'altra primitiva molto comune è quella dell'*atomic exchange* la quale si incarica di effettuare lo scambio di un valore nei registri con uno in memoria e viceversa. Tale metodo può essere utilizzato per implementare un semplice lock:

1. Prendiamo in considerazione un registro nel quale se il valore è 0 allora il lock è libero se invece il valore è 1 il lock non è disponibile

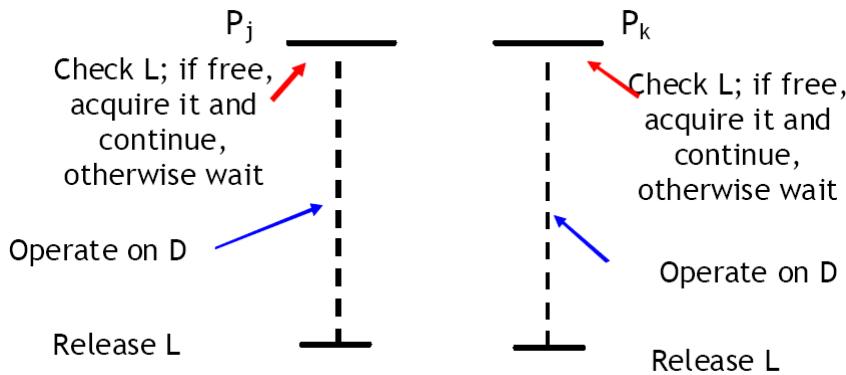


Figura 95: Esempio di due processi che sfruttano il lock

2. Un processo prova ad acquisire un lock su di una risorsa scambiando il valore in memoria con il valore 1 contenuto in un registro.
3. Il valore scambiato risulterà essere 1 se il lock è già stato scambiato, 0 altrimenti.

Due soluzioni più recenti per implementare dei meccanismi di lock sono *load linked* e la *store conditional*. Se un contenuto di una memoria puntato da una *load linked* cambia prima che la *store conditional* effettui la scrittura allora la scrittura fallisce; ovvero se tra una load e una store il processore effettua un *cambio di contesto* allora la store fallisce.

## 8.2 Dalle primitive di sincronizzazione ai metodi di sincronizzazione

Utilizzando le primitive di sincronizzazione è possibile costruire dei meccanismi che implementano la sincronizzazione a livello software. Il primo che analizziamo è lo *spin lock* un meccanismo che utilizza le operazioni atomiche. In questo meccanismo un processo tenta di acquisire un lock in modo continuo tramite un loop fino a quando non ha successo. Tale meccanismo è utilizzato soprattutto quando il programmatore si aspetta che i lock siano trattenuti per un breve periodo di tempo e siano richiesti con scarsa frequenza. Un esempio del funzionamento di questo meccanismo è mostrato in Figura 96. Un meccanismo di sincronizzazione ideale deve mantenere alcune caratteristiche per essere performante:

**Bassa latenza:** ovvero il tempo necessario per acquisire un lock da parte di un processo quando questo è libero e nessun altro processo sta tentando di acquisirlo deve essere basso.

**Basso traffico:** se molti processi stanno tentando di acquisire uno stesso lock contemporaneamente essi lo devono ottenere in sequenza generando un traffico sul bus di dimensioni ridotte.

### Scalabilità

### Basso costo di memoria

Lo *spin lock* tuttavia non rispetta molte di queste caratteristiche, infatti la latenza è bassa solo nel caso in cui uno stesso processo acquisisce il lock molte volte in successione, in caso di competizione per l'acquisizione del lock si genera molto traffico sul bus e difficilmente scalabile con l'aumentare dei processi.

Un altro meccanismo di sincronizzazione è quello del *barrier* utilizzato soprattutto nel caso

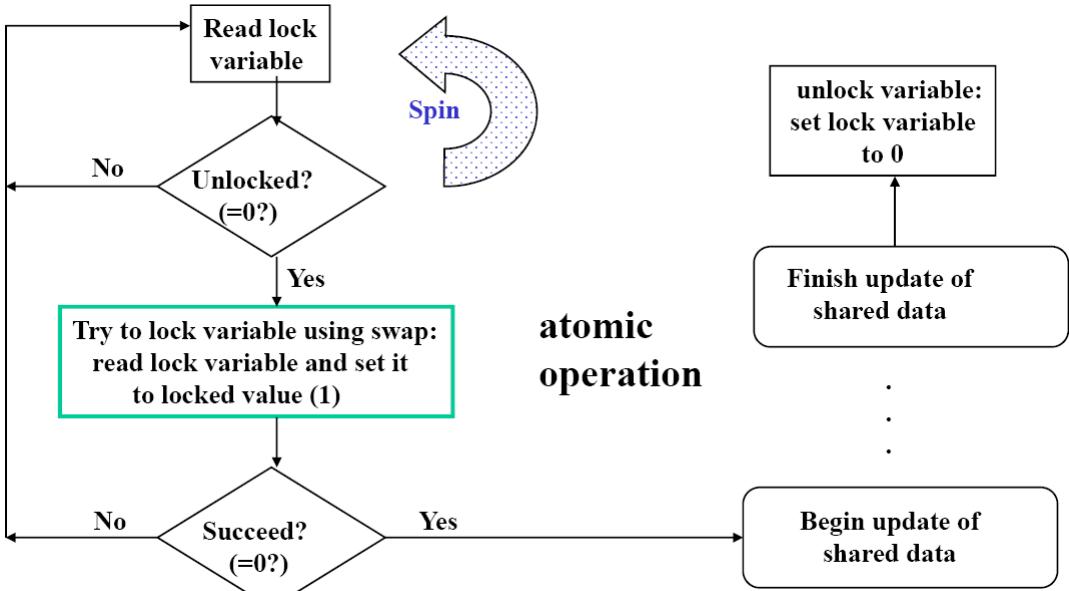


Figura 96: Diagramma di flusso dello spin lock

di loop paralleli. Tale meccanismo è di tipo globale e si applica ad un numero predefinito di processi  $p$ . Per implementare tale meccanismo si utilizzano i lock, i contatori condivisi e le flags. Per analizzare il funzionamento dell'algoritmo introduciamo l'esempio di un piccolo videogioco nel quale i frame vengono prima preparati e poi mostrati. Il codice di esempio che mostra il funzionamento è il seguente:

```

while (true) {
    frame.prepare();
    frame.display();
}

```

è possibile ottimizzare il processo dividendo il processo in due thread il primo che disegna il frame e il secondo che prepara il successivo come mostrato in Figura 97 Tuttavia possono

```

while (true) {
    if (phase) {
        frame[0].display();
    } else {
        frame[1].display();
    }
    phase = !phase;
}
while (true) {
    if (phase) {
        frame[1].prepare();
    } else {
        frame[0].prepare();
    }
    phase = !phase;
}

```

Figura 97: Esempio di parallelizzazione del codice precedente

sorgere alcuni problemi come un processo che si sveglia in ritardo e non prepara il frame in tempo come si mostra in Figura 98 per questo motivo è stato introdotto il meccanismo del barrier. Il meccanismo del barrier centralizzato mantiene il numero dei processi che arrivano alla barriera, per ogni processo che arriva un contatore viene incrementato, tale incremento

