

<h1>Pipeline ideale</h1> <p>A regime, una istruzione conclusa per ogni colpo di clock</p> <div><table><tr><td>F</td><td>D</td><td>E</td><td>M</td><td>W</td></tr><tr><td></td><td>F</td><td>D</td><td>E</td><td>M</td><td>W</td></tr><tr><td></td><td></td><td>F</td><td>D</td><td>E</td><td>M</td><td>W</td></tr><tr><td></td><td></td><td></td><td>F</td><td>D</td><td>E</td><td>M</td><td>W</td></tr><tr><td></td><td></td><td></td><td></td><td>F</td><td>D</td><td>E</td><td>M</td><td>W</td></tr></table></div> <div>Tra le diverse fasi sono presenti dei registri di pipeline che permettono di <u>separare i dati elaborati</u> da ciascuna istruzione.</div>		F	D	E	M	W		F	D	E	M	W			F	D	E	M	W				F	D	E	M	W					F	D	E	M	W	<h1>Hazard</h1> <p>Situazioni che impediscono ad un'istruzione di essere eseguita nel momento 'ideale'. Tre tipologie:</p> <ul style="list-style-type: none">• Strutturali -> Più istruzioni richiedono la stessa risorsa (es accesso a register file, ma c'è una sola porta);• Di dato -> Un'istruzione dipende da dati elaborati da istruzioni precedenti (RAW);• Di controllo -> Causati da branch e altre istruzioni che alterano il PC.		<h1>Stalli</h1> <p>Sono uno dei modi di gestire gli hazard: in generale, l'istruzione stallata e le successive attendono, le precedenti continuano l'esecuzione. Possono essere utilizzati per gestire hazard strutturali e di dato.</p>		<p>Il forwarding prende i dati provenienti dagli output di</p> <ul style="list-style-type: none">• ALU• memoria <p>e destinati agli input di</p> <ul style="list-style-type: none">• ALU• memoria• unità di rilevamento dello zero, per i salti <p>Quindi dagli stati EX/MEM e MEM/WB agli stati IF/ID, ID/EX ed EX/MEM</p>	
F	D	E	M	W																																						
	F	D	E	M	W																																					
		F	D	E	M	W																																				
			F	D	E	M	W																																			
				F	D	E	M	W																																		
<h1>Control hazard</h1> <p>Dipendono dal fatto che il PC possa essere modificato da istruzioni di branch dopo il fetch di una o più istruzioni successive. Possibili soluzioni:</p> <ul style="list-style-type: none">• Stallare la pipeline (soluzione più semplice).• Predict untaken: supporre che il branch non venga mai preso. Se viene preso, tutte le azioni eseguite nel frattempo vengono annullate.• Predict taken: assumere che il salto non venga mai preso.• Branch delay slot. <p>Il MIPS scrive il PC alla fine della decode, se il salto è preso.</p>		<h1>Branch delay slot</h1> <p>Siccome l'istruzione successiva al branch viene comunque eseguita, invece che annullarla la si sfrutta per eseguire un'istruzione che dovrebbe comunque essere eseguita, e che quindi non sarà annullata se il salto viene preso. Il compilatore ha responsabilità di sfruttare il Branch Delay Slot opportunamente.</p>		<h1>Operazioni multiciclo</h1> <p>Alcune unità, che effettuano operazioni complesse (es, quelle sui numeri floating point) sono multiciclo (la loro fase EX dura più cicli). Se una istruzione I1 deve <u>scrivere</u> in un registro che deve essere ancora scritto da un'altra istruzione I2 ancora in EXE, I2 stalla in decode ed entra in EXE solo quando ne è uscita I1. <u>Latency</u> -> Tempo tra istruzione che produce un dato e istruzione che lo usa <u>Initiation interval</u> -> Numero di cicli tra l'issuing di due istruzioni che utilizzano la stessa unità (minimo 1). Alcune unità (normalmente la divisione) NON sono nemmeno pipelined. <u>Hazard strutturali più frequenti</u> -> Aggiungere porte costa, quindi normalmente si sceglie di stallare l'istruzione.</p>		<h1>Nuovi hazard di dato</h1> <ul style="list-style-type: none">• RAW - Read After Write, o <i>Dipendenza</i> -> I2 ha bisogno del risultato di I1, che però non lo ha ancora calcolato;• WAW - Write After Write, o <i>Output Dipendenza</i> -> I2 (più breve) scrive in R1 <u>prima</u> di I1 (più lunga), quindi I1 sovrascrive il dato di I2;• WAR - Write After Read, o <i>Antidipendenza</i> -> I2 deve ancora leggere il dato scritto da I1, ed I2 è eseguita per prima.																																				
