

# Contents

<b>Sistemi Operativi Avanzati</b>	<b>1</b>
<b>Hardware Insights</b>	<b>2</b>
Introduzione	2
Scheduling e parallelismo nell'architettura di Von Newman	2
Velocità di computazione	3
Pipeline	3
Pipeline vs Sviluppo	5
Pipeline superscalare	5
Algoritmo di Robert Tomasulo	7
Dipendenza RAW	7
Dipendenza WAW e WAR	7
Architettura di riferimento di Tomasulo	8
Esempi di schemi di esecuzione	9
Organizzazione architetturale dei processori x86 OOO	10
Processori hyper-threaded	10
Gestione degli interrupt	12
Gestione delle eccezioni	12
Attacco Meltdown	12
Insights sui processori x86 e x86-64	13
Codice assembly dell'attacco Meltdown	15
Contromisure per l'attacco Meltdown	16
Insights sui branch	17
Predittori per i salti condizionali	18
Tournament Predictor	19
Predittore per salti indiretti	20
Attacco Spectre	20
Spectre V1	21
Spectre V2	21
Contromisure per gli attacchi Spectre	22
Retpoline - Return trampoline	22
IBRS (Indirect Branch Restricted Speculation)	23
IBPB (Indirect Branch Prediction Barrier)	23
Sanitizzazione	23
Introduzione	23
Come funziona	24
Come funziona - for dummies	24
Loop unrolling	24
Power wall	25
Symmetric Multiprocessors	25
Chip Multi Processor (CMP) o Multicore	26
Symmetric Multi Threading (SMT) o Hyperthreading	26
Non Uniform Memory Access (NUMA):	27
Cache Coherency	27
Protocolli CC Cache Coherency	28
Protocollo MSI (Modified-Shared-Invalid)	29
Protocollo MESI (Modified-Exclusive-Shared-Invalid):	29
<b>Protocollo MOESI (Modified-Owned-Exclusive-Shared-Invalid):</b>	<b>30</b>

## Sistemi Operativi Avanzati

- Versione originale: Matteo Fanfarillo
- Umile formattazione: Simone Festa

# Hardware Insights

## Introduzione

In generale, non è possibile dire che lo stato di un'applicazione sia semplicemente dato dal “puzzle” degli stati dei componenti software sviluppati. Di fatto, quando mandiamo in esercizio tali componenti software, stiamo generando dei cambi di stato al livello hardware che concorrono a determinare qual è lo stato effettivo del nostro sistema. Tra l'altro, alcuni dei cambi di stato al livello hardware possono non essere voluti o specificati dal programmatore. Il fatto che l'esecuzione di moduli software sviluppati a qualsiasi livello impatti sui moduli sottostanti (e quindi sull'hardware), come vedremo, può rappresentare un problema importante per quanto riguarda la sicurezza: talvolta, per ottenere un sistema più sicuro e performante a livello hardware, sarà necessario ristrutturare l'applicazione software.

Ma quali sono le entità che si frappongono tra ciò che viene specificato dal programmatore e ciò che realmente avviene nel sistema?

- Il **compilatore**: il programma compilato è un oggetto molto complesso che può interfacciarsi direttamente con l'hardware; il compilatore può decidere ad esempio di inserire particolari istruzioni macchina secondo uno specifico ordine in modo completamente trasparente al programmatore.
- Le **hardware run-time decisions**: una volta che è stato generato il flusso esatto delle istruzioni macchina da eseguire per mandare in esercizio l'applicazione, l'hardware in realtà può prendere delle decisioni a run-time su come gestire tale flusso. A parità di istruzioni macchina, processori di vendor diversi possono prendere delle decisioni a run-time differenti.
- La **disponibilità** (o l'assenza) **di specifiche features dell'hardware**.

In pratica, si ha una sorta di non-determinismo dell'hardware e, quindi, del software.

**Esempio:** Se implementiamo l'algoritmo del Panificio di Lamport senza l'ausilio di librerie di sistema (e quindi senza l'ausilio di spinlock o semafori), l'algoritmo a un certo punto della sua esecuzione si romperà. Infatti, le macchine moderne (a meno che non siano single core, dove non è assicurata consistenza globale) non garantiscono una visione consistente per tutti i thread di ciò che sta succedendo in memoria. Prima, ad ogni cpu veniva associato un flusso di esecuzione, cosa non vera per i sistemi moderni.

## Scheduling e parallelismo nell'architettura di Von Newman

L'architettura di calcolatore più semplice a cui siamo stati abituati a pensare è quella di **Von Newman**, ed è caratterizzata da:

- Un'unica CPU.
- Un'unica memoria.
- Un unico flusso di controllo (*fetch – execute – store*).
- Transizioni nell'hardware time-separated: le istruzioni devono essere eseguite tutte una alla volta.
- Stato della memoria ben definito all'inizio di ciascuna istruzione, la quale quindi deve poter vedere la memoria in uno stato coerente con l'esecuzione delle istruzioni precedenti.

La maniera moderna di pensare le architetture non è basata sull'idea di seguire il flusso di *esecuzione esattamente così com'è* stato codificato nel programma, bensì è basata sul concetto di **scheduling** (e.g. dell'utilizzo dei componenti hardware) che, alla fine della fiera, porterà a eseguire un flusso di esecuzione equivalente al flusso codificato nel programma (gli stati intermedi possono essere diversi, possono cambiare di run in run, non sono deterministici, ma influenzabili da fattori esterni come la temperatura). In particolare, lo scheduling dovrebbe consentire di eseguire più cose in **parallelo**, anche in contesti con un'unica CPU, un'unica memoria e un unico flusso di controllo. Per quanto riguarda lo scheduling, ne esistono diverse tipologie. A livello **hardware** abbiamo:

- Scheduling delle istruzioni all'interno di un singolo program flow.
- Scheduling delle istruzioni in program flow paralleli (**speculativi** = non so se l'esecuzione sia corretta o meno, ad esempio la branch condition non sempre lo è. Il processore sceglie come propagare l'update, non è imposto. Ho garantita l'equivalenza software, non ciò che avviene dentro).
- Propagazione dei valori tra i componenti hardware del sistema.

A livello software invece abbiamo:

- Scheduling dei thread da assegnare alle CPU / ai CPU-core.
- Scheduling delle attività da eseguire sull'hardware (e.g. gli interrupt).
- Supporti di sincronizzazione tra thread software-based.

Anche per quanto riguarda il **parallelismo**, ne esistono diverse tipologie:

- A livello hardware si parla di **ILP (Instruction Level Parallelism)**, che consiste nell'impiegare le risorse hardware in modo tale da eseguire contemporaneamente istruzioni macchina diverse. (ad esempio posso eseguire al tempo 't' sia A sia B, anche in un unico flusso).
- A livello software, invece, abbiamo il **TLP (Thread Level Parallelism)**, secondo cui un programma può essere pensato come la combinazione di molteplici flussi di esecuzione concorrenti. (Ad esempio ho 3 flussi su 3 processori, singolarmente sarebbero ILP, che cambiano).

## Velocità di computazione

È generalmente correlata alla velocità di un processore (espressa in GHz), anche se in realtà esistono istruzioni che possono richiedere un numero arbitrario di cicli di clock a causa di più possibili fattori:

- Possono essere istruzioni più onerose per loro natura.
- Possono dover richiedere a un certo punto una risorsa hardware tuttora occupata da un'altra istruzione, per cui devono rimanere in attesa.
- Possono esservi delle asimmetrie a livello hardware (e.g. un CPU-core può essere più veloce di un altro).
- I pattern per l'accesso ai dati influiscono a loro volta sulle prestazioni: ad esempio, se un dato viene memorizzato in cache, l'accesso a esso sarà più efficiente e viceversa.

Nel corso base di Sistemi Operativi abbiamo parlato di thread CPU-bound e di thread I/O-bound; introduciamo ora una terza categoria di thread (che, di fatto, è una sottocategoria dei CPU-bound): i **memory-bound**. Essi sono dei thread che utilizzano in maniera intensiva la CPU ma, mentre sono in esecuzione, utilizzano in maniera intensiva anche la memoria. I thread (o comunque i programmi) che presentano questa caratteristica possono rappresentare un problema dal punto di vista prestazionale: come riportato dal seguente grafico, il divario prestazionale tra processore e memoria aumenta sempre di più col tempo; questo fenomeno è detto **memory wall**.

Per evitare che i thread memory-bound sperimentino e causino ad altri thread un crollo delle prestazioni, sono necessari dei meccanismi avanzati ad-hoc al livello dell'hardware.

## Pipeline

È una **tecnica di scheduling e parallelismo hardware-based**. Infatti:

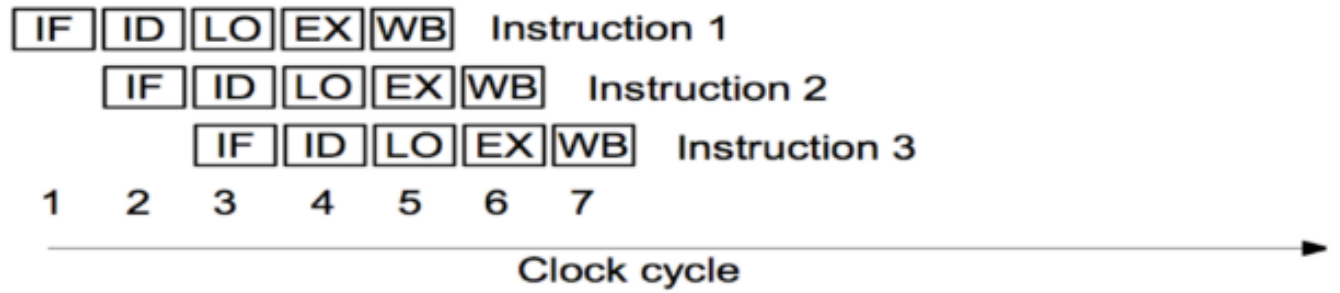
- Non prevede una separazione temporale tra le finestre di esecuzione delle diverse istruzioni: più istruzioni possono essere in esecuzione contemporaneamente (questo è *parallelismo*).
- Le istruzioni che vengono sequenzializzate dal programmatore non sono necessariamente eseguite secondo la stessa sequenza nell'hardware, anche per conseguire la proprietà di parallelismo (questo è *scheduling*).

In ogni caso, ci sono spesso e volentieri coppie di istruzioni ( $i1$ ,  $i2$ ) in cui  $i2$ , per essere eseguita, necessita del risultato ottenuto da  $i1$ . In tal caso, *la causalità deve essere preservata*, per cui  $i1$  e  $i2$  non possono essere schedate in modo del tutto arbitrario. Questo non è altro che un modello **data flow** per l'esecuzione dei programmi.

Ricordiamo che le fasi (stage) delle istruzioni sono:

- **IF (Instruction Fetch)**: caricamento dell'istruzione nel processore (richiede certamente un accesso in memoria).
- **ID (Instruction Decode)**: decodifica dell'istruzione, in cui viene stabilito ciò che deve effettivamente essere fatto per eseguire l'istruzione stessa.
- **LO (Load Operands)**: caricamento degli operandi richiesti per l'esecuzione dell'istruzione (potrebbe richiedere un accesso in memoria).
- **EX (Execute)**: esecuzione vera e propria dell'istruzione.
- **WB (Write Back)**: scrittura dell'output dell'istruzione su un registro o in memoria (potrebbe quindi richiedere un accesso in memoria).

Il fatto che un thread o un'applicazione sia memory-bound o meno dipende proprio dagli stage *LO* e *WB*. Inoltre, per permettere l'esecuzione di più istruzioni in parallelo, è necessario che ciascuno *stage coinvolga una componente hardware differente del processore*, in modo tale da non avere collisioni tra più istruzioni che vengono eseguite contemporaneamente. In particolare, un'astrazione di base della pipeline è quella riportata nella seguente figura.



In questo scenario, idealmente, sarebbe possibile completare un'istruzione per ogni ciclo di clock (anche se poi nella realtà non è esattamente così per i motivi esposti nel paragrafo "Velocità di computazione"). Supponiamo comunque di trovarci nello scenario ideale in cui ciascuno stage viene eseguito in un unico ciclo di clock, e supponiamo di voler fornire  $N$  risultati (1 per ogni istruzione), di avere  $L$  diversi stage per le istruzioni (gli stage sarebbero *IF*, *ID*, *LO*, *EX* e *WB* che devo attraversare) e di avere un ciclo di clock di durata pari a  $T$ .

- Senza pipeline si ha un ritardo pari a  $N \cdot L \cdot T$
- Con la pipeline si ha un ritardo pari a  $(N + L) \cdot T$

Lo speedup è dunque pari a  $(N \cdot L) / (N + L)$ , che tende a  $L$  per  $N$  tendente a infinito. Dal punto di vista delle prestazioni, sarebbe magnifico avere un  $L$  molto grande (ovvero molti stage diversi per le istruzioni, ovvero riesco a parallelizzare molto di più se aumento il numero di componenti parallelizzabili) ma ciò nella realtà non accade. Il caso estremo richiederebbe avere hardware infinito, visto che ogni stage è un pezzo hw! Nella realtà:

- I processori Pentium prevedevano 5 stage.
- I processori i3 / i5 / i7 prevedono 14 stage.
- I processori ARM-11 prevedono 8 stage.

Questa scelta è dovuta alla necessità di preservare la causalità tra le istruzioni. Più istruzioni si prendono in considerazione insieme, più è probabile che tra tali istruzioni ce ne siano alcune (e.g.  $i1$ ,  $i2$ ) legate da una relazione di causalità, e sappiamo che  $i2$ , per essere eseguita, deve attendere che il risultato di  $i1$  sia pronto; in tal caso, più  $L$  è grande, più la pipeline è lunga, più l'attesa di  $i2$  sarà lunga. Lo scenario di cui abbiamo appena parlato è una **data dependency**. Oltre a questa esiste anche la **control dependency**, che si può avere nel momento in cui c'è un salto condizionale il cui esito dipende dagli outcome delle istruzioni precedenti al salto. Per quanto concerne la data dependency tra le istruzioni  $i1$  e  $i2$ , abbiamo che lo stage **LO** di  $i2$  tipicamente non può precedere lo stage **WB** di  $i1$ . Analogamente, per quanto riguarda la control dependency tra le istruzioni  $i1$  e  $i2$  (dove  $i1$  è l'istruzione di salto condizionale), non si conosce l'esito del salto prima della fase **EX\*** di  $i1$ , per cui la fetch di  $i2$  non dovrebbe precedere lo stage **EX\*** di  $i1$ . Tutte queste condizioni possono comportare dei rallentamenti nell'esecuzione dell'applicazione rispetto al caso ideale di pipeline. Per gestire tali condizioni è possibile ricorrere a svariati meccanismi:

- **Stalli software** (compiler driven): possono essere aggiunti all'interno della pipeline con lo scopo di distanziare due istruzioni  $i1$ ,  $i2$  in modo tale che uno stage  $x$  (e.g. LO) di  $i2$  venga eseguito dopo uno stage  $y$  (e.g. WB) di  $i1$ .
- **Rischedulazione software** (compiler driven): se devono essere eseguite tre istruzioni  $i1$ ,  $i2$ ,  $i3$  tali per cui  $i2$  dipende da  $i1$  ma  $i3$  non dipende né da  $i1$  né da  $i2$ , il compilatore può frapporre  $i3$  tra  $i1$  e  $i2$  in modo tale da svolgere lavoro utile mentre  $i2$  è in attesa che si completi uno specifico stage di  $i1$ .
- **Propagazione hardware**: se l'istruzione  $i2$  dipende dall'istruzione  $i1$  e, ad esempio, lo stage LO di  $i2$  consiste nell'acquisire un valore prodotto dallo stage EX di  $i1$ , allora è possibile per  $i1$  propagare il valore a  $i2$  subito dopo la fase EX senza dover attendere il completamento della write-back.
- **Azzardi hardware supported**: contestualmente a un'istruzione di salto condizionale, il processore può provare a *indovinare* se il salto viene preso o meno per poi anticipare la fetch delle istruzioni successive di conseguenza. Dal punto di vista delle performance, una predizione errata equivale a un inserimento di stalli per attendere passivamente l'esito dell'istruzione di salto; di conseguenza, la tecnica degli azzardi è statisticamente conveniente da adottare.

- **Rischedulazione hardware:** è un meccanismo noto anche come **out-of-order pipeline (OOO)**, e prevede che le istruzioni vengano completate non necessariamente nel medesimo ordine con cui sono entrate all'interno della pipeline. In particolare, un'istruzione  $i_k$ , per accedere allo stage  $x$ , non deve necessariamente attendere che tutte le istruzioni a lei precedenti abbiano completato lo stage  $x$ . Si può dunque avere un meccanismo di **superamento** delle istruzioni all'interno della pipeline. Tale tecnica è vantaggiosa poiché permette all'istruzione  $i_k$  di essere completata prima e, quindi, di liberare un posto all'interno della pipeline. Tuttavia, è una tecnica che *richiede la possibilità da parte dell'hardware di ospitare più istruzioni contemporaneamente all'interno dello stesso stage* (altrimenti non sarebbe possibile effettuare il sorpasso): i processori con questa caratteristica sono detti **superscalari**. Inoltre, ovviamente, affinché si possa avere una out-of-order pipeline, *l'istruzione che effettua il sorpasso non deve dipendere dalle istruzioni che vengono sorpassate*. Non si hanno effetti reali sull'**Instruction Set Architecture** finché non si deve “mostrare” l'output.

### Esempio:

Se in pipeline su un thread ho  $A \rightarrow B$ , e nella pipeline B supera A, ma A è oggetto di *trap* (quindi non può fare quello che stava facendo, come dividere per 0, di conseguenza non potrò averla in ISA), cosa accade a B? Vedremo che in OOO si producono valori registrati in maniera “speculativa” che poi butterò, ma non posso eseguire una UNDO.