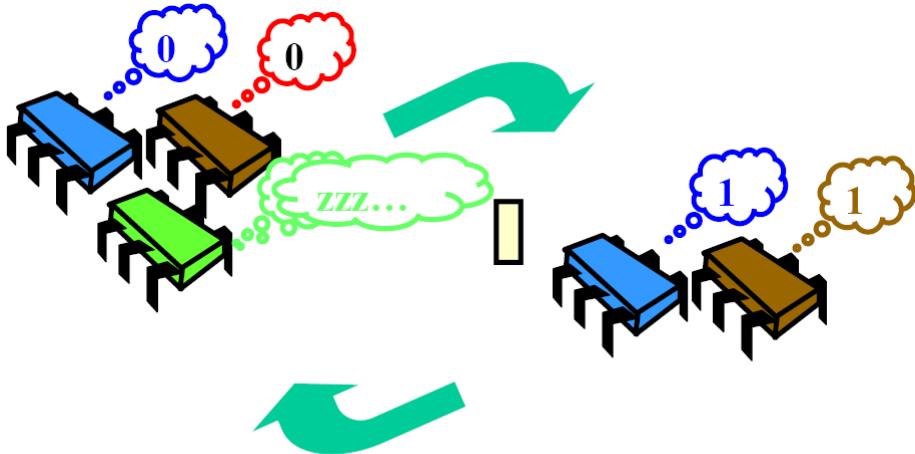


Figura 98:

deve essere atomico. Dopo ogni incremento il processo controlla se il numero dei processi  $p$  coincide con quello del contatore, in caso negativo attende, in caso affermativo comunica agli altri processi tramite dei flag che la loro attesa è finita. Un esempio di implementazione di un sistema di barrier è mostrato in Figura 99 Il problema di questo sistema è che risulta impossibile

```
public class Barrier {
    AtomicInteger count;
    int size;
    public Barrier(int n){
        count = AtomicInteger(n);
        size = n;
    }
    public void await() {
        if (count.getAndDecrement()==1) {
            count.set(size);
        } else {
            while (count.get() != 0);
        }
    }
}
```

Ristituta Silvana, Sotirios Xydis - Politecnico di Milano - 44 -

Figura 99: Implementazione di un sistema di barrier

riutilizzare il barrier, per risolvere tale inconveniente si è pensato ad un sistema denominato *sense-reversing barrier*. Ogni oggetto barrier contiene, diversamente dal caso precedente, un campo booleano che indica il senso dell'esecuzione corrente. Ogni thread mantiene traccia del senso corrente dell'esecuzione quando un thread raggiunge il barrier controlla se è l'ultimo thread e se lo è oltre a resettare il contatore inverte il senso dell'esecuzione, in caso contrario attende che il senso dell'esecuzione cambi; un esempio di questo meccanismo è mostrato in Figura 100

```
public class Barrier {  
    AtomicInteger count;  
    int size;  
    boolean sense = false;  
    ThreadLocal<boolean>...  
  
    public void await {  
        boolean mySense = threadSense.get();  
        if (count.getAndDecrement() == 1) {  
            count.set(size); sense = !mySense  
        } else {  
            while (sense != mySense) {}  
        }  
        threadSense.set(!mySense)}  
    }  
}
```

Figura 100: Implementazione di un sistema di barrier con senso di esecuzione

## 9 Introduzione ai multiprocessori

Negli ultimi quindici anni siamo entrati in una nuova era delle architetture dei calcolatori, i multiprocessori ora sono i più diffusi sia in ambito embedded che in quello general-purpose, in quanto forniscono elevate prestazioni, scalabilità ed affidabilità. Si parla di multiprocessori quando si hanno dei computer con più processori strettamente accoppiati i quali sono coordinati e controllati dal sistema operativo e condividono la memoria e lo spazio degli indirizzi. Si parla, invece, di multicores quando più core risiedono sullo stesso chip.

I multiprocessori comportano una serie di problematiche che spaziano dal metodo di collegamento a problemi di coordinamento; il primo problema che ci poniamo è come connettere i diversi processori. Esistono due soluzioni principali, la prima prevede l'utilizzo di un singolo bus come quello mostrato in Figura 101 dove il mezzo di connessione (il *bus*) è situato tra i processori e la memoria ed esso viene utilizzato ogni qualvolta sia necessario un accesso in memoria; questo meccanismo tuttavia comporta alcune problematiche prima tra tutte l'impossibilità di connettere una serie infinita di processori ma solamente un numero limitato causando problemi di *saturazione*. Una seconda tipologia, mostrata in Figura 102 di collegamento tra diversi pro-

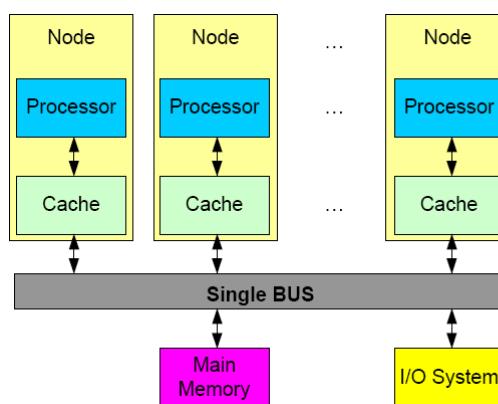


Figura 101: Architettura multiprocessore a singolo bus

cessori è quella che utilizza una *rete* nella quale ogni processore ha a disposizione una sua parte di memoria e la rete connette solamente i nodi i quali si scambiano comunicazione interprocesso. Per quanto riguarda la tipologia di connessione esistono diverse topologie che la rete può assu-

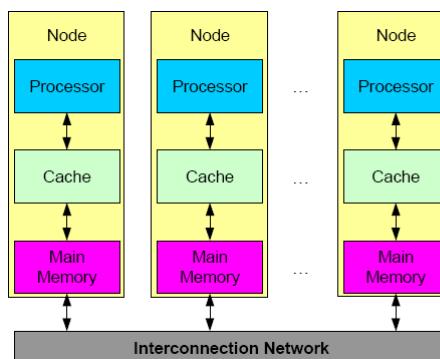


Figura 102: Architettura multiprocessore con connessione tramite rete

mere, la forma e le interconnessioni dipendono da un trade-off tra costi e prestazioni, una rete

completamente connessa è molto performante tuttavia è anche molto costosa. Alcuni esempi di topologia di reti sono:

**Single bus**

**Anello**

**Maglia**

**N-cubo**

**Completamente connessa**

Per descrivere una rete si possono utilizzare dei grafici come quello di Figura 103 che rappresenta una rete ad anello nei quali i *nodi*, rappresentati dai quadrati sono i nodi che comprendono processore e memoria, i cerchi rappresentano invece gli *switch* che connettono i nodi alla rete. Gli *archi* invece rappresentano le linee di comunicazione tra gli switch e sono sempre di tipo bidirezionale. Il costo di una rete è definito dal numero di *switch*, dal numero di archi che connettono i diversi switch, dalla lunghezza dei collegamenti.

Per analizzare le performance di una determinata tipologia di rete introduciamo due metriche

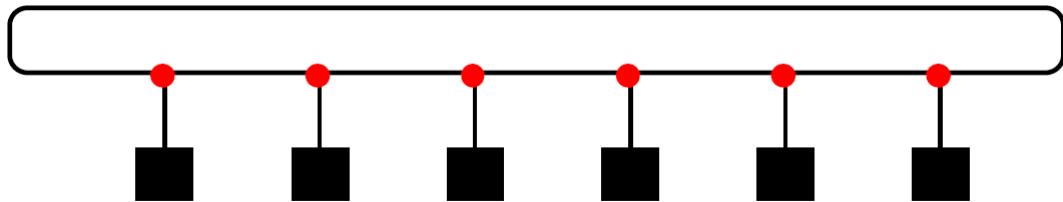


Figura 103: Esempio di rete ad anello

di riferimento, una che analizza il caso migliore ed una che analizza il caso peggiore. Il caso migliore è dato dal *total network bandwidth* ovvero la capacità massima di trasferimento che è data dalla banda di ogni link moltiplicata per il numero di link. Definendo:

$$P = \# \text{ di nodi}$$

$$b = \text{banda del singolo link}$$

Per una rete a singolo bus abbiamo che nel caso migliore la banda della rete è data dalla massima banda del link ( $1 \times b$ ); per una rete ad anello invece la massima banda è data dal numero di nodi moltiplicato per la banda di un link ( $P \times b$ ). Nel caso peggiore invece si considera la rete divisa in due parti ognuna delle quali con la metà dei nodi e si sommano le capacità dei link che attraversano l'immaginaria linea di separazione dei nodi; si ha così che per la rete ad anello la banda diventa ( $2 \times b$ ) mentre per la rete a singolo bus rimane invariata ( $1 \times b$ ). Nel caso di reti non simmetriche dobbiamo considerare la divisione che da le peggiori performance per la rete. Analizziamo ora una rete di tipo *completamente connessa* nel quale ogni nodo è connesso ad ogni altro nodo da un link bidirezionale; tale struttura ha un costo molto alto ma le sue performance sono ottime. La banda totale è data da  $\{(P \times (P - 1))/2\} \times b$  mentre il caso peggiore è dato da  $(P/2)^2 \times b$ .

