

Arquitectura de Computadoras

TEMA 4

Lanzamiento múltiple, Límites de ILP,
Multithreading

Abril 2016/20

- o Introducción: $CPI < 1$
- o Lanzamiento múltiple de instrucciones: Superescalar, VLIW
- o Superescalar simple
- o VLIW
- o Superescalar con planificación dinámica
- o Límites de ILP
- o Ejemplo: Implementaciones X86
- o Thread Level Parallelism y Multithreading

- o Bibliografía
 - o Capítulo 3 y 4 de [HePa07]
 - o Capítulos 4 , 6 y 7 de [SiFK97]

□ Introducción

- ¿ Por que limitar a una instrucción por ciclo?
- Objetivo: $CPI < 1$
- Lanzar y ejecutar simultáneamente múltiples instrucciones por ciclo
- ¿Tenemos recursos?
 - Más área de silicio disponible
 - Técnicas para resolver las dependencias de datos (*planificación*)
 - Técnicas para resolver las dependencias de control (*especulación*)

□ Alternativas

- Procesador Superescalar con planificación estática
- Procesador Superescalar con planificación dinámica+(especulación)
- Procesadores VLIW (very long instruction processor)
- ✓ Superescalar
 - ✓ Lanza de 1 a 8 instrucciones por ciclo
 - ✓ Reglas de ejecución
 - o Ejecución en orden-planificación estática
 - o Ejecución fuera de orden-planificación dinámica
- ✓ VLIW
 - ✓ Numero fijo de instrucciones por ciclo
 - ✓ Planificadas estáticamente por el compilador
 - ✓ EPIC (Explicitly Parallel Instruction Computing) Intel/HP

□ Alternativas

Tipo	Forma del issue	Detección de azares	Planificación	Ejemplos
Superescalar estático	Dinámico	HW	estática	Embedded MIPS, ARM
Superescalar dinámico	Dinámico	HW	dinámica	ninguno
Superescalar especulativo	Dinámico	HW	Dinámica con especulación	P4, Core2, Power5, SparcVI
VLIW	Estático	Básicamente SW	estática	TI C6x Itanium

❑ SUPERESCALAR

Grado 2
2 Vías

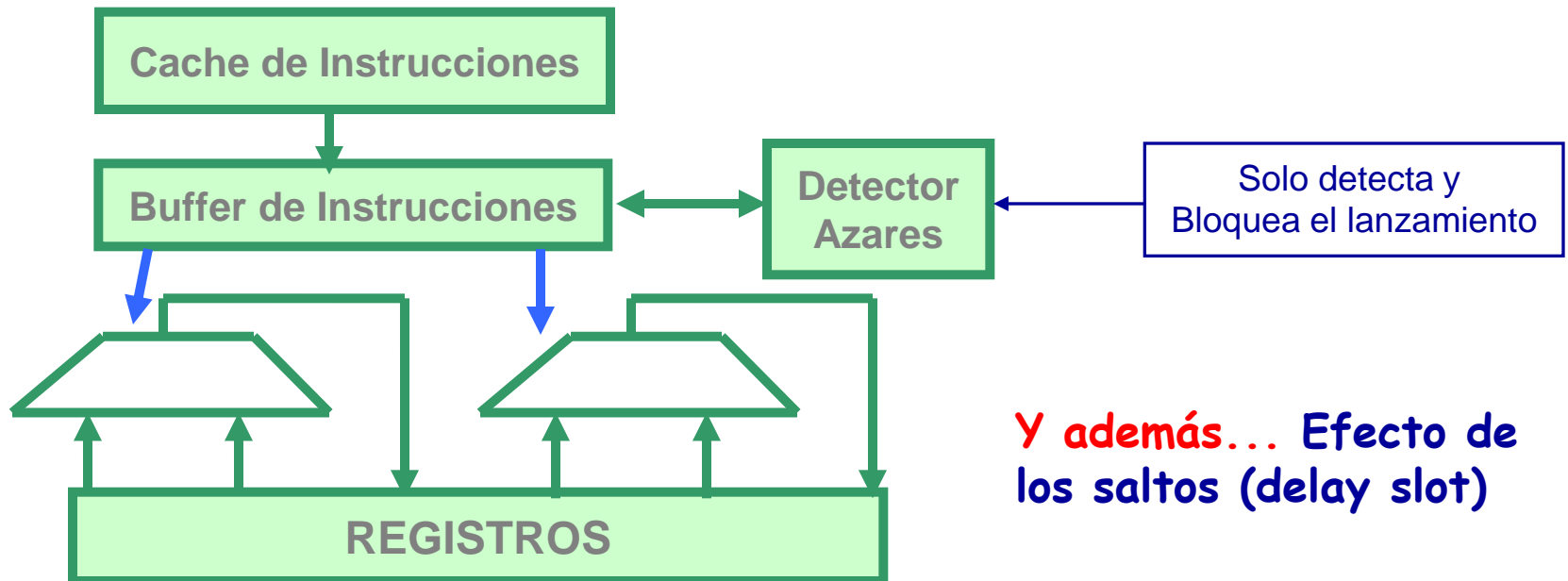


Instruction	1	2	3	4	5	6	7
i	IF	ID	EX	MEM	WB		
i+1	IF	ID	EX	MEM	WB		
i+2		IF	ID	EX	MEM	WB	
i+3		IF	ID	EX	MEM	WB	
i+4			IF	ID	EX	MEM	WB
i+5			IF	ID	EX	MEM	WB

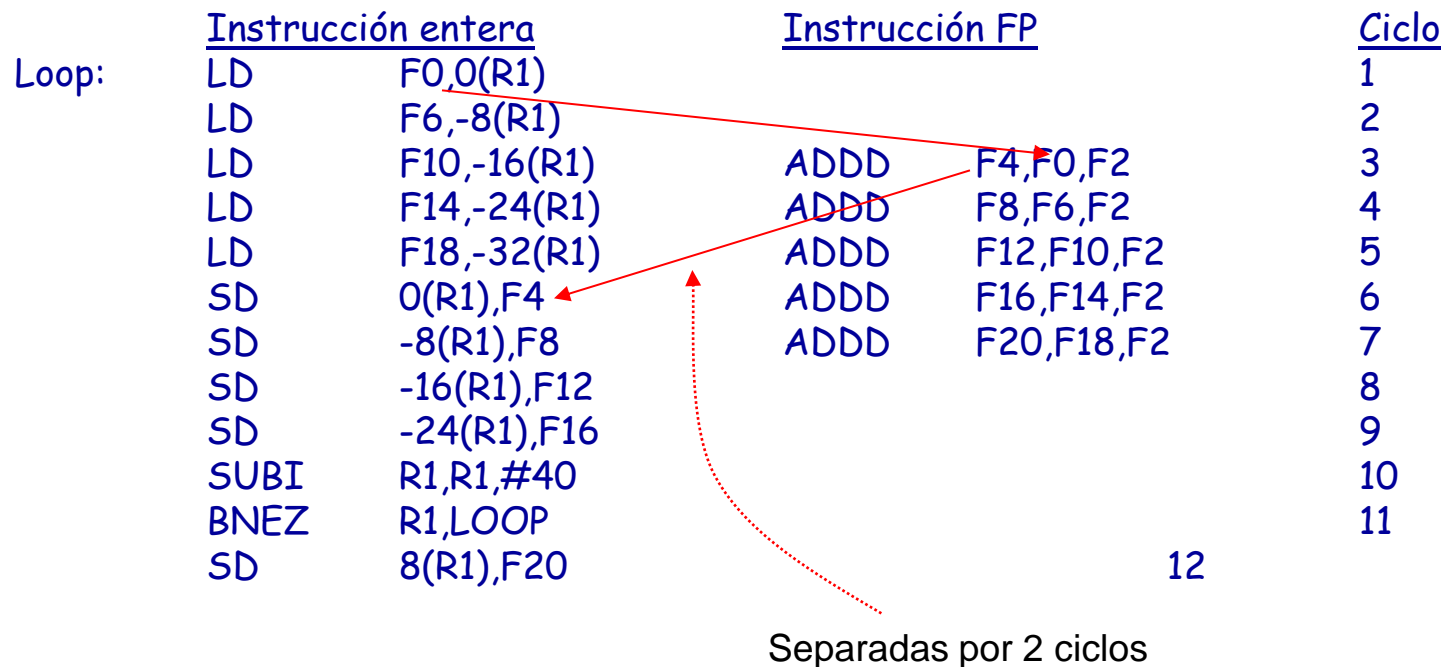
- **Duplicar todos los recursos:**
 - o Puertas bloque de Registros
 - o Fus
 - o Puertas de memoria,...
 - o Control
- **Ideal CPI= 0.5 se reduce por:**
 - o Saltos
 - o LOADs
 - o Dependencias verdaderas LDE
- **Necesita:**
 - o Predicción sofisticada
 - o Tratamiento de LOAD; Cargas especulativas, técnicas de prebúsqueda
- **Más presión sobre la memoria**
- **Efecto incremental de los riesgos**
- **Se puede reducir complejidad con limitaciones (Un acceso a memoria por ciclo)**

❑ SUPERESCALAR Simple (estático, en orden)

- Regla de lanzamiento: Una instrucción FP+ una instrucción de cualquier otro tipo
- Buscar y decodificar dos instrucciones por ciclo (64 bits)
Ordenamiento y decodificación
Se analizan en orden. Sólo se lanza la 2ª si se ha lanzado la 1ª (conflictos)
- Unidades funcionales segmentadas (una ope. por ciclo) ó múltiples (división, raíz), más puertas en el bloque de registros
- Lanzamiento simple, recursos no conflictivos (diferentes reg y UF,..), excepto
Conflictos de recursos; load, store, move FP → más puertas en el bloque de reg.
Conflictos de datos LDE → más distancia entre instrucciones.



□ SUPERESCALAR Simple (estático, en orden)



- Desarrollo para ejecución superescalar: se desarrolla una iteración más.
12 ciclos, 2.4 ciclos por iteración
- El código máquina está compacto en la memoria

❑ SUPERESCALAR Simple (estático, en orden)

➤ *Ventajas*

- No modifica código. Compatibilidad binaria
- No riesgos en ejecución

➤ *Desventajas*

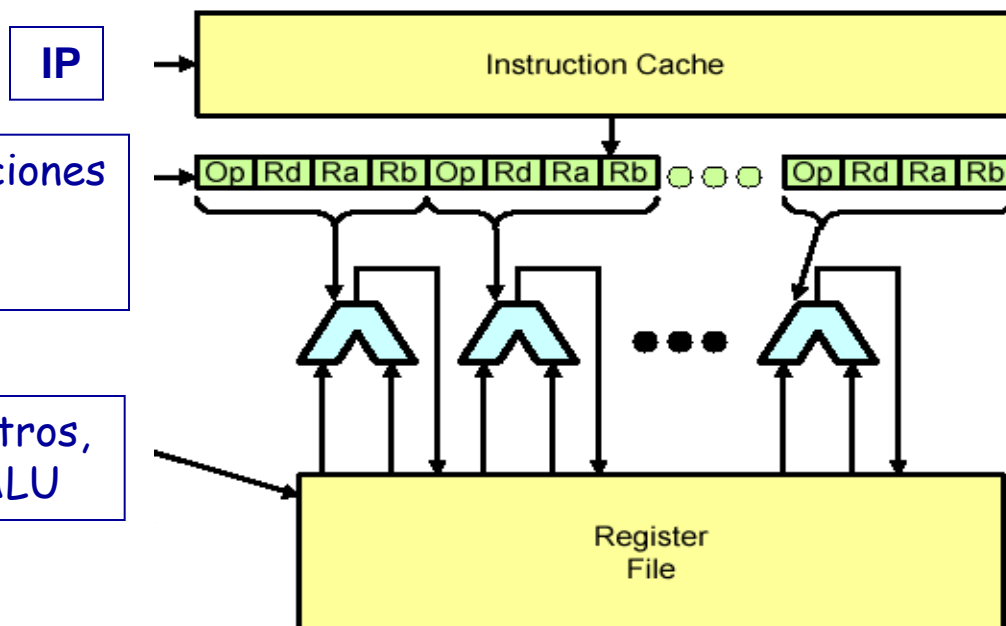
- Mezcla de instrucciones. Solo obtiene CPI de 0.5 en programas con 50 % de FP
- Bloqueos en el lanzamiento
- Planificación fija: No puede adaptarse a cambios en ejecución (Fallos de cache)
- Los códigos deben de ser replanificados para cada nueva implementación (eficiencia)

□ VLIW

- El análisis de dependencias en tiempo de compilación
- Muchas operaciones por instrucción (IA64 packet, Tramsmeta molecula)
- Todas las operaciones de una instrucción se ejecutan en paralelo
- Instrucciones con muchos bits
- Muchas operaciones vacías (NOP)

Instrucción: Incluye varias instrucciones convencionales de tres operandos una por ALU

Bloque de registros,
3 puertas por ALU



Más ILP: Lanzamiento múltiple

□ VLIW Ejemplo Tema3

```

LOOP    LD      F0,0(R1)
        ADDD    F4,F0,F2
        SD      0(R1),F4
        SUBI    R1,R1,#8
        BNEZ    R1,LOOP
  
```

- Aplicar técnicas conocidas para minimizar paradas

- Unrolling
- Renombrado de registros



- Latencias de uso: LD a ADD 1 ciclo, ADD a SD 2 ciclos
- Opción: desarrollar 4 iteraciones y planificar: 14 ciclos, 3.5 ciclos por iteración

```

LOOP:   LD      F0, 0(R1)
        LD      F6, -8(R1)
        LD      F10, -16(R1)
        LD      F14, -24(R1)
        ADDD    F4, F0, F2
        ADDD    F8, F6, F2
        ADDD    F12, F10, F2
        ADDD    F16, F14, F2
        SD      0(R1), F4
        SD      -8(R1), F8
        SD      -16(R1), F12
        SUBI    R1, R1, #32
        BNEZ    R1, LOOP
        SD      8(R1), F16; 8-32 = -24
  
```

Más ILP: Lanzamiento múltiple

□ VLIW

Loop unrolling en VLIW

```

LOOP:  LD F0,0(R1)           ; F0 = array element
        ADDD F4,F0,F2        ; add scalar in F2
        SD 0(R1),F4          ; store result
        SUBI R1,R1,#8        ; decrement pointer
        BNEZ R1, LOOP        ; branch if R1!=0
  
```

<u>Mem ref 1</u>	<u>Mem ref 2</u>	<u>FP op</u>	<u>FP op</u>	<u>Int op/branch</u>
LD F0,0(R1)	LD F6,-8(R1)			
LD F10,-16(R1)	LD F14,-24(R1)			
LD F18,-32(R1)	LD F22,-40(R1)	ADDD F4,F0,F2	ADDD F8,F6,F2	
LD F26,-48(R1)		ADDD F12,F10,F2	ADDD F16,F14,F2	
		ADDD F20,F18,F2	ADDD F24,F22,F2	
SD 0(R1),F4	SD -8(R1),F8	ADDD F28,F26,F2		
SD -16(R1),F12	SD -24(R1),F16			
SD -32(R1),F20	SD -40(R1),F24			SUBI R1,R1,#56
SD 8(R1),F28				BNEZ R1, LOOP

- ✓ 7 iteraciones en 9 ciclos: 1.3 ciclos por iteración
- ✓ 23 operaciones en 45 slots (<50% de ocupación)
- ✓ Muchos registros necesarios

□ VLIW

VENTAJAS

- Hardware muy simple
 - No detecta dependencias
 - Lógica de lanzamiento simple
- Puede explotar paralelismo a todo lo largo del programa

DESVENTAJAS

- Planificación estática; Muy sensible a fallos de cache
- Necesita desenrollado muy agresivo
- Bloque de registros muy complejo en área y tiempo de acceso
- Muchas NOP
 - Poca densidad de código
 - Capacidad y AB de la cache de instrucciones
- Compilador muy complejo
- No binario compatible
- Operación síncrona para todas las operaciones de una instrucción

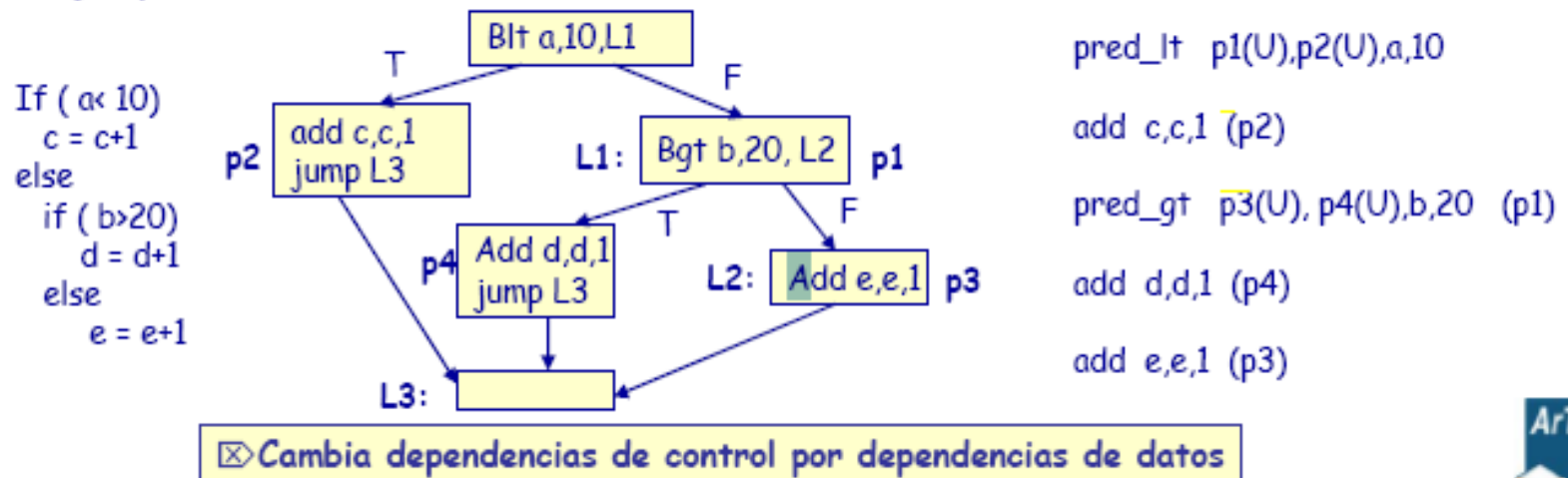
Ejecución Predicada

Ejecución con predicados

Idea básica

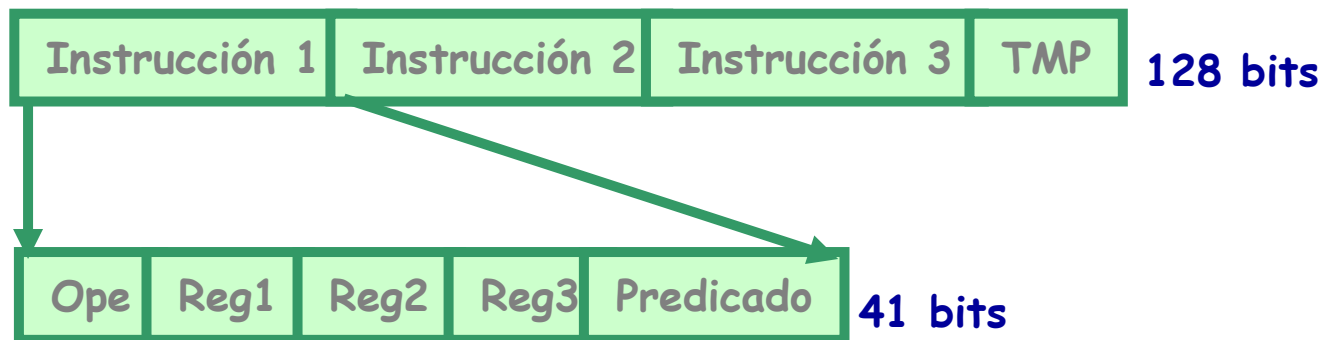
- Transformar todas las **instrucciones en condicionales**
- Una instrucción de ejecución condicional está formada por:
 - Una parte de condición, denominada **predicado o guarda**
 - Una parte de operación
- Ejecución predicada:
 - Si la condición es **cierta** \Rightarrow La instrucción se ejecuta
 - Si la condición es **falsa** \Rightarrow La instrucción se comporta como NOP
- Transforma dependencias de control en dependencias de datos y elimina fallos de predicción.