<h1>Gestire gli hazard di dato</h1> <ul style="list-style-type: none">• RAW: il forwarding aiuta, dopo il forwarding l'unica soluzione sono gli stalli.• WAW e WAR: si gestiscono mediante <u>register renaming</u> e la <u>modifica dell'ordine delle istruzioni</u>.		<h1>Instruction level parallelism</h1> <p>L'ILP è sfruttato dalla pipeline per eseguire più istruzioni in parallelo: <u>più parallelismo esiste, maggiori saranno le prestazioni</u>.</p> <p>Per migliorare l'ILP si può modificare l'ordine delle istruzioni mediante un approccio dinamico, basato sul processore, oppure statico, basato sul compilatore.</p> <p>Un basic block è una sequenza di istruzioni senza branch né in entrata né in uscita, eccetto quelli all'inizio e alla fine.</p>		<h1>Rescheduling</h1> <p>All'interno di un basic block è possibile implementare il rescheduling, la modifica dell'ordine delle istruzioni in modo da eliminare gli stalli.</p> <p>Tuttavia, nel MIPS i basic block sono normalmente tra 4 e 7 istruzioni: per migliorare il parallelismo si può dunque considerare il loop unrolling.</p>		<h1>Dipendenze</h1> <p>Esistono tra istruzioni e possono creare degli hazard, ma non sono la stessa cosa. Tipologie</p> <p>-> Di dato, su registri e su memoria (più difficile rilevare hazard, impossibile a compile time) -> Di nome (antidipendenza e output dipendenza), non c'è una dipendenza tra i dati ma questi vengono comunque alterati.</p>		<h1>Loop unrolling</h1> <p>Se sono necessarie N iterazioni, il corpo del loop viene opportunamente replicato K volte, rendendo così necessarie N ÷ K iterazioni.</p> <p>Pro -> Meno overhead per il controllo dell'iterazione; -> Corpo del loop più largo, quindi più possibilità di fare rescheduling;</p> <p>Contro -> Dimensione del codice maggiore.</p>																																		

<div><div>Esempio di RAW</div><div><div>ADD R1, R2, R3</div><div>SUB R4, R5, R1</div></div><div>La SUB ha bisogno del dato che la ADD salva in R1, ma il dato non è ancora pronto quando la SUB viene eseguita.</div></div>	<div><div>Esempio di WAR</div><div><div>DIV R6, R2, R4</div><div>SUB R4, R5, R1</div></div><div>La SUB potrebbe scrivere R4 <u>prima</u> che la DIV lo legga durante l'EXE, ad esempio perché la DIV impiega più colpi di clock della SUB.</div></div>	<div><div>Esempio di WAW</div><div><div>DIV R6, R2, R4</div><div>SUB R6, R5, R1</div></div><div>Se la SUB viene conclusa prima della DIV, R6 avrà un valore non consistente.</div></div>	<div><div>Dipendenze <i>di nome</i> e <i>di dato</i></div><div>Notare come in un RAW la dipendenza sia <u>di dato</u>, ossia ciò che manca alla SUB (in questo caso) sia il risultato della precedente ADD, mentre WAR e WAW non sarebbero presenti se si utilizzassero registri diversi: si tratta infatti di dipendenze <u>di nome</u>.</div></div>
-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<h2>Cache</h2> <p>La cache è una memoria piccola e ad alta velocità, posta tra il processore e la memoria principale. Sono normalmente a più livelli, ed organizzate in un array.</p> <div><div></div><div></div><div></div><div></div><div></div><div></div></div> <div>Validity bitTagData blocks</div>		<h2>Organizzazione della cache</h2> <p>La cache è organizzata in linee, ogni linea contiene un blocco di memoria che include alcune parole.</p> <p>Ogni linea è associata a un tag, indicante il blocco di memoria presente nella specifica linea di cache.</p>		<h2>Hit e miss</h2> <ul style="list-style-type: none">• Hit: il blocco richiesto è presente in cache;• Miss: il blocco richiesto NON è presente in cache, e viene caricato dalla memoria. <div><div>Tag (T)Index (I)Offset (O)</div><div>Comparator</div><div>Effective AddressValidity bit(s)</div><div>Hit / Miss</div></div>		<h2>Performance</h2> <ul style="list-style-type: none">• In caso di hit: impatto sempre positivo;• In caso di miss:<ul style="list-style-type: none">◦ Se la cache carica l'intero blocco e poi la fornisce, impatto decisamente negativo;◦ Se la cache accede la memoria e fornisce la word immediatamente (<i>load-through</i>), impatto leggermente negativo.	