Nel caso di topologia a maglia come quella di Figura 104 abbiamo che dati  $P$  nodi possiamo individuare un numero  $N = \sqrt{P}$  che è la dimensione della maglia il numero di canali di comunicazione è dato da  $N \times (N - 1)$  canali orizzontali e  $N \times (N - 1)$  canali verticali per ogni

switch interno abbiamo che il numero di link è pari a 5 (4 provenienti dagli altri switch e uno proveniente dal nodo) mentre per gli switch sul bordo abbiamo che il numero di link è pari a 3. La banda nel caso migliore è data da  $\{2 * N * (N - 1)\} \times b$  mentre nel caso peggiore è data da  $N \times B$ .

Nel caso in cui la topologia sia di tipo iper-cubica il numero di nodi determina la dimensione

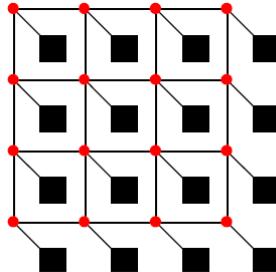


Figura 104: Esempio di topologia a maglia con  $N=4$

del cubo tramite la formula  $P = 2^N$  in Figura 105 vediamo la configurazione di un ipercubo con  $P = 16$  ogni nodo a un numero di vicini pari ad  $N$  e ogni switch ha esattamente  $N + 1$  link. La banda totale di questa configurazione è data dalla formula  $\{(N \times P)/2\} \times b$  mentre la banda di bisezione è data da  $2^{n-1} \times b$ .

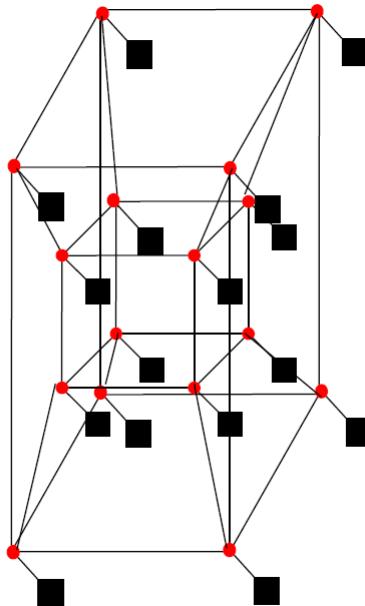


Figura 105: Esempio di topologia ad ipercubo.

## 9.1 Gestione della memoria

In questo paragrafo ci chiediamo come i diversi processi possono condividere i dati e come possono gestire la memoria. Partiamo con l'analizzare come i diversi processori condividono lo memoria, tale condivisione può avvenire in due modi tramite il *modello di memoria a spazio degli indirizzi*, il primo modo è avere un singolo spazio degli indirizzi condiviso mentre il secondo

è avere più spazi degli indirizzi logicamente separati come mostrato in Figura 106. Nel caso di

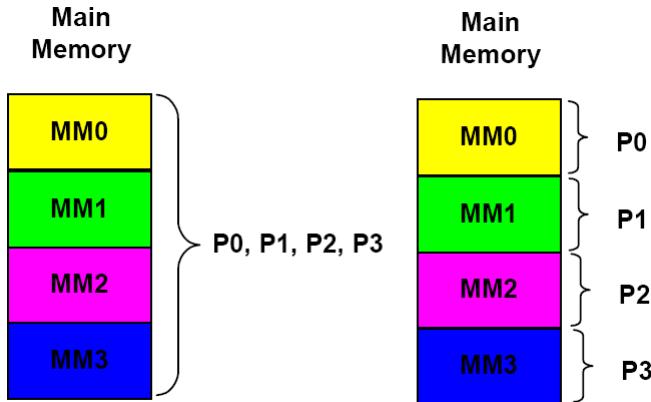


Figura 106: Esempio di singolo spazio degli indirizzi condiviso e spazio degli indirizzi multiplo.

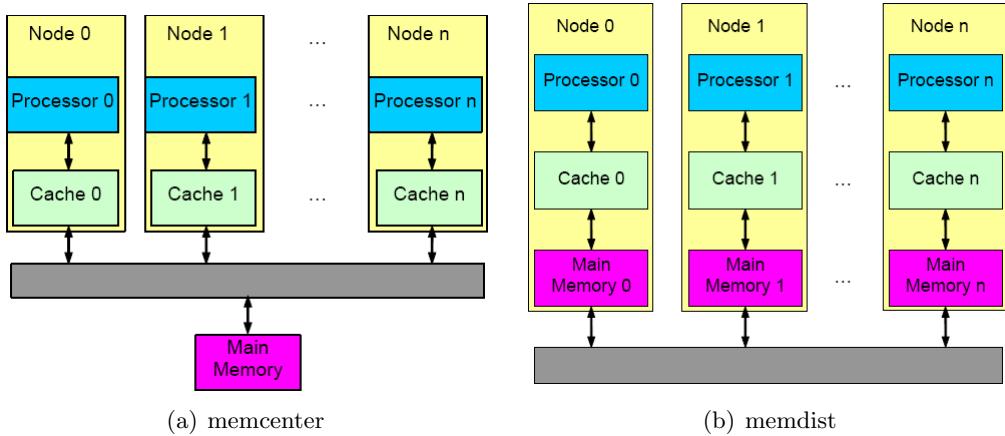
singolo spazio degli indirizzi condiviso un qualsiasi processo può fare riferimento a qualsiasi area di memoria, gli indirizzi sono condivisi tra i processi ed un indirizzo fisico su due processi diversi fa riferimento alla stessa area di memoria. La comunicazione tra processi e thread avviene tramite lo spazio di memoria condiviso, la gestione delle comunicazione è implicita e avviene tramite l'utilizzo delle operazioni di load e di store. Una memoria condivisa non implica necessariamente che esista un'unica memoria centrale. Questo tipo di modello tuttavia comporta problemi di coerenza della cache.

Nel caso di spazio degli indirizzi multiplo i processi comunicano tra di loro tramite le primitive *send* e *receive*; lo spazio di memoria è logicamente diviso e quindi due processi non potranno mai fare riferimento alla stessa area di memoria. La gestione della comunicazione tra i processi è gestita in modo esplicito tramite l'utilizzo delle primitive *send* e *receive*, la memoria di un processo non può essere acceduta da un altro processo se esso non è supportato dal protocollo software evitando così problemi di cache coherence.

Un altro problema riguardante la gestione della memoria nei sistemi multiprocessore è quello di dove posizionare fisicamente la memoria. Esistono due possibilità, la prima, mostrata in Figura 107(a) prevede di posizionare la memoria all'esterno dei nodi in un punto accessibile da tutti i processori. Nel secondo caso, mostrato in Figura 107(b) prevede di suddividere la memoria su ogni singolo nodo.

Nel caso di memoria centrale si parla di **UMA** (*Uniform Memory Access*) il tempo di accesso alla memoria è uniforme per tutti i processori indipendentemente da quale sia il processo e da quale area di memoria esso richieda. Nel caso invece di memoria condivisa si parla di **NUMA** (*Non Uniform Memory Access*) in questo caso il tempo di accesso dipende dalla posizione di memoria che si richiede e da posizione del processore.