Ejemplo



□EPIC: Explicitly Parallel Instruction Computing IA64

- Instrucciones de 128 bits
 - Operaciones de tres operandos
 - TMP codifica dependencias entre las operaciones
 - 128 registros enteros (64bits), 128 registros FP (82bits)
 - Ejecución predicada. 64 registros de predicado de 1 bit
 - Cargas especulativas
 - Hw para chequeo de dependencias



Primera implementación Itanium (2001), 6 operaciones por ciclo, 10 etapas, 800Mhz

Segunda implementación Itanium2 (2005), 6 operaciones por ciclo, 8 etapas, 1,66Ghz

Más ILP: Lanzamiento múltiple

□ SUPERESCALAR con Planificación Dinámica. Fuera de orden

➤ Un Diseño Simple

- Estaciones de reserva separadas para enteros (+reg) y PF (+reg)
- Lanzar dos instrucciones en orden (ciclo de lanzamiento: partir en dos subciclos)
- Solo FP load causan dependencias entre instrucciones enteras y PF
 - Reemplazar buffer de load con cola. Las lecturas se hacen en orden
 - Ejecución Load: "check" dirección en cola de escritura para evitar LDE
 - Ejecución Store: "check" dirección en cola de lecturas para evitar EDL

➤ Rendimiento del procesador

<u>Iteración</u> <u>no.</u>	<u>Instrucción</u>	<u>Lanzada</u>	<u>Ejecutada</u> (número de ciclo)	<u>Escribe resultado</u>
1	LD F0,0(R1)	1	2	4
1	ADDD F4,F0,F2	1	5	8 ←
1	SD 0(R1),F4	2	9	
1	SUBI R1,R1,#8	3	4	5
1	BNEZ R1,LOOP	4	6	
2	LD F0,0(R1)	5	6	8 ←
2	ADDD F4,F0,F2	5	9	12
2	SD 0(R1),F4	6	13	
2	SUBI R1,R1,#8	7	8	9
2	BNEZ R1,LOOP	8	10	

4 ciclos por
iteración

❑ Más paralelismo

Procesador superescalar con:
Planificación dinámica
Predicción de saltos agresiva



Especulación

➤ **Especulación:**

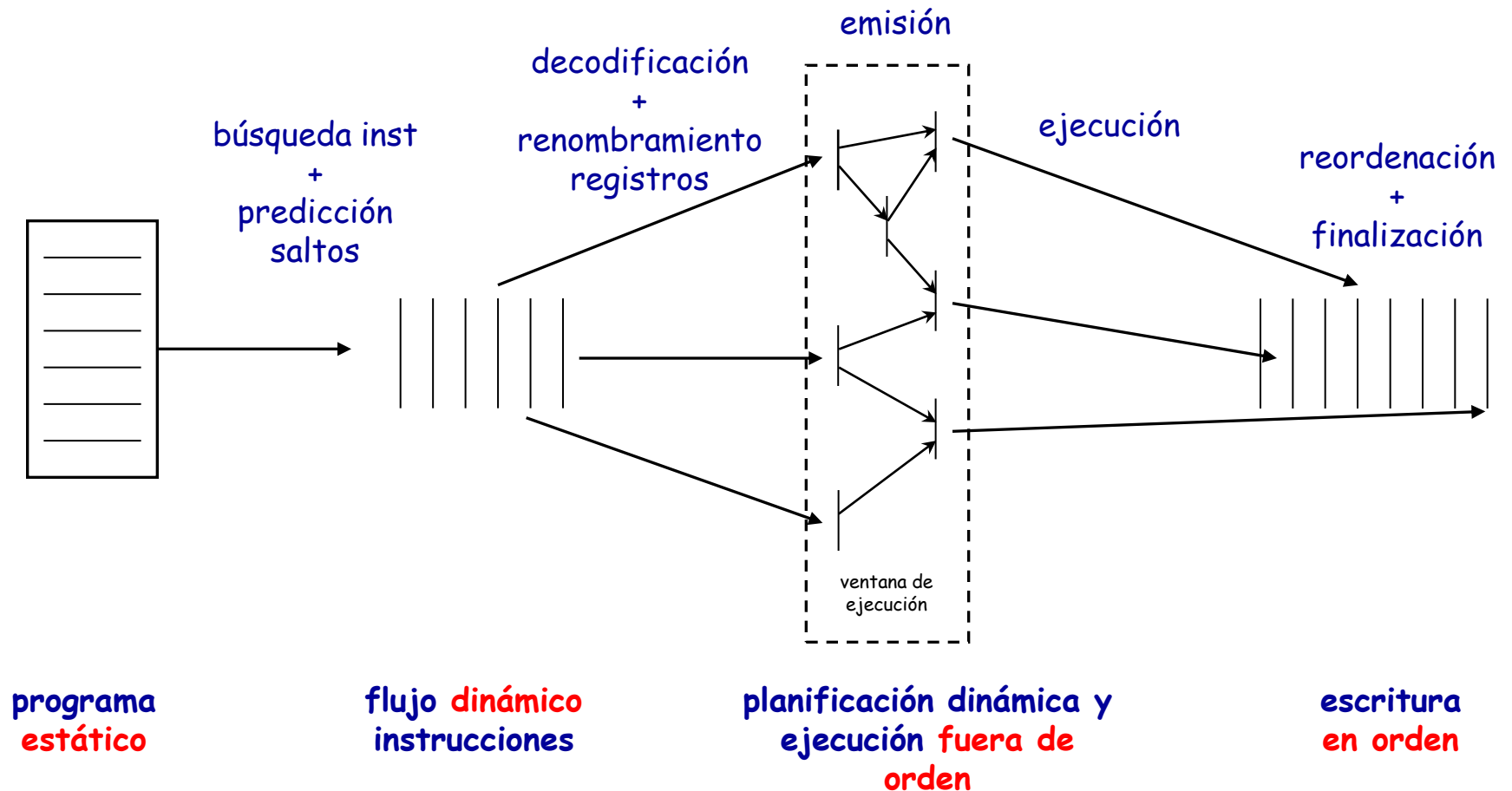
- Permite lanzar instrucciones dependientes de una predicción de un salto sin consecuencias
 - HW "undo"
 - Tratamiento de excepciones
- Una instrucción deja de ser especulada cuando su salto es resuelto (inst. *COMMIT*)
- Ejecución fuera de orden, finalización en orden

➤ **Características**

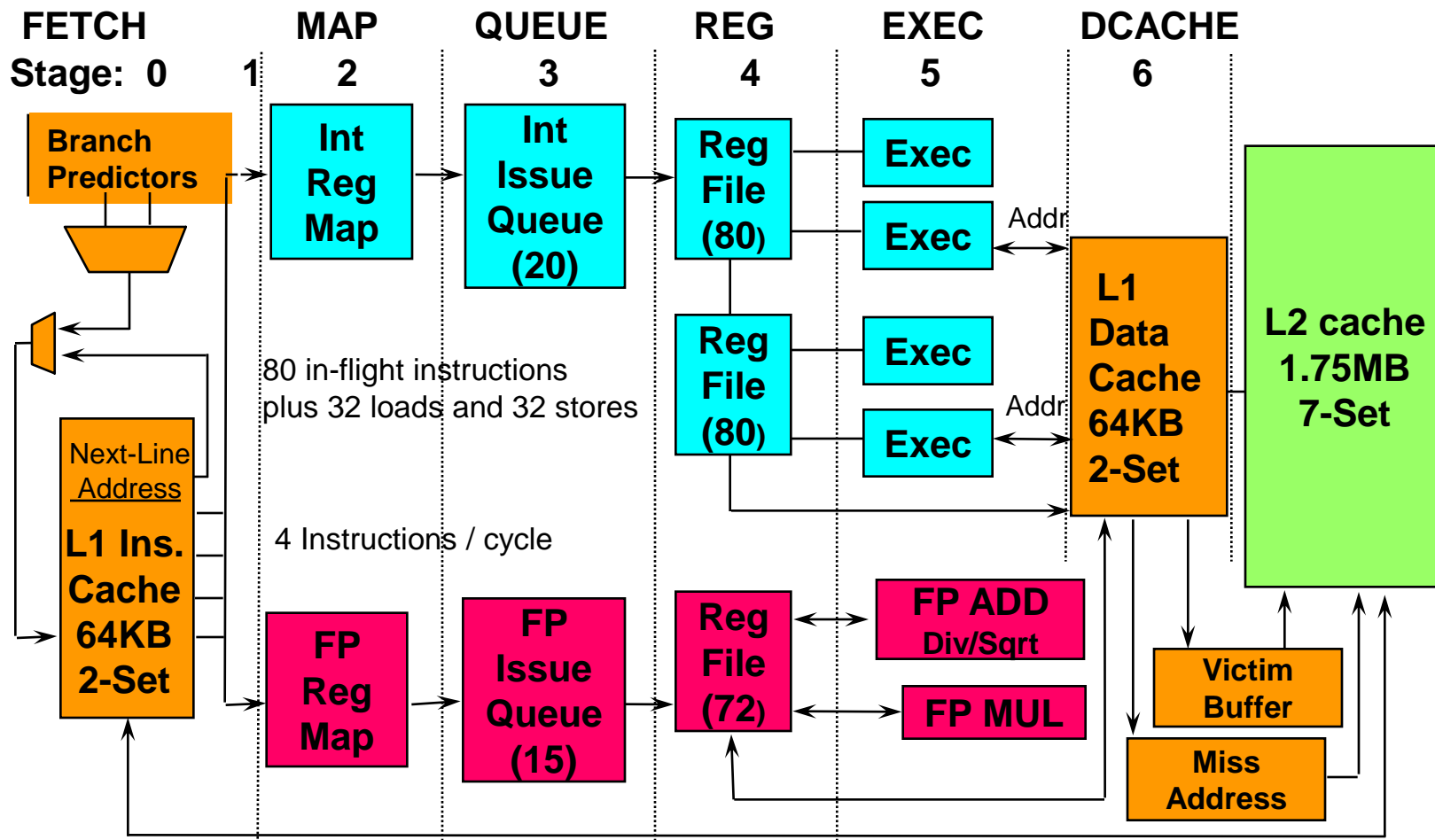
- Determina el orden de ejecución en el momento del lanzamiento
- Planifica con conocimiento de las latencias variables (fallos de cache)
- Compatibilidad binaria (mejor recompilar)
- >>> Complejidad HW

❑ SUPERESCALAR con Planificación Dinámica y Especulación

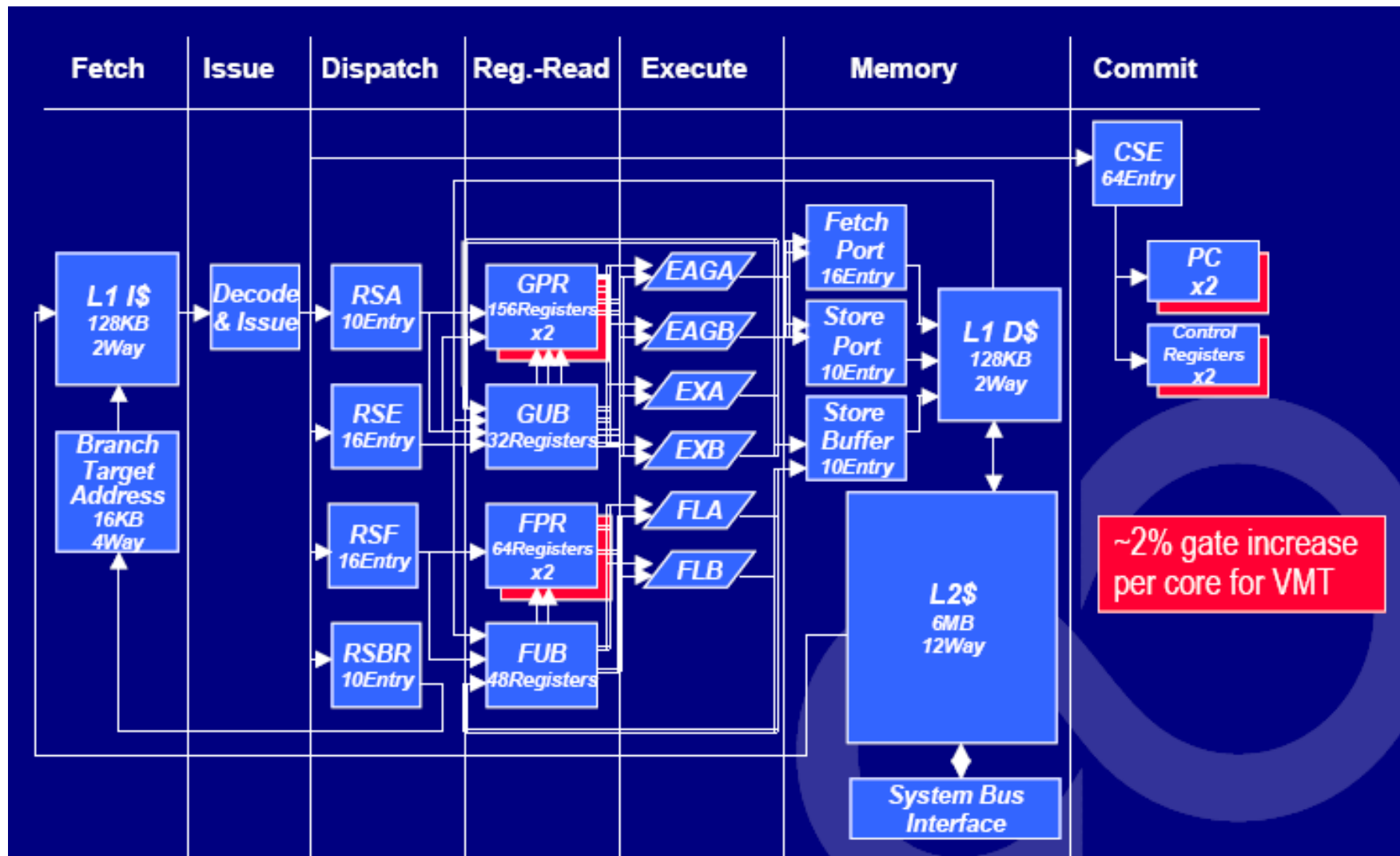
Ejecución fuera de orden. Finalización en orden



□ EV7 ALPHA 21364 Core (2003)



❑ SPARC64 VI (2006/7)



```
❑ #define N 9984
double x[N+8], y[N+8], u[N]
loop () {
    register int i;
    double q;
    for (i=0; i<N; i++) {
        q = u[i] * y[i];
        y[i] = x[i] + q;
        x[i] = q - u[i] * x[i];
    }
}
```

Operaciones por iteración:

- o 3 lecturas (u,y,x)
- o 2 escrituras (y,x)
- o 2 multiplicaciones
- o 1 suma
- o 1 resta

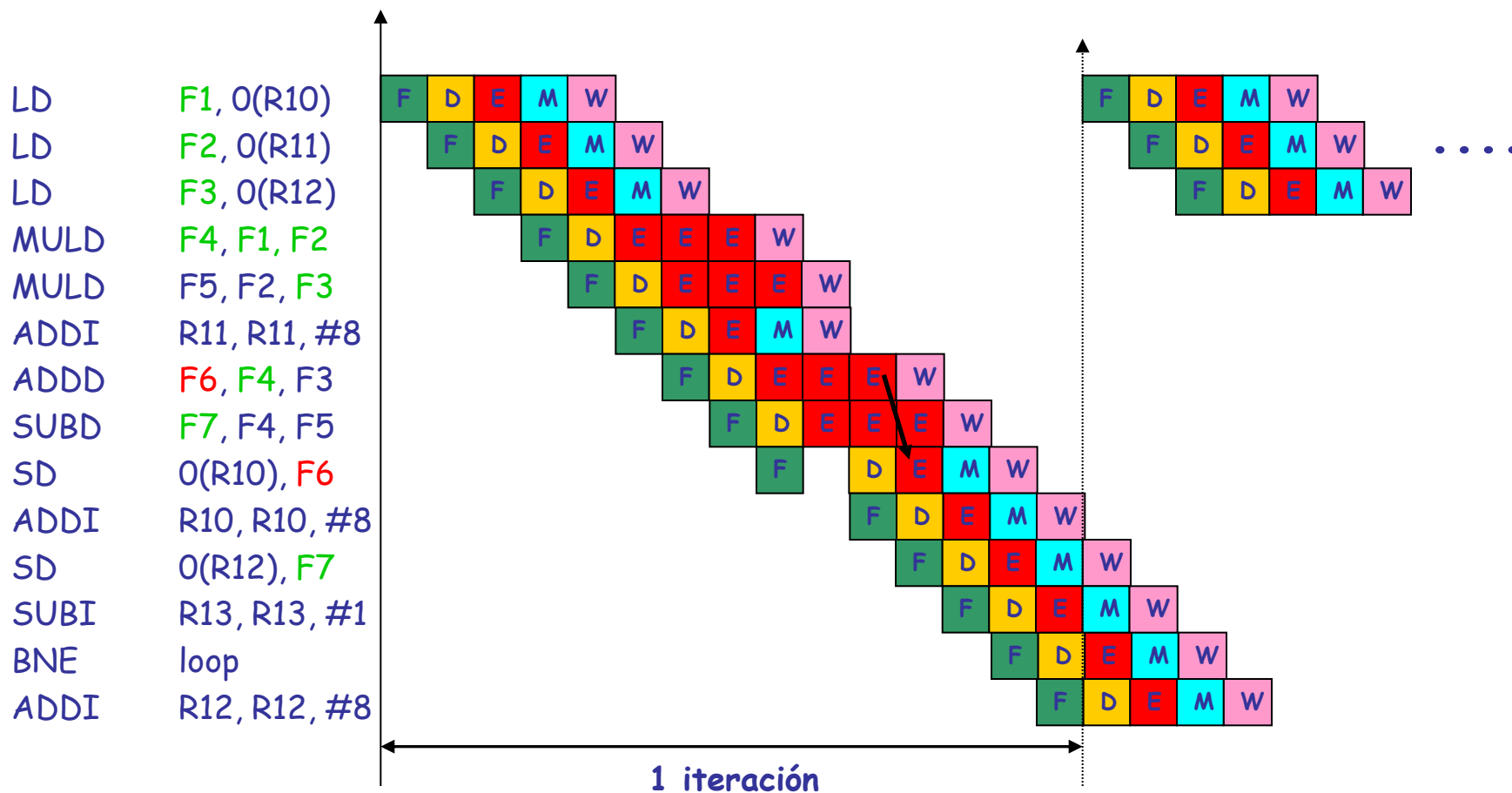
```

LD      R10, @y[0] ; dirección inicial del vector y a R10
LD      R11, @u[0] ; Idem u
LD      R12, @x[0] ; Idem x
LD      R13, N      ; numero de componentes a R13
loop:   LD      F1, 0(R10) ; elemento de y a F1
LD      F2, 0(R11) ; elemento de u a F1
LD      F3, 0(R12) ; elemento de x a F1
MULD    F4, F1, F2 ; q = u x y
MULD    F5, F2, F3 ; u x x
ADDI    R11, R11, #8 ; actualiza puntero de u
ADDD    F6, F4, F3 ; y = x + q
SUBD    F7, F4, F5 ; x = q - u x x
SD      0(R10), F6 ; almacena y
ADDI    R10, R10, #8 ; actualiza puntero de y
SD      0(R12), F7 ; almacena x
SUBI    R13, R13, #1 ; actualiza N
BNE     loop
ADDI    R12, R12, #8 ; delay slot actualiza puntero de x

```

- Operaciones en coma flotante 3 ciclo de latencia
- Operaciones enteras 1 ciclo de latencia
- Operaciones segmentadas
- Una unidad de cada tipo
- Saltos efectivos en etapa Decode (D)

