

Calcolatori Elettronici L-A

ISA DLX: Implementazione Tramite Struttura “Pipelined”

DLX pipelined 1

Principio del Pipelining (1)

Il pipelining è oggi la principale tecnica di base impiegata per rendere “veloce” una CPU

L’idea alla base del pipelining è generale, e trova applicazione in molteplici settori dell’industria (linee di produzione, oleodotti ...)

Un sistema, S , deve eseguire N volte un’attività A :



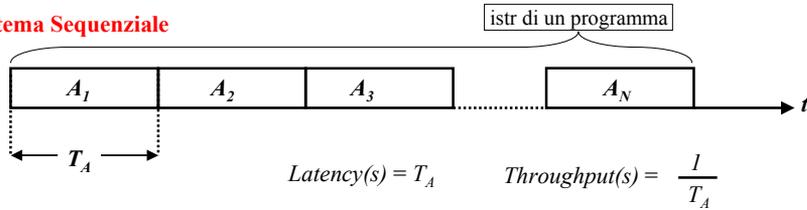
Def.: *Latency*: (o latenza [sec]) T_A tempo che intercorre fra l’inizio ed il completamento dell’attività A .

Def.: *Throughput*: (produttività [Hz]) frequenza con cui vengono completate le attività ($\frac{1}{T_A}$)

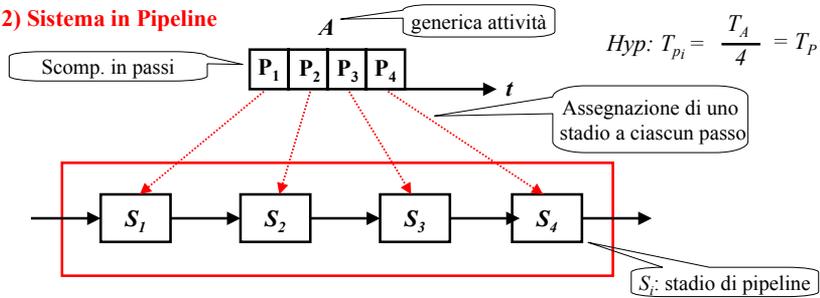
DLX pipelined 2

Principio del Pipelining (2)

1) Sistema Sequenziale

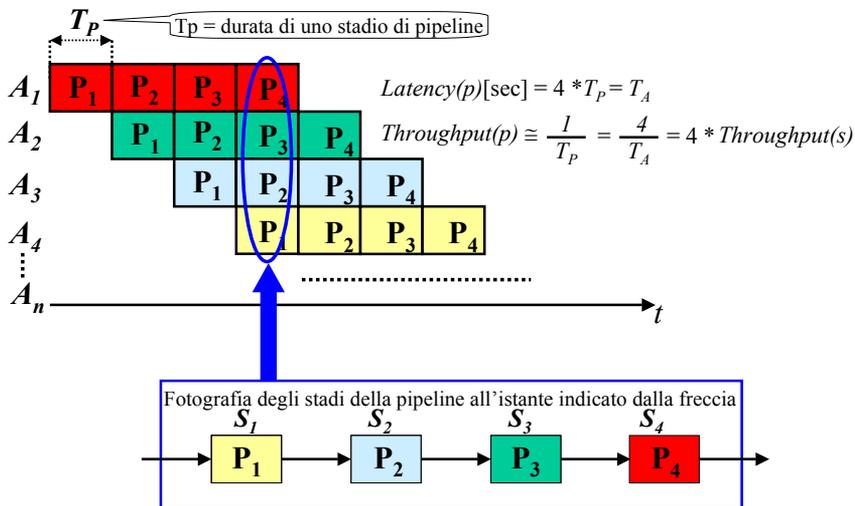


2) Sistema in Pipeline



DLX pipelined 3

Principio del Pipelining (3)



DLX pipelined 4

Principio del Pipelining (4)

Il pipelining *non* riduce il tempo necessario al completamento di una singola attività:

$$Latency(p) [sec] = Latency(s)$$

Il pipelining *incrementa* il *Throughput*, moltiplicandolo per un fattore pari al numero degli stadi (K):

$$Throughput(p) [1/sec] = K * Throughput(s)$$

Ciò porta ad una riduzione dello stesso fattore *del tempo totale di esecuzione* di una sequenza di N attività $T_N = [sec]$:

$$T_N = \frac{N}{Throughput} \rightarrow T_N(s) = \frac{N}{Throughput(s)}, \quad T_N(p) = \frac{N}{Throughput(p)}$$

$$\Rightarrow Speedup(p) = \frac{T_N(s)}{T_N(p)} = \frac{Throughput(p)}{Throughput(s)} = K$$