Nella maggior parte dei sistemi multiprocessore attuali troviamo una struttura con singolo spazio degli indirizzi e una memoria centrale che permette un accesso a memoria uniforme, tali processori sono anche chiamati *Symmetric Multiprocessor (SMPs)*. La maggior parte dei processori *multicore* attualmente esistenti sono di tipo SMPs con un ridotto numero di cores ( $\leq 8$ ) tipicamente esiste un livello di cache condiviso e uno o più livelli di cache privati per ogni core. Quando il numero di core invece aumenta ogni risorsa centralizzata del sistema diventa un collo di bottiglia. Per incrementare la banda di comunicazione si utilizzano bus multipli o reti di comunicazione dove le memorie condivise sono configurate come banchi di memoria multipli. Un esempio di questa architettura è mostrato in Figura 107 dove ogni processore condivide l'in-



terro spazio di memoria tuttavia il tempo di accesso dipende dalla posizione alla quale si vuole accedere (*NUMA*); l'accesso al banco di memoria connesso direttamente al processore è molto più veloce rispetto agli altri. Negli ultimi anni si sono affermati una nuova tipologia di compu-

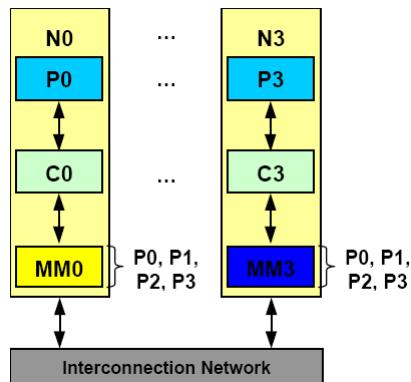


Figura 107: Architettura di un sistema multiprocessore a memoria distribuita

ter multiprocessore, i *cluster* composti da computer individuali con spazi di memoria privati e l'interconnessione avviene tramite il passaggio di messaggi.

Esistono alcuni limiti al parallelismo tra processi, il primo è quello dovuto alla quantità di parallelismo che si può estrapolare tramite il compilatore. Il secondo limite dipende molto dai costi di comunicazione infatti, un accesso a memoria tra cores sullo stesso chip è dell'ordine dei 35-50 cicli mentre si aggira tra i 100 e i 500 cicli per accessi tra cores su chip differenti. Diventa così molto importante meccanismi di cache, in quanto riducono il tempo medio di accesso, inoltre riducono i problemi di contesa in quanto permettono la copia dei dati, tuttavia questo porta a problemi di *coerenza* della cache.

Per sopperire ai problemi di coerenza della cache esistono diversi meccanismi i principali sono:

**Migrazione:** che consiste nello spostamento dei dati che permette un accesso molto rapido.

**Replicazione:** consiste nel creare molte copie dei dati riducendo così la contesa su di essi.

Per assicurare la coerenza si utilizzano dei *protocolli di coerenza*, possiamo distinguere due classi di protocolli di coerenza:

**Snooping protocol:** nei quali ogni core tiene traccia dello stato dei dati su ogni blocco, un controllore (*snoop*) sul bus controlla le richieste dirette ad altre cache.

**Directory-based protocol:** lo stato di ogni blocco di memoria è mantenuto in un punto denominato *directory* che può essere un unico punto nel caso di SMP o multiplo in caso di DSM.

Analizzeremo i protocolli di coerenza più in dettaglio nel paragrafo seguente.

L'ultima problematica che ci poniamo è quella di come sincronizzare e coordinare i diversi processori. Esistono diversi modelli di comunicazione alcuni già visti come la *memoria condivisa* tramite cui i processori comunicano utilizzando uno spazio degli indirizzi condiviso, tale modello è di facile implementazione e adatto a un ridotto numero di macchine in quanto facile da implementare con bassa latenza e permette di utilizzare i controlli di coerenza dell'hardware. Il secondo modello è quello del *passaggio di messaggi* in questo caso è adatto a processori con meccanismi di memoria privati richiede un hardware minore e si concentra sull'ottimizzazione dei costi sulle operazioni non locali. L'ultimo modello è il *data parallel* in questo caso si adatta ad operazioni che possono essere eseguite in parallelo su un numero molto grande di strutture dati come gli *array*. In questo caso un controllore distribuisce i dati sugli altri processori, i dati sono distribuiti su tutte le memorie. Il principio del **SIMD** (*Single Instruction Multiple Data*) ha portato allo sviluppo della programmazione di dati parallela nella quale tutti i processori eseguono uno stesso programma.

## 9.2 Il problema della coerenza della cache

Come abbiamo visto nel paragrafo precedente il fatto di avere della memoria condivisa pur avendo della cache dedicata ad ogni processore può portare ad avere problemi di coerenza in quanto un dato in una cache può essere replicato più volte, in quanto questo meccanismo permette di ridurre le contese dovute a letture su dati condivisi.

Un meccanismo per ridurre i problemi di coerenza della cache è quello di forzare l'accesso a memoria centrale per i dati condivisi, questo però riduce notevolmente le prestazioni. I problemi di coerenza si possono verificare solamente quando si svolgono delle *read* o delle *write*, tuttavia, l'accesso in lettura a più copie non crea problemi di coerenza, ma il processore deve avere accesso esclusivo quando accede in scrittura. Un processore deve avere la copia più recente di un dato prima di accedervi in lettura così ogni processore deve avere il nuovo valore quando si effettua una scrittura.

Per evitare i problemi di coerenza bisogna effettuare due operazioni, la prima è tenere traccia di tutte le copie cahce di un particolare dato condiviso e la seconda è che una scrittura su di un dato condiviso deve o *invalidare* oppure *aggiornare* tutte le altre copie condivise. La soluzione sono i *protocolli di coerenza*, esistono due classi di protocolli, gli *snooping protocols* e *directory-based protocol*. Nel caso di protocolli di *snooping* ogni processore tiene traccia dello stato delle condivisione di ogni blocco di memoria ad esso associato, un controllore sul bus denominato *snoop* verifica eventuali richieste indirizzate ad altre cache. Un esempio di questo meccanismo è rappresentato in Figura 108 Ad ogni richiesta sul bus lo *snoop* controlla se è in possesso di una copia del blocco richiesto e risponde di conseguenza, ogni cache che ha una copia del blocco condiviso ha anche una copia del suo stato così facendo non esiste un'entità centrale nel quale lo stato è mantenuto. Le richieste per dei dati condivisi vengono inviati a tutti i processori, le richieste vengono inviate in *broadcast* fino a quando non il processore non riesce a ricostruire le informazioni. Questo tipo di meccanismo è molto adatto per sistemi a memoria centrale condivisa e in particolare per multiprocessori di piccole dimensioni con **single snoopy bus**. Tuttavia questa tecnica presenta anche alcuni problemi, ogni qualvolta si cerca un blocco

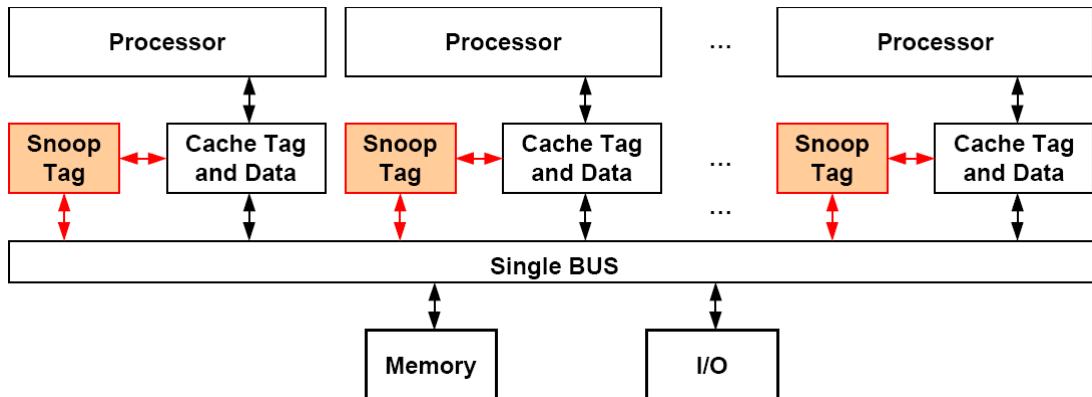


Figura 108: Struttura di uno snooping protocol

condiviso è necessario controllare i *tag* della cache e questo comporta delle interferenze con le normali operazioni del processore in quanto la cache risulta non disponibile. Per evitare questo problema si possono duplicare i tag in modo da poter effettuare le attività di *snooping* per fare ciò è sufficiente aggiungere alla parte dei tag una porta in lettura.