□ Ejecución escalar con una unidad entera y una de PF



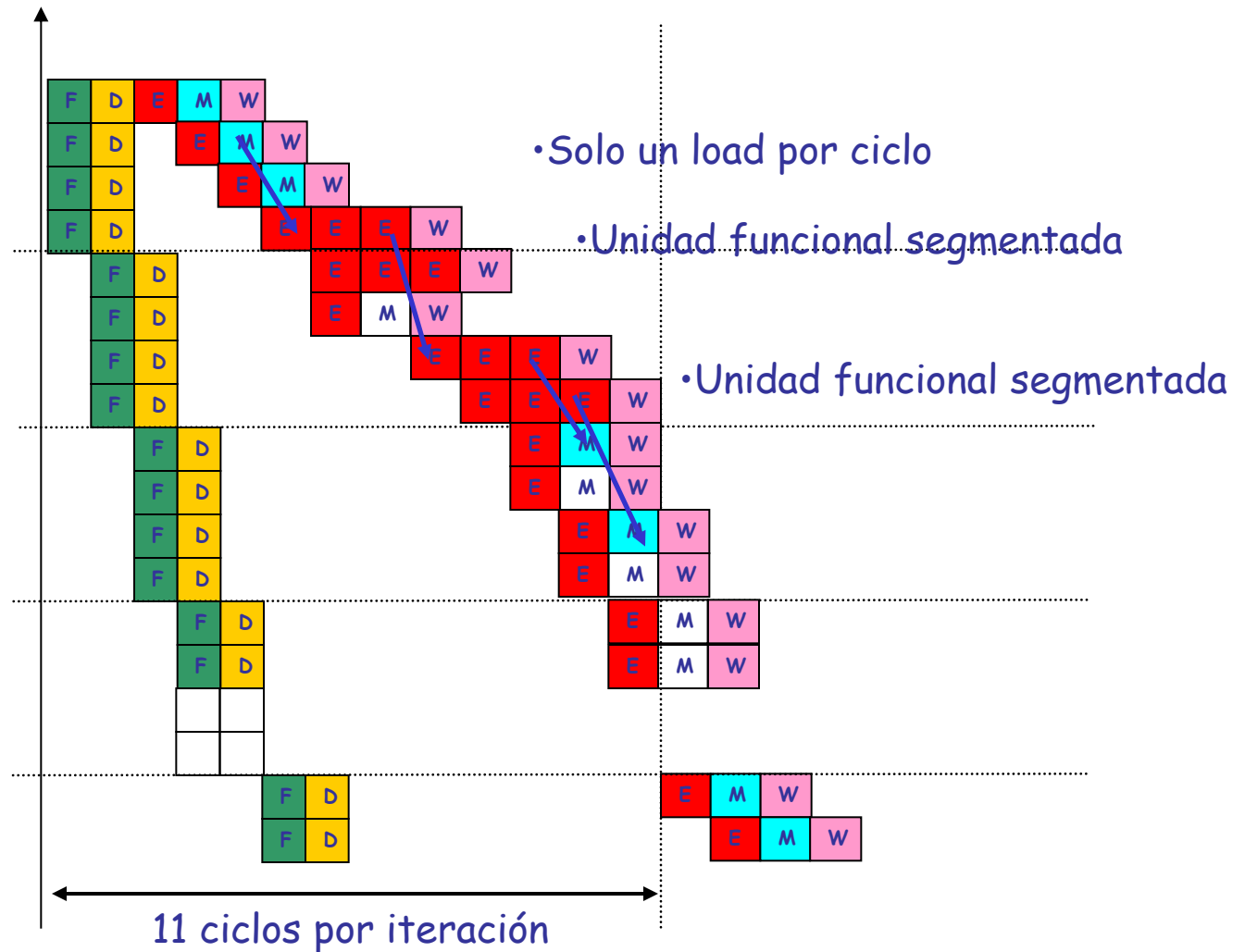
•Una iteración de 14 instrucciones cada 15 ciclos

Modelos de Ejecución

- Superescalar de 4 vías
- Ejecución en orden

```

LD      F1, 0(R10)
LD      F2, 0(R11)
LD      F3, 0(R12)
MULD    F4, F1, F2
MULD    F5, F2, F3
ADDI    R11, R11, #8
ADD    F6, F4, F3
SUBD    F7, F4, F5
SD      0(R10), F6
ADDI    R10, R10, #8
SD      0(R12), F7
SUBI    R13, R13, #1
BNE     loop
ADDI    R12, R12, #8
  
```

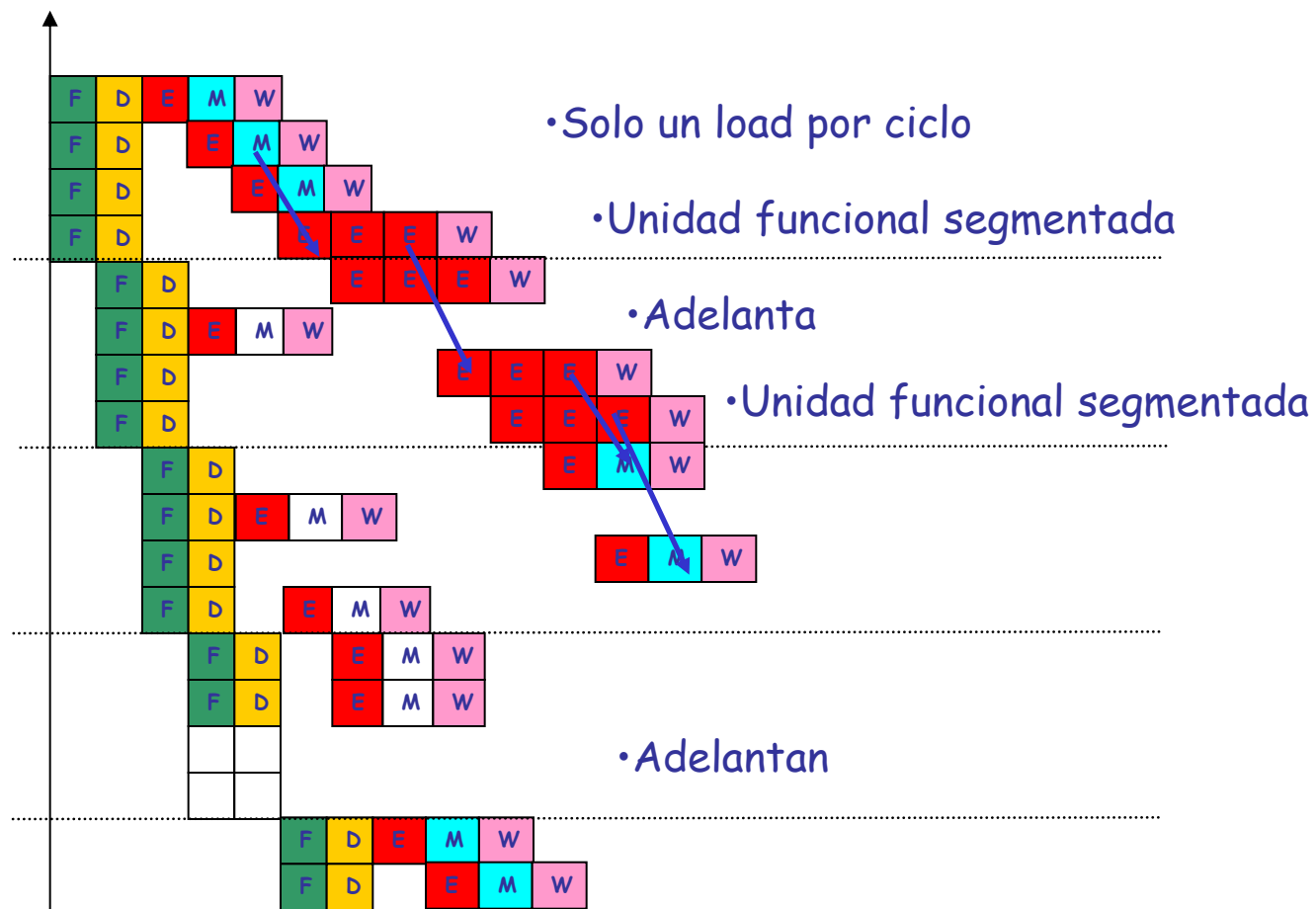


Modelos de Ejecución

- Superescalar de 4 vías
- Ejecución fuera de orden
- Predicción de saltos

```

LD      F1, 0(R10)
LD      F2, 0(R11)
LD      F3, 0(R12)
MULD    F4, F1, F2
MULD    F5, F2, F3
ADDI    R11, R11, #8
ADDD    F6, F4, F3
SUBD    F7, F4, F5
SD      0(R10), F6
ADDI    R10, R10, #8
SD      0(R12), F7
SUBI    R13, R13, #1
BNE     loop
ADDI    R12, R12, #8
  
```



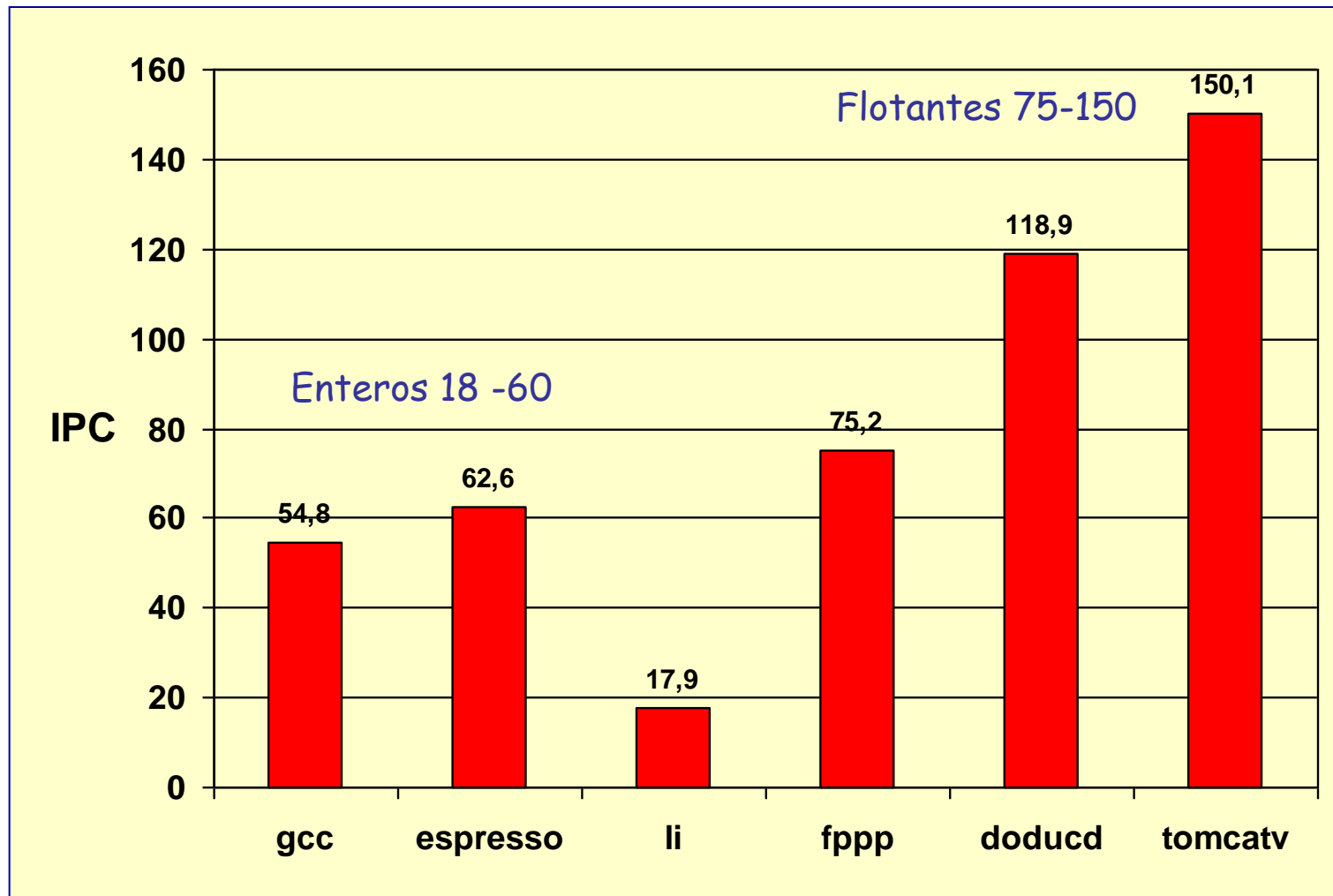
- ❑ El lanzamiento múltiple permite mejorar el rendimiento sin afectar al modelo de programación
- ❑ En los últimos años se ha mantenido el mismo ancho superescalar que tenían los diseños del 1995
- ❑ El gap entre rendimiento pico y rendimiento obtenido crece
- ❑ ¿Cuanto ILP hay en las aplicaciones?
- ❑ ¿Necesitamos nuevos mecanismos HW/SW para explotarlo?
 - o Extensiones multimedia:
 - o Intel MMX,SSE,SSE2,SSE3, SSE4
 - o Motorola AltiVec, Sparc, SGI, HP

- ❑ ¿Cuanto ILP hay en las aplicaciones?
- ❑ Supongamos un procesador superescalar fuera de orden con especulación y con recursos ilimitados
 - o Infinitos registros para renombrado
 - o Predicción perfecta de saltos
 - o Caches perfectas
 - o Lanzamiento no limitado
 - o Desambiguación de memoria perfecta

❑ Modelo versus procesador real

	Modelo	Power 5
Instrucciones lanzadas por ciclo	Infinito	4
Ventana de instrucciones	Infinita	200
Registros para renombrado	Infinitos	48 integer + 40 Fl. Pt.
Predicción de saltos	Perfecta	2% to 6% de fallos de predicción (Tournament Branch Predictor)
Cache	Perfecta	64KI, 32KD, 1.92MB L2, 36 MB L3
Análisis de Memory Alias	Perfecto	??

□ Límite superior: Recursos ilimitados



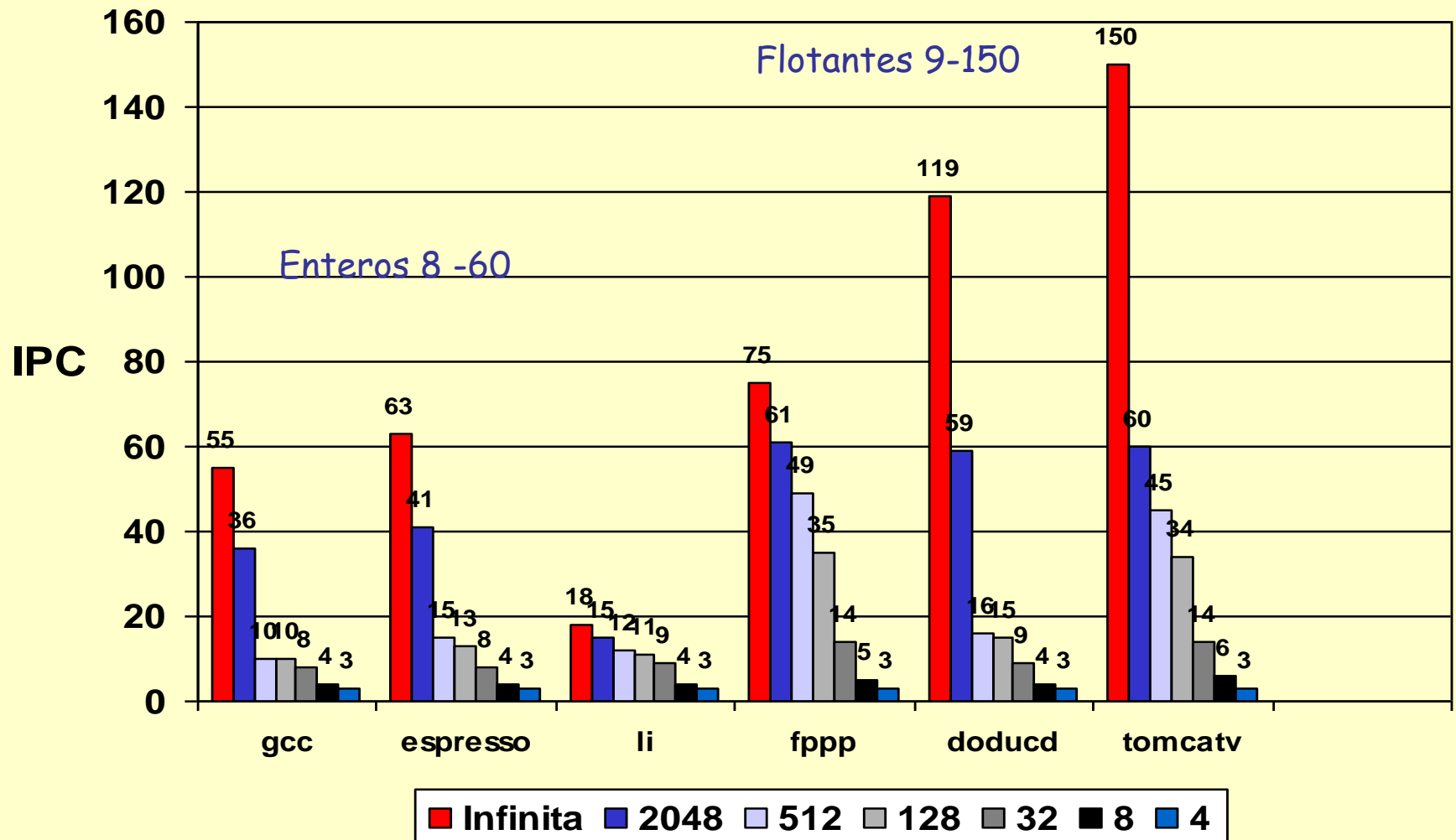
Algunas aplicaciones tienen poco paralelismo (Li)

❑ Modelo versus procesador real

	Nuevo Modelo	Modelo	Power 5
Instrucciones lanzadas por ciclo	Infinitas	Infinitas	4
Ventana de instrucciones	Infinito vs. 2K, 512, 128, 32, 8, 4	Infinita	200
Registros para renombrado	Infinitos	Infinitos	48 enteros + 40 Fl. Pt.
Predicción de saltos	Perfecta	Perfecta	Tournament
Cache	Perfecta	Perfecta	64KI, 32KD, 1.92MB L2, 36 MB L3
Análisis de Memory Alias	Perfecto	Perfecto	Perfecto

Límites del ILP

□ HW más real: Impacto del tamaño de la ventana de instrucciones



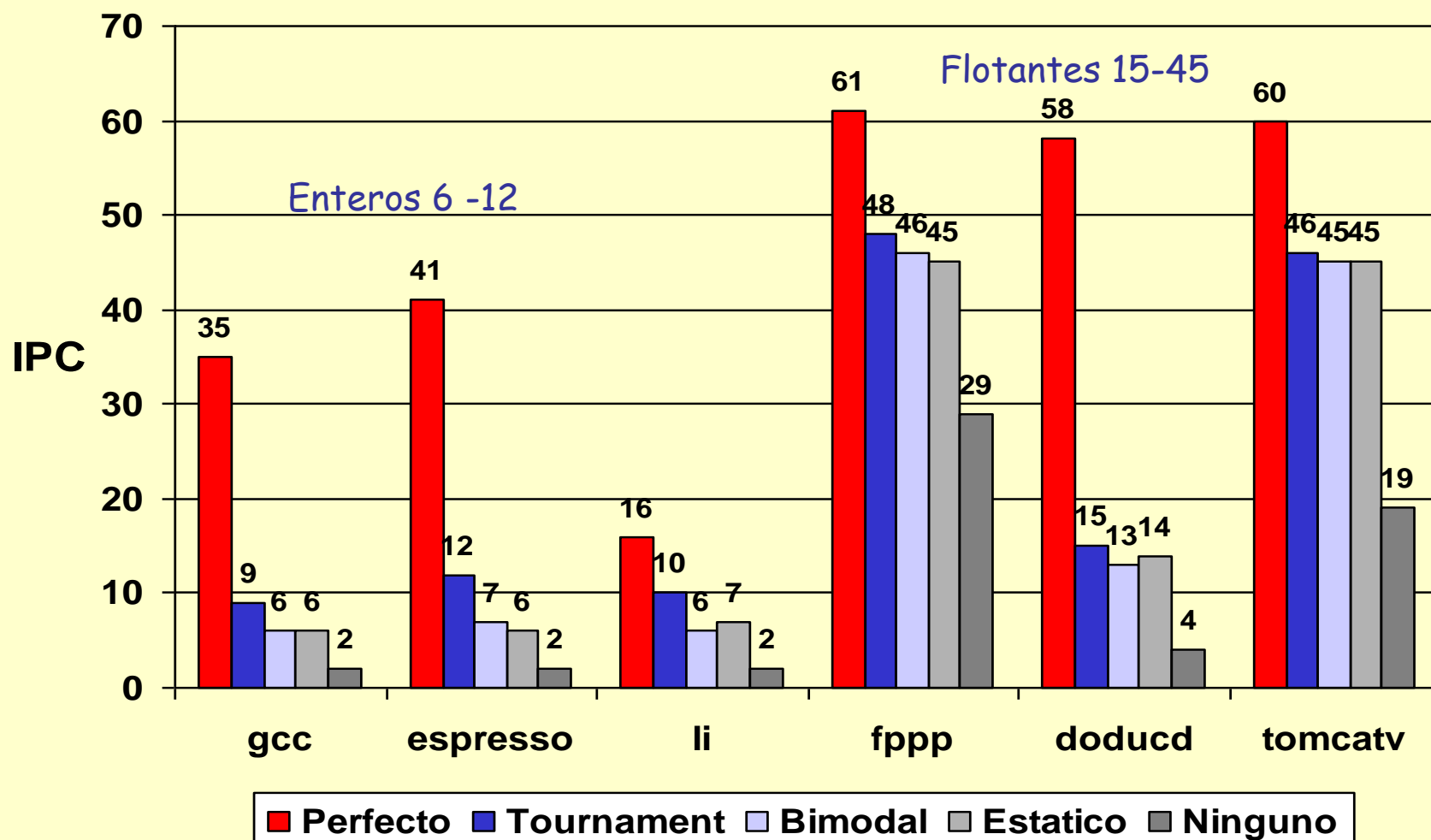
❑ Modelo versus procesador real

	Nuevo Modelo	Modelo	Power 5
Instrucciones lanzadas por ciclo	64 sin restricciones	Infinitas	4
Ventana de instrucciones	2048	Infinita	200
Registros para renombrado	Infinitos	Infinitos	48 enteros + 40 Fl. Pt.
Predicción de saltos	Perfecto vs. 8K Tournament vs. 512 2-bit vs. profile vs. ninguno	Perfecta	Tournament
Cache	Perfecta	Perfecta	64KI, 32KD, 1.92MB L2, 36 MB L3
Análisis de Memory Alias	Perfecto	Perfecto	Perfecto