DLX pipelined 5

Bilanciamento di una Pipeline

- Caso Ideale

$$T_p = T_{pi} = \frac{T_A}{K} \Rightarrow \text{Pipeline perfettamente bilanciata} \Rightarrow Speedup = K$$

- Caso Reale

$$T_p = \max(T_{P1}, T_{P2}, \dots, T_{PK}) \Rightarrow \text{Bilanciamento imperfetto} \Rightarrow Speedup < K$$

- Un Esempio

$T_A = 20 t$ (t : unità di tempo, es.: cicli di clock)

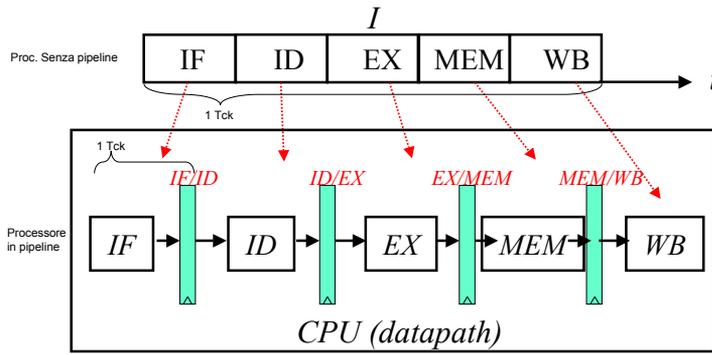
$$T_{P1} = 5 t, T_{P2} = 5 t, T_{P3} = 6 t, T_{P4} = 4 t \Rightarrow T_p = 6 t \quad (\text{è il più lento})$$

$$\Rightarrow Speedup(p) = \frac{T_A}{T_p} = \frac{20 t}{6 t} \cong 3.33 (< 4)$$

DLX pipelined 6

Pipelining in una CPU: DLX

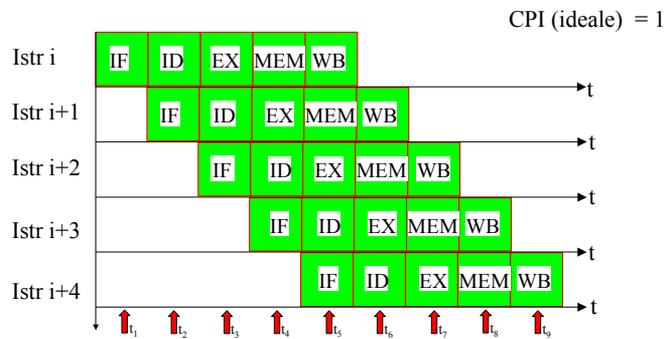
Attività: $A_1, A_2, A_3 \dots A_N$ \rightarrow Istruzioni: $I_1, I_2, I_3 \dots I_N$



Pipeline Cycle \rightarrow Clock Cycle \rightarrow Ritardo dello stadio più lento
 \downarrow
 CPI=1 (idealmente !)

DLX pipelined 7

Evoluzione della pipeline del DLX



Carico aggiuntivo introdotto dai Pipeline Registers:

$$T_{clk} = T_d + T_P + T_{su}$$

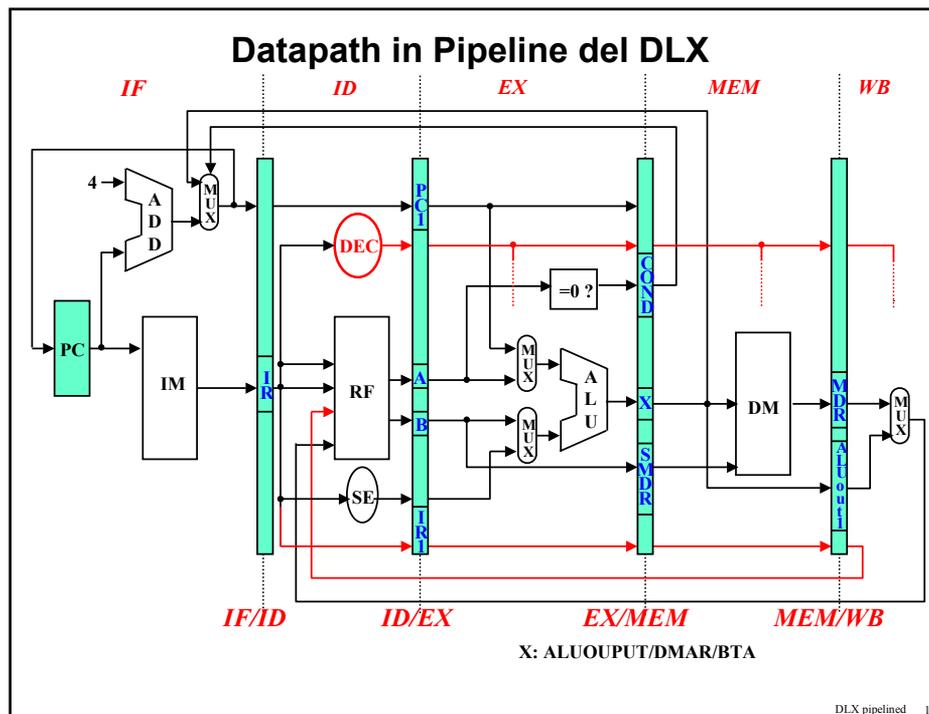
T_{clk} \rightarrow Clock Cycle
 T_d \rightarrow Ritardo registro a monte
 T_P \rightarrow Ritardo stadio più lento
 T_{su} \rightarrow Set-up registro a valle

DLX pipelined 8

Requisiti per l'implementazione in pipeline

- 1) Ogni stadio deve essere attivo in ogni ciclo di clock.
- 2) E' necessario incrementare il PC in IF (invece che in ID).
- 3) E' necessario introdurre un ADDER ($PC \leftarrow PC+4$) nello stadio IF.
- 4) Sono necessari due MDR (che chiameremo MDR e SMDR) per gestire il caso di una LOAD seguita immediatamente da una STORE (sovrapposizione di WB su registro e MEM su ram).
- 5) In ogni ciclo di clock devono poter essere eseguiti 2 accessi alla memoria (IF, MEM) riferiti a Instruction Memory (IM) e Data Memory (DM) (architettura "Harvard").
- 6) Il clock della CPU è determinato dallo stadio più lento: IM, DM devono essere delle memorie cache (on-chip)!
- 7) I Pipeline Registers trasportano sia dati sia informazioni di controllo (il controller è "distribuito" fra gli stadi della pipeline).

DLX pipelined 9



Esecuzione in pipeline di un'istruzione "ALU"

IF	$IR \leftarrow M[PC]; \quad PC \leftarrow PC+4$
ID	$A \leftarrow RS1; B \leftarrow RS2; PC1 \leftarrow PC;$ $IR1 \leftarrow IR$ (per gli stadi successivi)
EX	$ALUOUTPUT \leftarrow A \text{ op } B$ oppure $ALUOUTPUT \leftarrow A \text{ op } IR_{0..15} \text{ ## } (IR_{15})^{16}$
MEM	$ALUout1 \leftarrow ALUOUTPUT$ ("parcheggio" in attesa di WB)
WB	$RD \leftarrow ALUout1$

DLX pipelined 11

Esecuzione in pipeline di un'istruzione "MEM" (Load o Store)

IF	$IR \leftarrow M[PC]; \quad PC \leftarrow PC+4$
ID	$A \leftarrow RS1; B \leftarrow RS2; PC1 \leftarrow PC; IR1 \leftarrow IR$
EX	$DMAR \leftarrow A \text{ op } IR_{0..15} \text{ ## } (IR_{15})^{16}$ $SMDR \leftarrow B$
MEM	$MDR \leftarrow M[DMAR]$ (LOAD) oppure $M[DMAR] \leftarrow SMDR$ (STORE)
WB	$RD \leftarrow MDR$ (LOAD)

DLX pipelined 12

Esecuzione in pipeline di un'istruzione "BRANCH"

IF	$IR \leftarrow M[PC]; \quad PC \leftarrow PC+4$
ID	$A \leftarrow RS1; B \leftarrow RS2; PC1 \leftarrow PC; IR1 \leftarrow IR$
EX	$BTA \leftarrow PC1 + IR_{0..15} \#\# (IR_{15})^{16}$ $Cond \leftarrow A \text{ op } 0$
MEM	if (Cond) $PC \leftarrow BTA$
WB	-----

BTA = BRANCH TARGET ADDRESS
indirizzo del salto (calcolato in ID)

DLX pipelined 13

Alee nelle Pipeline

Si verifica una situazione di "Alea" (hazard) quando in un determinato ciclo di clock un'istruzione presente in uno stadio della pipeline non può essere eseguita in quel clock.