*Nota:* le istruzioni tra due thread sono indipendenti, dipendono solo se usano info condivise, ma ciò è di interesse al programmatore, non al processore. L'ordine della sorgente (programma) non è detto coincida con l'ordine del compilatore, tuttavia abbiamo la sicurezza che il data flow sia compatibile. A livello hardware, se ho operazione  $x, y, z$  può capitare quindi di avere  $z, x, y$ ; ma ciò è realizzato dinamicamente dall'hardware. Ho sorpasso se non ho dipendenza.

### Pipeline vs Sviluppo

Come abbiamo visto, i programmatori non hanno il diretto controllo del comportamento di un processore (trattasi di microcodice) ma, comunque sia, il modo con cui viene scritto il software può impattare sulle prestazioni effettive della pipeline. A livello ISA vedo il set di istruzioni e risorse usabili, ma non vedo la pipeline. **Esempio:** Esempi di ISA sono *ADD, COMPARE, JUMP*. Consideriamo le seguenti istruzioni C:

- 1)  $a = *++p$
- 2)  $a = *p++$

L'istruzione 1 prevede che prima debba essere incrementato il puntatore  $p$  affinché referenzi la entry successiva e solo dopo si possa accedere alla entry appena referenziata; di conseguenza, l'accesso dipende dall'incremento del puntatore. Al contrario, l'istruzione 2 prevede che prima si debba accedere alla entry attualmente referenziata dal puntatore  $p$  e solo poi si debba incrementare  $p$ ; stavolta, l'accesso non dipende dall'incremento del puntatore. Questo vuol dire che c'è dipendenza. Consideriamo due programmi, di cui il primo prevede l'istruzione 1 all'interno di un loop e il secondo prevede l'istruzione 2 all'interno di un loop. Nel secondo c'è indipendenza. La differenza prestazionale tra i due programmi è abbastanza significativa (può raggiungere tranquillamente il 20/25%).

Per giunta, esistono delle istruzioni macchina che, da sole, hanno degli effetti devastanti sulle prestazioni del programma. Un esempio è **cpuid**, che ha lo scopo di restituire l'id del processore (o del CPU-core o dell'hyperthread) che ha processato l'istruzione stessa. Ma, oltre a questo, effettua anche lo **squash** della pipeline: in particolare, nel momento in cui **cpuid** viene realmente eseguita, la pipeline viene svuotata e le altre istruzioni al suo interno vengono buttate. Questo non deve necessariamente rappresentare uno svantaggio: avere istruzioni pendenti all'interno della pipeline significa dire che tali istruzioni possono essere schedate dinamicamente dall'hardware secondo regole non meglio identificate, e ciò può impattare non solo sulla correttezza, ma anche sulla sicurezza del sistema. Più precisamente, quello di flushare la pipeline è un concetto legato al termine “**serializzazione**”. Di fatto, **cpuid** è detta **istruzione serializzante**, poiché riporta la pipeline a lavorare secondo uno schema sequenziale. Non solo: **cpuid**, come tutte le istruzioni serializzanti, garantisce che qualunque modifica apportata a flag, registri e memoria da parte delle istruzioni precedenti sia completata (finalizzata) **prima** che una qualsiasi istruzione a lei successiva venga fetchata ed eseguita. *Ciò implica anche che cpuid non può superare alcuna istruzione davanti a lei nella pipeline.* (e.g: Perché le istruzioni successive devono ancora essere fetchate ed eseguite, quindi come potrei superarle? Ciò che viene dopo questa istruzione serializzante è come se andasse in stallo finché non completo questa istruzione serializzante, e garantisce anche che le istruzioni precedenti vengano viste da quelle successive.)

Se avessi SOLO istruzioni serializzanti, il sistema sarebbe molto più lento.

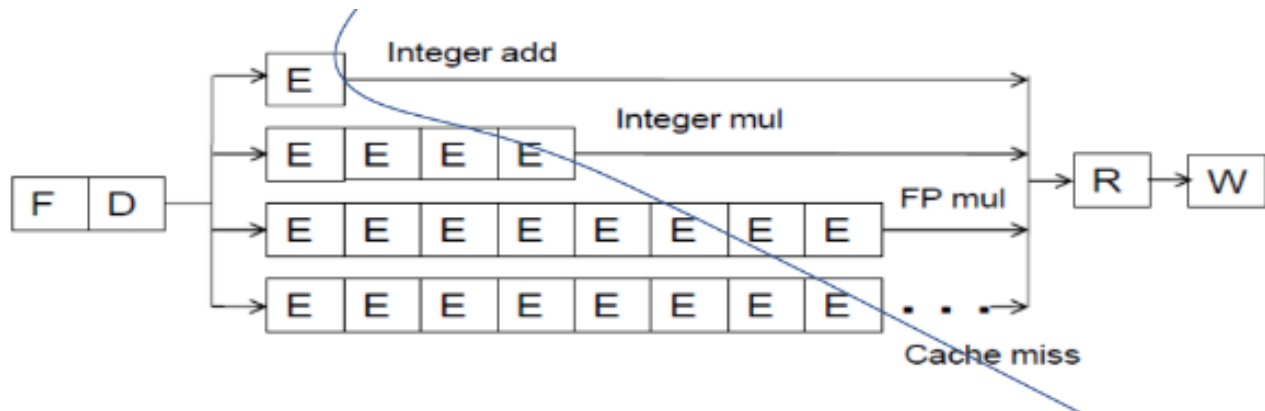
### Pipeline superscalare

È un insieme di molteplici pipeline che operano simultaneamente all'interno del processore. La si può ottenere aggiungendo *ridondanza alle risorse hardware* in modo tale che più componenti distinti siano adibite a uno stesso stage della pipeline.

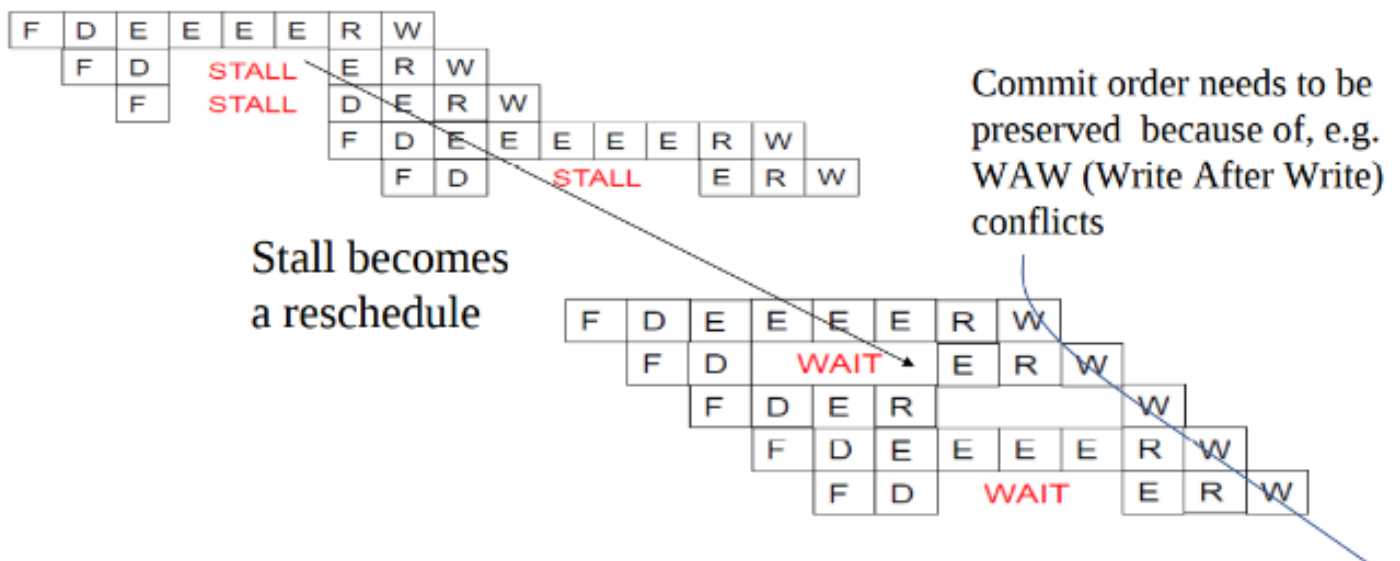
Come detto in precedenza, questa possibilità permette anche di adottare il modello OOO, la cui idea di base consiste in:

- Effettuare il **commit** (o **retire** o **finalizzazione**) delle istruzioni esattamente nel medesimo ordine in cui sono *entrate nella pipeline*, indipendentemente dagli eventuali sorpassi avvenuti all'interno della pipeline. Ciò significa che le scritture dei risultati delle istruzioni in memoria, su un registro qualunque o su un registro di stato *devono avvenire esattamente nell'ordine prestabilito*.
- **Processare le istruzioni indipendenti** (sia sui dati che sulle risorse hardware) **il prima possibile**, dove un'istruzione indipendente sulle risorse hardware è un'istruzione che non ha bisogno di attendere che un qualche componente si liberi per processare un certo stage  $x$ .

L'immagine riportata di seguito mostra molto bene come lo stage EX delle istruzioni possa avere una durata variabile in termini di numero di cicli di clock, ed è a partire da tale presupposto che risulta utile avere un out-of-order pipeline.



La seguente altra figura mostra invece la differenza prestazionale che si può avere tra il modello non OOO e il modello OOO:



: Con una pipeline che ha ridondanza hardware (come i pc moderni), possiamo ad esempio parlare di “*Pipeline parallele a X canali*”, che è come avere una autostrada a X corsie. Con Out-Of-Order Pipeline si evitano “grosse latenze”, poiché se “libero una corsia perchè sono veloce” anche gli altri andranno più veloci (liberando risorse utili ad altri). Lo stallo impatta su tutti invece. Nella O.O.O, se posso andare avanti, lo faccio, non mi importa che devo aspettare il commit, se posso eseguire qualcosa lo eseguo. Per questo parliamo di WAIT e non più di STALL. Nella WAIT, appena ho un risultato (outcome di una “corsia”) lo fornisco a chi lo necessita, non mi serve aspettare che si liberi tutta la corsia.

Ora, poiché tra l’**emission** (= iniezione all’interno della pipeline, normalmente di più di una istruzione) e il retire (= **commit**) delle istruzioni si ha una fase di esecuzione in cui le istruzioni stesse possono sorpassarsi a vicenda, si va incontro

al cosiddetto problema delle **eccezioni imprecise (OFFENDING)**: supponiamo di avere un'istruzione  $i2$  che ha superato un'istruzione  $i1$ , e assumiamo che  $i1$ , durante lo stage EX, generi un'eccezione (*e.g.*: perché magari si tratta di una divisione per 0); a questo punto, anche se il risultato di  $i2$  non è stato ancora committato e, quindi, esposto a livello di ISA,  $i2$  può aver comunque toccato delle risorse non esposte a livello di ISA ed effettuato dunque dei cambi di stato (in particolare cambi di **stato micro-architetturale**). Quindi con offending instructions non esponiamo su ISA, ma potrei cambiare lo stato hardware, il che è usabile da malintenzionati. Questo perché, lavorare in un contesto Out-Of-Order permette il superamento di istruzioni (**SPECULAZIONE**) che lascia delle tracce. Di conseguenza, l'eccezione sollevata da  $i1$  vede uno stato interno dell'hardware che ingloba anche delle attività relative a  $i2$ . Un problema analogo lo si ha anche a parti invertite: se l'istruzione che genera l'eccezione è  $i2$ , l'eccezione vedrà uno stato interno dell'hardware che non comprende il risultato prodotto da  $i1$ . Questo è un problema dello stato e delle informazioni all'interno del processore (*ma non esposte nell'ISA*). Peggio ancora:  $i2$  potrebbe sollevare un'eccezione che in realtà, secondo il program flow, non sarebbe dovuta mai esistere, magari perché  $i1$  è a sua volta un'istruzione che solleva un'eccezione o un'istruzione di salto. In realtà, si possono avere eccezioni imprecise anche in assenza di sorpassi: se l'istruzione  $i2$  viene dopo l'istruzione  $i1$  che genera un'eccezione,  $i2$  può aver comunque attraversato degli stage e, quindi, può aver modificato lo stato interno dell'hardware. Come vedremo, tutto questo rappresenta un grave problema per la sicurezza dei sistemi: **Meltdown** è solo il primo di una valanga di attacchi che hanno sfruttato tale vulnerabilità.

## Algoritmo di Robert Tomasulo

Secondo Robert Tomasulo, in uno scenario di utilizzo di una pipeline speculativa out-of-order (dove speculativa = caratterizzata dall'esecuzione di istruzioni che potrebbero servire solo in un secondo momento), se consideriamo due istruzioni A, B tali che  $A \rightarrow B$  nell'ordine di programma, dobbiamo stare attenti nell'evitare i seguenti tre tipi di azzardo:

- **RAW (Read After Write)**: B deve leggere un dato R necessariamente dopo che A lo ha aggiornato.
- **WAW (Write After Write)**: B deve scrivere su un dato R necessariamente dopo che A lo ha aggiornato.
- **WAR (Write After Read)**: B deve scrivere su un dato R necessariamente dopo che A ne ha letto il valore precedente (altrimenti vorrebbe dire che A è in grado di leggere dal futuro).

### Dipendenza RAW

Qui è necessario bloccare l'istruzione B e tenere traccia di quando il dato che deve essere letto da B sarà disponibile (disponibile non vuole dire necessariamente committato).

### Dipendenza WAW e WAR

Qui è necessario adottare una tecnica nota come **register renaming**, secondo cui al programmatore vengono esposti dei registri logici, e ciascuno di questi registri logici (che sono di fatto dei multi-registri) ingloba un insieme di registri fisici non visibili al livello dell'ISA. Quindi noi non scriviamo MAI sul registro esposto in ISA, ma su delle 'copie numerate' o alias, cioè registri multivalore. Ci ritroviamo dunque nella seguente situazione:

### Registro logico R



Qui i vari registri fisici rappresentano diverse versioni del medesimo registro logico  $R$ . Nel momento in cui all'interno della pipeline entra un'istruzione che vuole scrivere sul registro logico  $R$ , si considera il primo registro fisico  $R_k$  all'interno del quale viene effettivamente memorizzato il valore (quindi scrivo sempre sul primo tag libero, e leggo dall'ultimo tag in pipeline. Se A scrive in TAG 0, e B scrive in TAG 1, leggo da TAG1). Tipicamente, per la selezione del registro fisico, si segue un approccio round-robin; se però a un certo punto tutti i registri fisici di  $R$  sono stati sovrascritti e nessuna delle istruzioni che ha eseguito la scrittura è andata in commit, per un'eventuale altra scrittura su  $R$  bisognerà attendere.

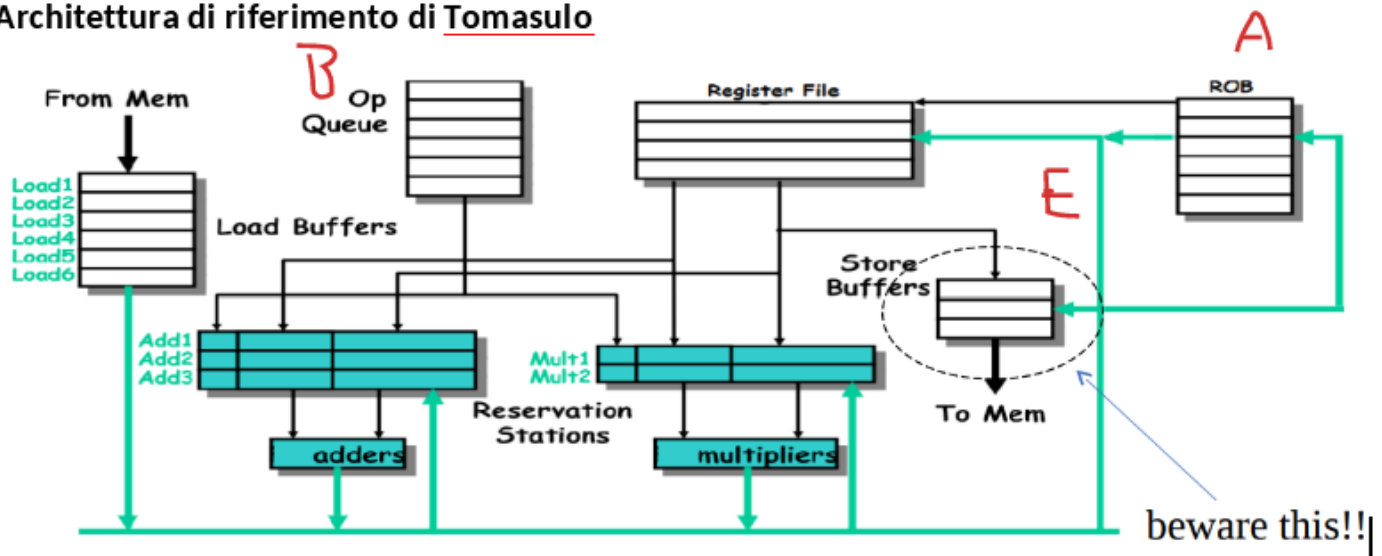
Posso superare solo dopo essere entrato nella pipeline (nelle condizioni discusse nelle pagine precedenti), ma non posso superare prima di entrare in pipeline.

In pratica, un *renamed register* materializza il concetto di speculatività di una pipeline: quando vengono effettuate delle write, vengono scritti dei valori che non necessariamente verranno considerati come validi, ma che comunque potranno essere letti dalle istruzioni successive (che sono state introdotte nella pipeline speculativamente).

Relativamente al registro *R*, vengono memorizzati anche dei metadati di gestione di *R* che indicano qual è il tag contenente la versione committata; la versione committata è la versione del valore di *R* scritto dall'ultima istruzione andata in commit che ha toccato proprio *R*. In questi registri abbiamo il "futuro" che avverrà, perchè non li ho ancora committati.