<h2>Ubicazione della cache</h2> <p>La cache, di norma, è collocata tra la CPU e il bus. Questo consente di ridurre la pressione sul bus, e supporta architetture multiprocessore.</p> <p>La cache può essere unica, oppure separata (Architettura Harvard) per istruzioni e per i dati. La cache delle istruzioni è più facilmente gestibile, in quanto viene solo letta.</p>		<h2>Dimensione della cache</h2> <p>Una cache, quanto più è grande, tanto più:</p> <ul style="list-style-type: none">• I costi salgono;• Le performance migliorano;• La cache diventa più lenta. <p>Le dimensioni, normalmente, variano tra alcuni Kb ad alcuni MB.</p>		<h2>Mapping</h2> <p>Il meccanismo di <i>mapping</i> definisce, dato in ingresso un blocco di memoria, quale linea della cache debba essere scritta. Esistono diversi metodi:</p> <ul style="list-style-type: none">• Direct mapping;• Set associative;• Full associative.			
<h2>Direct mapped</h2> <p>Ogni blocco di memoria <i>i</i> è associato in maniera statica ad una linea <i>k</i> della cache, mediante l'equazione</p> $k = i \bmod N$ <p>Ad esempio, in una cache grande 4 linee, se la memoria ha 16 blocchi una linea avrà associati 4 blocchi.</p>		<h2>Full Associative</h2> <p>Ogni blocco della memoria può essere posizionato in un qualsiasi blocco della cache.</p>		<h2>Set Associative</h2> <p>Le linee della cache sono divise in S set, ognuno contenente W linee. Ogni blocco di memoria <i>i</i> è associato ad un set <i>k</i>, calcolato con la formula</p> $k = i \bmod S$ <p>Il blocco <i>i</i> può essere piazzato in qualsiasi delle W linee del set <i>k</i>. Una cache con W linee in ogni set è chiamata cache a <i>W-vie</i>.</p>			
Meccanismo semplice da implementare, è sufficiente leggere i bit meno significativi dell'indirizzo di memoria per trovare il corrispondente in cache		Possono crearsi dei conflitti se due locazioni di memoria associate alla stessa linea sono accedute spesso		Massima flessibilità	L'hardware diventa più complesso, in quanto dovrà cercare un blocco libero su tutta la cache	Bilanciamento tra flessibilità nella collocazione dei blocchi e ricerca di un blocco libero	
<h2>Replacing algorithm</h2> <p>Il <i>replacing algorithm</i> definisce, per i casi <i>full associative</i> e <i>set associative</i>, quale linea della cache usare per scrivere i dati provenienti dalla memoria.</p> <ul style="list-style-type: none">• LRU (Least Recently Used);• FIFO (First in First Out);• LFU (Least Frequently Used);• Random.		<h2>Write back e Write through</h2> <p>Quando la cache viene scritta dal processore, occorre decidere come aggiornare la memoria principale. Nel caso del Write back un bit detto <i>dirty bit</i> sarà settato, e al momento della sovrascrittura del blocco se tale bit è a 1 il blocco preesistente in cache verrà copiato in memoria. In questo caso la sostituzione sarà più lenta, ci potrebbero essere problemi di coerenza in sistemi multiprocessore e, in caso di guasti, potrebbe non essere possibile ripristinare correttamente la memoria.</p> <p>L'altra opzione è il Write through: i dati scritti in cache sono anche trascritti nella memoria principale. Questo introduce, in teoria, una perdita di efficienza, in pratica però, siccome le letture sono molte più delle scritture, il problema dunque è in realtà più limitato.</p>					