Límites del ILP

□ HW más real: Impacto de los saltos

Ventana de 2048 y lanza 64 por ciclo



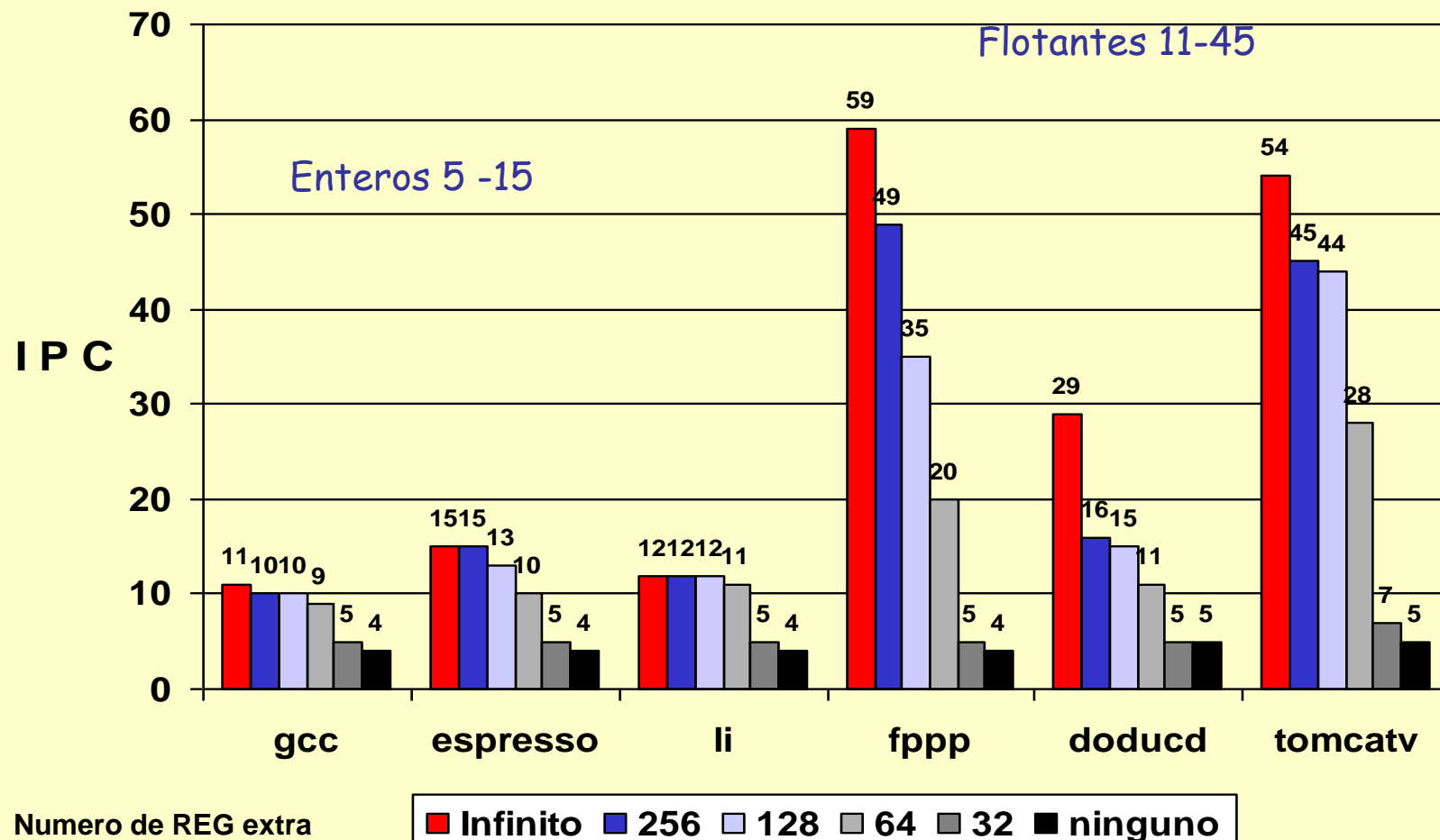
❑ Modelo versus procesador real

	Nuevo Modelo	Modelo	Power 5
Instrucciones lanzadas por ciclo	64 sin restricciones	Infinitas	4
Ventana de instrucciones	2048	Infinita	200
Registros para renombrado	Infinito v. 256, 128, 64, 32, ninguno	Infinitos	48 enteros + 40 Fl. Pt.
Predicción de saltos	8K Tournament (hibrido)	Perfecta	Tournament
Cache	Perfecta	Perfecta	64KI, 32KD, 1.92MB L2, 36 MB L3
Análisis de Memory Alias	Perfecto	Perfecto	Perfecto

Límites del ILP

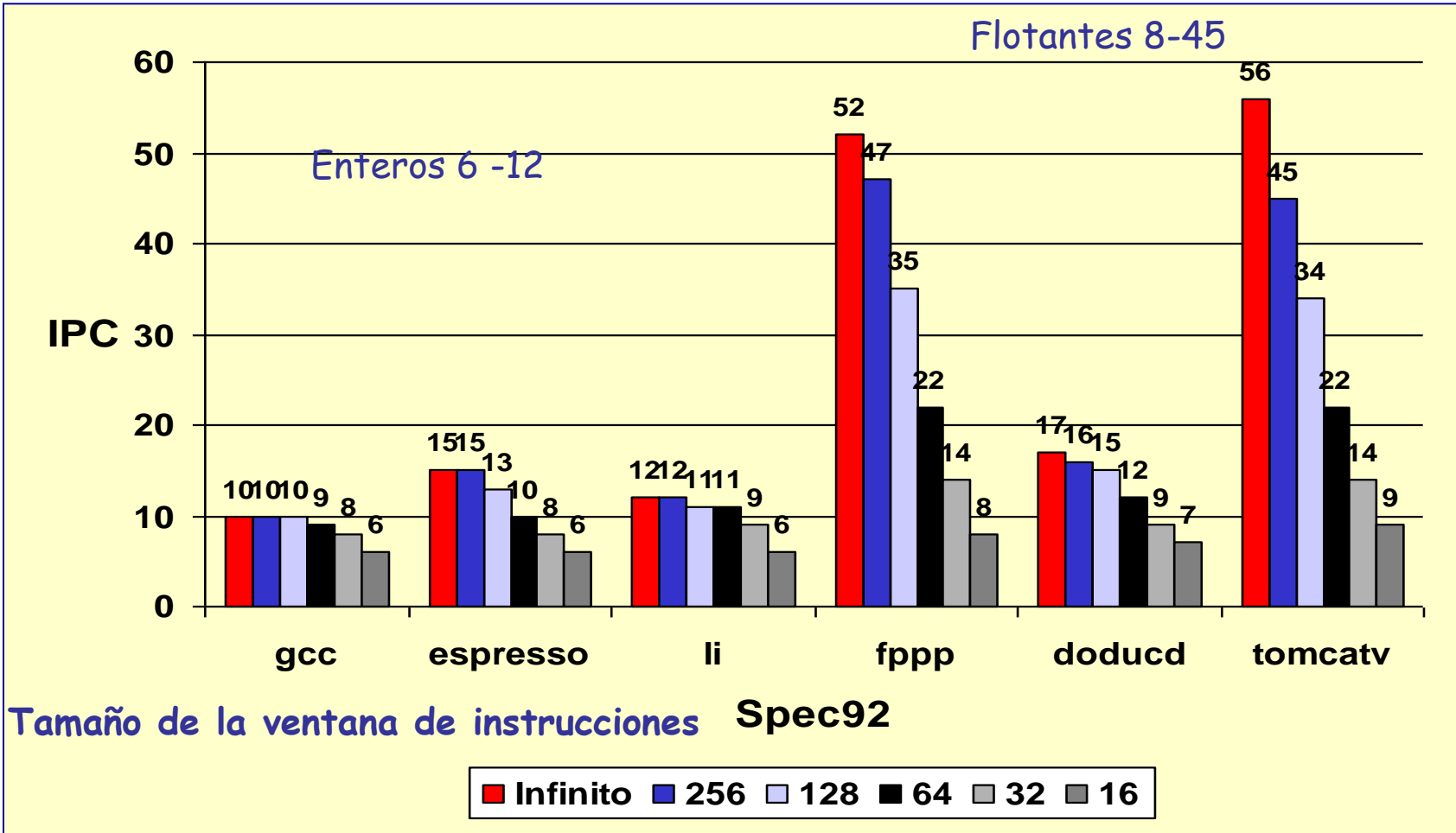
□ HW más real: Impacto del numero de registros

Ventana de 2048 y lanza 64 por ciclo, predictor híbrido 8K entradas



□ HW realizable

64 instrucciones por ciclo sin restricciones, predictor híbrido,
predictor de retornos de 16 entradas, Desambiguación perfecta,
64 registros enteros y 64 Fp extra



❑ Modelo versus procesador real

	Nuevo Modelo	Modelo	Power 5
Instrucciones lanzadas por ciclo	64 (sin restricciones)	Infinitas	4
Ventana de instrucciones	Infinito vs. 256, 128, 32, 16	Infinita	200
Registros para renombrado	64 Int + 64 FP	Infinitos	48 enteros + 40 Fl. Pt.
Predicción de saltos	1K 2-bit	Perfecta	Tournament
Cache	Perfecto	Perfecta	64KI, 32KD, 1.92MB L2, 36 MB L3
Análisis de Memory Alias	HW disambiguation	Perfecto	Perfecto

Ejemplos

Max Issue:
8 instrucciones por ciclo

Cache
L1 a L3 Itanium

AB memoria
75 GB/s

Tec 65nm, 10M

Mtrans 1700

Área 596mm²

Ventana de instrucciones
200 (Power 5 , 5+)

Etapas Pipe:
31 (Pentium NB)

Processor	AMD Opton 2222SE	AMD Opton 8222SE	Intel Pentium Extm Ed 965	Intel Core 2 X6800	Intel Xeon 5160	Intel Xeon 5355
Bit-width	32/64-bit	32/64-bit	32/64-bit	32/64-bit	32/64-bit	32/64-bit
Cores/chip x Threads/core	2 x 1	2 x 1	2 x 2	2 x 1	2 x 1	4 x 1
Clock Rate	3.00GHz	3.00GHz	3.73GHz	2.93GHz	3.00GHz	2.66GHz
Cache: L1-L2-L3 - I/D or Unified	2 x 64K/64K - 2 x 1M - N/A	2 x 64K/64K - 2 x 1M - N/A	2 x 16K/12K - 2 x 2M - NA	2 x 32K/32K - 2 x 4M - NA	2 x 32K/32K - 4M - NA	4 x 32K/32K - 8M - NA
Execution Rate/Core	3 x86 instr	3 x86 instr	3 uOPs	4 uOPs	5 uOPs	5 uOPs
Pipeline Stages	12/17 stages	12/17 stages	31 stages	14 stages	14 stages	14 stages
Out of Order	72ROPs	72ROPs	128 uOPs	96 uOPs	126 uOPs	126 uOPs
Memory B/W	8.5GB/s	8.5GB/s	8.5GB/s	8.5GB/s	10.5GB/s	10.5GB/s
Package	LGA-1207	LGA-1207	LGA-775	LGA-775	LGA-771	LGA-771
IC Process	90nm 9M SOI	90nm 9M SOI	65nm 8M	65nm 8M	65nm 8M	65nm 8M
Die Size	230mm ²	230mm ²	162mm ²	143mm ²	144mm ²	286mm ²
Transistors	227 million	227 million	376 million	291 million	291 million	582 million
List Price (intro)	\$873	\$2,149	\$1,207	\$999	\$851	\$1,172
Power (Max)	120W(MTP)	120W(MTP)	130W(TDP)	75W	80W	120W
Availability	2Q07	2Q07	2Q06	3Q06	2Q06	4Q06
Scalability	1-2 chips	2-8 chips	1 chip	1 chip	1-2 chips	1-2 chips
SPECint/fp2006 [cores]	13.5 / 14.3 [2]	11.2 / 13.0 [8]	11.7 / 12.7 [2]	18.5 / 16.8 [2]	18.9 / 17.1 [4]	17.3 / 16.2 [8]
SPECint/fp2006_rate [cores]	50.8 / 49.4 [4]	96.1 / 93.0 [8]	23.1 / 21.7 [2]	31.1 / 26.8 [2]	60.0 / 44.1 [4]	92.5 / 58.9 [8]
Architecture Status	Active	Active	Inactive	Active	Active	Active
Processor	Intel Itanium 2 9050	IBM Power6	IBM Power5+	Fujitsu SPARC64 VI	Sun UltraSPARC IV+	Sun UltraSPARC T1
Bit-width	64-bit	64-bit	64-bit	64-bit	64-bit	64-bit
Cores/chip x Threads/core	2 x 2	2 x 2	2 x 2	2 x 2	2 x 1	8 x 4
Clock Rate	1.60GHz	4.70GHz	2.20GHz	2.40GHz	1.95GHz	1.20GHz
Cache: L1-L2-L3 I/D or Unified	2 x 16K/16K - 1M/256K - 12M(on)	2 x 64K/64K - 2 x 4M - 32(off)	2 x 64K/32K 1.92M - 36M(off)	2 x 128K/128K - 5M - NA	2 x 64K/64K - 2M - 32M(off)	8 x 16K/8K - 3M - N/A
Execution Rate/Core	6 issue	7 issue	5 issue	4 issue	4 issue	1 issue
Pipeline Stages	8 stages	13 stages	15 stages	15 stages	14 stages	6 stages
Out of Order	None	"limited"	200 instr	64	None	None
Memory B/W	8.5GB/s	75GB/s	12.8GB/s	8GB/s	4.8GB/s	25.6GB/s
Package	mPGA-700	N/A	MCM-5370 pins	412 I/O pins	FC-LGA 1368	1933 pins
IC Process	90nm 7M	65nm 10m	90nm 10m	90nm 10M	90nm 9M	90nm 9M
Die Size	596mm ²	341mm ²	245mm ²	421mm ²	335mm ²	378mm ²
Transistors	1.72 billion	790 million	276 million	540 million	295 million	279 million
List Price (intro)	\$3,692	N/A	N/A	N/A	N/A	N/A
Power (Max)	104W	~100W	100W	120W	90W	70W
Availability	3Q06	2Q07	4Q05	2Q07	3Q06	4Q05
Scalability	1-64 chips	2-32 chips	1-32 chips	4-64 chips	1-72 chips	1 chip
SPECint/fp2006 [cores]	14.5 / 17.3 [2]	17.8 / 18.7 [1]	10.5 / 12.9 [1]	9.7 / 11.1 [32]	N/A	N/A
SPECint/fp2006_rate [cores]	1534 / 1671 [128]	410 / 379 [16]	197 / 229 [16]	1111 / 1160 [128]	1120 / N/A [144]	N/A
Architecture Status	Active	Active	Inactive	Active	Inactive	Active

All SPEC scores are base.

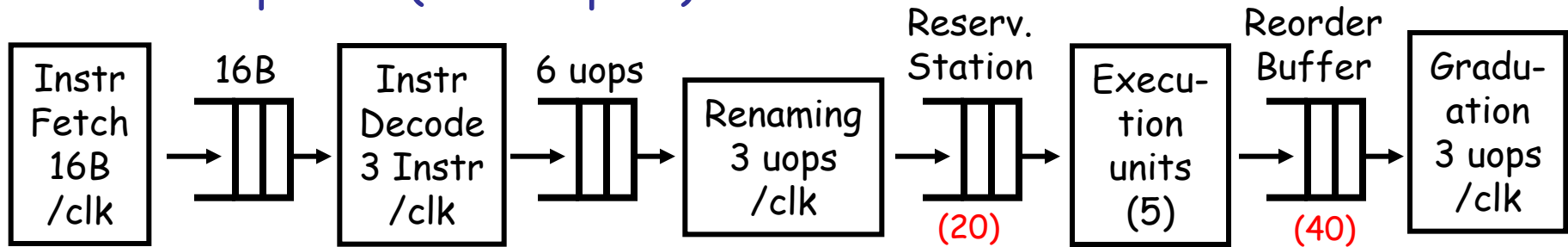
**Un ejemplo:
Implementaciones X86
P6, Netburst(Pentium4)**

❑ P6 (implementacion de la arquitectura IA-32 usada en Pentium Pro, II, III)

Modelo	Año	Clock	L1	L2
Pentium Pro	1995	100-200Mhz	8+8 KB	126-1024KB
Pentium II	1998	233-450Mhz	16+16 KB	256-512KB
Celeron	1999	500-900Mhz	16+16 KB	128KB
Pentium III	1999	450-1100Mhz	16+16 KB	256-512KB
PentiumIII Xeon	1999	700-900Mhz	16+16 KB	1MB-2MB

- ¿ Como segmentar un ISA con instrucciones entre 1 y 17 bytes?
- o El P6 traduce la instrucciones originales a microoperaciones de 72 bits (similar al DLX)
 - o Cada instrucción original se traduce a una secuencia de 1 a 4 microoperaciones. Pero ...
 - La instrucciones más complejas son traducidas por una secuencia almacenada en una memoria ROM(8Kx72)(microprograma)
 - o Tiene un pipeline de 14 etapas
 - o Ejecución fuera de orden especulativa, con renombrado de registros y ROB

□ P6 Pipeline (14 etapas)