## Architettura di riferimento di Tomasulo

### Architettura di riferimento di Tomasulo



- **A):** I metadati associati alle istruzioni fetchate vengono memorizzati all'interno del **ROB (Re-Order Buffer)**, che ci aiuta a scrivere le istruzioni nello store buffer, in quanto posso andare in memoria solo se sono committed, ma io non posso aspettare sempre questa fase prima di utilizzare i dati. Tali metadati possono essere:
    - L'ordine con cui le istruzioni dovranno andare in commit.
    - Che cosa dovrebbero fare le istruzioni nella loro esecuzione speculativa.
    - Qual è l'alias (l'istanza fisica) dei registri che dovranno essere usati da ciascuna istruzione; ad esempio, se un'istruzione A deve scrivere su un certo registro R e un'istruzione B successiva deve leggere dallo stesso registro R con una dipendenza RAW, è chiaro che A e B dovranno utilizzare il medesimo alias.
  - **B):** Le operazioni vere e proprie da eseguire (quindi gli OP code delle istruzioni) vengono memorizzate all'interno dell'**OP queue** e, in base alla loro tipologia, verranno date in input a un particolare componente di processamento (add, mult e così via). Nel momento in cui un'istruzione è pronta per essere processata, devono essere recuperati i metadati associati a essa e, in particolare, gli alias dei registri da usare. A tale scopo, c'è un bus comune (detto **Common Data Bus – CDB**) che mette in comunicazione i componenti di processamento col ROB. Se un alias non è pronto per essere utilizzato (perché magari bisogna attendere che venga sovrascritto da un'istruzione precedente), l'istruzione viene posta in attesa all'interno del relativo componente di processamento. I componenti di processamento sono detti anche **reservation station**. *Cosa è Reservation Station?* Per ogni componente in grado di eseguire uno specifico calcolo, si ha una coda/buffer/corsia dove posso mettere varie istruzioni, e le operazioni identificate dai registri usati (sorgenti). Non si ha sorpasso per utilizzare tali componenti in queste code, al massimo se ho due istruzioni dipendenti cerco di metterle nella stessa reservation station.
  - Si fa uso di una **cache** per leggere in modo più efficiente i dati provenienti dalla memoria.
  - Nel momento in cui viene **committata** un'istruzione, l'eventuale alias aggiornato da tale istruzione viene installato all'interno del **register file** come valore valido, ovvero come **valore esposto all'interno dell'ISA**.
  - **E):** Se l'istruzione committata prevede una scrittura in memoria, il valore di output, prima di essere riportato in memoria, viene inserito nello **store buffer** in modo tale da velocizzare il completamento dell'istruzione. Tuttavia, ciò implica che il valore non sarà in memoria per un po' di tempo, il che può rappresentare un problema dal punto di vista della consistenza.
- Nota:* Supponiamo che un'istruzione sia in fase di ritiro, quindi è totalmente conclusa. In questo istante, il dato



dovrebbe passare da Store Buffer alla memoria, ma in realtà non è istantaneo, passa un piccolo lasso di tempo. Ciò crea problemi con  $> 1$  thread, perchè se pesco un dato in memoria potrei non vedere alcuni dati. Per questo l'algoritmo della “*pasticceria*” non funziona nei processori moderni, suppone che tale tempo sia nullo. Esistono però delle istruzioni per controllare lo stato dello Store Buffer. Questo approccio è l'unico per lavorare con O.O.O, ovvero non esistono altre tecniche per affrontare le O.O.O pipeline.

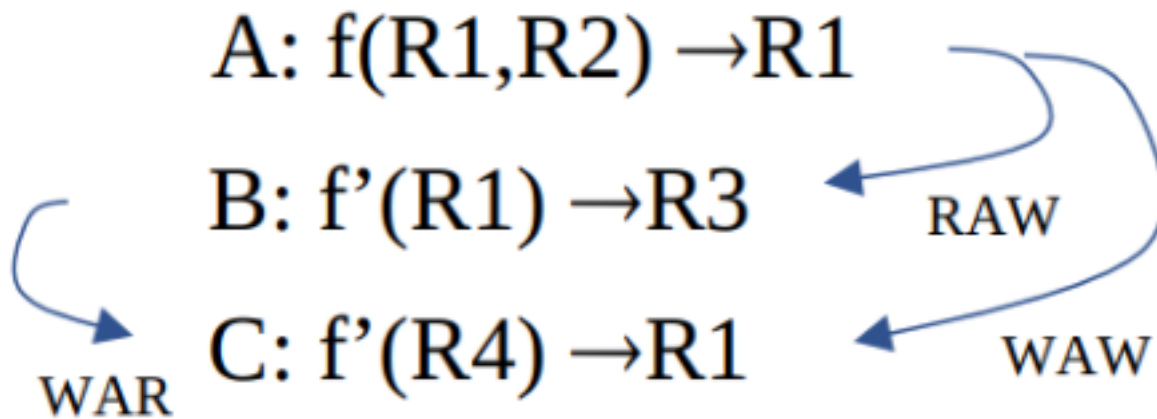
### Esempi di schemi di esecuzione

**Esempio 1** Supponiamo di avere tre istruzioni  $A, B, C$  tra cui  $B$  e  $C$  hanno una stessa latenza  $d$ , mentre  $A$  ha una latenza  $d' > d$ .

Assumiamo inoltre che:

- $A$  esegua l'operazione  $f(R1, R2)$  e riporti l'output sul registro  $R1$ . (in particolare, scrive su un alias di  $R1$ ).
- $B$  esegua l'operazione  $f'(R1)$  e riporti l'output sul registro  $R3$ . (legge da STESSO alias di  $R1$ )
- $C$  esegua l'operazione  $f'(R4)$  e riporti l'output sul registro  $R1$ . (scrive su UN ALTRO alias di  $R1$ ).

Ci ritroviamo dunque nel seguente scenario:



È abbastanza evidente che  $B$  debba leggere **lo stesso alias** scritto da  $A$ , mentre  $C$  potrà riportare il valore di output su un **alias differente**. In tal modo è possibile **processare  $C$  parallelamente ad  $A$**  :

In fondo, anche se  $C$  genera il suo output prima di  $A$ , le dipendenze che legano le tre istruzioni non vengono violate. Infatti,  $B$  sarà comunque in grado di leggere dall'alias sovrascritto da  $A$  e, ovviamente, rimarrà comunque possibile mandare in commit  $C$  solo dopo il completamento di  $A$  e  $B$ . Il vantaggio che porta questo approccio è la riduzione del tempo necessario per il completamento di  $C$ .

**Esempio 2** Vediamo con un secondo esempio riportato qui di seguito come viene associata una entry del ROB a ciascun alias differente:

<b>Before:</b> add <span style="color: blue;">r3</span> , r3, 4 add <span style="color: green;">r4</span> , r7, r3 add <span style="color: magenta;">r3</span> , r2, r7	<b>after</b>	add <span style="color: blue;">rob6</span> , r3, 4 add <span style="color: green;">rob7</span> , r7, <span style="color: blue;">rob6</span> add <span style="color: magenta;">rob8</span> , r2, r7
---	--------------	--

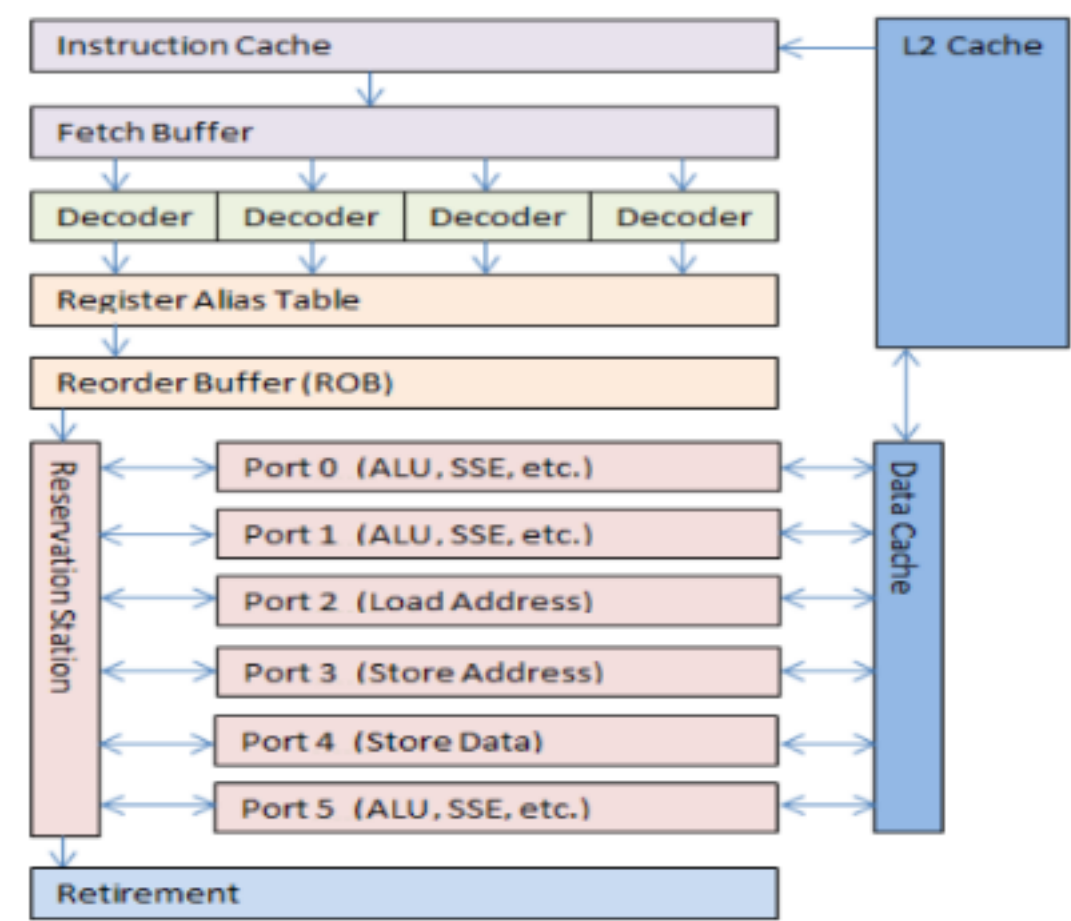
E' importante ricordare che, anche se è possibile eseguire in parallelo, non posso pensare di aumentare all'infinito le prestazioni, perchè ad un certo punto eccedo il numero di istruzioni che posso gestire.

Questo porta ad un sistema scarico, ovvero non occupo corsie perchè ci sono caricamenti di elementi dalla memoria, oppure eseguo una predizione/azzardo sbagliata che mi porta a svuotare etc. . .

E' come fare un'autostrada a 20 corsie: è molto difficile saturarla il “giusto”.

Che soluzione è possibile adottare? Introduciamo i **Processori HYPERTHREAD**.

## Organizzazione architetturale dei processori x86 OOO



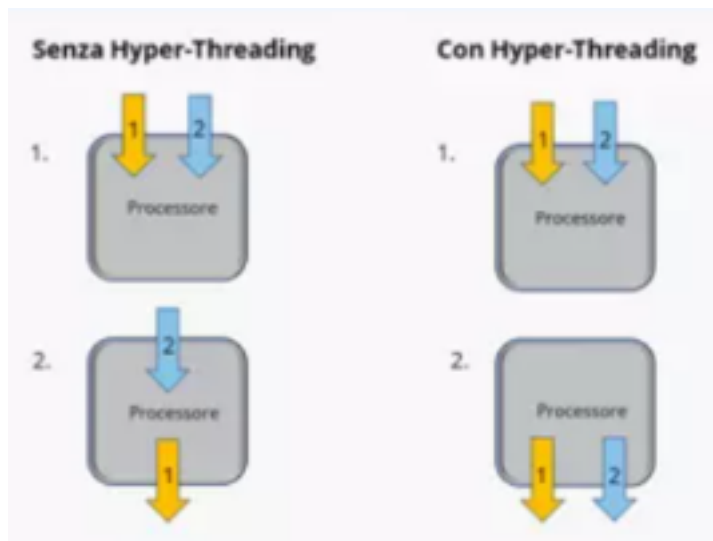
### Processori hyper-threaded

Come abbiamo osservato all'inizio della trattazione, le architetture moderne sono soggette al **memory wall**. In particolare, si ha una latenza di esecuzione non solo quando si devono *recuperare i dati dalla memoria*, ma anche quando devono essere estratte le istruzioni stesse. Dunque, si può anche avere una pipeline molto efficiente con un **ILP** elevato, ma il processore non arriva mai a processare in parallelo tutte le istruzioni possibili a *causa dei ritardi causati dalla memoria*, per cui la potenza e le risorse del processore risultano sprecate. Per ovviare all'inconveniente, si può assegnare una *stessa reservation station a più program flow diversi*. Da qui nascono i processori **hyper-threaded**, che hanno un *unico core fisico (i.e. un'unica information station)* che permette di eseguire le reali operazioni dettate dalle istruzioni; d'altra parte, la porzione di processore che memorizza le istruzioni da fetchare, decodifica le istruzioni per uno specifico flusso (Decode) e tiene traccia dei registri logici dell'ISA (mediante Register Alias Table) viene replicata e viene definita **hyperthread**. Una CPU può contenere uno o più core. Un core è un'unità all'interno della CPU che realizza l'effettiva esecuzione. E' un "oggetto engine". Con l'Hyper-Threading, un microprocessore fisico si comporta come se avesse due core logici, ma virtuali. In questo modo si consente a un unico processore l'esecuzione di più thread contemporaneamente. Questo procedimento aumenta le prestazioni della CPU e migliora l'utilizzo del computer.

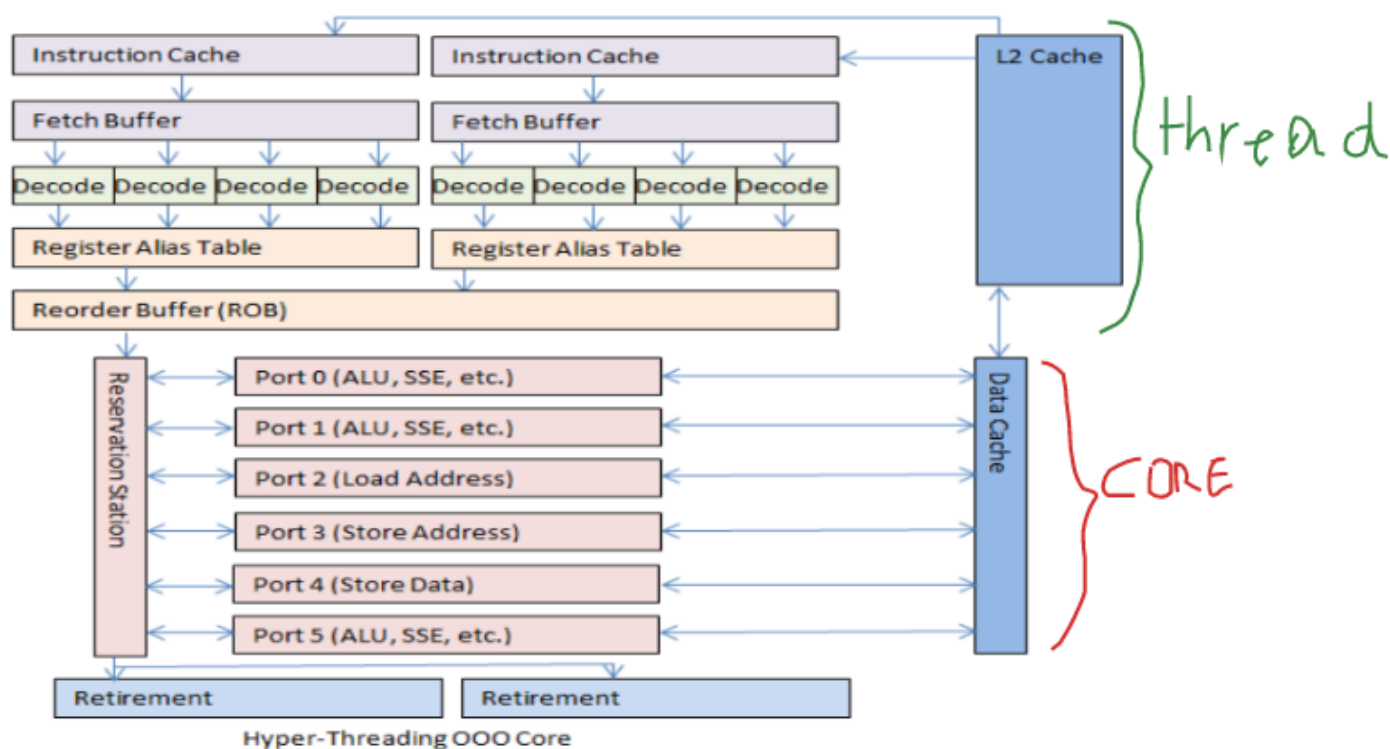
---

**Esempio:** ho una bocca (core) e due mani con cui prendo il cibo (thread). Per migliorare la prestazione devo puntare ad aumentare i thread (più mani = più cibo in entrata). Se aumentassi il numero di bocche, non migliorerei, perchè sempre quel cibo ingerisco, anche se finisco prima devo aspettare la prossima forchettata.

---



In tal modo, è possibile eseguire più **workflow in parallelo su un unico core fisico**, e l'architettura risultante è riportata nella pagina seguente:



Tale architettura risulta molto vantaggiosa per la maggior parte dei workload, mentre risulta un po' più avversa nel caso in cui tutti i thread che girano sullo **stesso core sono CPU-bound** ma non memory-bound. Infatti, in quest'ultimo caso è più probabile che si verifichi un conflitto tra i thread in esecuzione nell'utilizzo delle risorse della reservation station. Un processore hyper-threaded viene tipicamente marcato con l'indicazione su quanti *core* (= "motori") e quanti *hyperthread* (= "thread fisici") possiede.

Esempio:

Un processore a due core e quattro thread fisici viene marcato come **2C/4T**.

In un contesto senza Hyperthread avremmo 2 core e 2 thread (1 per ogni core).

In un contesto con Hyperthread possiamo avere 2 core e 4 thread (2 thread per ogni core).

Se ho un caso 2C/4T, vuol dire che il mio programma non può generare più di 4 thread?