<h2>Cache coherence</h2> <p>In sistemi multiprocessore, la memoria è comune a tutti i processori <u>mentre la cache è, di norma, locale ad ogni singolo processore</u>: questo provoca potenziali problemi di coerenza.</p> <p>Per questo motivo esiste il validity bit: ogni volta che un certo blocco di memoria principale è scritto da un processore, eventuali linee di cache di altri processori vengono segnate come obsolete settando il validity bit a 0. Se a 0, tutti gli accessi alla cache producono un miss.</p>	<h2>Livelli di cache</h2> <p>La cache può essere divisa in più livelli: la CPU cercherà il dato nel primo livello, quindi passerà al secondo e così via, copiando se necessario le informazioni nel livello più basso. I livelli successivi al primo possono essere condivisi tra più core.</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<h2>Branch prediction</h2> <p>Predire i branch è fondamentale per evitare stalli dovuti a dipendenze di controllo. Le tecniche di predizione possono essere divise in:</p> <ul style="list-style-type: none"> • Tecniche statiche: gestite dal compilatore, che si affida ad una analisi preliminare del codice; • Tecniche dinamiche: gestite dall'hardware, che si affida al comportamento in tempo reale del codice. 	<h2>Tecniche statiche</h2> <p>Il compilatore può predire il comportamento dei branch in diversi modi:</p> <ul style="list-style-type: none"> • Predire sempre i branch come taken; • Predire in base alla direzione del branch <ul style="list-style-type: none"> ◦ In avanti, normalmente untaken; ◦ All'indietro, normalmente taken. • Predire in base alle informazioni provenienti da precedenti esecuzioni, identificando una serie di dati in ingresso, eseguendo il programma un certo numero di volte salvando le statistiche circa tali esecuzioni, che saranno poi utilizzate come statistiche. 	<h2>Tecniche dinamiche</h2> <p>Le tecniche dinamiche si basano sull'hardware, e possono essere di diversi tipi:</p> <ul style="list-style-type: none"> • Branch history table (BHT); • Two-level prediction schemes; • Branch-target buffer; • ... <p>Esse si basano sul principio di località: come esiste quella temporale (una certa porzione di memoria acceduta a tempo T sarà probabilmente acceduta anche poco dopo T) e spaziale (una certa porzione di memoria acceduta alla posizione P, le posizioni intorno a P saranno probabilmente accedute poco dopo), <u>la località di salto è il comportamento normalmente assunto ad esempio in caso di <i>for</i> e <i>while</i>, che potrebbero eseguire la stessa porzione di codice molte volte</u>. Questo comportamento è dunque prevedibile, e se è prevedibile sarà anche ottimizzabile. La località non è normalmente spaziale, in quanto le seppur poche possibilità potrebbero essere collocate distanti le une dalle altre.</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<h2>Tecniche dinamiche - Funzionamento</h2> <pre> graph LR A[Predizione nella fase di decode] --> B[Issue dell'istruzione predetta] B -- "In base alla tecnica scelta" --> C{Predizione corretta?