- **Alee Strutturali** - Una risorsa è condivisa fra due stadi della pipeline: le istruzioni correntemente in tali stadi non possono essere eseguite simultaneamente.
- **Alee di Dato** - Sono dovute a dipendenze fra le istruzioni. Ad esempio un'istruzione tenta di leggere un registro prima che l'istruzione precedente l'abbia scritto (RAW).
- **Alee di Controllo** - Le istruzioni che seguono un branch dipendono dal risultato del branch (taken/not taken).

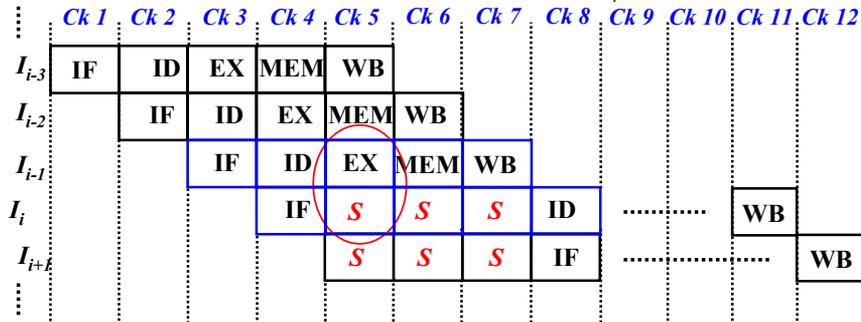


L'istruzione che non può essere eseguita viene bloccata ("stallo della pipeline"), insieme a tutte quelle che la seguono, mentre le istruzioni che la precedono avanzano normalmente (così da rimuovere la causa dell'alea). L'introduzione degli stalli è una possibile soluzione al problema delle alee

DLX pipelined 14

Alee e Stalli

Caso di un'alea di dato RAW nel DLX: $I_{i-1} = \text{ADD } R3, R1, R4$
 $I_i = \text{SUB } R7, R3, R5$



$$T_5 = 8 * CK = (5 + 3) * CK$$

$$T_N = N * 1 * CK$$

$$T_5 = 5 * (1 + 3/5) * CK$$

$$T_N = N * (1 + S) * CK$$

CPI ideale

Stalli per istruzione

CPI effettivo

generalizz

DLX pipelined 15

Impatto degli stalli sullo Speedup

s: non-pipelined
 p: pipelined

$$\text{Speedup} = \frac{T_N(s)}{T_N(p)} = \frac{N * CPI(s) * Ck(s)}{N * CPI(p) * Ck(p)} = \frac{CPI(s)}{1 + S} * \frac{Ck(s)}{Ck(p)}$$

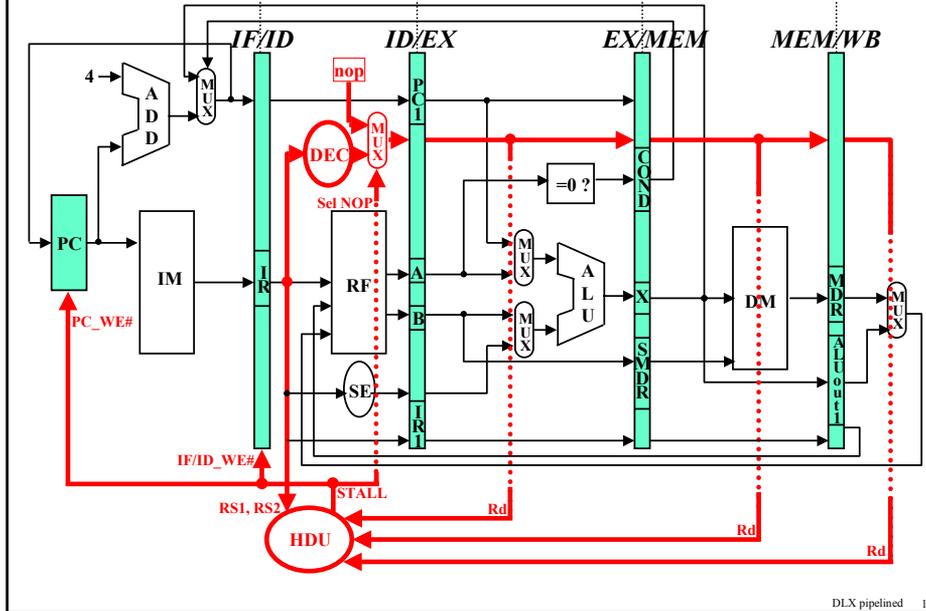
Ipotizzando che i clock siano identici: $\text{Speedup} = \frac{CPI(s)}{1 + S}$