No, vuol dire che ogni core può gestire al massimo due thread con un proprio set di registri aventi alias (quindi due workflow in parallelo su unico core), quindi in totale posso gestire al massimo quattro thread/workflow in contemporanea. Ci sarà lo scheduling che realizzerà il parallelismo, ma comunque ne verranno gestiti quattro insieme. I flussi sui thread possono essere speculativi.

## Gestione degli interrupt

Ne sono esempi: muovere mouse, premere sulla tastiera.

L'interrupt è un segnale *proveniente da un certo dispositivo hardware* che indica la necessità di **cambiare il flusso di esecuzione** (e.g. andando a eseguire un handler). Poiché è buona norma processare gli interrupt il prima possibile, quando ne occorre uno, si attende solo che una delle istruzioni correntemente in pipeline vada in commit; dopodiché si effettua lo squash (ovvero si svuota) la pipeline e si iniziano a fetchare le istruzioni dell'handler dell'interrupt. Ciò può portare a un rallentamento dell'esecuzione, ma non vale la pena memorizzare da qualche parte lo stato di esecuzione delle istruzioni che vengono buttate (richiederebbe hardware aggiuntivo, inoltre non ho certezza che dopo l'handler, ritornando al flusso di esecuzione, io parta dall'istruzione subito dopo); tra l'altro, anche mentre viene eseguito il gestore di un certo interrupt può subentrare un ulteriore interrupt, e questo può avvenire iterativamente un numero arbitrario di volte.

### Attenzione:

Buttare delle istruzioni dalla pipeline implica prevenire i loro effetti sull'ISA ma, se esse hanno sporcato lo stato micro-architetturale, quest'ultimo non può essere ripristinato: rimangono comunque gli effetti dell'esecuzione speculativa delle istruzioni che sono state buttate. A livello hardware non ho meccanismi di gestione priorità o trap, operazioni offending come una divisione per 0 non vengono svolte dal processore, ma passano ad un handler.

## Gestione delle eccezioni

Esse avvengono durante esecuzione di un programma, a livello software. Le eccezioni risultano più complesse da gestire rispetto agli interrupt poiché vengono sollevate dall'esecuzione di particolari istruzioni. Di fatto, un'istruzione A (cosiddetta **offending**) che solleva un'eccezione  $e_A$  può essere eseguita speculativamente all'interno della pipeline e, di conseguenza, potrebbe non esistere nel flusso di esecuzione definito dal programmatore. Ad esempio, poco prima dell'istruzione A potrebbe esserci un'istruzione B che solleva a sua volta un'eccezione  $e_B$ : in tal caso sarà  $e_B$  a esistere realmente e non  $e_A$ . Una Istruzione è offending se vuole usare qualcosa che non è usabile. Non genera una trap, perchè non so se arrivo in retire. Se l'istruzione successiva in fondo non viene committata, cancello tutto. A valle di queste considerazioni, un'eccezione viene presa in carico solo nel momento in cui l'istruzione che l'ha generata va in retire, anche se ci si può accorgere prima che l'istruzione sia offending. (Vedi sotto, etichettamento offending).

Ma quali sono gli stage in cui un'istruzione può rivelarsi offending?

- **Instruction Fetch / Memory stages**  
(MEM stage = stage in cui avviene l'accesso agli operandi): qui si può avere un page fault (ovvero un tentato accesso a un indirizzo logico di memoria che attualmente non ha un corrispettivo fisico), un accesso in memoria disallineato o una violazione delle protezioni applicate a una pagina di memoria (e.g. si tenta di accedere in scrittura a una locazione read-only).
- **Instruction Decode stage**: qui può emergere che l'istruzione in esercizio sia illegale.
- **Execution stage**: qui può essere sollevata un'eccezione aritmetica (e.g. divisione per zero).
- **Write-Back stage**: qui non possono essere sollevate eccezioni.

**Per etichettare un'istruzione come offending**, si imposta a 1 un apposito bit dei metadati. Quando tale istruzione va in retire, se il bit vale 1 allora il commit non viene effettuato, bensì viene attivato il gestore di eccezioni opportuno. A tal punto tutte le istruzioni successive a quella offending diventano *phantom* (fantasma).

Attenzione: Con le eccezioni si verifica lo stesso problema riscontrato con gli interrupt. In particolare, quando un'istruzione offending va in retire, si hanno altre istruzioni all'interno della pipeline che hanno già modificato lo stato micro-architetturale e, dunque, hanno lasciato tracce di informazione. Questa problematica lascia spazio al cosiddetto attacco **Meltdown**. Sto propagando "attività illecite" nella pipeline.

## Attacco Meltdown

Sappiamo bene che *ciascun processo ha il proprio address space* suddiviso in zone di memoria dedicate all'esecuzione *user mode* e zone di memoria dedicate all'esecuzione *kernel mode*. L'attacco ha come scopo ultimo quello di accedere a un'area di memoria situata nel kernel anche se non si hanno i permessi, magari per andare a leggere delle informazioni sensibili di un utente differente del sistema. Andando nel dettaglio, l'attacco viene eseguito nella maniera spiegata qui di seguito. Anzitutto si definisce un array A nella memoria user space e si svuota la cache tramite l'istruzione **cflush**. Dopodiché, mediante una **mov**, si preleva un **particolare byte**  $x$  (e.g. address) situato nel kernel e si memorizza il byte in un registro; naturalmente questa **mov** è un'istruzione offending. Si accede alla entry con **spiazzamento**  $x$  rispetto all'indirizzo base dell'array A (e.g. caricandone il contenuto in un registro) in modo da farla salire in cache: tale aggiornamento della cache rappresenta proprio il *side effect* lasciato sullo stato micro-architetturale da parte dell'istruzione successiva a quella



- **GPR:** registri general purpose.
- **SSE & SSE2:** registri vettoriali, che consentono a singole istruzioni di fare più cose in parallelo.
- **x87:** floating-point stack registers.
- **RIP:** program counter.

Vediamo alcune istruzioni fondamentali per il trasferimento di dati nell'architettura x86-64:

Istruz. (a 64 bit)	Variante a 32 bit	Variante a 16 bit	Variante a 8 bit	Descrizione	Esempio*
mov	movl	movw	movb	Trasferisce un valore su un registro destinazione.	movw \$5, %ax
push	pushl	pushw	\	Aggiunge un elemento sullo stack.	pushl %ebx
pop	popl	popw	\	Elimina un elemento dallo stack e lo salva su un registro.	popl %ebx

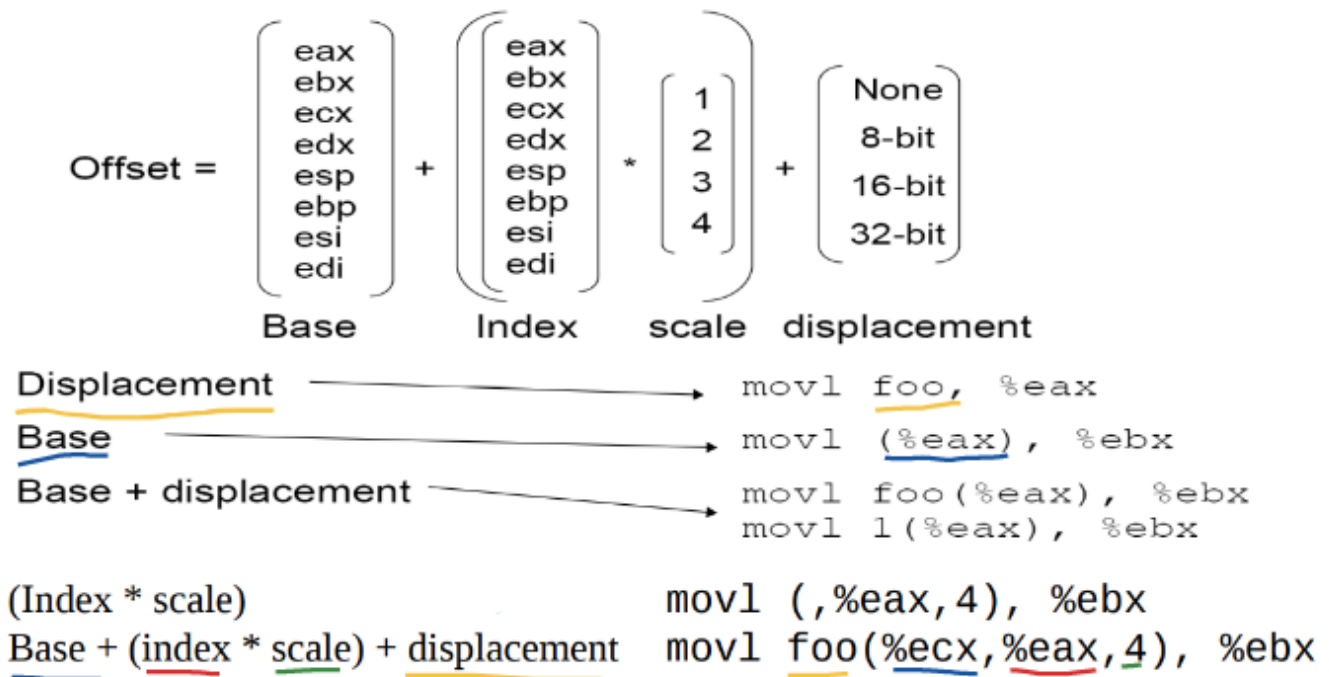
\*Negli esempi specificati in tabella si è adottata la sintassi AT&T che, a differenza di quella Intel, prevede che la sorgente venga posta prima della destinazione.

La parte più complessa però sta nello specificare la sorgente e/o la destinazione **in memoria logica** (e non su un registro). La figura riportata in seguito mostra il meccanismo utilizzato. In *ISA* non ho i nomi delle variabili, solo quelli dei registri, per arrivare ad una variabile di memoria devo usare l'indirizzo di tale variabile in memoria.

Nota: Le componenti *base*, *index*, *scale* sono sempre contenute nelle parentesi, ad esempio *movl (... , ... , ...)*

- **Displacement:** può ad esempio essere un indirizzo assoluto noto a tempo di compilazione. Possiamo vederlo come il punto dell'address space che specifica la sorgente (es: *foo* è *displacement diretto*)
- **Base:** può essere relativo al valore di un pointer.
- **\*\*Index\*scale\*\*:** può rappresentare uno spiazzamento rispetto all'indirizzo base; ad esempio, quando si itera su un array di interi, *scale* vale sempre 4 mentre *index* viene incrementato di 1 a ogni iterazione.

(nb: in *%ebx* è dove carico il tutto).



Vediamo ora le istruzioni logiche e aritmetiche:

Istruzione (a 64 bit)	Variante a 32 bit	Variante a 16 bit	Variante a 8 bit	Descrizione
and	andl	andw	andb	dest = source & dest
or	orl	orw	orb	dest = source    dest
xor	xorl	xorw	xorb	dest = source ^ dest
not	notl	notw	notb	dest = ^dest
sal	sall	salw	salb	dest = dest << source
sar	sarl	sarw	sarb	dest = dest >> source
add	addl	addw	addb	dest = source + dest
sub	subl	subw	subb	dest = dest - source
inc	incl	incw	incb	dest = dest + 1
dec	decl	decw	decb	dest = dest - 1
neg	negl	negw	negb	dest = ^dest
cmp	cmpl	cmpw	cmpb	Confronta due valori; se essi sono uguali, imposta un determinato registro di stato.

## Codice assembly dell'attacco Meltdown

Nella pagina seguente viene mostrato il codice assembly che permette di effettuare l'attacco Meltdown. La sintassi utilizzata è quella di Intel. Chiaramente il codice può essere incapsulato all'interno di un programma C.

```

; rcx = kernel address
; rbx = probe array (array A)
retry:
mov al, byte [rcx] //istr. offending. Prendo byte[rcx] e lo metto in 'al', ultimo byte del reg. 'eax'
shl rax, 0xc      //shlto rax di 0xc = 12, cioè 12 bit a 0 (verso sx shlto la parte meno significativa) lo

```



```
jz retry          //se il valore letto nel kernel è nullo, cioè rax=0 riprovare.
mov rbx, qword [rbx + rax] //qui si effettua l'accesso al probe array con spiazamento pari a rax
```

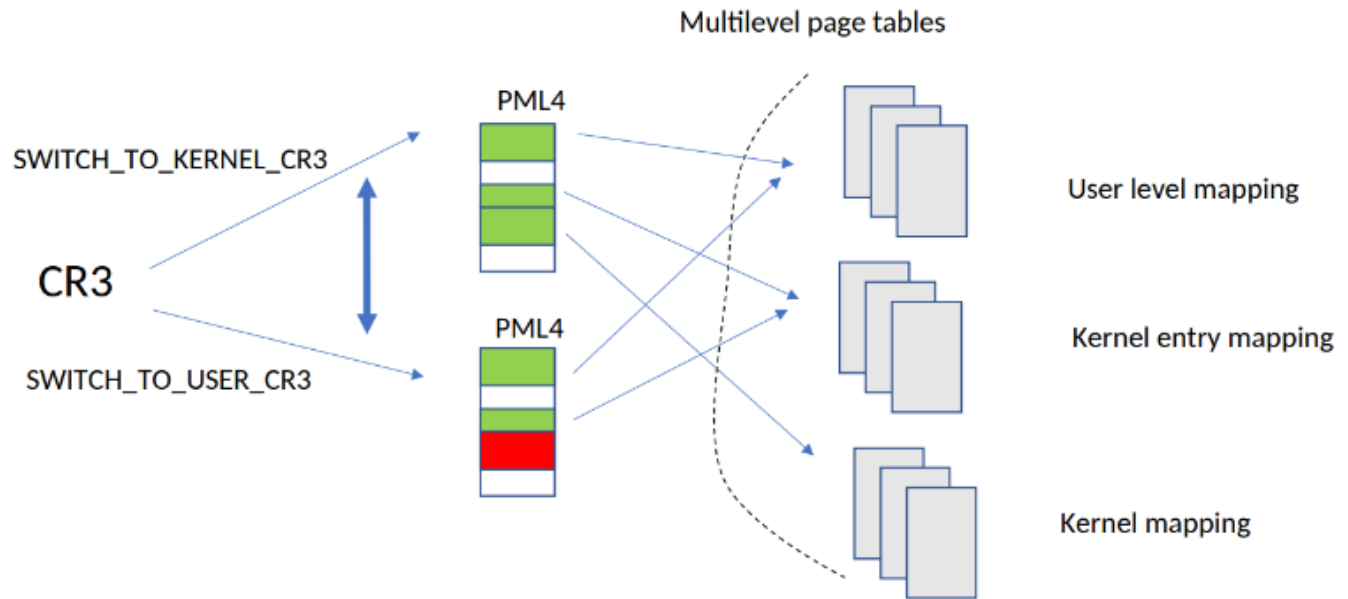
Vale la pena fare alcune osservazioni:

- 4096 byte è esattamente la dimensione di una pagina di memoria. Ma perché gli spiazamenti che si prendono in considerazione all'interno dell'array A sono esclusivamente la prima entry di ciascuna pagina di memoria? Per evitare problemi con la cache: di fatto, quando si accede a un certo valore in memoria, non è solo lui a essere caricato in cache, ma come minimo una **linea di cache**, che è lunga 64 byte. Leggiamo sempre lo 0-esimo byte, perché la cache lavora a blocchi e potrei trovare '64 byte' rapidi, e quindi non essere in grado di differenziare gli accessi. Per giunta, i processori tipicamente applicano il cosiddetto **pre-fetching**: nel momento in cui viene caricata in cache una linea, vengono conseguentemente caricate anche alcune linee adiacenti per il principio di località. Perciò, per evitare che i valori relativi a più di uno spiazamento in A vengano caricati in cache a seguito di un unico accesso, si prendono gli spiazamenti in modo tale che siano distanti tra loro, e una distanza di 4096 byte risulta essere sufficiente.
- Perché se nel kernel space viene letto il byte 0 si fa un nuovo tentativo? Perché un'area di memoria inizializzata a zero o è memoria non valida o è relativa a un terminatore di stringa, per cui si tratta di un'area di memoria non significativa. Resta comunque possibile ritentare ad accedere al medesimo byte all'interno del kernel space perché non è escluso che, in concorrenza, il kernel possa cambiare il valore di quel byte (da zero a un valore significativo e di interesse).