} C -- Si --> D[L'esecuzione continua] C -- No --> E[Flush della pipeline e issue dell'istruzione corretta] </pre>	<h2>Branch History Table (BHT)</h2> <p>Si tratta della tecnica più semplice: una piccola memoria</p> <ul style="list-style-type: none"> • indicizzata dalla parte meno significativa dell'indirizzo delle istruzioni, significativa poi solo per quelle di branch; • contenente, per ogni record, uno o più bit che tengono traccia se il branch sia stato o meno preso l'ultima volta che l'istruzione è stata eseguita. <p>L'algoritmo prevede:</p> <ul style="list-style-type: none"> • accesso alla BHT mediante la parte meno significativa dell'indirizzo dell'istruzione; • calcolo del PC sulla base dell'esito della predizione; • aggiornamento della BHT in base al risultato effettivo. <p>Più i record della BHT si riferiranno con precisione al branch di interesse e più la predizione sarà accurata, maggiore sarà la performance di questo metodo.</p>	<h2>Schemi a n bit</h2> <p>In questo caso, invece di avere un singolo bit nella BHT, il record è costituito da un contatore. Quando il contatore è a più di metà del valore massimo il salto è predetto come taken, altrimenti come untaken.</p> <p>Ad ogni aggiornamento il contatore è incrementato o decrementato in base all'esito effettivo del branch.</p>
-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

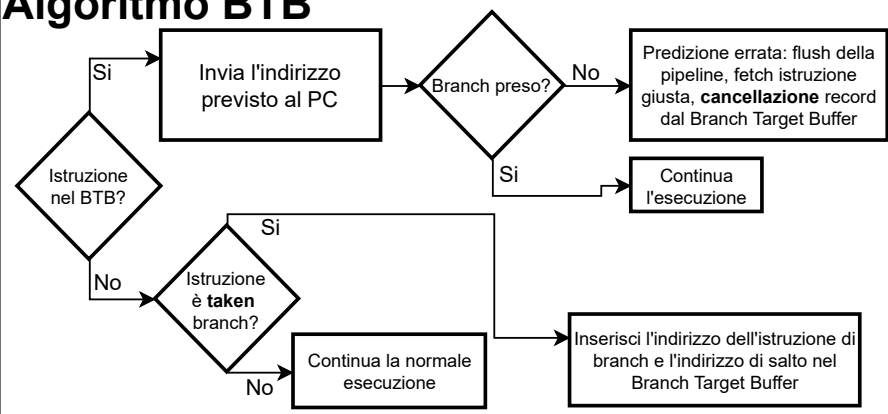
<h2>Schemi a due livelli (m,n), esempio (2,2)</h2> <p>Questi schemi utilizzano la cronologia degli ultimi m branch per la predizione; ogni predittore ha un contatore a n bit che indica la probabilità che il branch sia preso o no.</p> <p>Per ogni indirizzo di branch esiste una riga della tabella associata (più indirizzi per riga).</p> <ul style="list-style-type: none"> • Un registro a 2 bit (<i>shift register</i>) tiene traccia degli esiti degli ultimi 2 branch (in assoluto, lungo tutta l'esecuzione del codice): ad esempio, 10 indica che l'ultimo branch non è stato preso, il penultimo sì (shift verso sinistra); • La tabella dei predittori contiene 2² = 4 combinazioni di cronologia, con un predittore a 2 bit per ogni combinazione: 00, 01, 10 e 11. Nei 2 bit del predittore sono conteggiati gli esiti degli ultimi branch; • Supponiamo che la cronologia (shift register) iniziale sia 10: andiamo a cercare il predittore associato alla cronologia 10, e supponiamo che il contatore associato sia a 2 (10). In questo caso il predittore fa prendere il branch. • Se il branch viene preso effettivamente, il predittore associato alla cronologia 10 verrà incrementato (diventando 3), e la cronologia passerà a 01; diversamente la cronologia passerà a 00, e il contatore associato alla cronologia 10 decrementato (diventando 1). <p>Questo approccio sfrutta le dipendenze tra branch consecutivi che si possono trovare in caso di iterazioni; tuttavia, aumentare m e n non crea significativi vantaggi, aumentando però il costo hardware.</p>

Branch target buffer

Oltre a conoscere se un branch sarà o meno preso, **è necessario, per ridurre gli effetti negativi delle dipendenze di controllo, conoscere anche il nuovo valore del PC se il branch sarà preso**: il *branch target buffer* esiste per questo motivo.

- Si tratta di una tabella che immagazzina:
- Indirizzi delle istruzioni di branch;
 - Target dell'istruzione da caricare nel PC.