Possiamo distinguere due sotto categorie di snooping protocol in base a quello che succede quando si effettua un'operazione di write; le strategie sono due:

- Write-Invalidate Protocol
- Write-Update Protocol

Nel primo caso quando un processore scrive in un blocco condiviso esso invia un segnale di *invalidazione* sul bus che fa sì che tutte le copie nelle altre cache siano invalidate prima di aggiornare la propria copia locale, il processore è libero di aggiornare il suo dato fino a quando un'altra cache non richiederà il dato. All'arrivo di un segnale di invalidazione ogni cache controlla se ha una copia del dato e in caso affermativo lo marca come non valido. Questo schema permette di avere delle operazioni di lettura simultanee ma solamente una scrittura per volta. Il bus viene utilizzato solamente durante la prima operazione di write effettuata da un processore tutte le altre risultano trasparenti. Questo tipo di protocollo fornisce dei vantaggi paragonabili a quelli del protocollo *write-back* in termini di riduzione di utilizzo del bus.

Nel caso di *write-update*, invece, si ha che il processore distribuisce il nuovo valore sul bus e tutte le altre cache verificano la presenza del dato nella loro memoria e in caso affermativo aggiornano con il nuovo valore, in questo modo tutte le copie sono sempre aggiornate ma questo meccanismo richiede un continuo broadcast dei nuovi valori sul bus. Questo protocollo è simile al *write-through* in quanto tutte le scritture effettuano una comunicazione sul bus. Questo protocollo ha il vantaggio di tenere tutte le copie aggiornate e così facendo permette di ridurre la latenza.

La maggior parte dei prodotti commerciali attualmente utilizza una cache di tipo *write-back* per ridurre il traffico sul bus e uno snooping protocol di tipo *write-invalidate* l'unico problema che rimane da risolvere è la serializzazione delle scritture in quanto il bus è un punto di *arbitrarietà*. Per ogni blocco di memoria possiamo definire tre stati:

**Shared:** quando tutte le copie cache sono aggiornate al valore di memoria

**Modified:** quando esiste un'unica cache in cui è presente il valore di memoria

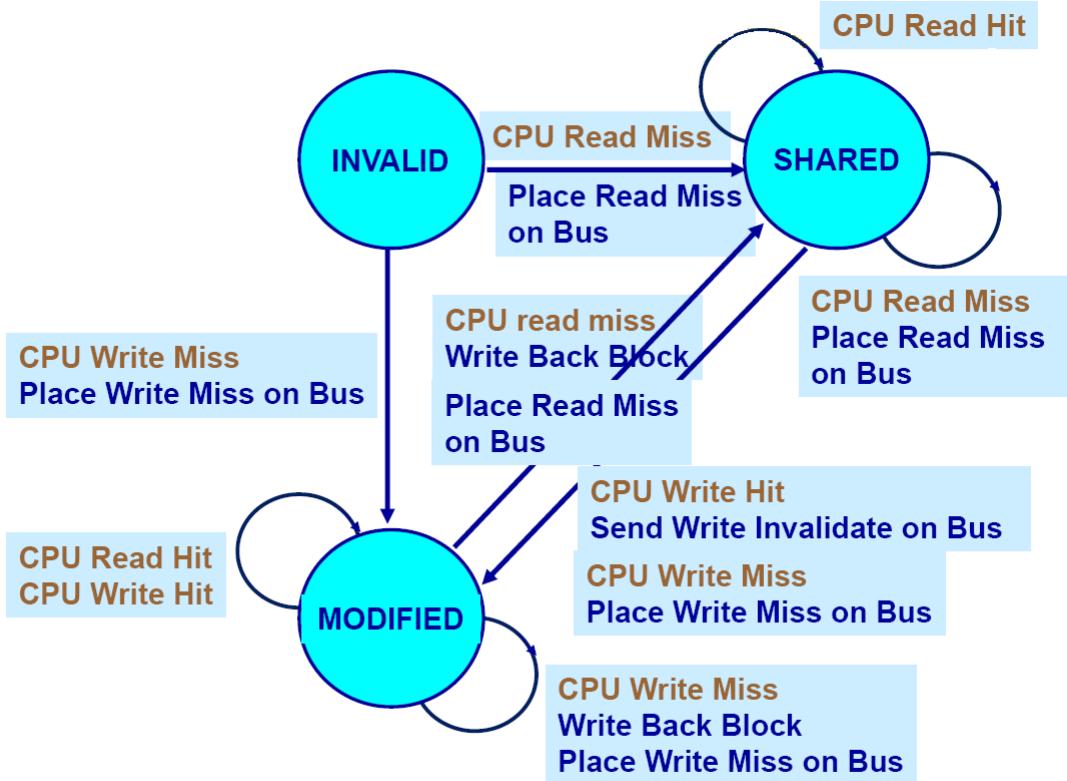


Figura 109: Macchina a stati finiti che rappresenta un blocco di cache

In nessuna cache.

Per i blocchi in cache possiamo distinguere tre stati:

**Shared:** il blocco è aggiornato e può essere letto

**Modified:** il blocco è l'unico aggiornato e può essere scritto

**Invalid:** il blocco non contiene dati validi

In Figura 109 si possono vedere i diversi stati di un blocco cache e le transizioni da uno stato all'altro. Una variazione di questo meccanismo prevede per i blocchi di cache un quarto stato, **Exclusive**: in questo caso il blocco non è ancora stato scritto ma è l'unica copia presente nelle cache; tale stato fa sì che una successiva scrittura non comporti un invio del segnale di invalidazione sul bus.

Analizziamo ora i protocolli *directory based*: lo stato di ogni blocco di memoria fisico è mantenuto in un unico spazio chiamato *directory*; ogni riga nella directory corrisponde a un blocco di memoria; per quei sistemi con memoria condivisa distribuita anche la directory è distribuita come nel caso di Figura 110 per evitare colli di bottiglia. Questo protocollo risulta più scalare rispetto a quello di snooping. Le directory mantengono informazioni riguardanti lo *stato* di ogni blocco che può essere:

**Uncached:** ovvero non esistono copie cache valide di quel blocco.

**Shared:** uno o più processori mantengono una copia cache aggiornata del dato.

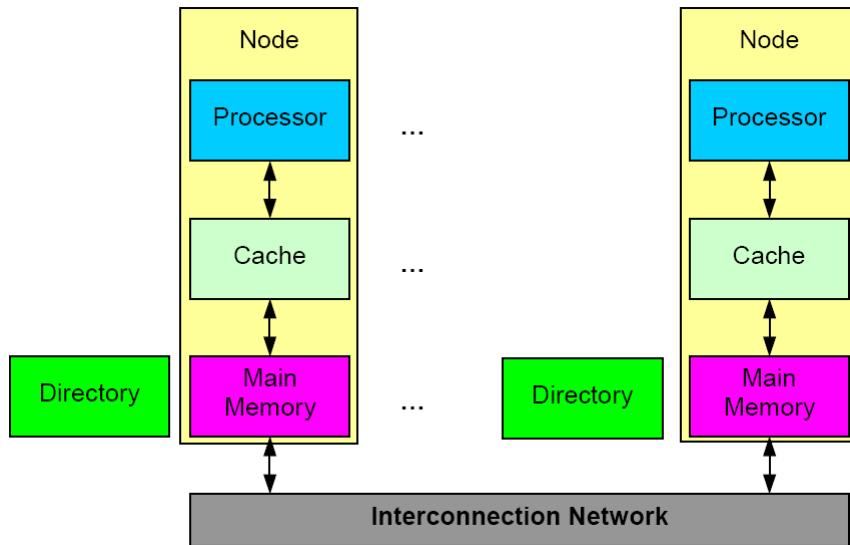


Figura 110: Esempio di directory base protocol distribuito

**Modified:** solamente un processore (*il possessore*) ha una copia aggiornata del dato, anche la memoria è obsoleta.

Oltre allo stato la directory mantiene un indice del processore o dei processori che hanno una copia del dato imponendo a 1 il campo corrispondente di un vettore di bit come mostrato nell'esempio di Figura 111 nel caso in cui il blocco risulti modificato un unico bit del vettore è a 1. Il protocollo directory based sfrutta i *messaggi* per comunicare attraverso i nodi (comunicazione

Block	Coherence State	Sharer / Owner Bits
B0	Uncached	- - - -
B1	Shared	1 0 1 0
B2	Shared	0 0 1 0
B3	Modified	0 1 0 0

Figura 111: Esempio di directory per quattro blocchi di memoria utilizzati da quattro processori

punto-a-punto) ed è necessaria un esplicita risposta dal destinatario questo permette di non avere un unico punto di arbitrarietà.