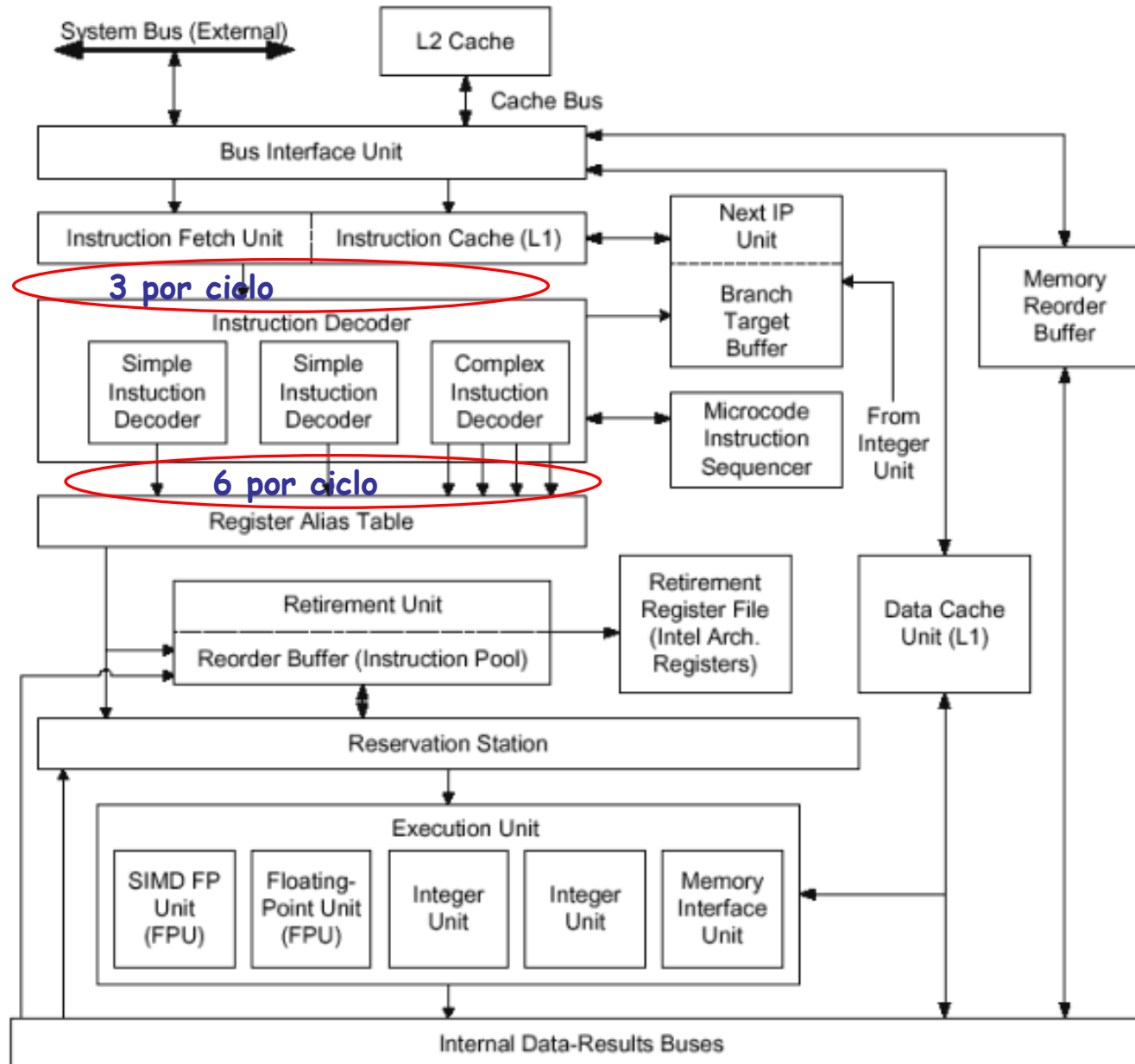


8 etapas para fetch, decodificación y issue en orden

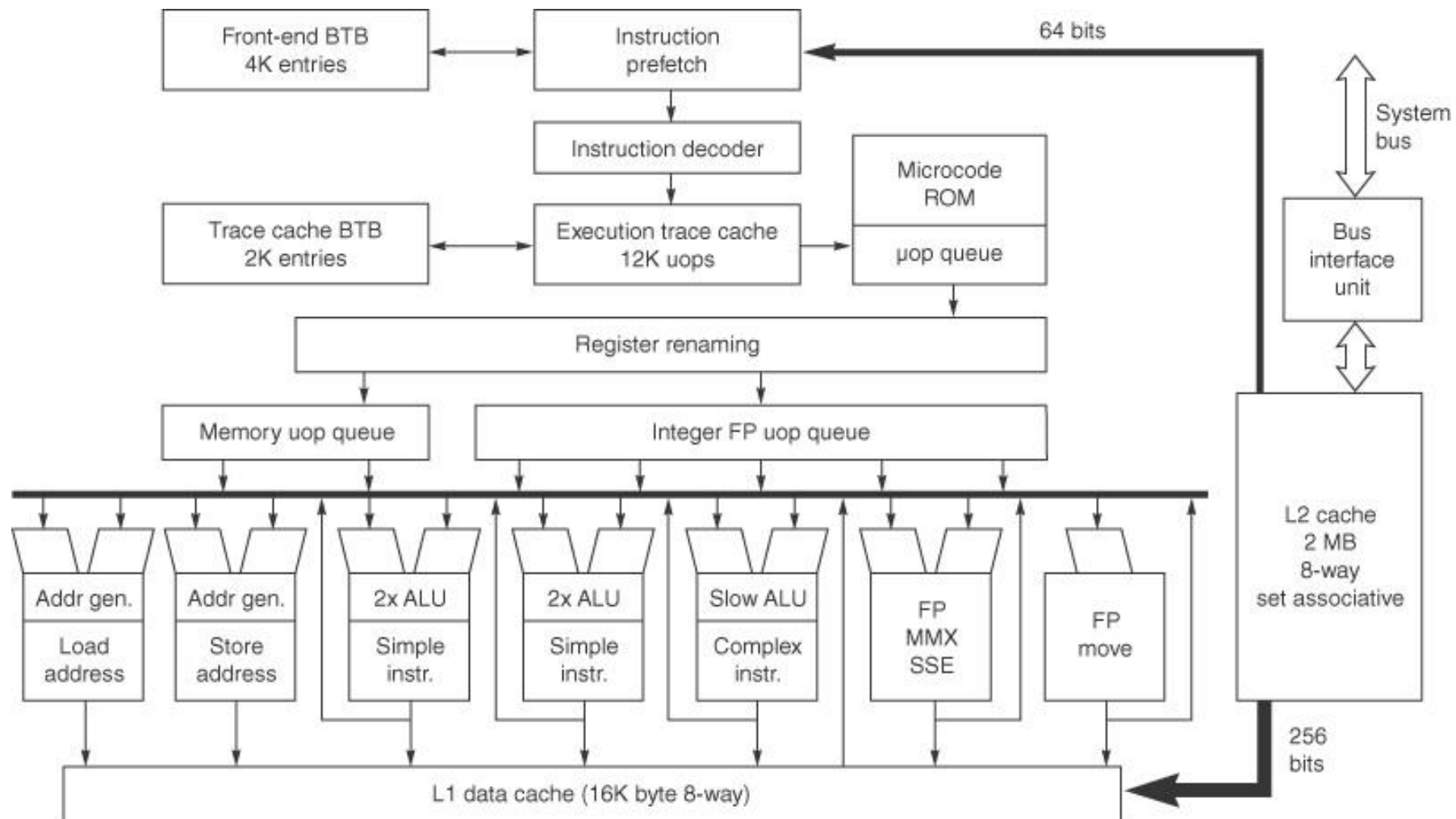
- o 1 ciclo para determinar la longitud de la instrucción 80x86 in + 2 más para generar las microoperaciones

- 3 etapas para ejecución fuera de orden en una de 5 unidades funcionales
- 3 etapas para la finalización de la instrucción (commit)

<u>Parameter</u>	<u>80x86</u>	<u>microops</u>
Max. instructions issued/clock	3	6
Max. instr. complete exec./clock		5
Max. instr. committed/clock		3
Window (Instrs in reorder buffer)	40	
Number of reservations stations	20	
Number of rename registers	40	
No. integer functional units (FUs)	2	
No. floating point FUs	1	
No. SIMD Fl. Pt. FUs	1	
No. memory Fus	1 load + 1 store	



Pentium 4 Microarchitecture



© 2007 Elsevier, Inc. All rights reserved.

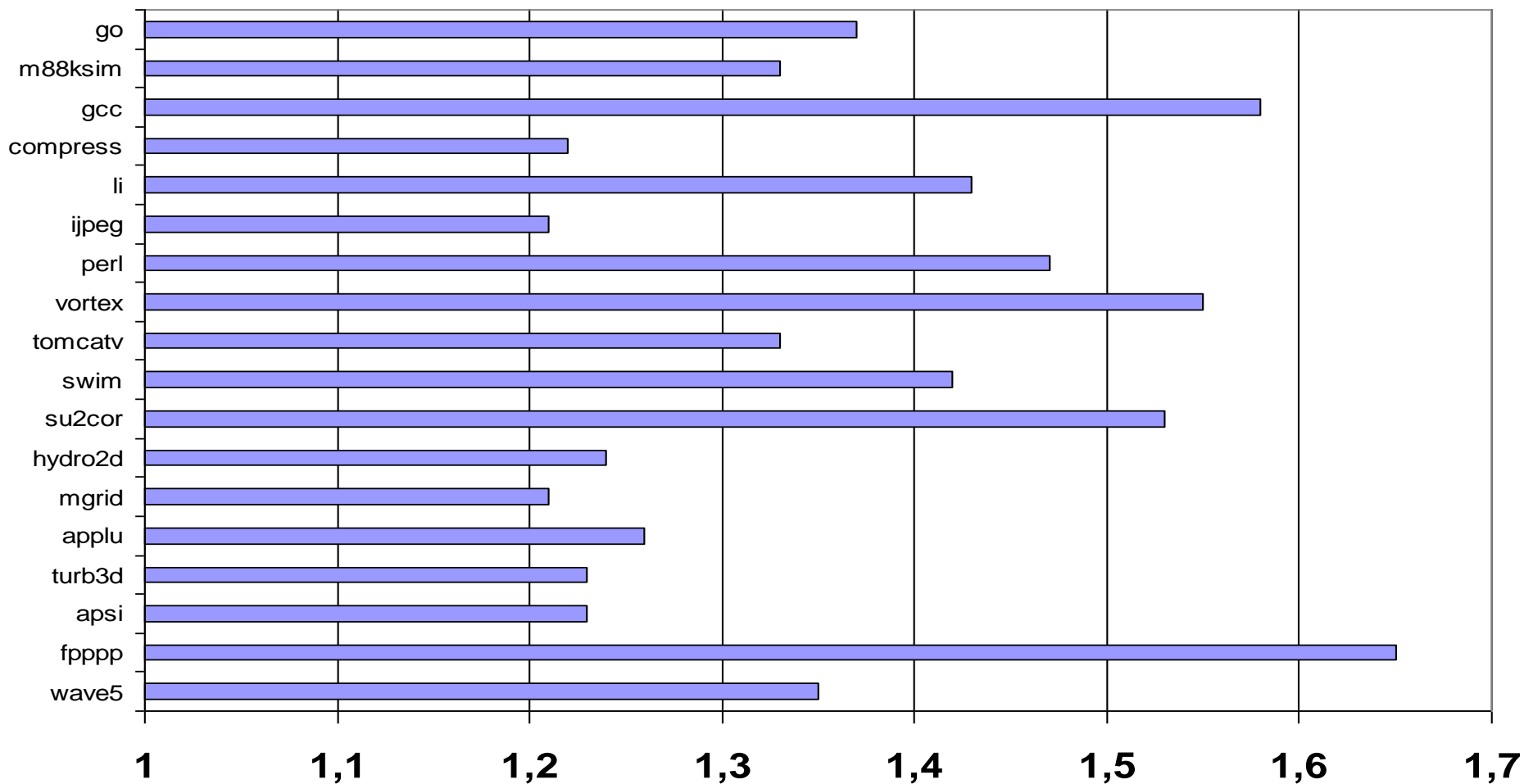
- ❑ BTB = Branch Target Buffer (branch predictor)
- ❑ I-TLB = Instruction TLB, Trace Cache = Instruction cache
- ❑ RF = Register File; AGU = Address Generation Unit
- ❑ "Double pumped ALU" means ALU clock rate 2X \Rightarrow 2X ALU F.U.s

Intel Netburst Microarchitecture

- ❑ Traduce instrucciones 80x86 a micro-ops (como P6)
- ❑ P4 tiene mejor predicción de saltos (x8) y más FU (7 versus 5)
- ❑ La Cache de instrucciones almacena micro-operations vs. 80x86 instrucciones.
"trace cache" (TC), BTB TC 2K entradas
 - o En caso de acierto elimina decodificación
- ❑ Nuevo bus de memoria: 400(800) MHz v. 133 MHz (RamBus, DDR, SDRAM)
(Bus@1066 Mhz)
- ❑ Caches
 - o Pentium III: L1I 16KB, L1D 16KB, L2 256 KB
 - o Pentium 4: L1I 12K uops, L1D 16 KB 8-way, L2 2MB 8-way
- ❑ Clock :
 - o Pentium III 1 GHz v. Pentium 4 1.5 GHz (3.8 Ghz)
 - o 14 etapas en pipeline vs. 24 etapas en pipeline (31 etapas)
- ❑ Instrucciones Multimedia: 128 bits vs. 64 bits => 144 instrucciones nuevas.
- ❑ Usa RAMBUS DRAM
 - o Más AB y misma latencia que SDRAM. Costo 2X-3X vs. SDRAM
- ❑ ALUs operan al doble del clock para operaciones simples
- ❑ Registros de renombrado: 40 vs. 128; Ventana: 40 v. 126
- ❑ BTB: 512 vs. 4096 entradas. Mejora 30% la tasa de malas predicciones

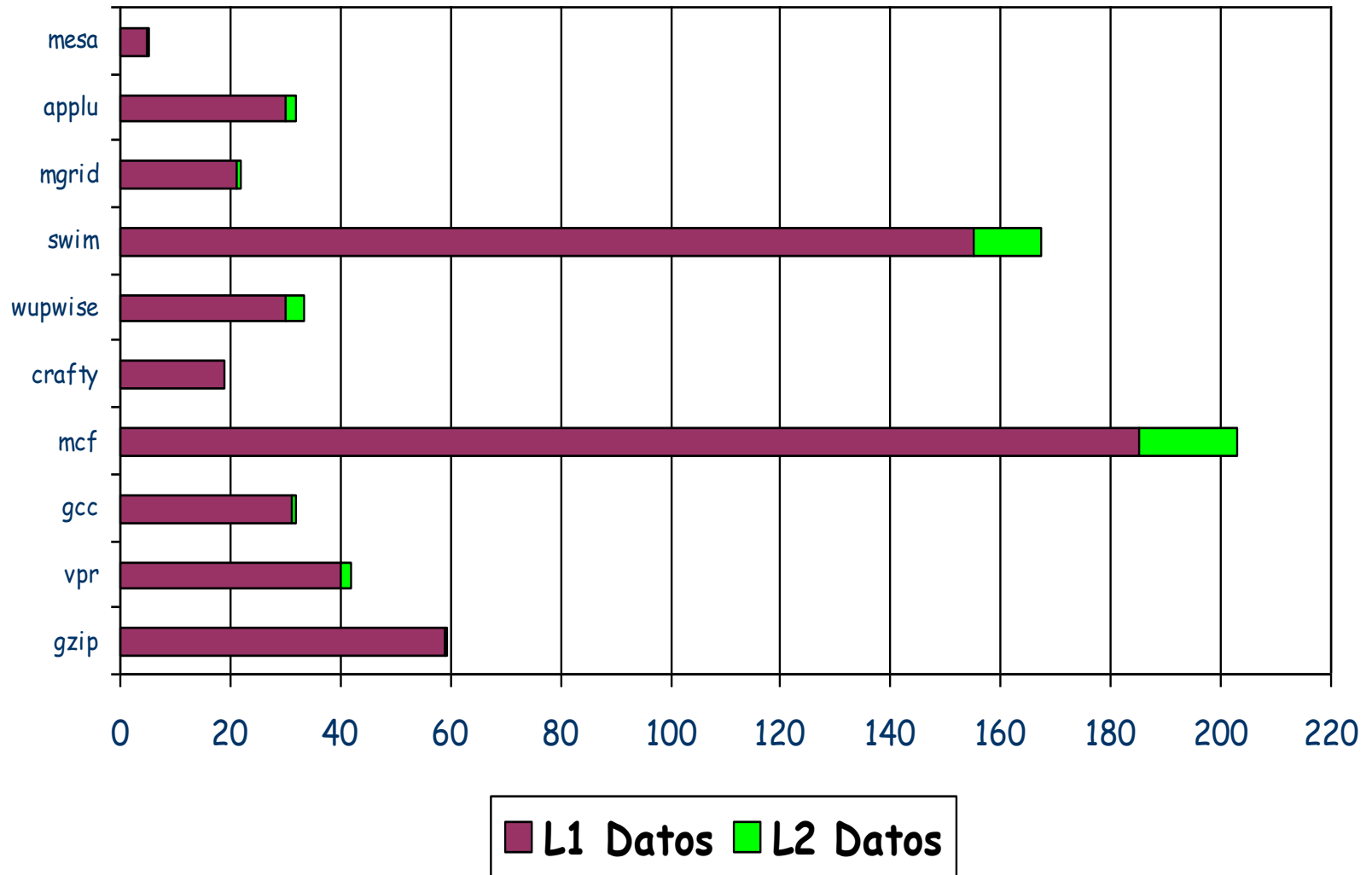
Netburst(Pentium 4) Rendimiento

Relación u-operaciones s/x86 instrucciones



1.2 to 1.6 uops per IA-32 instruction: 1.36 avg. (1.37 integer)

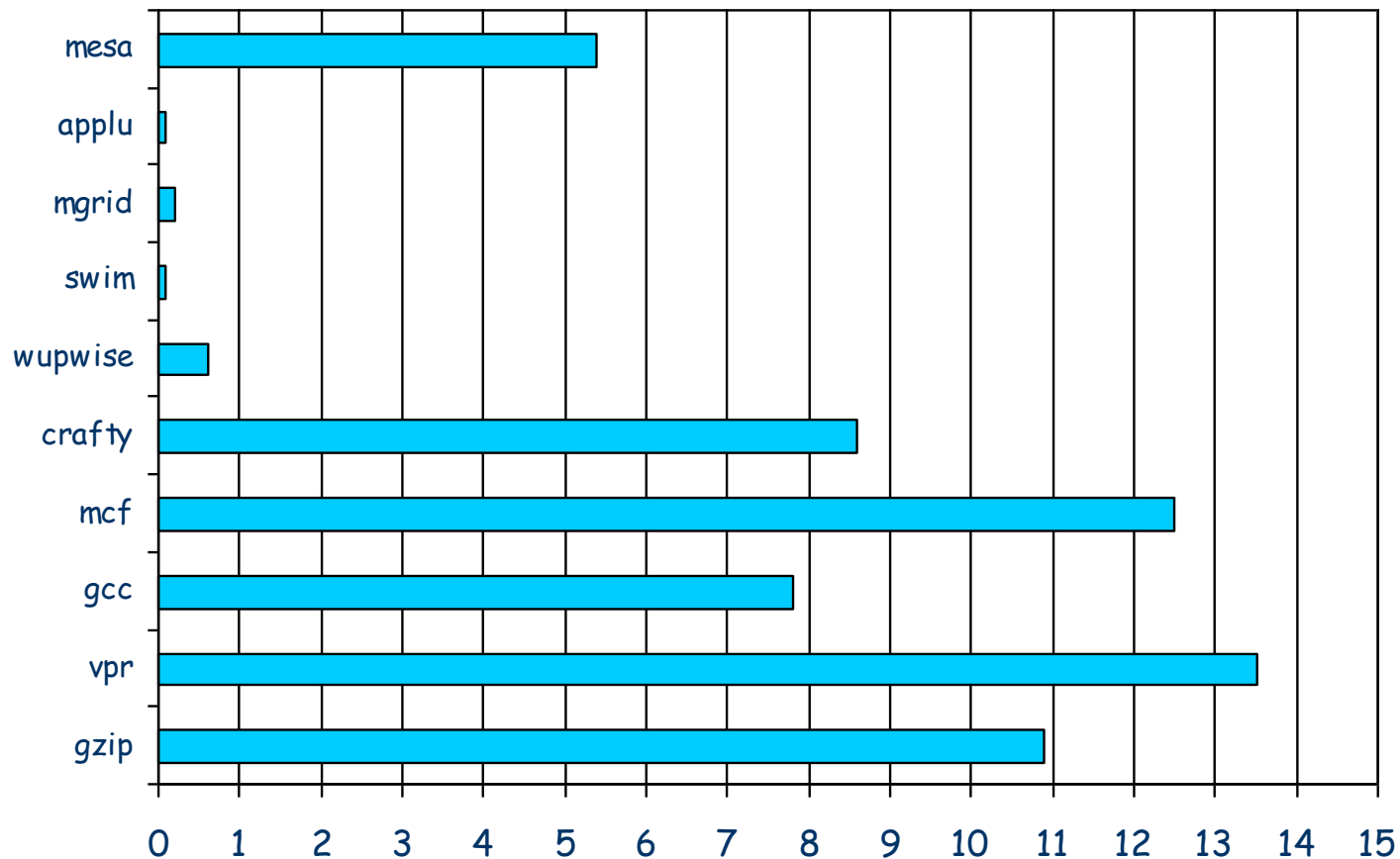
Fallos de cache



De 10 a 200 fallos por cada mil instrucciones

Netburst(Pentium4) Rendimiento

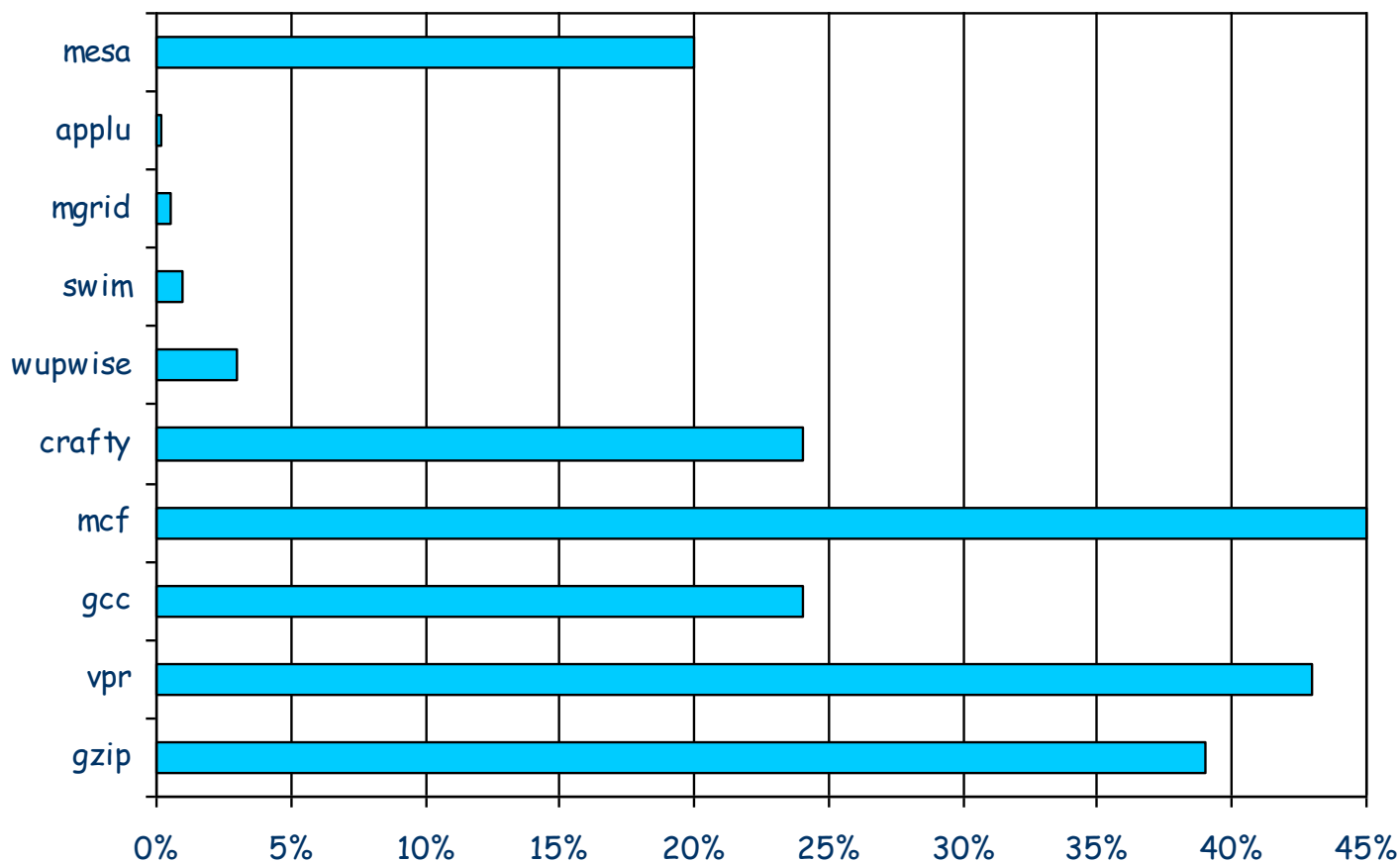
Comportamiento del predictor: Dos niveles, con información global y local
4K entradas