Algoritmo BTB



<h2>Dynamic scheduling</h2> <p>Il dynamic scheduling permette di riordinare le istruzioni in hardware, riducendo gli stalli.</p> <p>Il Dynamic Scheduling:</p> <ul style="list-style-type: none"> Identifica dipendenze non conosciute a compile time, ad esempio gli indirizzi di memoria; Semplifica il lavoro del compilatore; Permette al processore di tollerare ritardi non predicibili; Permette di eseguire lo stesso codice su diversi processori pipelined. 	<h2>Esempio</h2> <div> <div>DIV.D F0, F2, F4</div> <div>ADD.D F10, F0, F8</div> <div>SUB.D F12, F8, F14</div> </div> <p>In una situazione come questa, la DIV.D causa lo stallo della ADD.D per la dipendenza di dato, ma anche la SUB.D viene stallata, anche se non ha nessuna dipendenza con le istruzioni precedenti e potrebbe essere senza problemi eseguita.</p> <p>Miglioramento: Rimuovere il vincolo della <i>in-order execution</i> Rischio di hazard WAR, WAW ed eccezioni imprecise.</p>	<h2>Dynamic scheduling (II)</h2> <ul style="list-style-type: none"> Rende impossibile la gestione precisa delle eccezioni, in quanto, con I_e istruzione che ha sollevato l'eccezione <ul style="list-style-type: none"> Una istruzione prima di I_e potrebbe dover essere ancora eseguita; Una istruzione dopo I_e potrebbe essere già stata eseguita. Richiede di dividere la fase ID in due parti: <ul style="list-style-type: none"> <i>issue</i>: decode dell'istruzione e ricerca di hazard strutturali, <u>questa fase è eseguita in ordine</u>; <i>read operands</i>: attesa dell'assenza di hazard di dato, quindi lettura degli operandi: in questa fase le istruzioni possono sorpassarsi, entrando così nella fase di esecuzione <i>out-of-order</i>. <p>Esistono diversi approcci, tra cui lo scoreboarding e l'algoritmo di Tomasulo.</p>
<h2>Architettura di <u>Tomasulo</u></h2> <p>Introduce un'architettura nuova e più complessa, basata su:</p> <ul style="list-style-type: none"> Tracciamento della disponibilità degli operandi mediante Reservation station. Esse <ul style="list-style-type: none"> immagazzinano gli operandi in attesa che vengano utilizzati dalle istruzioni: essi vengono salvati nella RS non appena disponibili; implementano l'<i>issue</i>; identificano univocamente un'istruzione nella pipeline, in quanto le istruzioni in attesa di esecuzione designeranno la RS che fornirà loro l'operando in input. Register renaming, ogni volta che un dato entra in una RS per un'istruzione, il nome del registro di provenienza è mutato in quello della reservation station. <u>Questo permette di eliminare gli hazard WAW e WAR.</u> Common Data Bus (CDB): i dati sono passati direttamente alle unità funzionali invece che passare dai registri, in modo che tutte le unità in attesa di un operando lo carichino insieme. 	<h2>Issue</h2> <p>In questa fase le istruzioni sono prelevate dalla coda in modalità FIFO. Se non ci sono RS disponibili, l'istruzione stalla, altrimenti se c'è una RS disponibile per l'operazione richiesta l'istruzione vi viene inviata insieme agli operandi se disponibili, altrimenti viene registrata la FU responsabile per la loro generazione.</p>	<h2>Execution</h2> <p>Istruzioni aritmetiche</p> <p>Quando un operando appare sul CDB, esso è letto dalla RS, e non appena tutti gli operandi sono disponibili l'istruzione viene eseguita.</p> <p>Istruzioni load/store</p> <ol style="list-style-type: none"> Non appena il registro contenente l'indirizzo di base è disponibile, l'effective address è calcolato e salvato nel <i>load/store buffer</i>; In base al tipo di istruzione: <ol style="list-style-type: none"> Load: eseguita appena la memoria è disponibile; Store: attesa dell'operando da scrivere, quindi scrittura non appena la memoria è disponibile. <p>Nessuna istruzione può iniziare l'esecuzione prima che tutti i precedenti branch siano stati completati.</p>
<h2>Write result</h2> <p>In questa fase</p> <ul style="list-style-type: none"> i risultati sono, tramite il CDB <ul style="list-style-type: none"> scritti nei registri; portati alle FU che ne hanno bisogno. le istruzioni store effettuano la scrittura in memoria. <p>I risultati sono associati a degli identificatori interni alle RS, che funzionano come nomi di registro virtuale, implementando in questo modo il register renaming.</p>	<h2>Reservation station</h2> <div> <div> <div>Operazione da effettuare</div> <div>Op</div> </div> <div> <div>Valori degli operandi in input</div> <div>V_j V_k</div> </div> <div> <div>RS che produrranno gli operandi se non disponibili</div> <div>Q_j Q_k</div> </div> <div> <div>Solo per load/store, immagazzina prima l'immediato, poi l'effective address</div> <div>A</div> </div> <div> <div>Stato della RS e dell'unità funzionale corrispondente</div> <div>Busy</div> </div> </div> <p>Campo Q_i del register file: contiene il numero della RS contenente l'istruzione il risultato della quale sarà salvato in quel registro. Se nullo, nessuna istruzione sta attualmente lavorando al risultato di quel registro.</p>	<h2>Vantaggi e svantaggi - Load e store</h2> <p>Sono eliminati gli hazard WAW e WAR, e la logica di controllo degli altri hazard è distribuita. Tuttavia, questa architettura ha un'elevata complessità, e il Common Data Bus può rappresentare un collo di bottiglia.</p> <p>Non è necessario effettuare il loop unrolling, in quanto le istruzioni sono naturalmente eseguite in parallelo.</p> <p>Gli stalli WAR e RAW possono comunque avvenire cambiando l'ordine delle load e delle store. Per evitare hazard in questi casi:</p> <ul style="list-style-type: none"> In caso di load, lo store buffer viene controllato per trovare store ancora da eseguire sullo stesso indirizzo. Nel caso, la load attende; In caso di store, sia lo store che il load buffer vengono controllati alla ricerca di istruzioni che agiscono sullo stesso indirizzo. Nel caso, la store attende.