Come per lo snooping protocol anche nel directory based i blocchi in cache possono trovarsi nei tre stati *shared*, *modified* e *invalid*. Solitamente in un protocollo directory based sono coinvolti tre processori:

**Local node:** il nodo da cui la richiesta ha origine.

**Home node:** il nodo dove è immagazzinato il dato fisicamente in memoria.

**Remote node:** nodo nel quale è presente una copia cache del blocco.

Il *local node* (**L**) e *home node* (**H**) possono essere lo stesso nodo, in questo caso la comunicazione può essere di tipo *intra-nodo*. Anche il *remote node* (**R**) e **H** possono essere lo stesso nodo.

Quello che non può essere è che **R** ed **L** siano lo stesso nodo.

In Figura 112 è mostrato lo stato dei blocchi di memoria con relative richieste provenienti dalle cache nel caso di protocollo directory based.

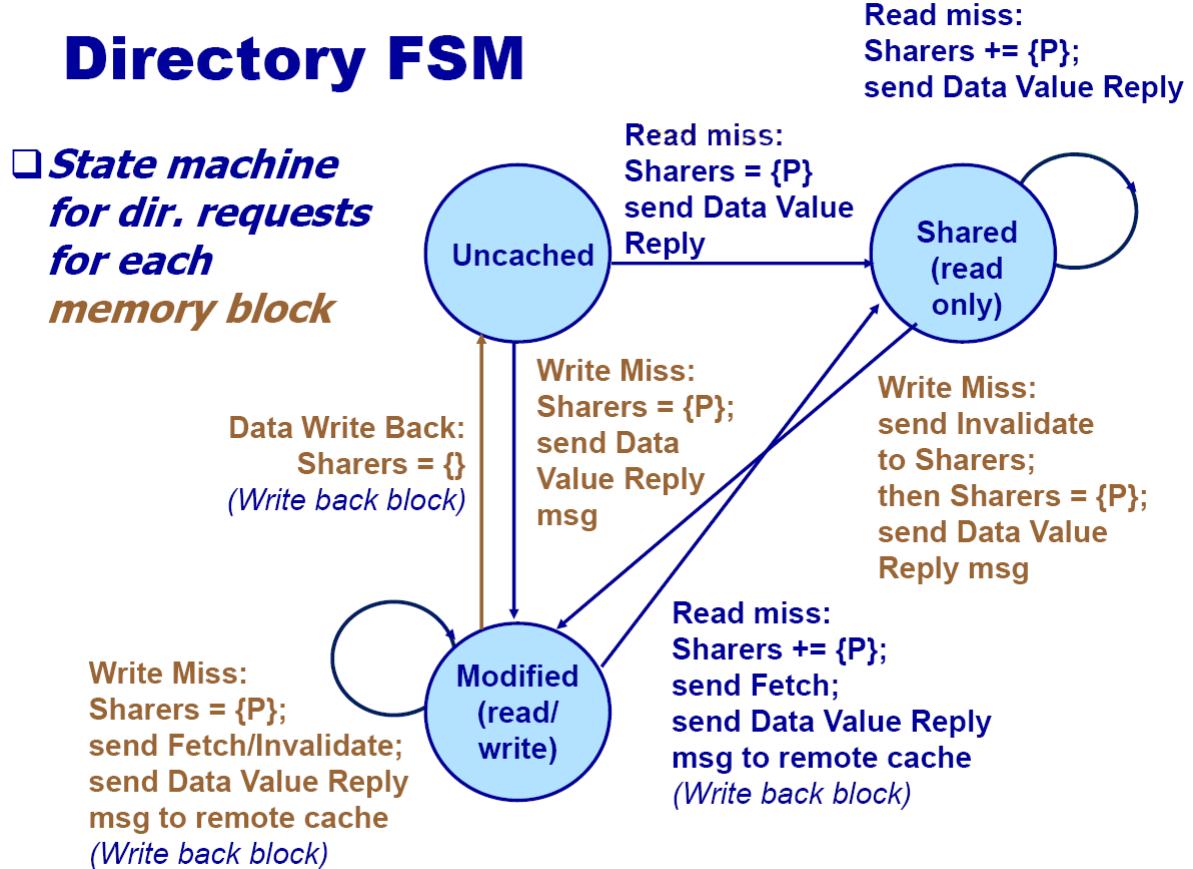


Figura 112: Macchina a stati finiti per i blocchi di memoria nel caso di protocollo directory based

## 10 Analisi delle performance

Nei capitoli precedenti abbiamo visto alcune tecniche per migliorare le performance dei processori, ora vogliamo valutare quantitativamente questo miglioramento analizzando caso per caso le diverse tecniche. Cosa significa tuttavia misurare un miglioramento? esistono diversi aspetti da considerare, ad esempio se consideriamo il costo dobbiamo tener presente il rapporto costo/prestazioni, se consideriamo il design di un sistema dobbiamo misurare il miglioramento delle prestazioni rispetto all'aumento di costo sia in termini monetari che di superficie di silicio. Prima di tutto però dobbiamo valutare che cosa significa *performance*, possiamo considerare due aspetti:

- Il **tempo di esecuzione** ovvero il tempo necessario per compiere un lavoro.
- Il **throughput** ovvero il numero di lavori completati nell'unità di tempo.

Infine dobbiamo avere un metodo di paragone per confrontare due soluzioni diverse, tale metodo è dato da:

$$X \text{ is } n\% \text{ piu veloce di } Y = \frac{\text{tempo esecuzione (y)}}{\text{tempo esecuzione (x)}} = 1 + \frac{n}{100}$$

### 10.1 Legge di Amdahl

Secondo la legge di Amdahl l'idea di base è quella velocizzare le operazioni più comuni, più precisamente

*Il miglioramento delle performance dato dall'uso di qualche modulo di esecuzione più veloce è limitato alla frazione di tempo nella quale questo modulo viene utilizzato*

Considerando la  $Frazione_E$  come la frazione di tempo di computazione nella macchina originale che può essere velocizzato e lo  $SpeedUp_E$  come il miglioramento di prestazioni dovuto al nuovo modulo.

$$SpeedUp(E) = \frac{ExTime \text{ w/o } E}{ExTime \text{ w/E}} = \frac{Performance \text{ w/E}}{Performance \text{ w/o } E}$$

Supponendo che il modulo  $E$  migliori una frazione  $F$  del lavoro di un fattore  $S$  mentre il resto del processo non ne viene influenzato allora abbiamo

$$ExTime(\text{with } E) = ((1 - F) + F/S) \times ExTime(\text{without } E)$$

$$SpeedUp(\text{with } E) = \frac{1}{(1 - F) + F/S}$$

### 10.2 Analisi delle performance in un processore pipelined

Nel caso di processore con pipeline abbiamo che viene incrementato il *throughput* delle istruzioni ma non viene ridotto il tempo di esecuzione della singola istruzione, anzi molte volte è leggermente incrementato a causa del bilanciamento degli stage della pipeline

$$IC = Instruction Count$$

$$\# Clock Cycle = IC + \# Stalli + 4$$

$$CPI = Clock \text{ per Istruzione} = \# Clock Cycle / IC = (IC + \# Stalls + 4) / IC$$

$$MIPS = \frac{f_{clock}}{(CPI * 10^6)}$$

Nel caso di un ciclo composto da **m** istruzioni il quale compie **n** iterazioni e richiede **k** stalli per iterazione abbiamo che per ogni iterazione:

$$\begin{aligned} IC_{per\_iter} &= m \\ \# Clock Cycle_{per\_iter} &= IC_{per\_iter} + \# Stalli_{per\_iter} + 4 \\ CPI_{per\_iter} &= (IC_{per\_iter} + \# Stalli_{per\_iter} + 4)/IC_{per\_iter} = (m + k + 4)/m \\ MIPS_{per\_iter} &= \frac{f_{clock}}{(CPI_{per\_iter} * 10^6)} \end{aligned}$$