### Contromisure per l'attacco Meltdown

- KASLR (Kernel Address Space Randomization): Quando si fa il setup del kernel del sistema operativo, le strutture dati di livello kernel possono cadere ovunque all'interno del kernel space: in particolare, si stabilisce randomicamente un *offset*  $F$  rispetto all'indirizzo base del kernel a partire da cui sono definite le varie strutture dati del kernel. In questo modo, si complica la vita all'attaccante poiché lui non conosce a priori l'indirizzo del kernel space in cui andare a effettuare l'attacco; tuttavia, con un approccio brute-force, può comunque fare in modo che, prima o poi, l'attacco vada a buon fine. Per questa ragione, KASLR non è la soluzione definitiva contro l'attacco Meltdown.
- Cash flush:  
Si effettua semplicemente un flush della cache ogni volta che il kernel prende il controllo od ogni volta che un thread viene rischedulato. Tuttavia, così facendo, si avrebbe un calo inaccettabile delle prestazioni: le cache sono condivise tra i vari hyperthread, per cui ciascun cache flush impatterebbe su tutti i thread in esecuzione all'interno dell'host.
- KAISER (Kernel Isolation in Linux):  
Quando un processo gira in modalità *user*, utilizza una page table  $PT_U$  all'interno della quale il mapping della maggior parte degli indirizzi del kernel space non è presente: ci sono esclusivamente quegli indirizzi della porzione del kernel che vengono utilizzati ad esempio per gestire gli interrupt (se non ci fossero almeno loro, non saremmo neanche in grado di gestire gli interrupt o comunque di trasferire il controllo al kernel). Dall'altro lato, quando un processo gira in modalità kernel, utilizza un'altra page table  $PT_K$  che contiene il mapping degli indirizzi di memoria mancanti in  $PT_U$ . Di fatto, gli indirizzi del kernel space il cui mapping è presente nella page table  $PT_U$  costituiscono la cosiddetta **entry zone** del kernel, che contiene quasi esclusivamente il codice che permette di effettuare lo switch delle page table ( $PT_U$  a  $PT_K$ ). Perché quest'ultima soluzione risulta funzionante? Nel momento in cui un processo che esegue in user mode incontra un'istruzione offending che vuole accedere a un indirizzo di memoria del kernel space, tipicamente non troverà tale indirizzo nella page table che ha a disposizione. In altre parole, si solleva un'eccezione di tipo **page fault**, che non consente al processo di continuare ad eseguire le istruzioni anche speculativamente. Inoltre, si tratta di una soluzione che, sì, ha dei costi prestazionali, ma mai quanto il cash flush; in particolare, il calo delle performance viene osservato dalla **TLB - Translation Lookaside Buffer** che, a ogni switch da user mode a kernel mode e viceversa, dovrà cachare le informazioni sulla nuova memory view su cui ci siamo spostati (che sia essa relativa alla  $PT_U$  o alla  $PT_K$ ). Entrando un po' più nel dettaglio sulle page table, se ne hanno di livelli differenti (page table di livello 1, page table di livello 2, e così via). Solo le **page table di livello 1** sono replicate tra l'esecuzione user mode e l'esecuzione kernel mode, mentre le altre sono a **istanza unica**; in particolare, al livello 1, la page table  $PT_K$  è in grado di raggiungere tutte le informazioni delle page table ai livelli inferiori, mentre la page table  $PT_U$  punta solo a un sottoinsieme di informazioni delle page table ai livelli inferiori. Da user possiamo comunque accedere a qualcosa di livello kernel, ma solo per i moduli kernel utili/indispensabili (compile time) per avere il minimo supporto necessario. KAISER sfrutta quindi questa tecnica di **PTI-Page Table Isolation**, disattivabile da GRUB. Il selettore CR3, ad ogni cambio, azzerla la TLB e genera cache miss. Però è un'operazione automatica ed abbastanza semplice.





## Insights sui branch

Type	Direction at fetch time	Number of possible next fetch addresses?	When is next fetch address resolved?
Conditional	Unknown	2	Execution (register dependent)
Unconditional	Always taken	1	Decode (PC + offset)
Call	Always taken	1	Decode (PC + offset)
Return	Always taken	Many	Execution (register dependent)
Indirect	Always taken	Many	Execution (register dependent)

Per i salti **condizionali** e **indiretti**, l'indirizzo da cui riprenderà l'esecuzione è noto soltanto a partire da un certo stage S della pipeline. Chiaramente, però, vogliamo che venga comunque eseguito del lavoro all'interno della pipeline prima che l'istruzione di salto raggiunga lo stage S, altrimenti il calo delle prestazioni sarebbe certo e significativo. A tal proposito, si introducono dei **dynamic predictor** (o **branch predictor**), che hanno lo scopo di predire quale sarà la destinazione verso cui si dovrebbe saltare con la più alta probabilità. Ciò permette di comporre in maniera speculativa un program flow all'interno della pipeline in funzione della predizione che è stata attuata dal predittore. La reale implementazione dei branch predictor è basata sui **Branch History Table (BHT)**, detti anche **Branch-Prediction Buffer (BPB)**, che contengono metadati che associano un'istruzione di salto all'ipotesi su dove tale istruzione salterà. E' una zona di cache hardware contenente [indirizzo istruzione, target da usare se jump], e l'indirizzo istruzione può anche essere il target stesso. E' un suggeritore, non so nulla dell'affidabilità. L'implementazione minimale (andando avanti verranno aggiunte altre componenti) per una BHT consiste in una cache indicizzata tramite i *bit meno significativi* di ogni istruzione di salto. Ogni qual volta viene fetchata un'istruzione di salto, si può avere una cache hit o una cache miss, dove la *cache hit* la si ha nel caso in cui sono disponibili sufficienti informazioni che permettano di effettuare una predizione sulla destinazione del salto. Queste informazioni consistono in particolari bit di stato: ciò che succede nel passato è rappresentativo di ciò che si prevede che accadrà in futuro. Nei salti **condizionali** ho due destinazioni, saprò quella corretta in fase di esecuzione, e quindi è register dependent, perchè dipende a runtime cosa ci sarà nel registro).

I salti **unconditional** e le **call**, dove salteremo è già noto nello stage di decode, salterò sempre in quel punto.

Anche le **return** vengono sempre eseguite, ma ho più destinazioni note nello stage di esecuzione.

Per i salti **indiretti** salterò sempre (come per le call), ma con più destinazioni. Il target potrebbe essere in un registro.

## Predittori per i salti condizionali

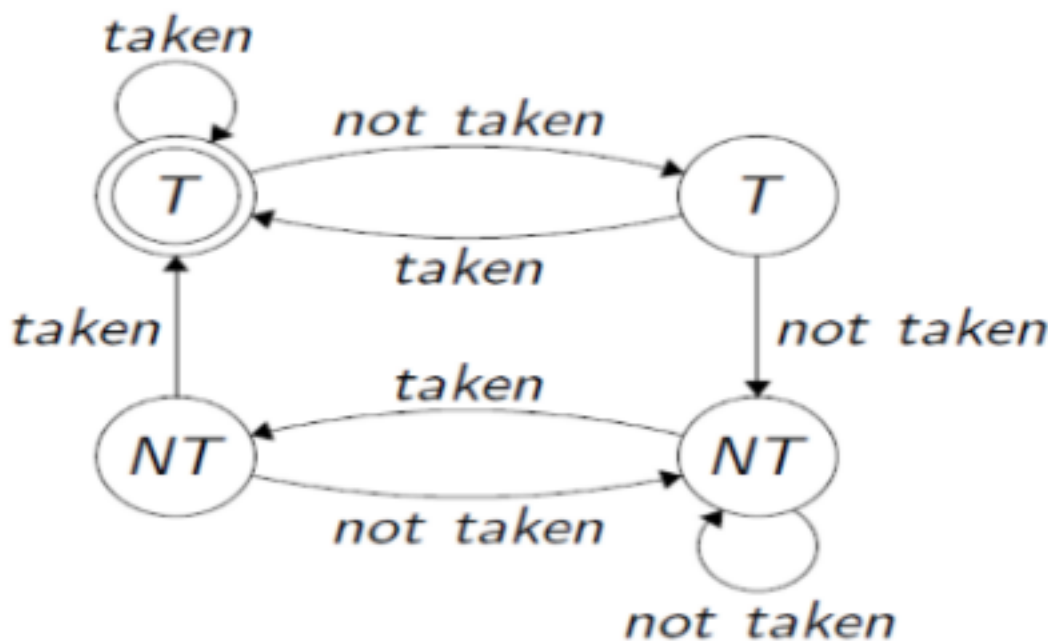
**Predittori coi bit di stato** Il modo più facile per implementare un predittore per i salti condizionali è sfruttando un **unico bit di stato**: se l'ultima istruzione di salto ha avuto come esito "taken" (ovvero ha realmente portato a effettuare il salto), allora il predittore prevederà che anche il prossimo salto avrà come esito "taken", e viceversa. Qui la storia passata è rappresentata esclusivamente dall'*ultima istruzione di salto*, per cui non si tratterebbe neanche di un vero e proprio storico.

**Esempio Caso semplice:**

ciclo for, sbaglio la predizione quando uscirò. Per questo motivo, sono stati introdotti anche i predittori a *due bit* di stato. Di fatto, questi predittori si comportano meglio negli scenari in cui si hanno molte istruzioni di salto con esito "taken" e poche istruzioni di salto con esito "not taken" (e viceversa). Lo scenario classico è quello dei *nested loop*. In particolare, la predizione cambia da "taken" a "not taken" (o viceversa) non più a seguito di un solo errore, ma a seguito di **due errori consecutivi del predittore**. La macchina a stati che rappresenta i predittori a due bit è la seguente, dove:

- T = "predico che il salto sarà taken"
- NT = "predico che il salto sarà not taken"

**Esempio doppio ciclo:** Nel ciclo interno itero, e la predizione mi dice che salterò. Quando finisco le iterazioni nel ciclo interno, la predizione sbaglia e dice che continuerò nel ciclo. Invece aumento l'iterazione sul ciclo esterno (1 errore). Però poi rientro effettivamente nel ciclo interno, quindi in realtà non devo cambiare predizione anche se ho fatto un errore, perchè stavolta ci rientro davvero dentro. *Solo se sbaglio 2 volte cambio previsione.*



Ricordiamo che, in caso di previsione sbagliata, in pipeline vengono caricate istruzioni che alla fine subiranno uno squash, che porterà al refill della cache. Quindi è importante effettuare la predizione con un buon tradeoff affidabilità/supporto hw.

Esistono anche predittori più sofisticati come il **Two-Level Correlated Predictor** e l'**Hybrid Local/Global Predictor** (noto anche come **Tournament Predictor**).

**Two-Level Correlated Predictor** Consideriamo il seguente codice:

```
if (aa == VAL) aa = 0;
if (bb == VAL) bb = 0;
if (aa != bb) {
    //do the work
}
```

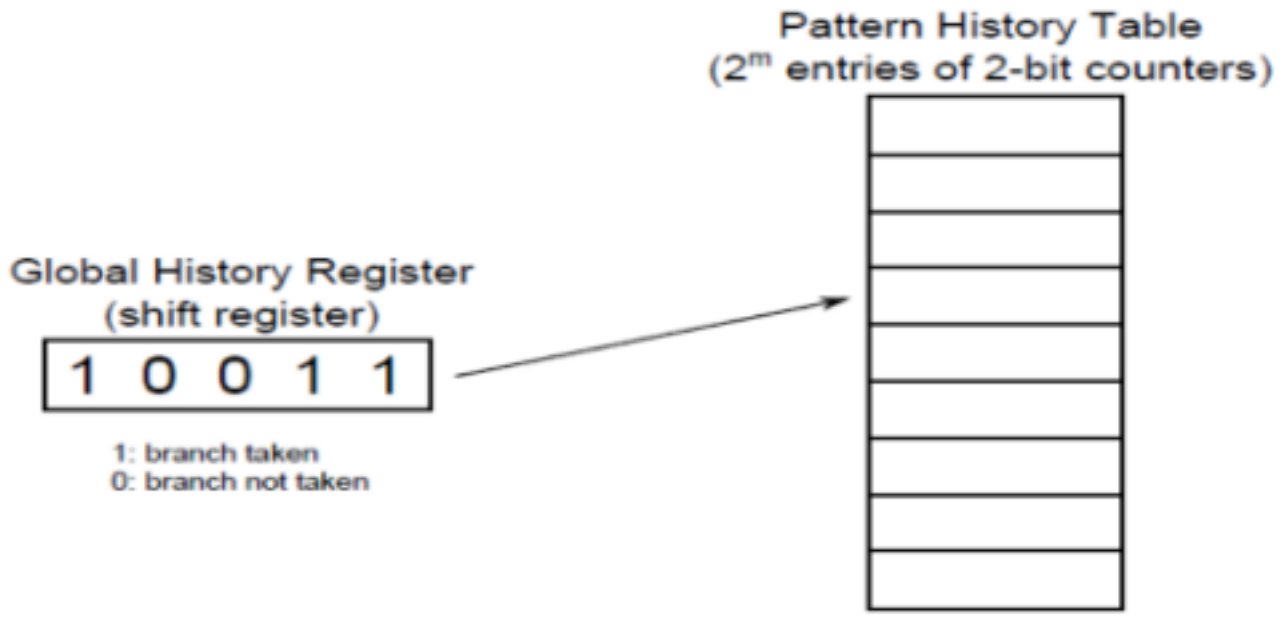
Ci piacerebbe avere un predittore tale per cui, se i branch relativi ai primi due if vengono presi, predica che il terzo branch (che dipende dai due precedenti) non venga preso. È quello che fa il **(m,n) Two-Level Correlated Predictor**, dove:

- m = *numero di salti precedenti*, il cui esito determina la predizione che verrà effettuata.

- $n$  = numero di bit del valore che indica quanti **errori** deve commettere il predittore prima di **cambiare** la sua **predizione** (esattamente come nei predittori con bit di stato).

#### Esempio:

Consideriamo il caso con  $m=5$ ,  $n=2$ . Si ha il seguente scenario:



Il **global history register** è composto da  $m$  bit, ciascuno dei quali indica l'esito di uno delle ultime  $m$  istruzioni di salto condizionale. Da destra ho i salti più recenti. Tale registro può assumere  $2^m$  valori diversi, e a ciascuno di essi è associato un predittore con bit di stato (= una macchina a stati) differente. Il predittore per la predizione corrente viene scelto sulla base dei risultati degli ultimi  $m$  branches, come è codificato nella  $2^m$  bitmask. Sostanzialmente usiamo il Global History Register come indice.

#### Un esempio semplice con ( $m = 2$ , $n = 2$ ).

Nel global register posso avere 4 casi: 00, 01, 10, 11, dove 0 vuol dire “not taken” e 1 “taken”. Per ogni casistica associo una entry nel Pattern History Table.

00 mappato nell'entry 0, 01 mappato nell'entry 1, 10 mappato nell'entry 2, 11 mappato nell'entry 3.

Il fatto che  $n = 2$  vuol dire che ogni entry mantiene 2 bit nell'array, ovvero include l'automa a stati visto precedentemente.

Se ho tre statement if, e salto sempre al terzo in modo condizionato, la mia branch sequence è 001001001001...

Se ( $m =$  indifferente,  $n = 2$ ) mantengo 4 entry con due bit ciascuno (00, 01, 10, 11).

- Entry 00: ho sbagliato due volte, allora il terzo salto sono sicuro.
- Entry 01: ho sbagliato una volta, resto su 0.
- Entry 10: ho sbagliato una volta, resto su 0.
- Entry 11: non capita perchè abbiamo detto di non saltare due volte consecutivamente.

Se avessi ( $m = 0$ ,  $n = 2$ ), avrei 0 elementi per identificare l'entry. Sarebbe un predittore basico.

#### Tournament Predictor

In realtà non è sempre vero che i salti condizionali siano correlati tra loro. Per questo motivo si introduce l'**Hybrid Local/Global Predictor**, che consiste in una “sfida” tra un **predittore locale** (come quelli con bit di stato) e un **predittore globale** (che si basa sulla storia passata e assume che i salti condizionali siano correlati). In pratica, si hanno entrambi questi predittori e in più una macchina a 4 stati che indica se correntemente conviene utilizzare il predittore locale oppure quello globale: se stiamo sfruttando il predittore locale, *switchiamo a quello globale dopo due errori consecutivi, e viceversa*. Ciò è realizzato mediante *Return Stack Buffer RSB*, tipo una semplice cache, con 32 entries (quindi 32 livelli di annidamento), cioè indirizzi associati al return point dell'istruzione call, in cui salvo nella prima entry libera tale indirizzo di ritorno.

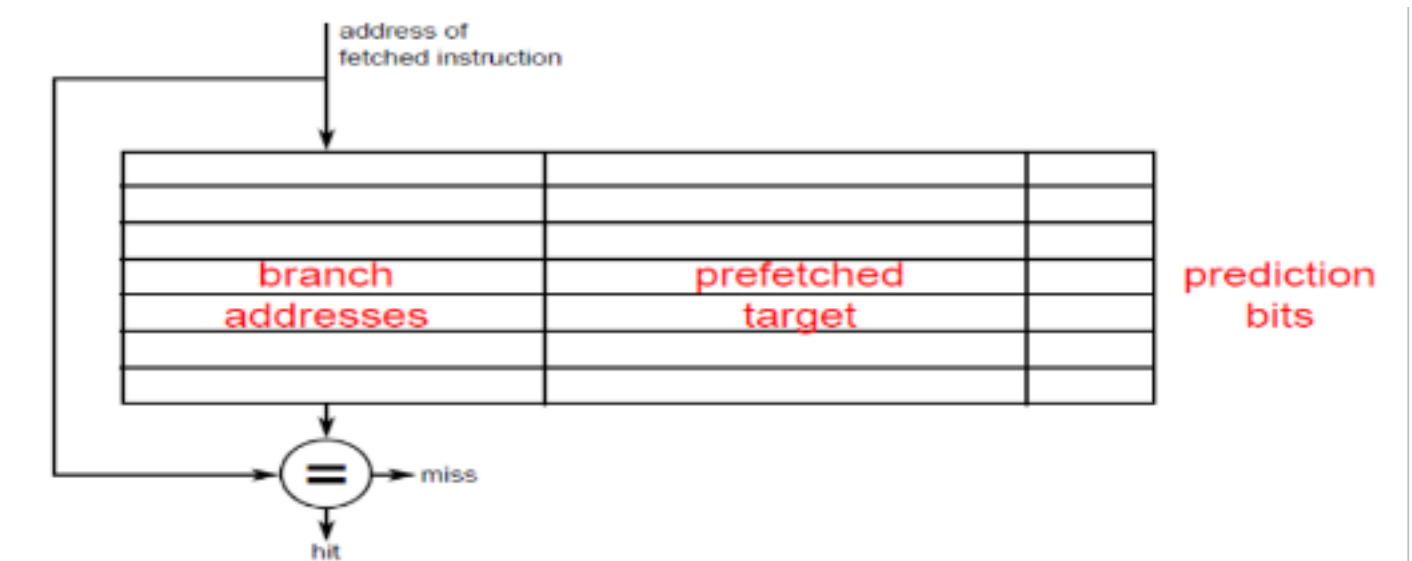
**Osservazione su RSB** Faccio call di una funzione, salvo indirizzo successivo alla call in RSB, quando esco dalla funzione il return punta all'indirizzo salvato nell'RSB precedentemente. RSB è una cache di tipo LIFO, contenuta nel processore ed è modificabile solo tramite *call* (inserimento valore) o *return* (prelievo valore). Non è esposta nell'ISA.

Questi switch e letture avvengono in elementi mantenuti in cache, allora non ho perdite di prestazioni.