<h2>Hardware based speculation</h2> <p>L'hardware-based speculation è una combinazione di tre idee:</p> <ul style="list-style-type: none">• Dynamic branch prediction;• Dynamic scheduling;• Speculation. <p>Un processore che supporta il branch prediction e il dynamic scheduling effettua il fetch e l'issue delle istruzioni considerando la predizione come sempre corretta. <u>Un processore che supporta la speculation effettua anche l'esecuzione delle istruzioni</u>. Le operazioni sono quindi eseguite non appena gli operandi sono disponibili.</p>		<h2>Architettura</h2> <p>L'architettura di base è quella di Tomasulo, estesa per supportare la speculazione con due step nell'esecuzione delle istruzioni:</p> <ul style="list-style-type: none">• Calcolo dei risultati, anche con bypass di altre istruzioni;• Aggiornamento di register file e memoria, effettuato quando l'istruzione non è più speculativa. Tale fase è detta <i>commit</i> ed è effettuata in ordine. <p>Il ReOrder Buffer (ROB) è la struttura preposta ad ospitare i risultati delle istruzioni fintantoché non hanno effettuato il commit. Ha registri virtuali aggiuntivi, ed integra lo store buffer dell'architettura base.</p> <p>I risultati già calcolati nell'architettura base sono letti dal register file, mentre <u>con la speculation sono letti dal ROB se l'istruzione non ha ancora fatto commit</u>, altrimenti dal register file.</p>	<h2>Struttura del ROB</h2> <p>Ogni record del ROB ha quattro campi:</p> <ul style="list-style-type: none">• <i>Instruction type</i>: branch, store, register;• <i>Destination</i>: numero del registro od indirizzo di memoria;• <i>Value</i>: risultato dell'istruzione che non ha ancora effettuato il commit;• <i>Ready</i>: indica se l'istruzione ha completato l'esecuzione. <p>Il ROB funziona in modo circolare.</p>
<h2>Fase di Issue o Dispatch</h2> <ul style="list-style-type: none">• L'istruzione è estratta dalla coda se ci sono:<ul style="list-style-type: none">◦ Una RS libera;◦ Uno slot libero nel ROB• Altrimenti, l'istruzione è stallata;• Gli operandi sono inviati alla RS se sono nel RF o nel ROB;• Il numero dello slot ROB per l'istruzione è inviato alla RS per contrassegnare l'istruzione e i suoi risultati, quando saranno scritti sul CDB.	<h2>Fase di Execute</h2> <ul style="list-style-type: none">• Non appena gli operandi sono pronti l'istruzione è eseguita (evita hazard RAW);<ul style="list-style-type: none">◦ Se non sono disponibili gli operandi, essi sono ricavati dal CDB non appena vengono prodotti. <p>La durata di questa fase dipende dalla tipologia dell'istruzione.</p>	<h2>Fase di Write Results e Fase di Commit</h2> <ul style="list-style-type: none">• Non appena i risultati sono pronti, essi vengono inviati sul CDB insieme al tag che li identifica, e vengono ricevuti dal ReOrder Buffer;• In questo modo, sono letti da ogni RS che li sta aspettando;• Viene contrassegnata come <i>disponibile</i> la relativa entry della RS. <p>Una volta nel ReOrder Buffer viene eseguita la <u>fase di Commit</u>:</p> <ul style="list-style-type: none">• Le istruzioni sono ordinate in base all'ordine originale;• <u>Quando l'istruzione raggiunge la testa del buffer</u>:<ul style="list-style-type: none">◦ Se il branch è stato predetto erroneamente, è effettuato il flush del buffer e l'esecuzione riprende con l'istruzione corretta successiva;◦ Altrimenti, i risultati sono scritti nel register file o in memoria• In ogni caso, il record del ROB è contrassegnato come libero.	