Analizzando asintoticamente le formule precedenti abbiamo che:

$$\begin{aligned} IC_{AS} &= m * n \\ \# Clock Cycle_{AS} &= IC_{AS} + \# Stalli_{AS} + 4 \\ CPI_{AS} &= \lim_{n \rightarrow \infty} (IC_{AS} + \# Stalli_{AS} + 4)/IC_{AS} \\ &= \lim_{n \rightarrow \infty} (m * n + k * n + 4)/m * n \\ &= (m + k)/m \\ MIPS_{AS} &= \frac{f_{clock}}{(CPI_{AS} * 10^6)} \end{aligned}$$

Il *CPI* ideale in un processore con pipeline dovrebbe essere uguale a 1 ma gli stalli causano un degradamento delle performance otteniamo così:

$$\begin{aligned} Ave. CPI &= CPI_{ideal} + Pipe Stall Cycle per Instruction \\ &= 1 + Pipe Stall Cycle per Instruction \end{aligned}$$

Il numero di stalli per istruzione dipende da *hazard strutturali* + *hazard sui dati* + *hazard di controllo* + *stalli per accedere alla memoria*. Possiamo misurare il miglioramento dato dalla pipeline come:

$$\begin{aligned} SpeedUp_{pipeline} &= \frac{Ave. Exec. Time Unpipelined}{Ave. Exec. Time Pipelined} \\ &= \frac{Ave. CPI Unp.}{Ave. CPI Pipe} \times \frac{Clock Cycle Unp.}{Clock Cycle Pipe} \end{aligned}$$

Ignorando l'overhead sul periodo di clock introdotto dalla pipeline e assumendo che gli stage siano perfettamente bilanciati allora il periodo di clock dei due processori può ritenersi uguale così l'ultima formula diventa:

$$SpeedUp_{pipeline} = \frac{Ave. CPI Unp.}{1 + Pipe Stall Cycle per Instruction}$$

Supponendo nel caso più semplice che le istruzioni richiedano tutte lo stesso numero di cicli per essere eseguite il quale è uguale al numero di stage della pipeline anche chiamato *pipeline depth* abbiamo che:

$$SpeedUp_{pipeline} = \frac{Pipeline Depth}{1 + Pipe Stall Cycle per Instruction}$$

Vediamo come nel caso non vi siano stalli si può aumentare il throughput del sistema semplicemente aumentando la profondità della pipeline.

Nel caso di salti abbiamo che lo speedup del sistema è dato da:

$$SpeedUp_{pipeline} = \frac{Pipeline Depth}{1 + Branch Freq. \times Branch Penalty}$$

### 10.3 Analisi delle performance nelle gerarchie di memoria

Prima di iniziare l'analisi delle performance reintroduciamo alcune definizioni già viste:

**Hit:** si ha quando un dato viene trovato in un blocco di memoria del livello più alto.

**Hit rate:** numero di accessi a memoria che trovano il dato rispetto al numero totale di accessi

$$\text{Hit Rate} = \frac{\#/\text{hits}}{\# \text{ memory accesses}}$$

**Hit time:** tempo per accedere ad un dato che si trova nel livello più alto della gerarchia compreso il tempo per decidere se esso si trova in tale livello.

**Miss:** si ha quando il dato deve essere recuperato da un livello più basso.

**Miss Rate:** numero di accessi a memoria che non trovano il dato nel livello più alto della gerarchia rispetto al numero di accessi totali.

$$\text{Miss Rate} = \frac{\#/\text{misses}}{\# \text{ memory accesses}}$$

$$\text{Miss Rate} + \text{Hit Rate} = 1$$

**Miss Penalty:** tempo necessario per accedere al livello più basso e rimpiazzare il blocco nel livello più alto

**Miss Time:**

$$\text{Miss Time} = \text{Hit Time} + \text{Miss Penalty}$$

$$\text{Hit Time} \ll \text{Miss Penalty}$$

Definiamo infine il tempo medio di accesso come:

$$\text{AMAT} = \text{Hit Rate} * \text{Hit Time} + \text{Miss Rate} * \text{Miss Time}$$

dalle formule precedenti possiamo semplificare la formula del tempo medio di accesso come :

$$\text{AMAT} = \text{Hit Time} + \text{Miss Rate} * \text{Miss Penalty}$$

Volendo valutare quale impatto hanno le gerarchie di memoria sul tempo di esecuzione di un programma possiamo definire il tempo di esecuzione come:

$$\text{CPU}_{\text{time}} = (\text{CPU exec cycles} + \text{Memory Stall Cycles}) \times T_{\text{CLK}}$$

dove:

$T_{\text{CLK}}$ : periodo di tempo del clock

$\text{CPU}_{\text{exec-cycles}}$ :  $IC \times CPI_{\text{exec}}$

$IC$ : Instruction count

$\text{Memory stall cycle}$ :  $IC \times \text{Miss per instr} \times \text{Miss penalty}$

A questo punto possiamo riscrivere la precedente equazione come:

$$CPU_{time} = IC \times (CPI_{exec} + Miss \text{ per } instr \times Miss \text{ penalty}) \times T_{CLK}$$

dove

$$Miss \text{ per } instr = Memory \text{ Access Per Instruction} \times Miss \text{ Rate}$$

$$CPU_{time} = IC \times (CPI_{exec} + MAPI \times Miss \text{ Rate} \times Miss \text{ penalty}) \times T_{CLK}$$

Nel caso ideale in cui abbiamo tutti hit la formula si riduce:

$$CPU_{time} = IC \times CPI_{exec} \times T_{CLK}$$

nel caso in cui invece prendiamo in considerazione un sistema senza cache:

$$CPU_{time} = IC \times (CPI_{exec} + MAPI \times Miss \text{ penalty}) \times T_{CLK}$$

Tenendo in considerazione tutti i tipi di stalli durante un'esecuzione abbiamo

$$CPU_{time} = IC \times (CPI_{exec} + Stalls \text{ per Instr.} + MAPI \times Miss \text{ Rate} \times Miss \text{ penalty}) \times T_{CLK}$$

Analizziamo ora i vantaggi di avere più livelli di cache, il tempo medio di accesso è dato da:

$$AMAT = Hit \text{ Time}_{L1} + Miss \text{ Rate}_{L1} * Miss \text{ Penalty}_{L1}$$

$$Miss \text{ Penalty}_{L1} = Hit \text{ Time}_{L2} + Miss \text{ Rate}_{L2} * Miss \text{ Penalty}_{L2}$$

Combinando queste due equazioni otteniamo:

$$AMAT = Hit \text{ Time}_{L1} + Miss \text{ Rate}_{L1} \times (Hit \text{ Time}_{L2} + Miss \text{ Rate}_{L2} \times Miss \text{ Penalty}_{L2})$$

Introducendo il concetto di *global miss* (vedi Capitolo 7)

$$Miss \text{ Rate}_{L1L2} = Miss \text{ Rate}_{L1} \times Miss \text{ Rate}_{L2}$$

Riscriviamo così l'equazione del tempo medio di accesso come:

$$AMAT = Hit \text{ Time}_{L1} + Miss \text{ Rate}_{L1} \times Hit \text{ Time}_{L2} + Miss \text{ Rate}_{L1L2} \times Miss \text{ Penalty}_{L2}$$

Il tempo di esecuzione diventa:

$$CPU_{time} = IC \times (CPI_{exec} + MAPI \times MR_{L1} \times HT_{L2} + MAPI \times MR_{L1L2} \times MP_{L2}) \times T_{CLK}$$