Ipotizzando che nell'implementazione sequenziale tutte le istruzioni impieghino un numero di clock pari al numero degli stadi della pipeline (indicato con K):

$$\text{Speedup} = \frac{K}{1 + S}$$

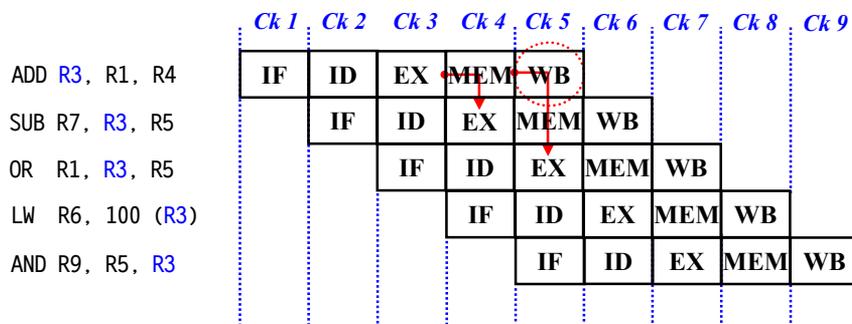
DLX pipelined 16

Circuito per l'introduzione degli stalli: Hazard Detection Unit per alee di dato



DLX pipelined 17

Forwarding



Se il RF ritorna il nuovo valore in caso di RD/WR sullo stesso registro non c'è alea fra la AND e la LW.

Il forwarding consente di eliminare quasi tutte le alee di tipo RAW della pipeline del DLX senza stallare la pipeline.

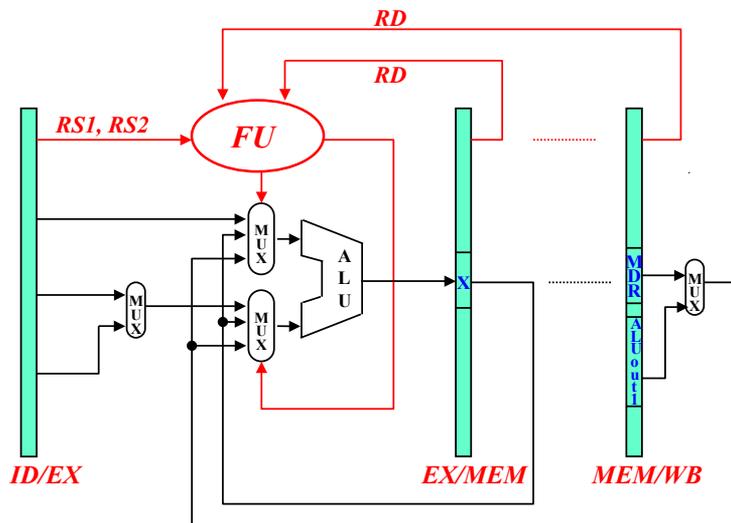
DLX pipelined 18

Esercizio

- Si disegni lo schema logico del Register File del DLX nell'ipotesi che esso debba restituire il nuovo valore in caso di lettura/scrittura simultanea sullo stesso registro.

DLX pipelined 19

Implementazione del Forwarding



DLX pipelined 20

Alea di dato dovuta alle istruzioni di LOAD

LW R1,32(R6)	IF	ID	EX	MEM	WB
ADD R4,R1,R7		IF	ID	EX	MEM
SUB R5,R1,R8			IF	ID	EX
AND R6,R1,R7				IF	ID

N.B. il dato richiesto dalla ADD è presente solo alla fine di MEM. L'Alea non può essere eliminata con il Forwarding!