<h2>Hazard</h2> <ul style="list-style-type: none">• WAW e WAR non possono avvenire, grazie al dynamic renaming;• Gli hazard RAW mediante la memoria sono prevenuti:<ul style="list-style-type: none">◦ Mantenendo l'ordine del programma nel calcolo dell'effective address delle load;◦ Evitando che una load possa iniziare l'esecuzione se esiste nel ROB una store con campo destinazione corrispondente al campo origine della load.	<h2>Istruzioni di Store</h2> <p>La scrittura in memoria è effettuata solo al commit, per questo motivo l'operando serve in realtà al momento del commit e non nella fase di write results.</p> <p>Per questo motivo, il ROB avrà un ulteriore campo che terrà traccia dell'origine del dato da scrivere, se l'istruzione è una store.</p>	<h2>Eccezioni</h2> <p>Le eccezioni sono eseguite non al momento in cui accadono, ma al momento del commit. Quando un'istruzione effettua il commit vengono, se del caso, anche eseguite le eccezioni, e viene effettuato il flush del buffer. Se l'istruzione che doveva eseguire un'eccezione è rimossa dal buffer, l'eccezione viene ignorata.</p> <p>Sono dunque supportare le eccezioni precise.</p>	<h2>Istruzioni dispendiose</h2> <p>La scrittura in memoria è effettuata solo al commit, per questo motivo l'operando serve in realtà al momento del commit e non nella fase di write results.</p> <p>Per questo motivo, il ROB avrà un ulteriore campo che terrà traccia dell'origine del dato da scrivere, se l'istruzione è una store.</p>

<h2>Multiple Issue Processors</h2> <p>Un CPI minore di 1 (quindi più di una istruzione, in media, è completata in un colpo di clock) può essere ottenuto effettuando l'<i>issue</i> di più istruzioni per ogni colpo di clock. Ci sono due tipologie di processori con questa possibilità:</p> <ul style="list-style-type: none">• <i>processori superscalari</i>, con scheduling statico o dinamico di istruzioni;• processori <i>VLIW - Very Long Instruction Word</i>. <p>In entrambi i casi, i processori hanno unità funzionali multiple.</p>	<h2>MIPS superscalare</h2> <p>Questa versione del MIPS implementa lo scheduling statico. Due istruzioni possono essere issued ad ogni colpo di clock se:</p> <ul style="list-style-type: none">• Una delle istruzioni è una load, store, branch od operazione aritmetica intera;• L'altra è un'operazione FP diversa da load e store. <p>Due istruzioni occupano 64 bit, e il fetch e decode sono effettuati ad ogni colpo di clock. Sono allineate a 64 bit e costituiscono un <i>issue packet</i>.</p> <p>I vecchi processori superscalari richiedevano una struttura fissa dell'issue packet (prima l'istruzione intera, poi la FP): oggi questa struttura rigida normalmente non esiste più. Ad ogni colpo di clock sono lette due istruzioni. Se le istruzioni appartengono a blocchi di cache separati, normalmente viene fatto il fetch solo di una. Normalmente è consentito un solo branch all'interno di un <i>issue packet</i>.</p>	<h2>Contention di registri FP, RAW</h2> <p>Quando la prima istruzione è una load, store o move FP si può creare una 'contesa' per la porta del registro FP. Le possibili soluzioni sono forzare la prima istruzione ad essere eseguita in autonomia, oppure aggiungere una porta al register file FP.</p> <p>Se la prima istruzione è una load, store o move FP e la seconda legge il suo risultato, è possibile un hazard RAW. In questo caso la seconda istruzione deve essere ritardata di un colpo di clock.</p>
<h2>Multiple issue dynamic scheduling</h2> <p>Questo tipo di processori adottano un'architettura simile a quella di Tomasulo. Per rendere l'implementazione più semplice, le istruzioni non sono mai issued fuori ordine.</p> <p>Il CDB può essere un collo di bottiglia, per questo motivo la performance può essere migliorata duplicandolo.</p>		