Fallos de predicción por cada mil instrucciones, 2% en FP, 6% en INT

Netburst(Pentium4) Rendimiento

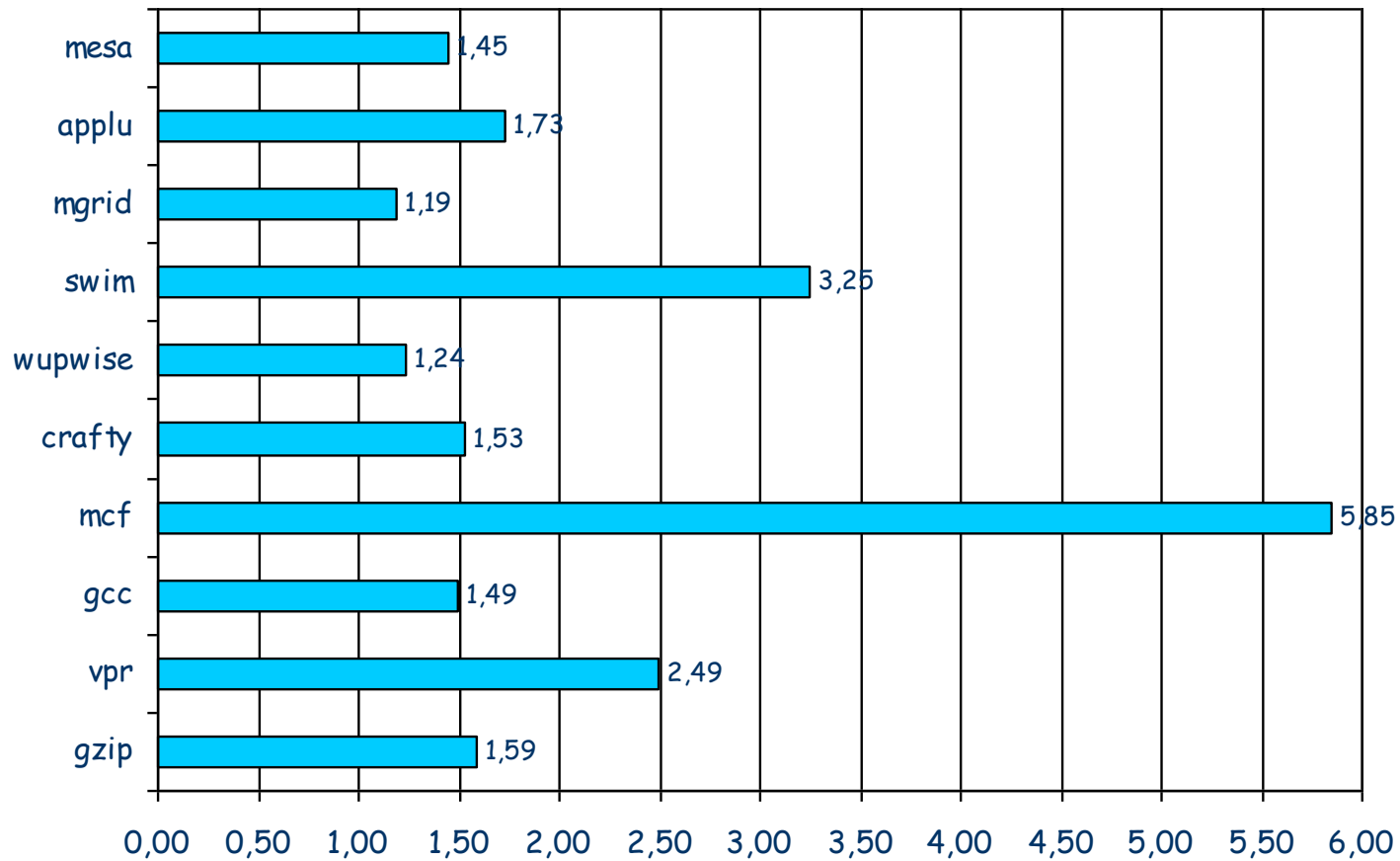
Especulación: porcentaje de μ -operaciones que no finalizan (commit)



Del 1% al 45 % de las instrucciones no finalizan

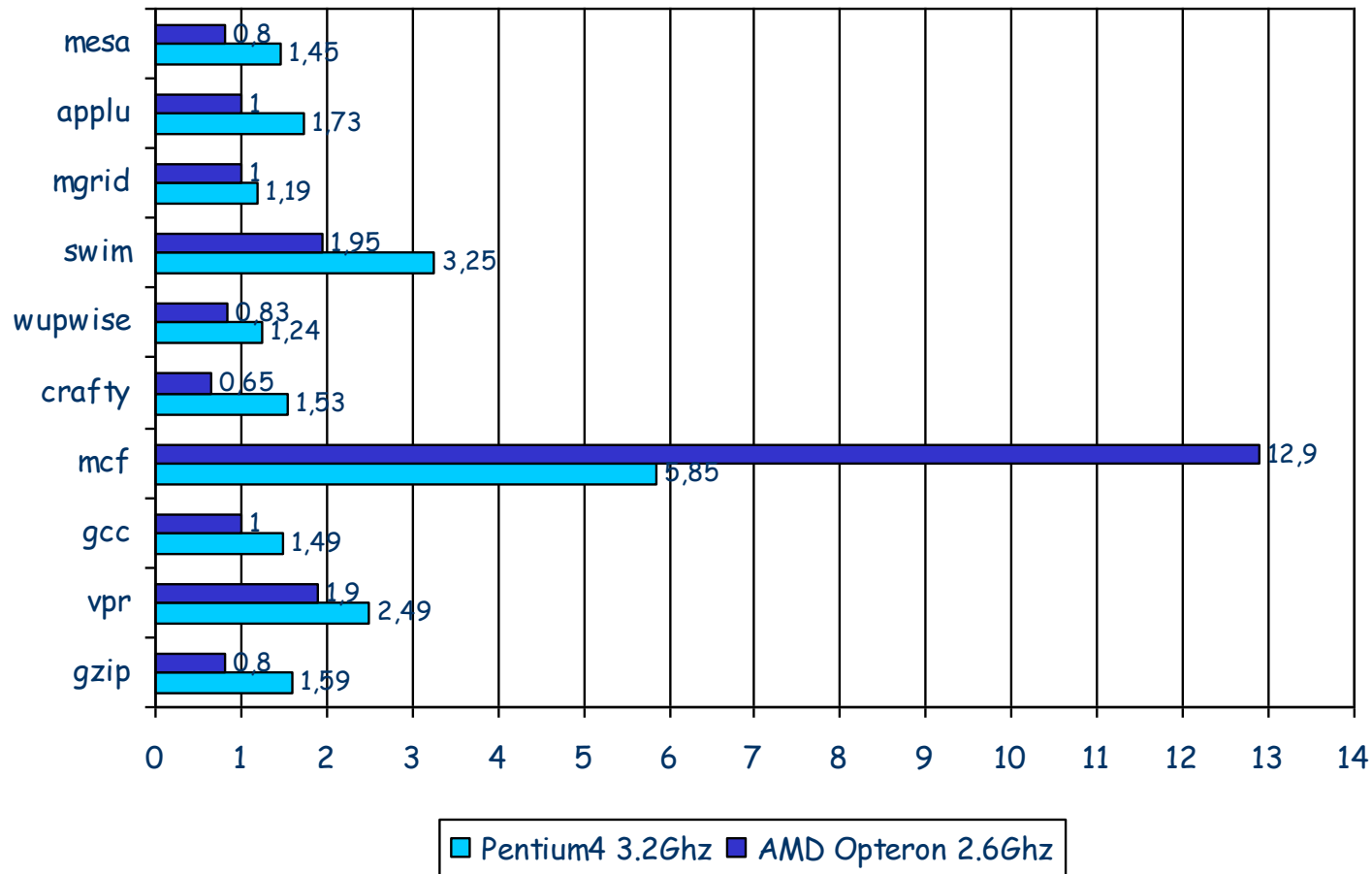
Netburst(Pentium4) Rendimiento

CPI real obtenido



Netburst(Pentium4) versus AMD

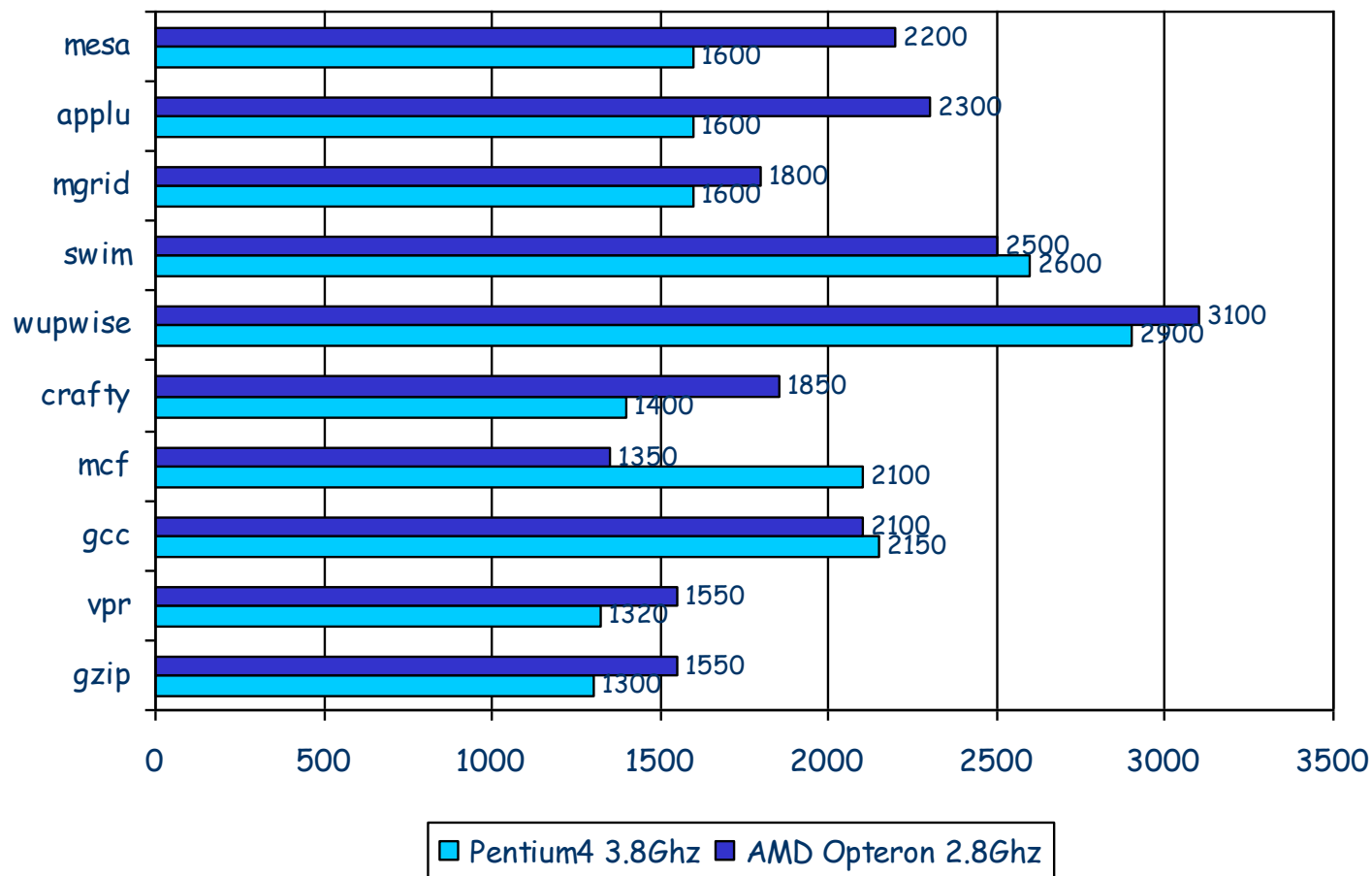
Diferencia fundamental; Pentium pipe profundo para alta frecuencia
CPI real obtenido



AMD CPI menor en un factor de 1,27. ¿ Se puede compensar con la mayor frecuencia?

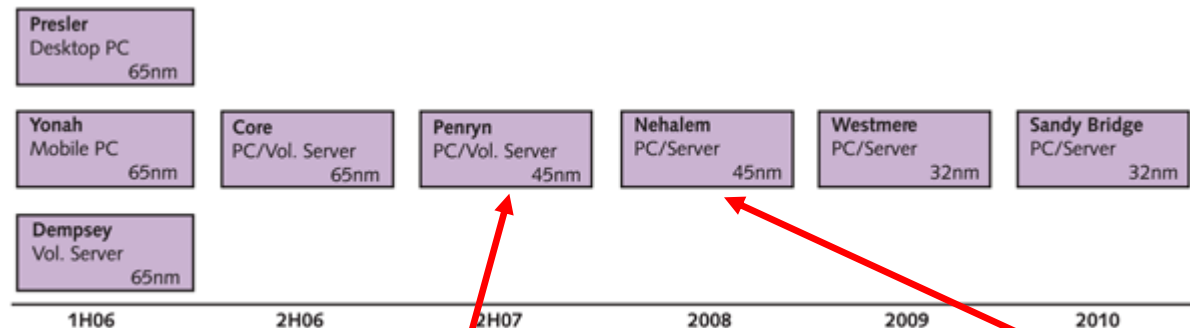
Netburst(Pentium4) versus AMD

Rendimiento en SPEC2000



AMD 1,08 más rendimiento. La mayor frecuencia no compensa el mayor CPI

Core, Core2,...



	Penryn	Core 2 Duo/Quad
Number of cores	2-4	2-4
Max Operating Frequency	>3.0GHz	2.93GHz
Max FSB Frequency	1.6GHz	1.06GHz
Max L2 cache	6-12MB	4-8MB
C-states	C0-C4 + Enhanced Deeper Sleep + Deep Power Down	C0-C4+ Enhanced Deeper Sleep
SSE Support	SSE4	SSE3
Super Shuffle Engine	Yes	No
Process	45nm (High-k /metal gate)	65nm
TDP	34-130W	34-130W
Dynamic Acceleration Support	Enhanced	Standard
Virtualization	Enhanced	Standard
Division	4bit/cycle (Radix 16)	2bit/cycle
Availability	2H07	Now

Features	Nehalem
Number of cores	2-8+
Max cores per die	Not specified
Multithreading	2 threads per core up to 16+ threads
Memory	Multilevel shared cache Integrated buffered/unbuffered memory controller
Power management	Unspecified enhancements
Graphics	Integrated
FSB	New CSI interface (up to four for servers)
PCIe	Not specified

- ❑ ¿Cómo superar los límites de este estudio?
- ❑ Mejoras en los compiladores e ISAs
- ❑ Eliminar riesgos EDL y EDE en la memoria
- ❑ Eliminar riesgos LDE en registros y memoria. Predicción
- ❑ Buscar otros paralelismos (thread)

❑ Buscar paralelismo de más de un thread

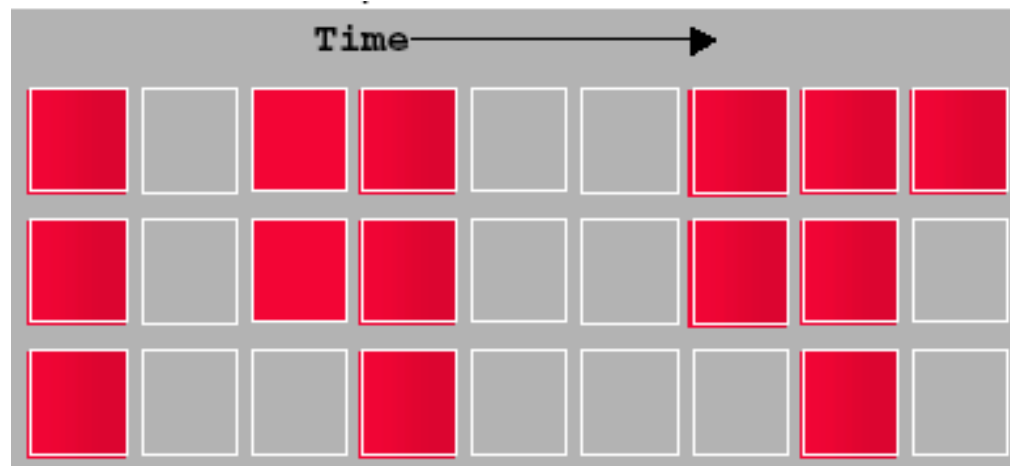
- o Hay mucho paralelismo en algunas aplicaciones (Bases de datos, códigos científicos)
- o Thread Level Parallelism
 - o Thread: proceso con sus propias instrucciones y datos
 - Cada thread puede ser parte de un programa paralelo de múltiples procesos, o un programa independiente.
 - Cada thread tiene todo el estado (instrucciones, datos, PC, register state, ...) necesario para permitir su ejecución
 - Arquitecturas(multiprocesadores, MultiThreading y multi/many cores)
- o Data Level Parallelism: Operaciones idénticas sobre grandes volúmenes de datos (extensiones multimedia y arquitecturas vectoriales)

Thread Level Parallelism (TLP) versus ILP

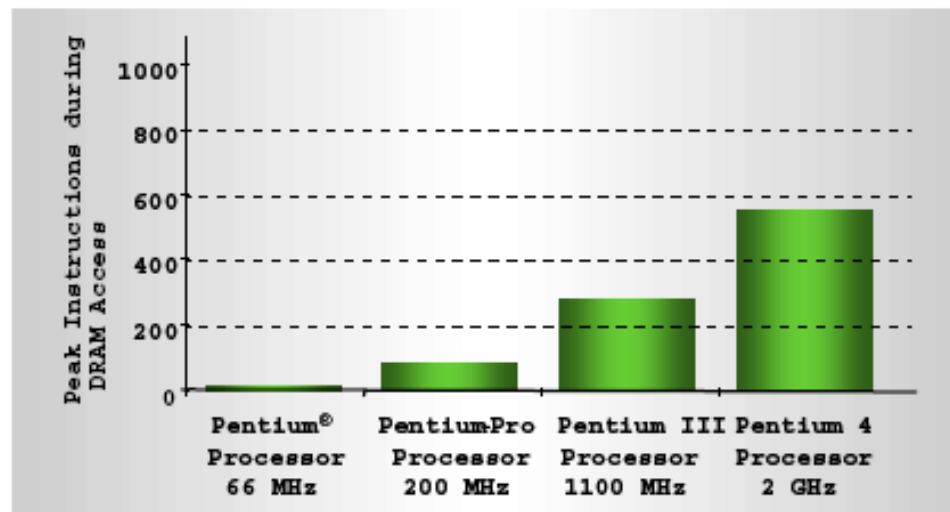
- ❑ ILP explota paralelismo implícito dentro de un segmento de código lineal o un bucle
- ❑ TLP representa el uso de múltiples thread que son inherentemente paralelos.
- ❑ Objetivo: Usar múltiples streams de instrucciones para mejorar;
 - o Throughput de computadores que ejecutan muchos programas diferentes.
 - o Reducir el tiempo de ejecución de un programa multi-threaded
- ❑ TLP puede más eficaz en coste que ILP

¿ Por que multithreading ?