➔ E' necessario stallare la pipeline

LW R1,32(R6)	IF	ID	EX	MEM	WB	
ADD R4,R1,R7		IF	ID	Stall	EX	MEM
SUB R5,R1,R8			IF	Stall	ID	EX
AND R6,R1,R7				Stall	IF	ID

DLX pipelined 21

Delayed load

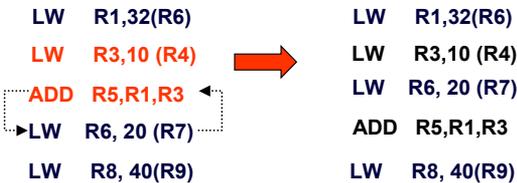
In diverse CPU RISC l'alea associata alla LOAD non è gestita in HW stallando la pipeline ma è gestita via SW dal compilatore (*delayed load*):

Istruzione LOAD

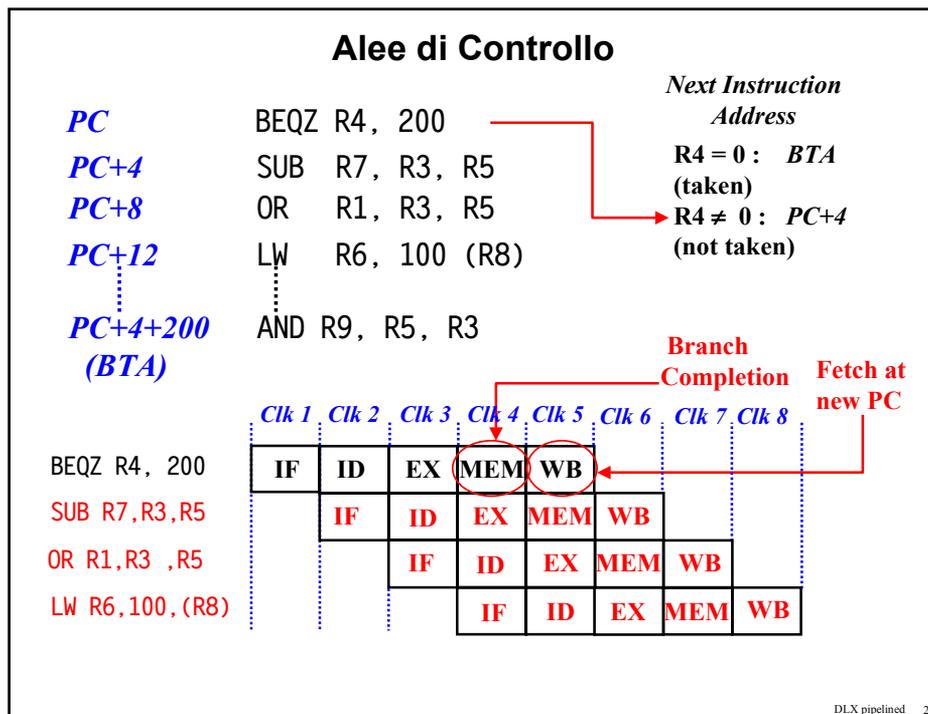
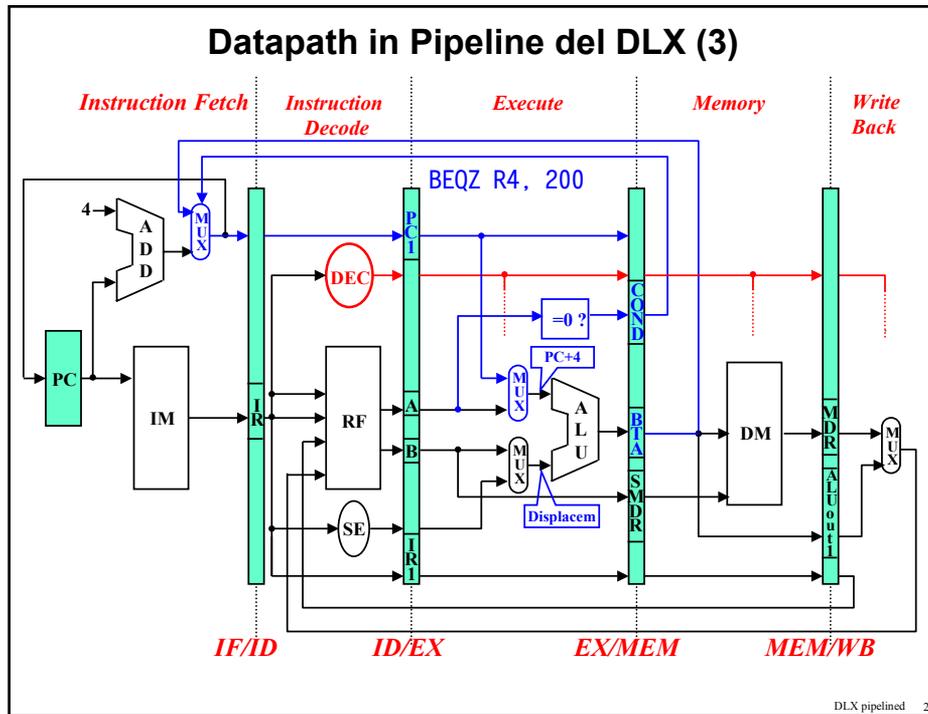