## Elenco delle figure

1	Esempio di istruzione ALU di tipo R-Format . . . . .	5
2	Divisione delle informazioni in un registro di un'istruzione ALU di tipo diretto . . . . .	5
3	Struttura di un'istruzione tipo load/store . . . . .	5
4	Struttura di un'istruzione tipo branch condizionato . . . . .	6
5	Divisione delle informazioni in un registro di un'istruzione tipo branch incondizionato . . . . .	6
6	Divisione dei registri nei diversi casi di operazione . . . . .	7
7	Cicli eseguiti da ogni operazione . . . . .	8
8	Latenza delle diverse operazioni . . . . .	9
9	Esempio di implementazione di MIPS . . . . .	9
10	Hardware necessario per realizzare l'Instruction Fetch . . . . .	9
11	Hardware che implementa un'istruzione di tipo aritmetico . . . . .	10
12	Hardware che implementa un'istruzione di tipo load . . . . .	10
13	Hardware che implementa un'istruzione di tipo store . . . . .	11
14	MIPS a singolo ciclo con logica di controllo . . . . .	12
15	Confronto tra esecuzione sequenziale e pipelined . . . . .	13
16	Fasi della pipeline necessarie ad ogni istruzione . . . . .	13
17	Schema di un MIPS con pipeline . . . . .	15
18	Esempio di data hazard . . . . .	16
19	Esempio di uso <code>nop</code> . . . . .	17
20	Esempio di uso degli stalli . . . . .	17
21	Esempio di forwarding path . . . . .	18
22	Schema MIPS con forwarding . . . . .	20
23	Forwarding path tra due fasi di memorizzazione successive. . . . .	21
24	Fase di memorizzazione multyciclo che porta alla creazione di dipendenze di tipo WAW . . . . .	21
25	Fase di esecuzione multyciclo che porta alla creazione di dipendenze di tipo WAW . . . . .	21
26	Esempio di istruzione di branch . . . . .	22
27	Suddivisione dell'esecuzione di un'istruzione di salto nelle varie fasi di una pipe . . . . .	22
28	Esempio di esecuzione di un'istruzione di salto . . . . .	23
29	Hardware aggiuntivo per risolvere i problemi di controllo . . . . .	24
30	Esempio di penalità dovuto ad una predizione sbagliata . . . . .	25
31	Esempio di utilizzo del branch delay slot . . . . .	26
32	Esempio di selezione dell'istruzione <i>From Target</i> . . . . .	27
33	Esempio di uso della tecnica <i>From Fall-Through</i> . . . . .	27
34	Esempio di <i>Branch History Table</i> . . . . .	28
35	Macchina a stati finiti per una BHT a 2 bit . . . . .	29
36	Esempio di salti correlati . . . . .	29
37	Esempio di predittore correlato di tipo (1,1) . . . . .	30
38	Esempio di predittore correlato (2,2) . . . . .	30
39	Comparazione delle performance per predittore non correlato e predittore corellato	31
40	Esempio di predittore adattativo . . . . .	32
41	Esempio di predittore GShare . . . . .	32
42	Esempio di record di un Branch Target Buffer . . . . .	32
43	Struttura di un Branch Target Buffer . . . . .	33
44	Esempi di dipendenza dei dati . . . . .	34

45	Esecuzione di istruzione in una pipeline dual-issue . . . . .	35
46	Schema hardware per una pipeline dual-issue con una unità ALU/BR e una unità load/store . . . . .	36
47	Confronto di prestazioni tra architetture . . . . .	38
48	Struttura di un processore superscalare . . . . .	38
49	Tabella delle dipendenze in uno scheduler dinamico . . . . .	39
50	Scheduler statico nel caso di processori VLIW . . . . .	39
51	Architettura di un sistema <i>Scoreboard</i> . . . . .	40
52	Architettura per l'algoritmo di Tomasulo . . . . .	42
53	Architettura di un'unità funzionale . . . . .	43
54	Algoritmo di Tomasulo all'istante 14 . . . . .	46
55	Associazione tra RF e RF fisico . . . . .	47
56	Meccanismo di register renaming . . . . .	47
57	Architettura di una pipeline per VLIW . . . . .	50
58	Esempio di utilizzo di Very Long Instruction Word con queste FU: 2 LD/ST, 1 Int e 1 Int/Branch . . . . .	51
59	Stage della pipeline di un processore Itanium . . . . .	52
60	Esempio di grafico delle dipendenze tra le istruzioni di un blocco base . . . . .	53
61	Esempio di <i>critical path</i> . . . . .	53
62	Esempio di scheduling List-Based con utilizzo di Reservation Table . . . . .	54
63	Esempio di unrolling con fattore 4 . . . . .	55
64	Loop con istruzioni indipendenti ad ogni iterazione . . . . .	57
65	Svolgimento delle diverse iterazioni . . . . .	57
66	Nuovo codice che implementa il loop . . . . .	58
67	Realizzazione di un superblocco . . . . .	59
68	Struttura dell'algoritmo di Tomasulo con Reorder Buffer . . . . .	61
69	Struttura di un Reorder Buffer . . . . .	62
70	Esempio di esistenza dei thread . . . . .	65
71	Esempio di processore super scalare 47 e di processore super scalare con multithreading di tip coarse-grained 71(b) . . . . .	66
72	Esempio di fine-grained MT 72(a) e di simultaneous MT 72(b) . . . . .	66
73	Esempio di gerarchia di memoria . . . . .	68
74	Sistema di memoria in un sistema server e in un device mobile. . . . .	69
75	Architettura cache di un processore Intel Core i7 . . . . .	70
76	Struttura di un record di cache . . . . .	71
77	Esempio di piazzamento di blocchi nel caso di cache di tipo direct map . . . . .	72
78	Indirizzamento in un sistema di cache direct map . . . . .	72
79	Campi di una cache completamente associativa . . . . .	73
80	Esempio di una cache completamente associativa . . . . .	73
81	Esempio di una cache associativa a due vie . . . . .	74
82	Indirizzamento in una cache associativa a due vie . . . . .	74
83	Problema del <i>block placement</i> nei tre meccanismi di cache . . . . .	75
84	Utilizzo del write buffer . . . . .	76
85	Incremento di performance dovuto al pre-fetching di due blocchi quando avviene un miss . . . . .	79
86	Esempio di merging array . . . . .	80
87	Esempio di loop interchange . . . . .	80
88	Esempio di loop fusion . . . . .	81

89	Esempio di utilizzo del blocking . . . . .	81
90	Tempi di hit per dimensione e livello di associatività di una cache. . . . .	82
91	Principali tecniche di incremento delle performance . . . . .	83
92	Principali tecniche di incremento delle performance . . . . .	84
93	Traduzione da indirizzo virtuale ad uno fisico . . . . .	84
94	Esempio di page table . . . . .	85
95	Esempio di due processi che sfruttano il lock . . . . .	87
96	Diagramma di flusso dello spin lock . . . . .	88
97	Esempio di parallelizzazione del codice precedente . . . . .	88
98	. . . . .	89
99	Implementazione di un sistema di barrier . . . . .	89
100	Implementazione di un sistema di barrier con senso di esecuzione . . . . .	90
101	Architettura multiprocessore a singolo bus . . . . .	91
102	Architettura multiprocessore con connessione tramite rete . . . . .	91
103	Esempio di rete ad anello . . . . .	92
104	Esempio di topologia a maglia con N=4 . . . . .	93
105	Esempio di topologia ad ipercubo. . . . .	93
106	Esempio di singolo spazio degli indirizzi condiviso e spazio degli indirizzi multiplo. . . . .	94
107	Architettura di un sistema multiprocessore a memoria distribuita . . . . .	95
108	Struttura di uno snooping protocol . . . . .	97
109	Macchina a stati finiti che rappresenta un blocco di cache . . . . .	98
110	Esempio di directory base protocol distribuito . . . . .	99
111	Esempio di directory per quattro blocchi di memoria utilizzati da quattro processori . . . . .	99
112	Macchina a stati finiti per i blocchi di memoria nel caso di protocollo directory based . . . . .	100