### Predittore per salti indiretti

I salti indiretti sono tali per cui l'indirizzo di destinazione viene caricato in un registro. Se l'istruzione deve saltare, questo è funzione del risultato delle istruzioni precedenti, ovvero di cosa queste hanno scritto nei registri, che verranno usati per saltare.

Quindi, il valore di tali registri può *variare sempre nel tempo*, il che rende la predizione più complessa e con più side effects generati. Per questo motivo, le destinazioni possibili possono essere molteplici, e la predizione può risultare più farraginosa. I predittori per i salti indiretti sono dotati della seguente struttura dati:



Ciascuna entry della tabella è relativa all'indirizzo di una particolare istruzione di salto (**branch address**). Ogni branch address è associato a un **prefetched target** (che è il *presunto* indirizzo destinazione del salto, se c'è cache miss non posso prelevare, oppure prendo l'istruzione successiva) e a dei **prediction bit** (che, come al solito, determinano numero di errori consecutivi da commettere prima di cambiare la predizione, ovvero prima di cambiare il prefetched target). In particolare, la prima volta che si incontra una certa istruzione di salto indiretto, si effettua il training, in cui il prefetched target verrà inizializzato all'effettivo indirizzo destinazione del salto. In queste tabelle riporto solo i bit meno significativi, per questioni di scalabilità. Tutto questo sta nel core, e se ho hyperthreading, la predizione viene fatta da una comune componente hardware per più thread. Ovvero se ho due thread (supponiamo facciano stesse cose) su due hyperthread e 1 core, questo core comune può essere sfruttato dal thread 2 per avere informazioni sui salti ereditate da thread 1, per ottenere un boost delle performance. A livello di sicurezza è un po' problematico: se thread1 seguisse un comportamento tale da suggerire al thread 2 di fare salti (anche speculativi) che in realtà non dovrebbe fare?

### Attacco Spectre

Chiaramente anche la predizione dei target dei salti condizionali può portare il processore a riempire la pipeline con istruzioni eseguite in modo speculativo. In particolare, se la predizione è scorretta, le istruzioni successive al salto dovranno poi essere buttate. Ma anche qui, come nel caso delle eccezioni, le istruzioni considerate in maniera speculativa che poi vengono scartate lasciano dei side effect nello stato micro-architetturale. (Solito discorso: un salto speculativo non installa nulla nell'ISA, ma lascia tracce a livello micro architetturale). Da qui nasce la possibilità di compiere una famiglia di attacchi denominati **Spectre** che sfruttano un covert-channel.

Questi attacchi sono anche più gravi di Meltdown poiché:

- **Non richiedono l'esecuzione di un'istruzione offending.**
- È possibile effettuare un **training scorretto** del branch predictor inserendo nel codice di livello user degli appositi branch / delle istruzioni di salto condizionale.

L'unica soluzione per evitare questi attacchi sarebbe quella di non fare predizioni, ma i sistemi diventerebbero molto più lenti (con o senza hyperthread), a livello di marketing sarebbe una mossa controproducente!

## Spectre V1

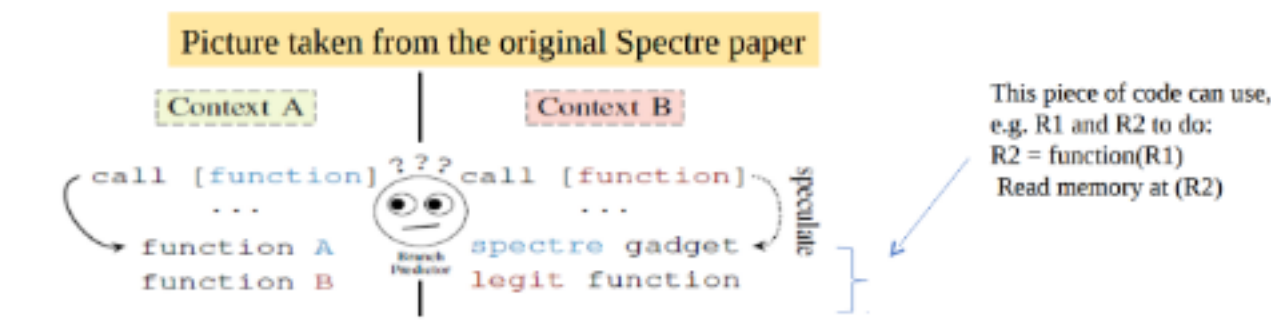
È un attacco che inizialmente effettua il flush della cache (come Meltdown) per poi sfruttare i salti condizionali. In particolare, prevede l'utilizzo del seguente blocco di codice:

```
if (x < array1_size)
    y = array2[array1[x] * 4096]
//il prodotto per 4096 corrisponde a uno shift a sx di 12 posizioni
```

- *array1* è un array che può trovarsi in qualunque punto dell'address space, e *array1[x]* è un particolare valore che verrà moltiplicato per 4096 per ottenere un indice di pagina, che verrà utilizzato per accedere ad *array2*.
- In ogni caso, il valore *x* viene selezionato in modo tale che sia molto elevato (maggiore di *array1\_size*), cosicché il branch non venga realmente preso, mentre magari il predittore prevedeva il contrario; inoltre, *x* può essere tale che l'accesso ad *array2* porti a reperire (speculativamente) un'informazione segreta (*y*), ad esempio all'interno del kernel space.
- Quel che si ottiene è che il valore *y* viene caricato all'interno della cache in modo speculativo. A questo punto, come in Meltdown, si effettuano gli accessi alla prima entry di ogni pagina di *array2* finché non si giunge a quella caricata in cache (di cui si osserva un tempo di accesso ridotto).
- Quindi ciò che carico è in funzione sia dell'array sia di *x*, se salto molto lontano potrei andare in una zona kernel. Lo scopo è far sì che questa condizione venga eseguita speculativamente, quindi il predittore deve essere indotto a credere che tale flusso di esecuzione sarà quello realmente eseguito. **Non devo eseguirlo realmente!**
- No trap perchè sono sempre in user mode e tutto dipende dall'esito dell'IF. Aver *patchato* Meltdown non implica risolvere anche Spectre. La soluzione di Meltdown risiedeva nel fatto che tale attacco creasse una trap poichè si passava da user a kernel, e c'era il cambio del puntatore alla page table. Qui sono sempre in user mode.

## Spectre V2

Con più thread passo alla versione v2! È un attacco che sfrutta i salti multi-target (indiretti) ed è di tipo **cross-context**. Ciò vuol dire che l'attaccante è in grado di inferire (aggiungere) delle informazioni sensibili da un contesto di esecuzione diverso dal suo: supponiamo di avere due thread A, B che girano sul medesimo hyperthread in modalità processor sharing o su due hyperthread relativi allo stesso core. Il thread A può portare avanti delle attività che vanno a cambiare lo stato del branch predictor all'interno del CPU-core. Di conseguenza, B può osservare il nuovo stato del branch predictor all'interno del medesimo CPU-core e, quindi, le attività che lui eseguirà in modo speculativo (in funzione dello stato del branch predictor) vengono *praticamente decise, e poi osservate mediante side effect, dal thread A*. Tale side effect, come al solito, può consistere nel caricamento di un'informazione sensibile all'interno della cache attraverso il suo accesso in memoria. L'efficacia dell'attacco diventa particolarmente evidente quando l'attaccante si basa su un contesto che utilizza una libreria condivisa (shared library) col contesto del thread B. Qui, infatti, le pagine di memoria della shared library vista dal thread B e le pagine di memorie della shared library vista dal thread A mappano esattamente sugli stessi indirizzi fisici. Perciò è ovvio che qualunque side effect la vittima lasci speculativamente su tali locazioni fisiche di memoria sia direttamente visibile all'attaccante.



Con riferimento alla figura, il **gadget** è un blocco di codice che è definito nell'address space della vittima e viene sfruttato dall'attaccante affinché la vittima lo esegua in modo speculativo a seguito di un errore del branch predictor; in tal modo, nello stato micro-architetturale, la vittima lascia i side effect desiderati dall'attaccante. Nell'esempio mostrato nella figura

vengono utilizzati due registri (*R1* e *R2*), di cui *R1* viene utilizzato per calcolare *R2* con una particolare funzione, mentre *R2* viene usato per memorizzare l'indirizzo verso cui accedere in memoria. In realtà, sarebbe sufficiente l'utilizzo di un solo registro *R1*.

NB: l'utilizzo del medesimo CPU-core tra thread A e thread B è richiesto solo per effettuare un training errato sul branch predictor. Dopodiché, poiché la cache è condivisa tra più CPU-core, gli accessi in memoria effettuati per stabilire quali dati sono saliti in cache possono essere effettuati anche in un momento in cui il thread vittima (B) gira in un CPU-core differente. Tra l'altro, per portare a termine l'attacco, è possibile anche utilizzare una macchina virtuale su cui possono girare delle applicazioni che, in qualche modo, vanno a cambiare lo stato del branch predictor: in fondo l'hardware sottostante viene utilizzato anche per eseguire le macchine virtuali!

## Contromisure per gli attacchi Spectre

### Retpoline - Return trampoline

Al posto di invocare l'istruzione di salto indiretto (che sia essa una jump o una call), si esegue un blocco di istruzioni funzionalmente equivalente che non consente di effettuare una mis-prediction (ovvero un training scorretto del branch predictor) per portare a compimento l'attacco Spectre. Una versione semplificata del blocco di istruzioni è riportata di seguito. Non faccio più salti indiretti, solo diretti!

Invece di *Jump su R*, eseguo *Return su R*, non c'è più predizione. Il *RET* usa il *Return Stack Buffer* per la predizione (di tipo *Lifo*), controllabile via software, a differenza del predittore di Branch Indiretto.

```
push target_address
1: call retpoline_target
    //put here whatever you would like (with no side effects)
    jmp 1b
retpoline_target:
    lea 8(%rsp), %rsp //we do not simply add 8 to RSP since FLAGS registry should not be modified
    ret    //this will jump to target_address
```

Vediamo nel dettaglio cosa fa questo blocco di istruzioni:

- Carica l'indirizzo destinazione del salto (**target\_address**) sullo **stack**. (ora è in cima allo stack)
- Effettua una call di **retpoline\_target** (per cui viene caricato sullo **stack** anche l'indirizzo dell'istruzione successiva alla call). (Graficamente, nell'address space avremo *target\_address* e sopra di lui, l'indirizzo successivo alla call. Questo perché chiamata *retpoline\_target*, prevediamo che al suo *return* ripartiamo dall'istruzione successiva ad essa). **rsp** a 64 bit, quindi 8 byte, per questo faccio quello shift.
- In **retpoline\_target** si **elimina** l'indirizzo dell'istruzione successiva alla call e, tramite la return, si salta verso l'indirizzo che si trova in cima allo **stack** (ovvero *target\_address*). Facendo lo shift di 8 byte, cancelliamo il punto di ritorno della chiamata *retpoline\_target*. Dopo ciò, in cima allo stack abbiamo proprio *target\_address*.
- Essendoci una **call** (non è register dependent, mi manda verso la funzione chiamata), il predittore non deve essere soggetto a training: prevede che, a seguito della return, si salti verso l'istruzione successiva alla call. Per questo motivo, dopo la call devono essere eseguite delle istruzioni innocue, come una jump verso l'istruzione stessa di call. La *return* vede qualcosa nell' *rsb* che non può essere bypassata. La *ret* esegue speculativamente ciò che gli dice *rsb*. Tuttavia, se la CPU specula, l'RSB la porta ad eseguire *jmp 1b*, intrappolandolo in un loop. Successivamente, la CPU realizza che il valore in RSB è diverso da quello nello stack (*target\_address*), e stoppa la speculazione. **O vado in target\_address, o resto in un loop**. RSB ha senso per singolo processo, mentre il predittore si muove tra più flussi.



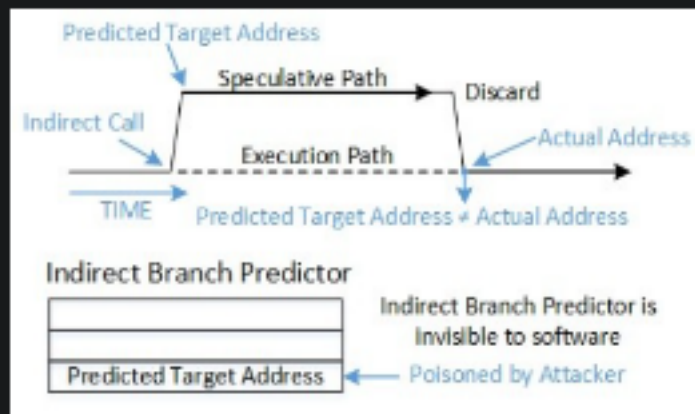


Figure 1: Speculative Execution without retpoline

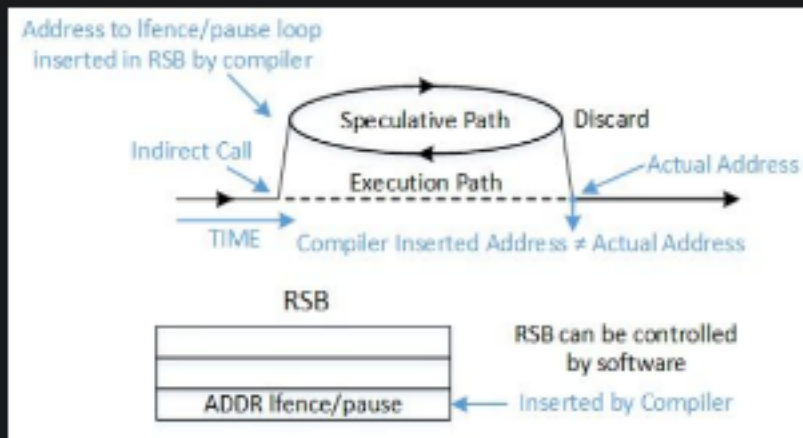


Figure 2: Speculative Execution with retpoline

## IBRS (Indirect Branch Restricted Speculation)

È possibile aggiornare il registro **MSR** (Model Specific Register, è un registro di controllo nella CPU) per creare un'enclave di esecuzione (= uno spazio di esecuzione chiuso entro determinati confini) dove la storia non è influenzata da cosa avviene al di fuori dell'enclave stessa. In pratica, a livello hardware siamo in grado di utilizzare la branch prediction in maniera differenziata a seconda se stiamo lavorando a livello user ( $MSR=0$ ) o a livello kernel ( $MSR=1$ ): nel primo caso il branch predictor dei salti indiretti si basa sulla storia passata della sola esecuzione a livello user, mentre nel secondo caso si basa sulla storia passata della sola esecuzione a livello kernel. Sostanzialmente, a tempo  $t$  decido che lo stato del predittore non dovrà più essere usato per un certo tempo. E' un guscio che mi protegge dall'esterno. Utile se ho app che lavorano a livelli diversi (user e kernel), se ad esempio volessi che le scelte kernel non venissero influenzate dall'user.

## IBPB (Indirect Branch Prediction Barrier)

Anche qui ci si basa sull'utilizzo del MSR. In particolare, nell'istante in cui si va a scrivere all'interno di tale registro, stiamo cancellando tutto ciò che il branch predictor dei salti indiretti ha imparato finora, creando così una sorta di barriera. Al tempo  $t$  resettiamo il branch predictor.

**IBRS** e **IBPB** sono operazioni software, spesso si usa la syscall `prcrcl()` per lavorare col Processor Control. Se lavoriamo con `exec()` (programma che chiama un altro programma) funzionano uguale, perchè facenti parte dello stesso programma iniziale.

## Sanitizzazione

### Introduzione

Per Spectre V1 abbiamo visto la possibilità di riusare le patch di Meltdown. Spectre V1 mi permetterebbe di leggere dati kernel? Sì, perchè la predizione a livello kernel potrebbe essere basata su predizioni livello user.

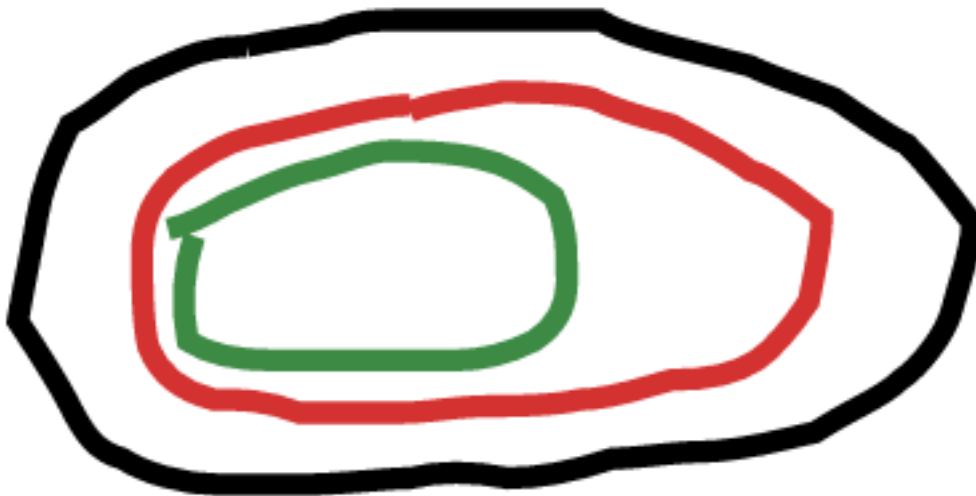
A livello kernel abbiamo una tabella, che chiamiamo **driver**, in cui per ogni indice ho un servizio che richiama una *system call*. L'indice è passato dall'user, ma è compatibile con la taglia della tabella? Dipende, con l'indice facciamo un salto verso

una entry, ma se eccedo la dimensione? **Speculativamente** starei usando un indice errato come pointer per una zona codice Kernel, perchè se le precedenti `call` erano corrette, il predittore si sarà settato in un certo modo. (Se ha sempre saltato, speculativamente lo rifarà). Soluzione? Sanitizzazione.