Istruzione Successiva

Il compilatore cerca di riempire il delay-slot con un'istruzione "utile" (caso peggiore: NOP).

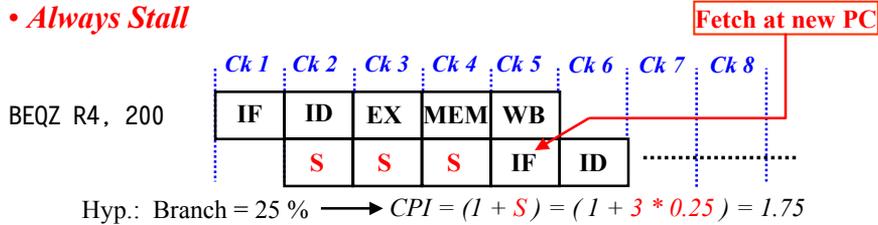


DLX pipelined 22

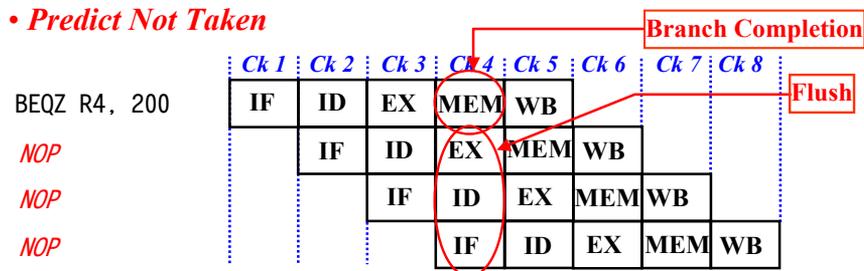


Gestione delle Alee di Controllo (1)

• Always Stall



• Predict Not Taken

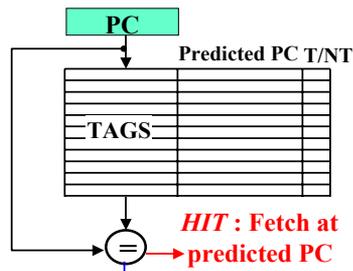


DLX pipelined 25

Gestione delle Alee di Controllo (2)

Hyp.: Branch = 25 % Taken = 65 % $\rightarrow CPI = (1 + S) = (1 + 3 * 0.65 * 0.25) = 1.48$

• Dynamic Prediction: Branch Target Buffer



HIT Predizione Corretta : 0 stalli
 Predizione Errata : (3 + 1) stalli

MISS Predict Not Taken
 Taken: (3+1) stalli
 Not Taken: 1 stallo

H: Hit Rate, E: % di Predizioni Errate
 B: % di Branch
 T/NT: % di Branch Taken/Not Taken

$$S = ((H * E * 4 + (1-H) * (T * 4 + NT * 1)) * B$$

$$S = ((0.9 * 0.1 * 4 + 0.1 * (0.65 * 4 + 0.35 * 1)) * 0.25 = 0.16 \quad CPI = (1 + S) = 1.16$$

DLX pipelined 26

Delayed branch

- Similmente al caso della LOAD, in diverse CPU RISC l'alea associata alle istruzioni di BRANCH è gestita via SW dal compilatore (*delayed branch*):

Istruzione LOAD

delay slot

delay slot

delay slot

Istruzione Successiva

Il compilatore cerca di riempire i delay-slot con istruzioni "utili" (caso peggiore: NOP).

- Introducendo opportune modifiche al datapath è possibile anticipare l'esecuzione del branch fino a ridurre la penalità ad un solo clock (un solo delay slot, come nel caso della LOAD).