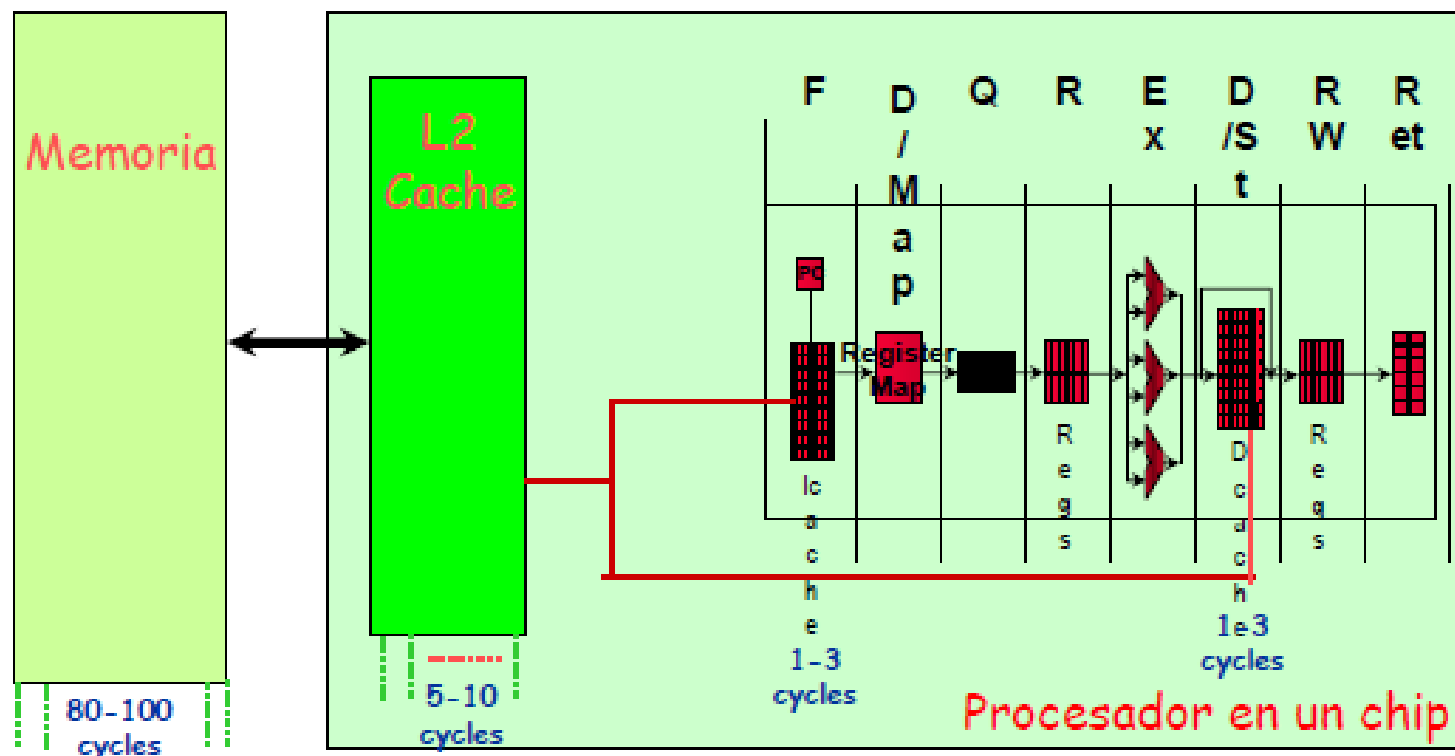
☐ Procesador superescalar



☐ La latencia de memoria crece.
¿ Como soportarla?



❑ Aspectos tecnológicos: Acceso a Memoria



❑ Tiempo de servicio de un fallo de cache: evolución

- o 21064 (200Mhz) 340ns, $340/5=68$ ciclos, $68 \times 2 = 136$ instrucciones
- o 21164 (300Mhz) 266ns, $266/3.3=80$, $80 \times 4 = 320$ instrucciones
- o 21264 (1Ghz) 180ns, $180/1=180$, $180 \times 6 = 1080$ instrucciones

La era del multi core y el multi thread

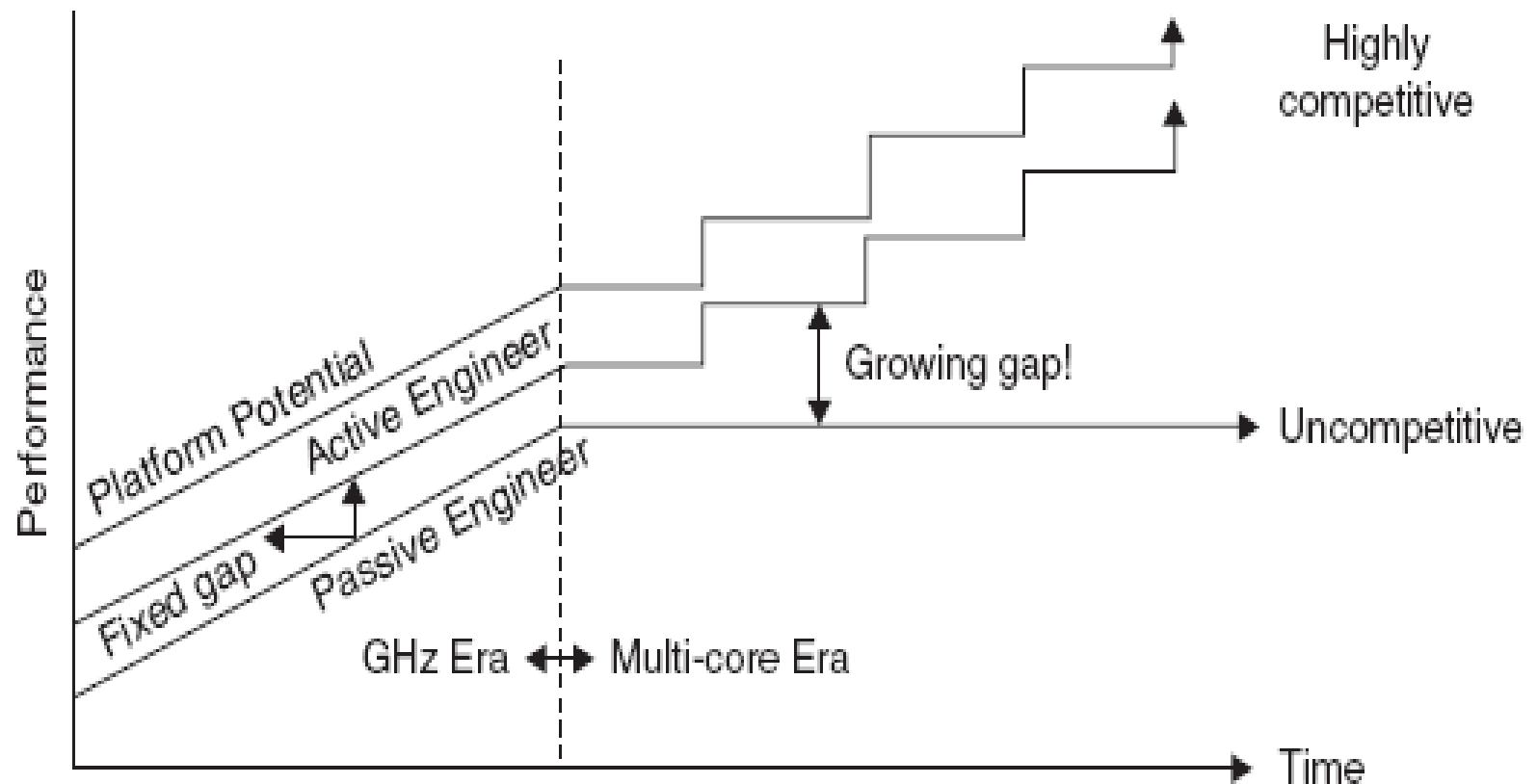
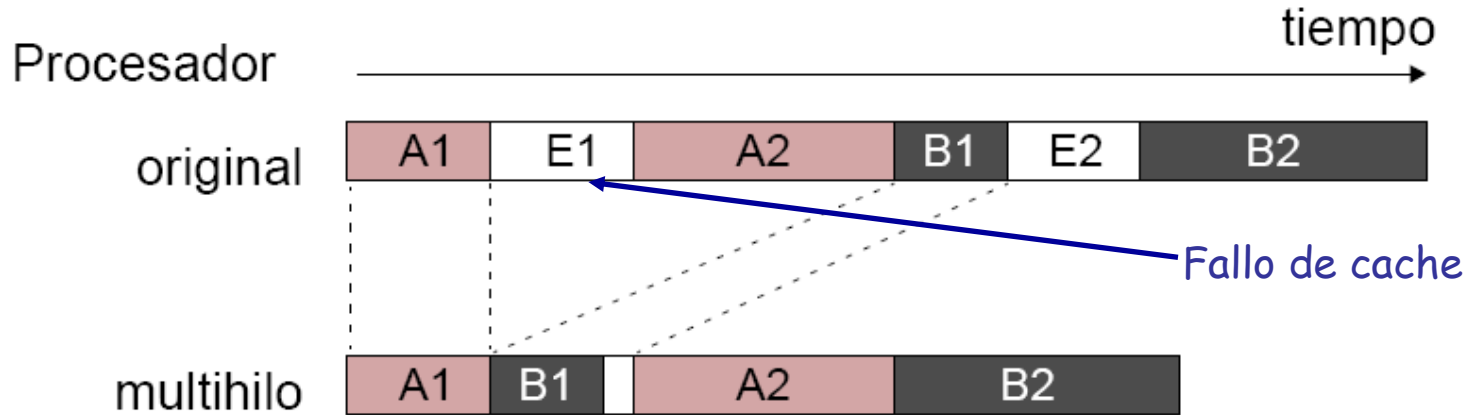


Figure 1.3: Taking advantage of multi-core processors

- ❑ La era Giga hertz comenzó en el año 2000 con la introducción procesadores con frecuencias mayor a 1 GHz y terminó en 2005 con la introducción de procesadores multi-core.
- ❑ Entrando en la era de procesador multi-core muestra una nueva tendencia sustitución de mayor frecuencia por multi core

❑ Multithreading



- Incrementar el trabajo procesado por unidad de tiempo
- Si los hilos son del mismo trabajo se reduce el tiempo de ejecución

La técnica multithreading no es ideal y se producen pérdidas de rendimiento. Por ejemplo, un programa puede ver incrementado su tiempo de ejecución aunque el número de programas ejecutados por unidad de tiempo sea mayor cuando se utiliza multithreading.

❑ Ejecución Multithreaded

- o Multithreading: múltiples threads comparten los recursos del procesador
 - o El procesador debe mantener el estado de cada thread e.g., una copia de bloque de registros, un PC separado, tablas de páginas separadas.
 - o La memoria compartida ya soporta múltiples procesos.
 - o HW para conmutación de thread muy rápido. Mucho mas rápido que entre procesos.
- o ¿Cuándo conmutar?
 - o Cada ciclo conmutar de thread (grano fino)
 - o Cuando un thread debe parar (por ejemplo fallo de cache)
- o HEP (1978), Alewife , M-Machine , Tera-Computer

□ Multithreading de Grano Fino

- o Conmuta entre threads en cada instrucción, entrelazando la ejecución de los diferentes thread.
- o Generalmente en modo "round-robin", los threads bloqueados se saltan
- o La CPU debe ser capaz de conmutar de thread cada ciclo.
- o Ventaja; puede ocultar stalls de alta y baja latencia, cuando un thread esta bloqueado los otros usan los recursos.
- o Desventaja; retarda la ejecución de cada thread individual, ya que un thread sin stall es retrasado por reparto de recursos (ciclos) entre threads
- o Ejemplo Niagara y Niagara 2 (SUN)

❑ Multithreading Grano Grueso

- o Conmuta entre threads solo en caso de largos stalls, como fallos de cache L2
- o Ventajas
 - o No necesita conmutación entre thread muy rápida.
 - o No retarda cada thread, la conmutación solo se produce cuando un thread no puede avanzar.
- o Desventajas; no elimina perdidas por stalls cortos. La conmutación es costosa en ciclos.
 - o Como CPU lanza instrucciones de un nuevo thread, el pipeline debe ser vaciado.
 - o El nuevo thread debe llenar el pipe antes las instrucciones empiecen a completarse.
- o Ejemplos; IBM AS/400, Montecio (Itanium 9000), Spacr64 VI

❑ Simultaneous Multi-threading

Motivación: Recursos no usados en un procesador superescalar

Un thread, 8 unidades

Ciclo M M FX FX FP FP BR CC

1								
2								
3								
4								
5								
6								
7								
8								
9								

Dos threads, 8 unidades

Ciclo M M FX FX FP FP BR CC

1								
2								
3								
4								
5								
6								
7								
8								
9								

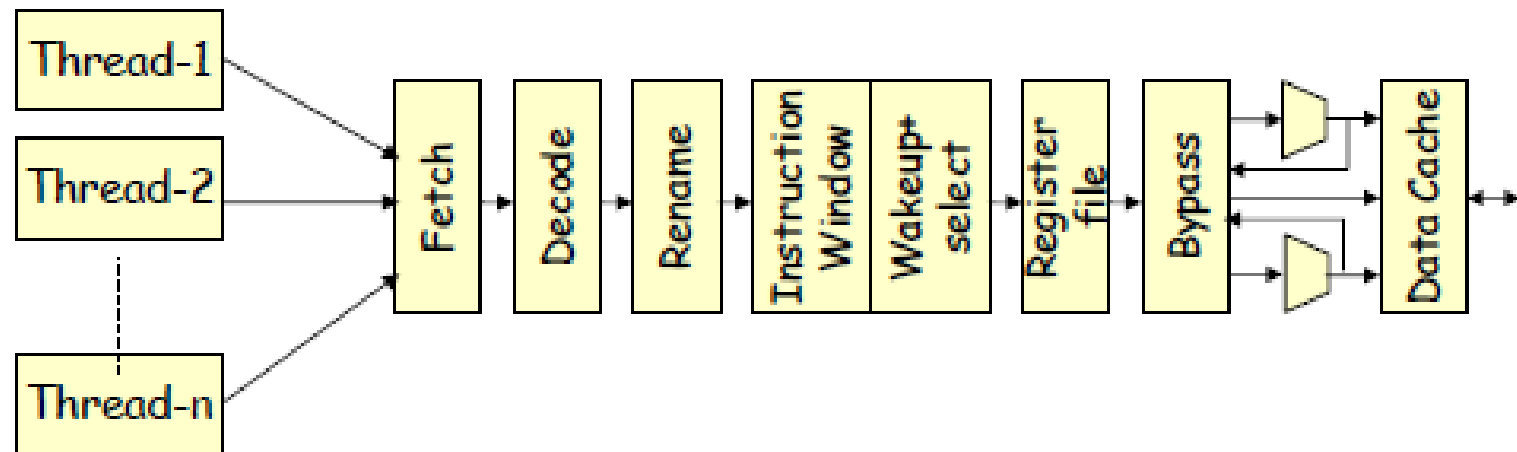
M = Load/Store, FX = Fixed Point, FP = Floating Point, BR = Branch, CC = Condition Codes

❑ Simultaneous Multithreading (SMT)

- ❑ Simultaneous multithreading (SMT): dentro de un procesador superescalar fuera de orden ya hay mecanismos Hw para soportar la ejecución de más de un thread
 - o Gran numero de registros físicos donde poder mapear los registros arquitectónicos de los diferentes threads
 - o El renombrado de registros proporciona un identificador único para los operandos de una instrucción, por tanto instrucciones de diferentes thread se pueden mezclar sin confundir sus operados
 - o La ejecución fuera de orden permite una utilización eficaz de los recursos.
- ❑ Solo necesitamos sumar una tabla de renombrado por thread y PC separados
 - o Commit independiente se soporta con un ROB por thread (Lógico o físico)
- ❑ Ojo conflictos en la jerarquía de memoria
- ❑ Ejemplos; Pentium4, Power5 y 6, Nehalem (2008)

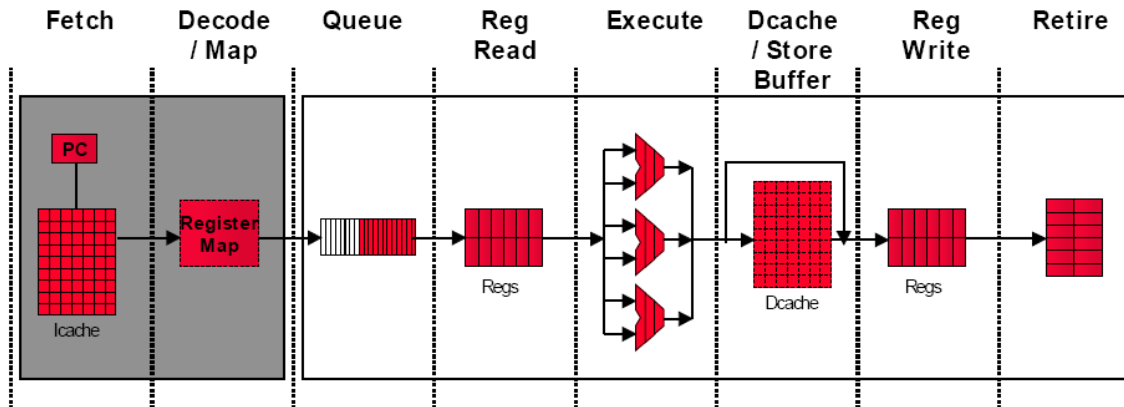
❑ Simultaneous multithreading

- o Diferentes "Threads" se ejecutan concurrentemente Los "Threads" pueden pertenecer al mismo o a diferentes procesos HEP (1978), Alewife , M-Machine , Tera-Computer Pentium 4, EV8(21464), Power5