### Come funziona

È una contromisura per i salti condizionali. In particolare, ciascuna condizione `if` che coinvolge una certa variabile `V` prevede dei valori per `V` ammissibili e dei valori per `V` non ammissibili. Quello che si fa è ridurre all'osso l'insieme dei valori non ammissibili, facendo sì che tutti gli altri valori non ammissibili non possano in alcun modo essere assunti da `V`. Questa è una tecnica estremamente utile per gli indici delle tabelle. Infatti, supponendo che per una tabella gli indici ammissibili vadano da 0 a  $j - 1$ , se per qualche motivo l'indice dovesse assumere un qualunque valore maggiore di  $j - 1$ , viene forzato ad assumere il valore  $j$ . In pratica, abbiamo imposto che  $j$  sia l'unico valore non ammissibile che può essere materializzato. A questo punto, se l'indice vale  $j$ , viene imposto ad esempio un accesso in memoria innocuo (i.e. alla prima locazione di memoria subito dopo la tabella vengono inserite delle informazioni non sensibili in modo tale che, anche speculativamente, viene acceduta o una entry della tabella o l'unica locazione di memoria ammissibile e innocua al di fuori della tabella).

### Come funziona - for dummies



Normalmente avremmo “indici OK” (verde) ed “indici NON OK” (rosso). Però se usassi questi indici NON OK, potrei andare in zone delicate. Introduco la “zona nera”, più ampia. Tutto quello che cade lì dentro lo butto nella zona rossa, che sarà la più piccola possibile, per limitare i danni.

Come lo faccio? Applicando una maschera di bit.

Perchè devo differenziare zona rossa e nera?

Perchè se lascio solo zona rossa posso andare in zone brutte. Se questa è piccola (ad esempio un indirizzo non critico) e faccio sì che tutti gli altri indici *NON OK* vadano lì dentro, limito i danni.

Perchè non voglio che applicando questa maschera ad un indice NON OK, questa mi porti ad una zona verde. Quindi prima applico la maschera, e poi faccio il controllo.

### Loop unrolling

Ricordiamo che i salti sono azzardi, sbagliare salto compromette performance (squash pipelin) e sicurezza. E se riducessimo i salti?

```
int s=0;
for (int i=0; i<16; i++)
{ s+=i;}
```

Questo ciclo si traduce nella seguente sequenza di istruzioni assembly (dove la sintassi adottata è AT&T):



400545:	8b 45 fc	mov	-0x4(%rbp), %eax
400548:	01 45 f8	add	%eax, -0x8(%rbp)
40054b:	83 45 fc 01	addl	\$0x1, -0x4(%rbp)
40054f:	83 7d fc 0f	cmpl	\$0xf, -0x4(%rbp)
400553:	7e f0	jle	400545 <main+0x18>

Come si può notare, a ogni iterazione del ciclo vengono eseguite cinque istruzioni, di cui tre di controllo (in blu scuro) e solo due di lavoro effettivo, in cui viene incrementata la variabile `s` (in nero): in pratica, in questo particolare esempio, abbiamo un overhead di computazione pari ai 3/5 delle istruzioni totali, il che porta inevitabilmente a un degrado delle prestazioni. Per questo motivo, corre in aiuto il **loop unrolling**, che è una tecnica per “srotolare” i cicli: da una parte si riduce il numero di iterazioni che devono essere eseguite, dall’altra, all’interno di ciascuna iterazione, si esegue il lavoro utile più volte. In tal modo, si riduce il numero di volte in cui devono essere eseguite le istruzioni di controllo. Per effettuare l’unroll in modo automatico, è possibile ricorrere alle seguenti direttive di C:

```
#pragma GCC push_options
#pragma GCC optimize ("unroll-loops")
//region to unroll
#pragma GCC pop_options
```

È anche possibile specificare esplicitamente il fattore di unroll mediante `#pragma unroll (N)`. Ma affinché queste direttive siano effettivamente attive, è necessario compilare il file C col `flag -O`.

**Attenzione:** il loop unrolling porta alla necessità di utilizzare un maggior numero di registri per eseguire tutte le operazioni della medesima iterazione. Inoltre, porta ad avere le istruzioni macchina del ciclo su un range di indirizzi più ampio, per cui si ha una minore località. Per questi motivi, non è una tecnica che può essere sfruttata in modo spregiudicato. Per sfruttarlo bene, bisogna capire architettura e obiettivi!

## Power wall

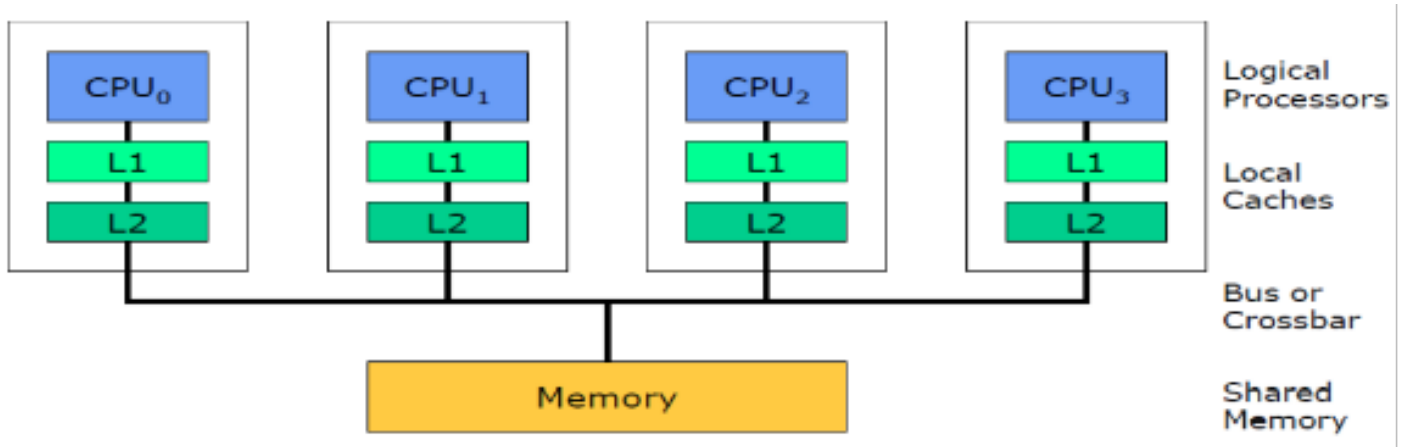
Per aumentare le prestazioni dei processori, è possibile seguire due approcci:

1. Aumentare la frequenza del clock.
2. Incrementare i componenti hardware a supporto del processore.

Per quanto riguarda la strategia 1, c’è una limitazione: la dissipazione massima che un processore può tollerare senza che si bruci è 130 W. Ma la dissipazione è pari a  $V \cdot V \cdot F$ , dove  $V$  è il voltaggio ed  $F$  è la frequenza del clock.  $V$  non può essere al di sotto di una certa soglia minima, altrimenti i componenti hardware non funzionerebbero proprio. Di conseguenza, esiste un upper-bound per  $F$  che è stato già raggiunto. Questa limitazione è detta **power wall**. A causa del power wall, ad oggi siamo obbligati a seguire la strategia 2 per rendere più efficienti i processori. In particolare, nel tempo, sono state adottate le soluzioni architetturali descritte di seguito.

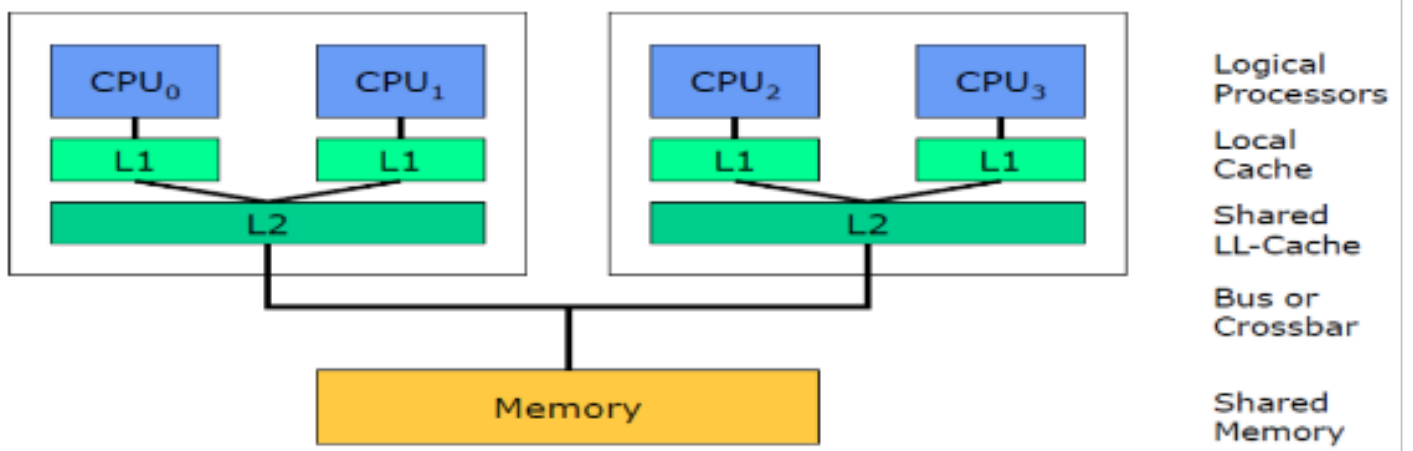
## Symmetric Multiprocessors

Si hanno semplicemente molteplici processori, ciascuno dei quali, per accedere alla memoria, impiega la stessa quantità di tempo. Abbiamo per ogni core un thread.



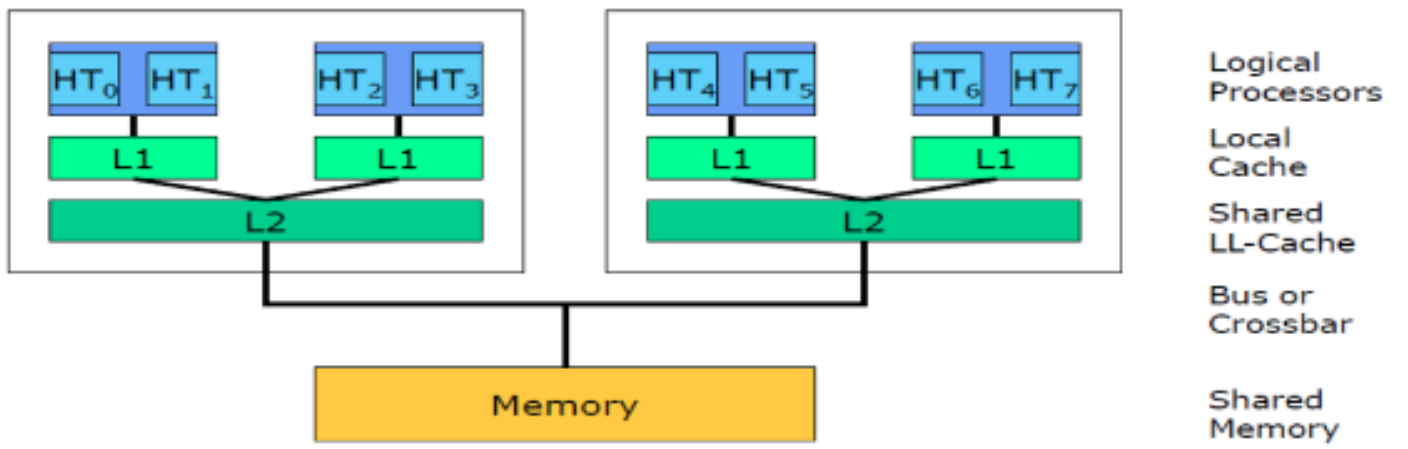
### Chip Multi Processor (CMP) o Multicore

All'interno di ciascun processore si hanno più motori (più core). Un hypethread per core.



### Symmetric Multi Threading (SMT) o Hyperthreading

All'interno di ciascun core possono esserci **più hyperthread distinti**. Dal punto di vista del sistema operativo, è esattamente come se ci fosse un numero di processori pari al numero totale di hyperthread. Il problema qui è che la memoria rappresenta un collo di bottiglia importante, anche perché viene acceduta in modo concorrente da tutti gli hyperthread. Dobbiamo vederla come architettura distribuita. La memoria è esposta in ISA (ma non le componenti cache  $L_1, L_2 \dots$ ).



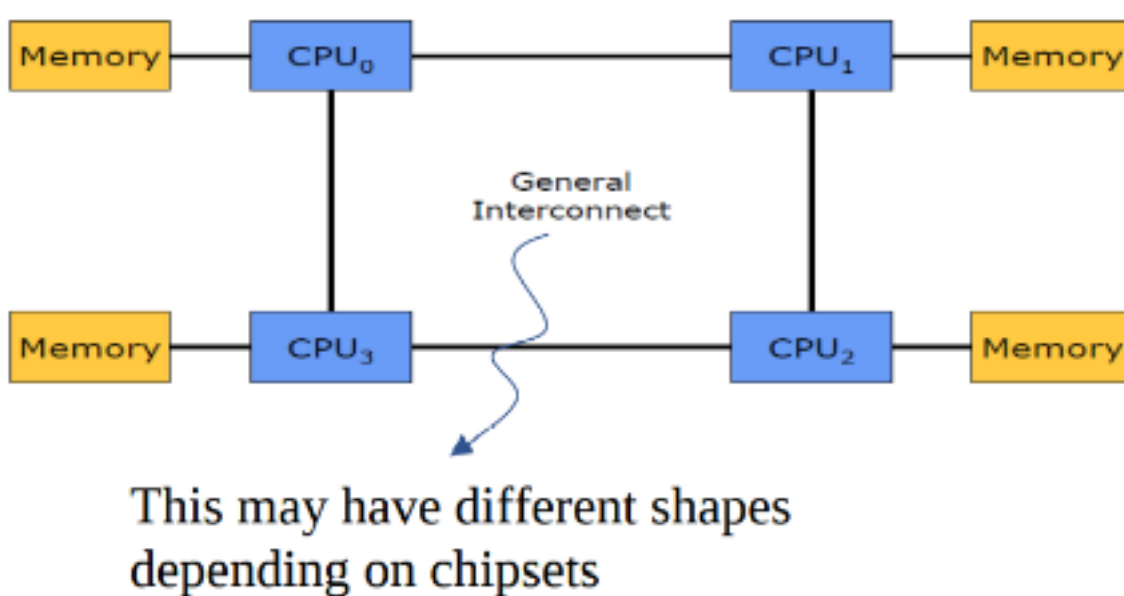
Successivamente si è passati all'**UMA**: memoria unica ad accesso uniforme, cioè stesse componenti hardware accedute in parallelo. Poi **NUMA**.

## Non Uniform Memory Access (NUMA):

Qui la memoria è divisa in banchi, ciascuno dei quali è direttamente collegato a uno specifico core (o a un insieme di core). Se un thread che gira in CPU-core  $i$  accede al banco di memoria a esso vicino, l'interazione con la memoria risulta molto efficiente, mentre se deve accedere a un altro banco di memoria, impiegherà più tempo. Tale soluzione risulta vantaggiosa nel momento in cui si riesce a fare in modo che ciascun CPU-core acceda prevalentemente al *proprio* banco di memoria (ovvero effettui prevalentemente degli accessi locali). Ciascuna **coppia** (CPU-core  $i$ , banco di memoria  $i$ ) o (insieme  $i$  di CPU-core, banco di memoria  $i$ ) costituisce un **nodo NUMA**. I thread possono leggere da indirizzi diversi (quello che abbiamo detto prima), tuttavia l'interconnect è time shared, si genera traffico, che deve essere ben gestito.

### Osservazione:

- Un accesso locale (CPU verso memoria adiacente) richiede un tempo pari a  $50(hit)+200(miss)$  cicli, quindi è anche facile capire, osservando il tempo, se si ha hit o miss.
- Se l'accesso è NON locale e SENZA TRAFFICO, si passa a  $200(hit)+300(miss)$ . Tempi totalmente diversi, che possono creare **problemi di coerenza nella cache**.



**NB:** ISA definisce come il processore interpreta e esegue le istruzioni, ma non specifica i dettagli sulla gerarchia di memoria o la presenza delle cache.

## Cache Coherency

Com'è possibile osservare, nelle architetture di memoria elencate precedentemente la cache è **replicata**: da qui sorge il problema della **cache coherency**, secondo cui bisogna stabilire qual è il *valore corretto da restituire a un thread che effettua una determinata lettura in memoria*. Disponiamo infatti di un processore e di una zona di memorizzazione (cache e RAM). Osservare solo l'interfaccia memoria-processore ci fornisce una visione limitata. Se una istruzione scrive una *mov*, questa passa prima per lo *store buffer*, non viene subito esposto nell'ISA. Se un altro programma dovesse leggere, non è detto che il valore sia già stato scritto. La coerenza viene **definita** in tre punti fondamentali:

- **Causal consistency**: se un processore  $p_1$  scrive un valore  $v$  nella locazione di memoria  $X$  e poi effettua una lettura da  $X$ , dovrà leggere il valore  $v$ , a meno che non ci siano altri processori che nel frattempo hanno finalizzato delle scritture concorrenti su  $X$  (coerenza di tipo RAW – Read After Write).
- **Avoidance of staleness**: se un processore  $p_1$  scrive un valore  $v$  nella locazione di memoria  $X$  (memoria qui intesa come RAM, visibile su ISA) e *dopo una quantità sufficiente di tempo* un altro processore  $p_2$  effettua una lettura da  $X$ ,  $p_2$  dovrà leggere il valore  $v$ , a meno che non ci siano altri processori che nel frattempo hanno finalizzato delle scritture concorrenti su  $X$ .
- **Avoidance of inversion of memory updates**: le varie scritture su uno stesso dato  $X$  devono essere viste da tutti i processori nello stesso ordine. Se ho una sequenza di scrittura dalla cache verso la RAM, l'ordine di scrittura in cache deve essere riproposto anche nella RAM. Queste proprietà valgono sulla singola locazione.

Le due principali tecniche utilizzate per mantenere la consistenza, ricordando che la copia master è in memoria, e le copie in cache, sono:

- **Write through cache:** Prima aggiorniamo in cache e poi in memoria. Questo porta al seguente problema: una CPU potrebbe leggere due volte lo stesso dato dalla cache, e tra le due letture un'altra CPU può aver toccato il dato e aggiornato in memoria. Quindi per la CPU in esame, abbiamo una incoerenza tra quello che c'è in cache e quello che c'è in memoria.
- **Write back cache:** Aggiornare la cache e procedere alla memorizzazione in un secondo momento che decidiamo noi. Il problema è che lasciar decidere a noi potrebbe portare a scritture in memoria non nello stesso ordine di come sono state effettivamente eseguite.

## Protocolli CC Cache Coherency

Permettono di conseguire la cache coherency e sono il risultato delle scelte di:

- Un insieme di transazioni (e.g. aggiornamento di una replica della cache) supportate dal cache system distribuito.
- Un insieme di stati per i cache block (= unità minime della memoria che possono essere portate all'interno della cache).
- Un insieme di eventi che capitano all'interno del cache system distribuito.
- Un insieme di transizioni di stato.

Il design dei protocolli CC dipende da svariati fattori come:

- La topologia dei componenti (e.g. single bus, gerarchica, ring-based).
- Le primitive di comunicazione (i.e. unicast, multicast, broadcast).
- Le cache policy (e.g. write-back, write-through).
- Le feature della gerarchia di memoria (e.g. numero di livelli della cache, inclusiveness).

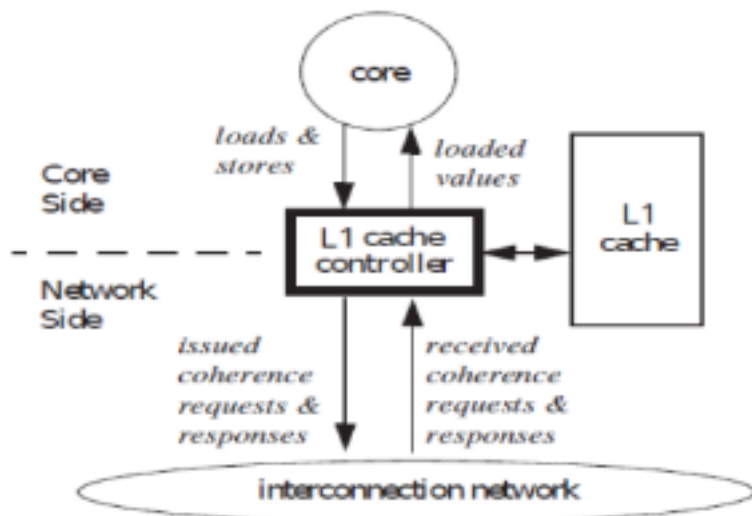
Implementazioni di cache coherency differenti possono presentare prestazioni diverse in termini di:

- **Latenza:** tempo per completare una singola transazione.
- **Throughput:** numero di transazioni completate per unità di tempo; aumenta se si fa in modo che scritture su aree di memoria diverse (e sufficientemente distanti) da parte di thread differenti avvengano in modo concorrente.
- **Overhead spaziale:** numero di bit richiesti per mantenere lo stato di un cache block.

Esistono due famiglie principali di protocolli CC:

- **Invalidate protocols:** quando un CPU-core scrive su una copia di un blocco, qualsiasi altra copia del medesimo blocco (i.e. della medesima informazione) viene *invalidata*: soltanto chi ha scritto l'ultima informazione ha la versione aggiornata del blocco. In tal caso, si utilizza poca banda ma si ha una latenza elevata per gli accessi in lettura: infatti, un CPU-core che si trova vicino a una replica invalidata, per andare a leggere il dato, ha la necessità di sfruttare l'unica replica valida, che è lontana. Sostanzialmente se aggiorni una linea di cache, spariscono tutte le altre repliche. Leggeranno tutti da me, perchè ho l'unica copia valida.
- **Update protocols:** quando un CPU-core scrive su una copia di un blocco, la nuova informazione viene propagata su tutte le altre copie del medesimo blocco. In tal caso, si ha una bassa latenza per gli accessi in lettura ma utilizza molta più banda.

Per poter invalidare / modificare una replica di un blocco che è stato aggiornato, si ricorre al cosiddetto **Snooping cache**: le componenti della cache e della memoria sono agganciati a un **broadcast medium** (= interconnection network, è un canale di comunicazione broadcast) tramite un controller, che ha la responsabilità di osservare le transizioni di stato degli altri componenti e di far reagire il componente locale di conseguenza (appunto invalidando o modificando la replica locale del blocco aggiornato). Naturalmente, il controller utilizza il broadcast medium anche per trasmettere gli aggiornamenti che avvengono in un blocco di cache locale. Mediante il broadcast medium siamo in grado di serializzare le transizioni di stato. *Inoltre, una transizione di stato non può occorrere finché il broadcast medium non viene acquisito in uso dal controller.*



Nel mondo reale viene adottato lo *Snooping cache* accoppiato con un protocollo basato sull'invalidazione. Vediamo dunque nel dettaglio alcuni di questi protocolli basati sull'invalidazione. Ad esempio, una componente di cache  $L_1$  ha una replica  $r$ , vuole eseguire operazione su tale replica mediante il richiamo di broadcast medium, in cui annuncio alle altre componenti che, chi ha una copia della replica  $r$ , deve aggiornare. Nel mentre, nessuno può eseguire altri aggiornamenti se ho io il broadcast medium. Ciò crea un problema di scalabilità, se lavoro con tanti dati e repliche.

**Osservazione:** devo per forza aggiornare le repliche? E se non mi servissero? Ad esempio, un thread su  $CPU_0$  con cache  $L_1$  vi esegue un aggiornamento su linea di memoria e comunica in maniera distribuita di cancellare quella linea di cache. Ma se quella linea di cache fosse già invalida per altri? Si mantiene traccia di queste operazioni mediante:

### Protocollo MSI (Modified-Shared-Invalid)

È un protocollo in cui qualsiasi transazione di scrittura su una **copia di un cache block** invalida tutte le altre copie del cache block. Quando dobbiamo servire delle transazioni di lettura:

- Se la policy della cache è **write-through**, recuperiamo semplicemente l'ultima copia aggiornata dalla memoria.
- Se la policy della cache è **write-back**, recuperiamo l'ultima copia aggiornata o dalla memoria o da un altro componente di caching.

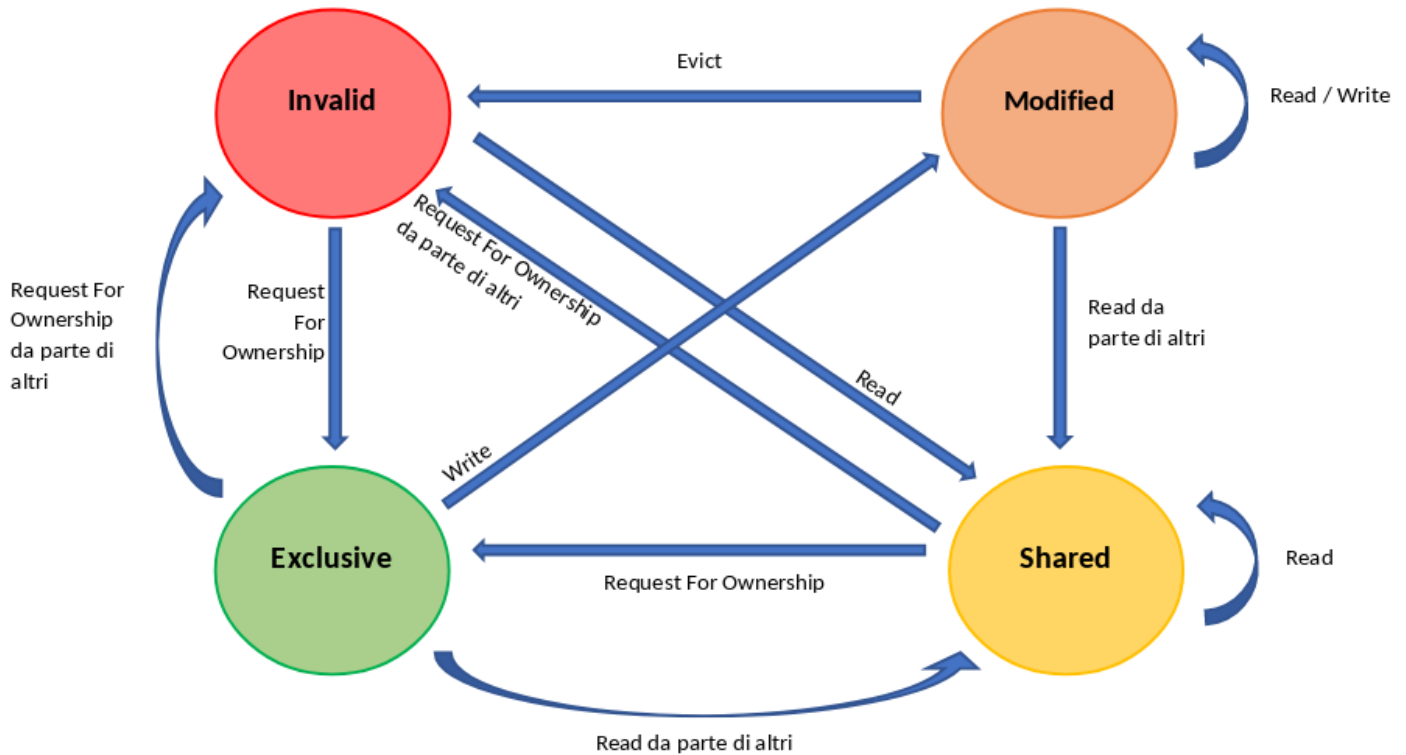
In questo protocollo è necessario tenere traccia dello stato della copia di un blocco:

- Si trova nello stato *modified* se è appena stata sovrascritta.
- Si trova nello stato *shared* se esistono dei CPU-core che stanno leggendo altre copie del medesimo blocco (per cui esistono più copie valide di tale blocco).
- Si trova nello stato *invalid* se un'altra copia del medesimo blocco è appena stata sovrascritta.

Il **problema** di MSI risiede nel fatto che richiede l'invio di un messaggio di invalidazione in broadcast (tramite il broadcast medium) ogni volta che una copia di un blocco viene modificata (e questo anche in caso di scritture successive sulla medesima copia). Questo può portare a impegnare massivamente il broadcast medium anche quando non ce n'è bisogno, perché magari tutte le altre copie erano nello stato invalid già da prima. Per ovviare a tale inconveniente, si ricorre ai protocolli descritti successivamente.

### Protocollo MESI (Modified-Exclusive-Shared-Invalid):

Rispetto a MSI, prevede uno stato in più, che è **exclusive**. In pratica, un CPU-core  $i$  che vuole apportare delle modifiche alla propria copia di un blocco deve richiedere l'**ownership** (l'esclusività) di quel blocco; questa è l'unica occasione in cui CPU-core  $i$  utilizza il broadcast medium per comunicare agli altri nodi che devono invalidare la propria copia del blocco. Una volta che la copia è nello stato *exclusive*, CPU-core  $i$  può modificarla tutte le volte che vuole senza dover comunicare più nulla a nessuno, finché un altro CPU-core non richiede di effettuare una lettura da un'altra copia dello stesso cache block (momento in cui la copia) passa allo stato *shared*). L'automa raffigurato di seguito descrive in modo completo il funzionamento del protocollo MESI. Schematizzando:

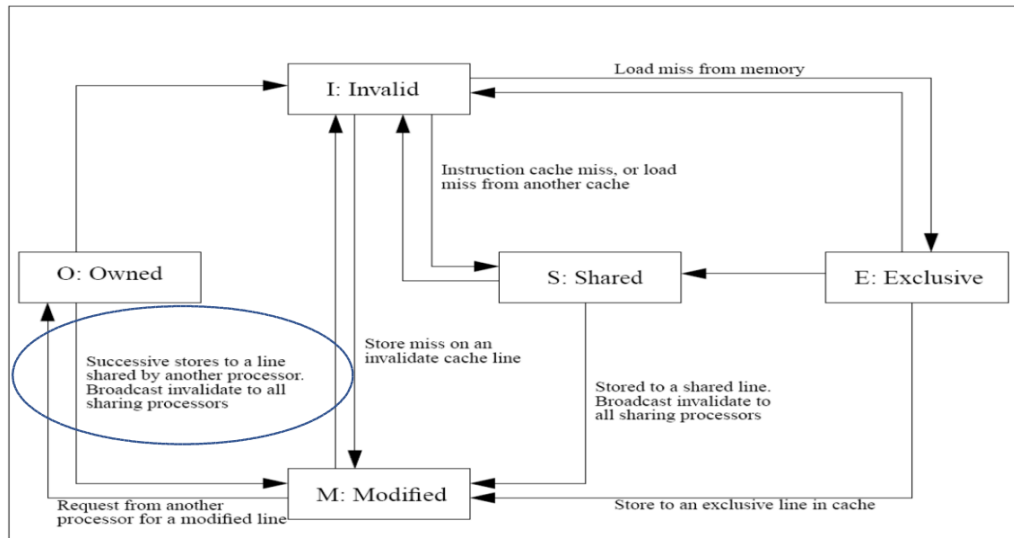


Ogni volta che si esce dallo stato modified, si effettua il *write back* in memoria delle nuove informazioni che sono state scritte. Se volessi solamente leggere? Chi ha l'informazione si trova in *exclusive* e poi passa a *shared*. Potrei introdurre un quinto stato per evitare queste due transizioni. La linea *exclusive* è a mio uso esclusivo. Solo lo stato Invalid crea interazioni distribuite. Lo stato Modified è importante per cache write back, ma non devo informare tutti.

### Protocollo MOESI (Modified-Owned-Exclusive-Shared-Invalid):

Rispetto a MESI, prevede un ulteriore stato in più, che è **owned**. In particolare, una copia di un cache block passa dallo stato modified allo stato *owned* (anziché shared) nel momento in cui stava per essere sovrascritta ma un'altra copia dello stesso blocco viene letta. *Owned* significa che la copia è tuttora in fase di aggiornamento ma, nel frattempo, altre copie dello stesso blocco vengono lette: se devono occorrere nuovi aggiornamenti quando si è nello stato *owned*, non ci si deve preoccupare di chiedere nuovamente l'esclusività del blocco, bensì si passa direttamente allo stato modified, invalidando tutte le altre copie. Di seguito è rappresentato l'automa completo del protocollo MOESI.

No need to pass through "exclusive" again



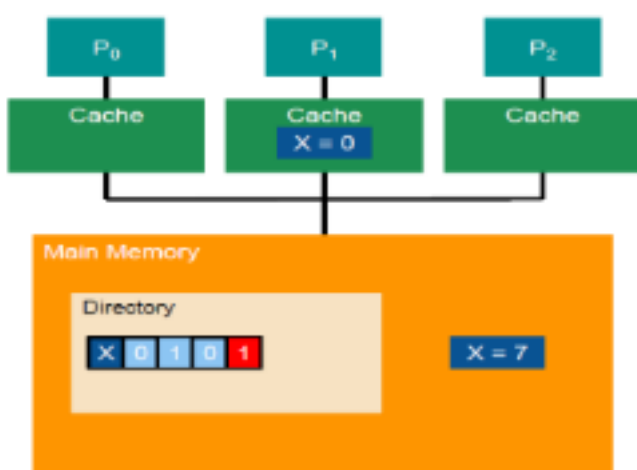
Nello stato **Owned** io ho l'uso esclusivo di un dato. Se mi chiedono di condividerlo, gli altri **non** possono scriverci. Tale stato ci permette di transitare tra più stati, in quanto non devo riacquisire l'esclusività.

**Protocolli directory based:** In realtà, i protocolli basati su snooping cache non scalano poiché le transazioni portano a una comunicazione broadcast. Per questo motivo sono stati introdotti i protocolli directory based, in cui i vari aggiornamenti possono essere point-to-point tra i vari CPU-core / processori. In particolare, supponiamo di avere un sistema con  $N$  processori  $P_0, P_1, \dots, P_{N-1}$ . Per ciascun blocco di memoria si mantiene una directory entry, composta da:

- $N$  bit di presenza (l' $i$ -esimo bit è impostato a 1 se il blocco si trova nella cache di  $P_i$ ).
- 1 dirty bit, che indica se il blocco è stato **modificato** senza che gli aggiornamenti siano stati riportati in memoria.

Se dirty bit = 0 e ho più possessori, linea cache condivisa.

Dirty bit = 1 se qualcuno lo scrive, però per passare a condiviso deve ritornare a 0. E' una machera che mi dice chi può fare cosa.



- **Caso 1:**  $P_0$  vuole leggere il blocco  $X$ , ha un cache miss e il dirty bit è pari a 0.  $X$  viene letto dalla memoria, il bit **presence**[0] viene settato a 1 e i dati vengono consegnati al lettore.
- **Caso 2:**  $P_0$  vuole leggere il blocco  $X$ , ha un cache miss e il dirty bit è pari a 1.  $X$  viene richiamato dal processore  $P_j$  tale che **presence**[ $j$ ]=1. Dopodiché il dirty bit viene settato a 0, il blocco viene condiviso e il bit **presence**[0] viene settato a 1. Infine, i dati vengono consegnati al lettore e propagati in memoria.
- **Caso 3:**  $P_0$  vuole scrivere sul blocco  $X$ , ha un cache miss e il dirty bit è pari a 0.  $P_0$  invia un messaggio di invalidazione non a tutti gli altri processori, bensì solo a quelli col bit **presence**[ $j$ ] pari a 1. Dopodiché tali **presence**[ $j$ ] vengono settati a 0, e **presence**[0] e il dirty bit vengono impostati a 1.
- **Caso 4:**  $P_0$  vuole scrivere sul blocco  $X$ , ha un cache miss e il dirty bit è pari a 1.  $X$  viene richiamato dal processore  $P_j$  tale che **presence**[ $j$ ]=1. Dopodiché **presence**[ $j$ ] viene settato a 0 e **presence**[0] viene impostato a 1.