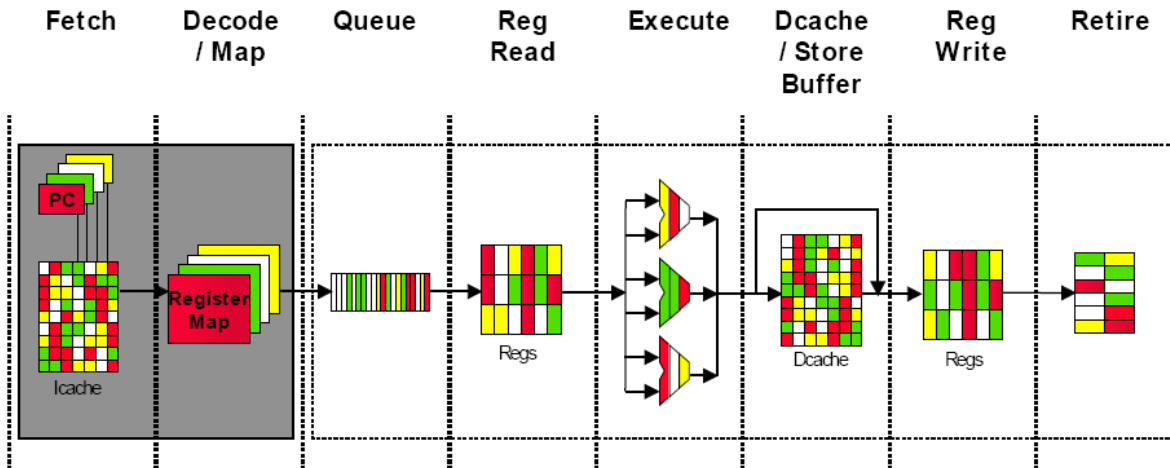
□ Simultaneous multithreading

- Un solo hilo: un flujo de instrucciones



✓ todos los recursos utilizados por un hilo

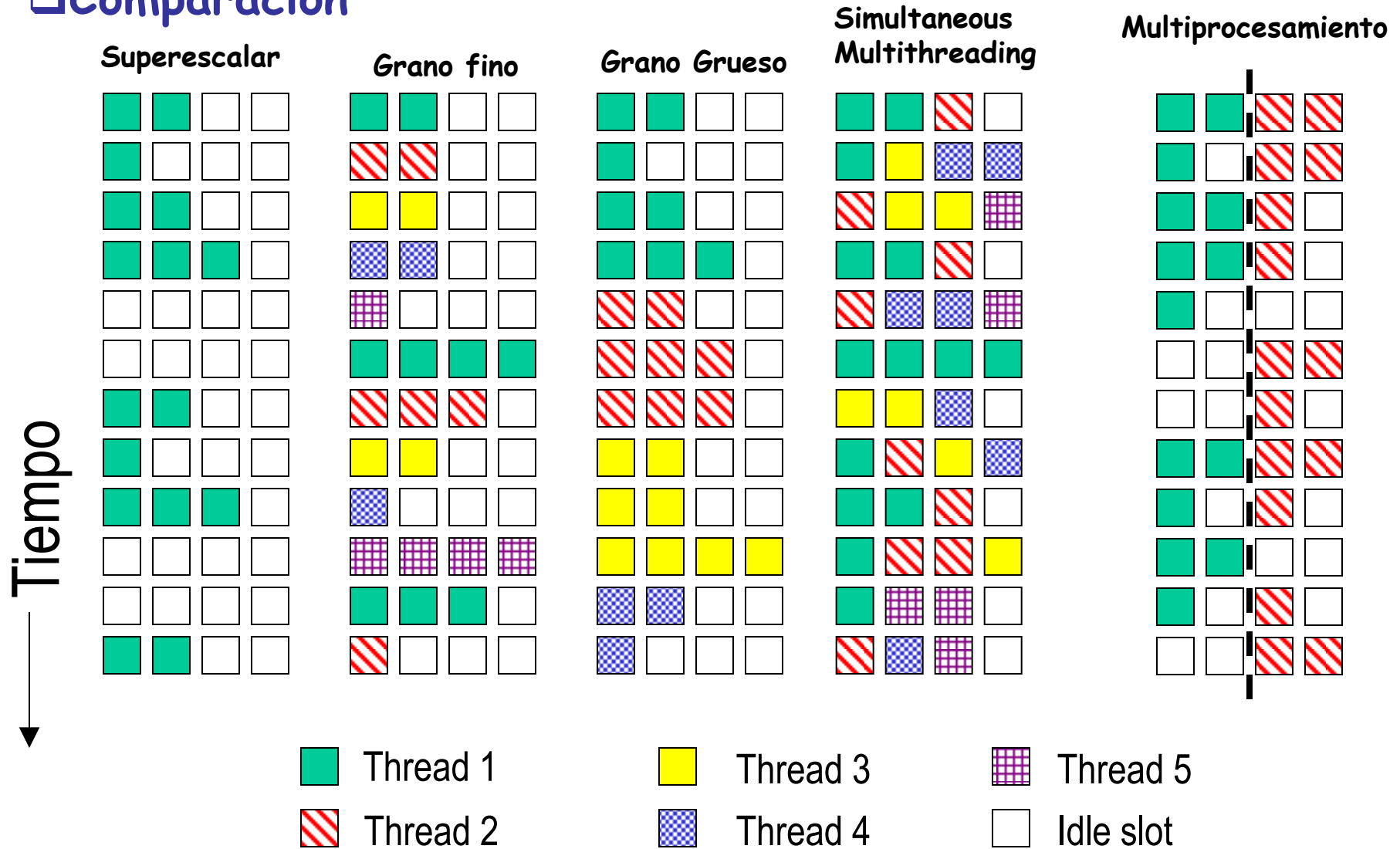
- Multihilo: varios flujos de instrucciones



✓ recursos para distinguir el estado de los hilos

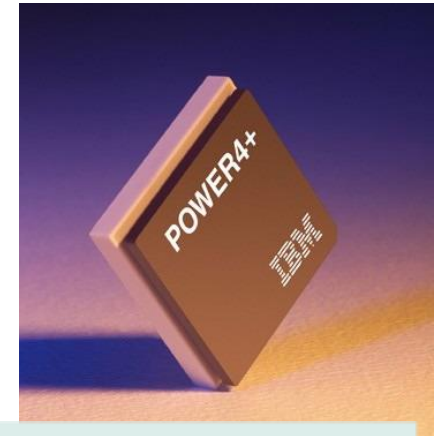
✓ los otros recursos se pueden compartir

Comparación



❑ Power 4

Predecesor Single-threaded del Power 5.
8 unidades de ejecución fuera de orden



Branch redirects

Instruction fetch

IF

IC

BP

D0

D1

D2

D3

Xfer

GD

Instruction crack and
group formation

MP

ISS

RF

EX

BR

WB

Xfer

MP

ISS

RF

EX

LD/ST

DC

Fmt

WB

Xfer

MP

ISS

RF

EX

FX

WB

Xfer

MP

ISS

RF

EX

FP

F6

WB

Xfer

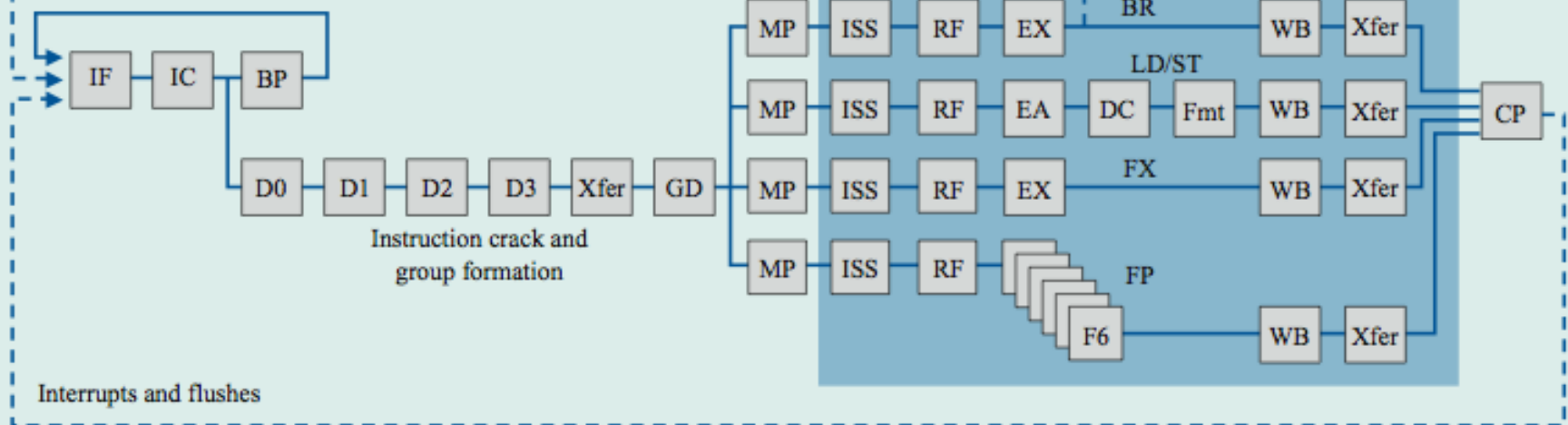
CP

Interrupts and flushes

Power 4

Branch redirects

Instruction fetch

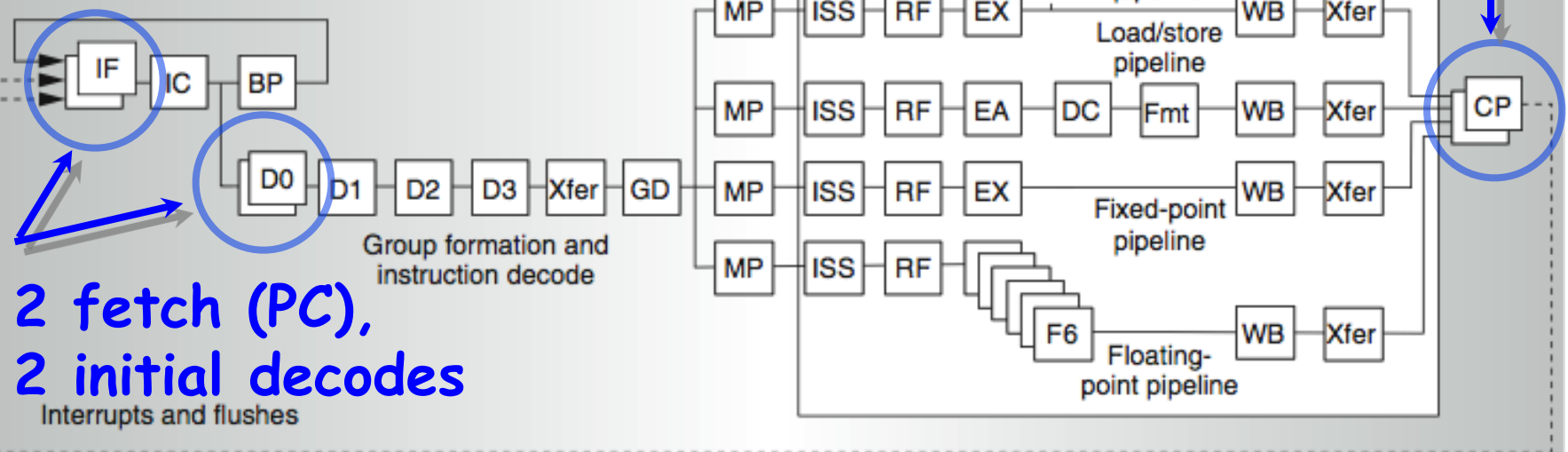


2 commits

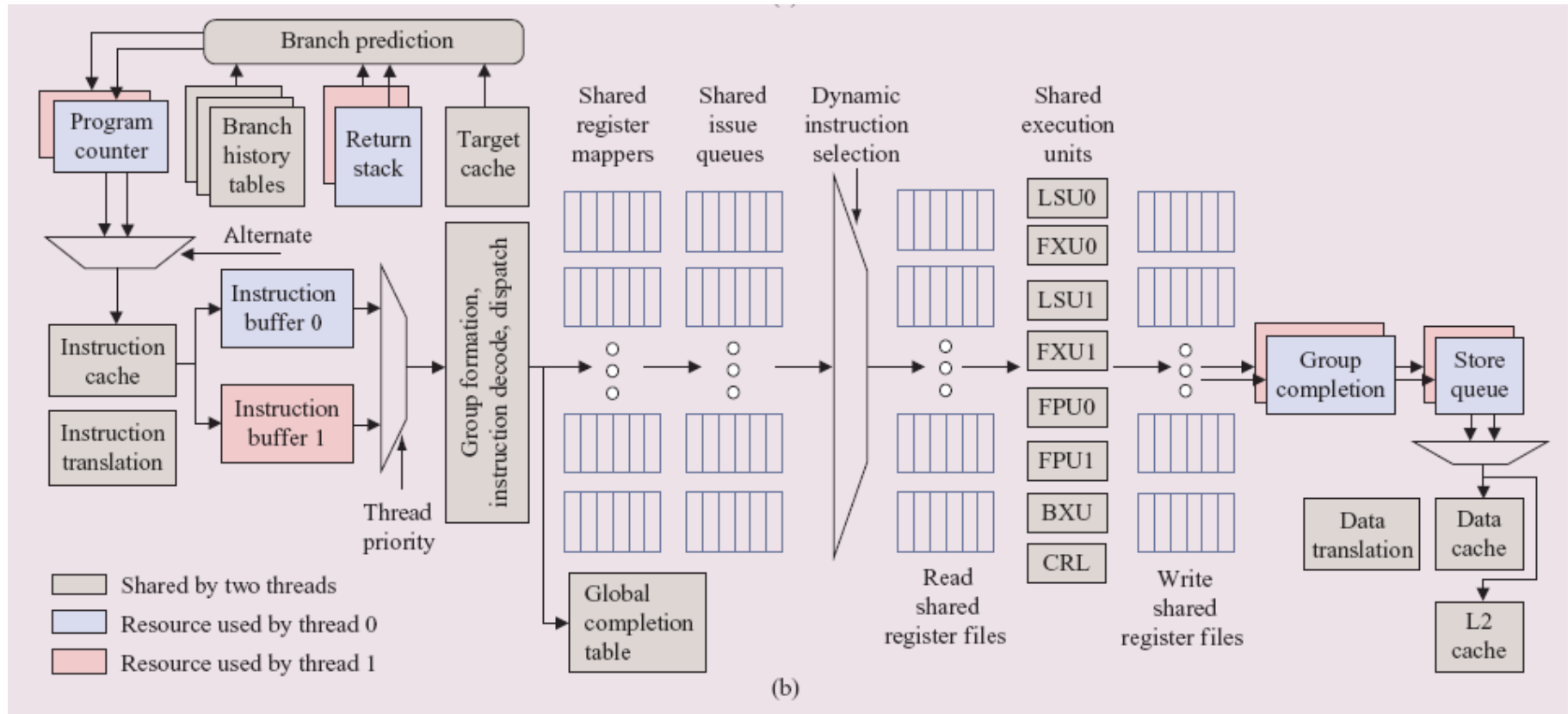
Power 5

Branch redirects

Instruction fetch



❑ Power 5



¿Por que solo 2 threads? Con 4, los recursos compartidos (registros físicos , cache, AB a memoria) son un cuello de botella.

□ Power 5

Balaneo de la carga dinámica

1- Monitorizar

cola de fallos (load en L2)
entradas ocupadas en ROB (GCT)

2 - Quitarle recursos;

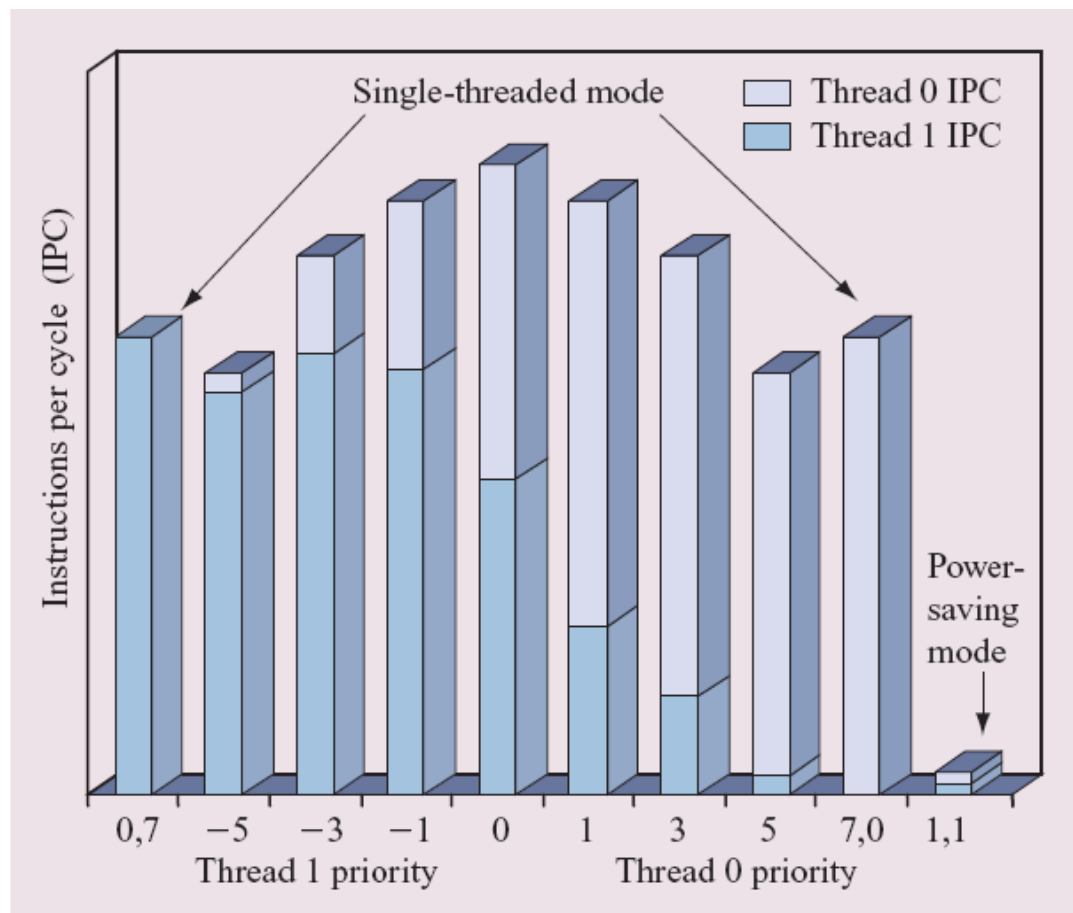
Reducir prioridad del hilo,
Inhibir decodificación (L2 miss)
Eliminar instrucciones desde emisión y
parar decodificación

3- Ajustar prioridad del hilo (hardware/software)

Baja en espera activa
Alta en tiempo real

8 niveles de prioridad

Da mas ciclos de decodificación al de mas
prioridad

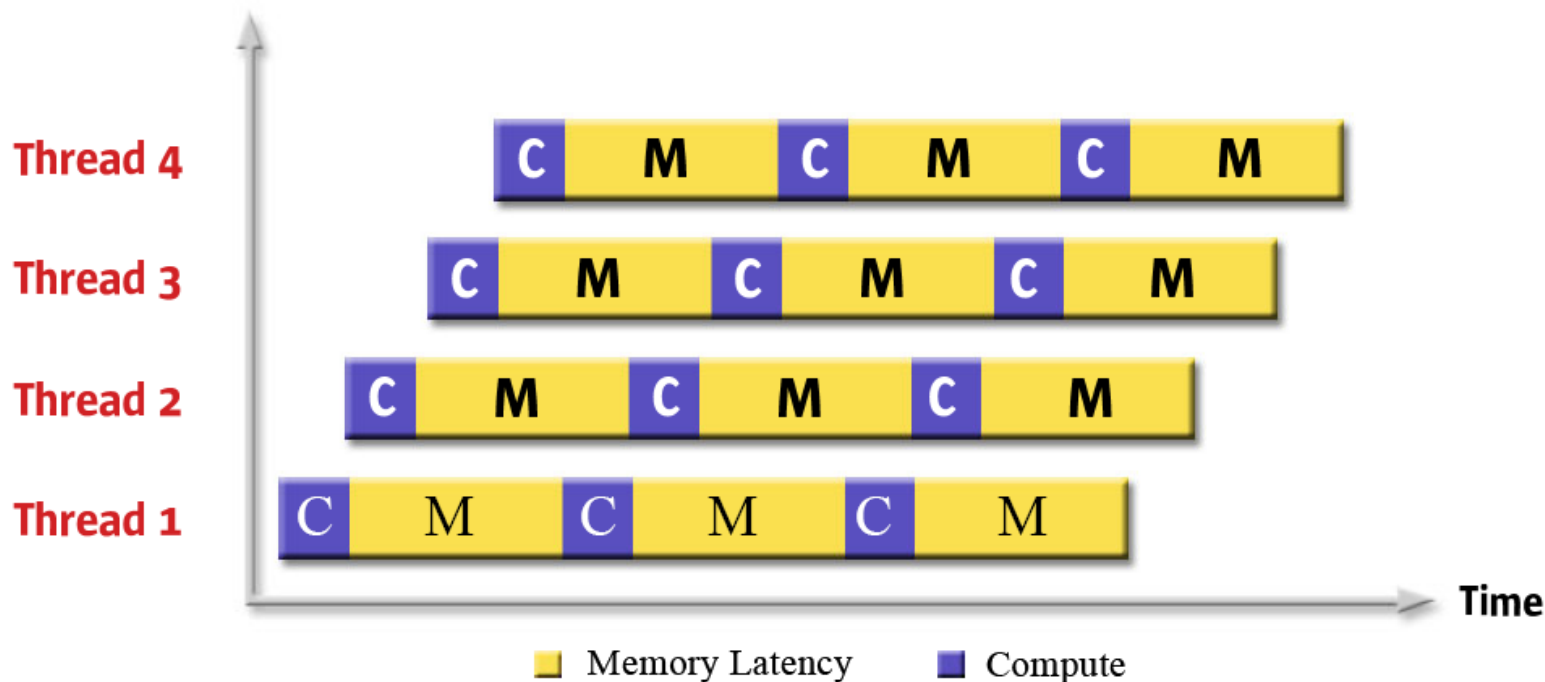


❑ Cambios en Power 5 para soportar SMT

- ❑ Incrementar asociatividad de la L1 de instrucciones y del TLB
- ❑ Una cola de load/stores por thread
- ❑ Incremento de tamaño de la L2 (1.92 vs. 1.44 MB) y L3
- ❑ Un buffer de prebúsqueda separado por thread
- ❑ Incrementar el número de registros físicos de 152 a 240
- ❑ Incrementar el tamaño de las colas de emisión
- ❑ El Power5 core es 24% mayor que el del Power4 para soportar SMT

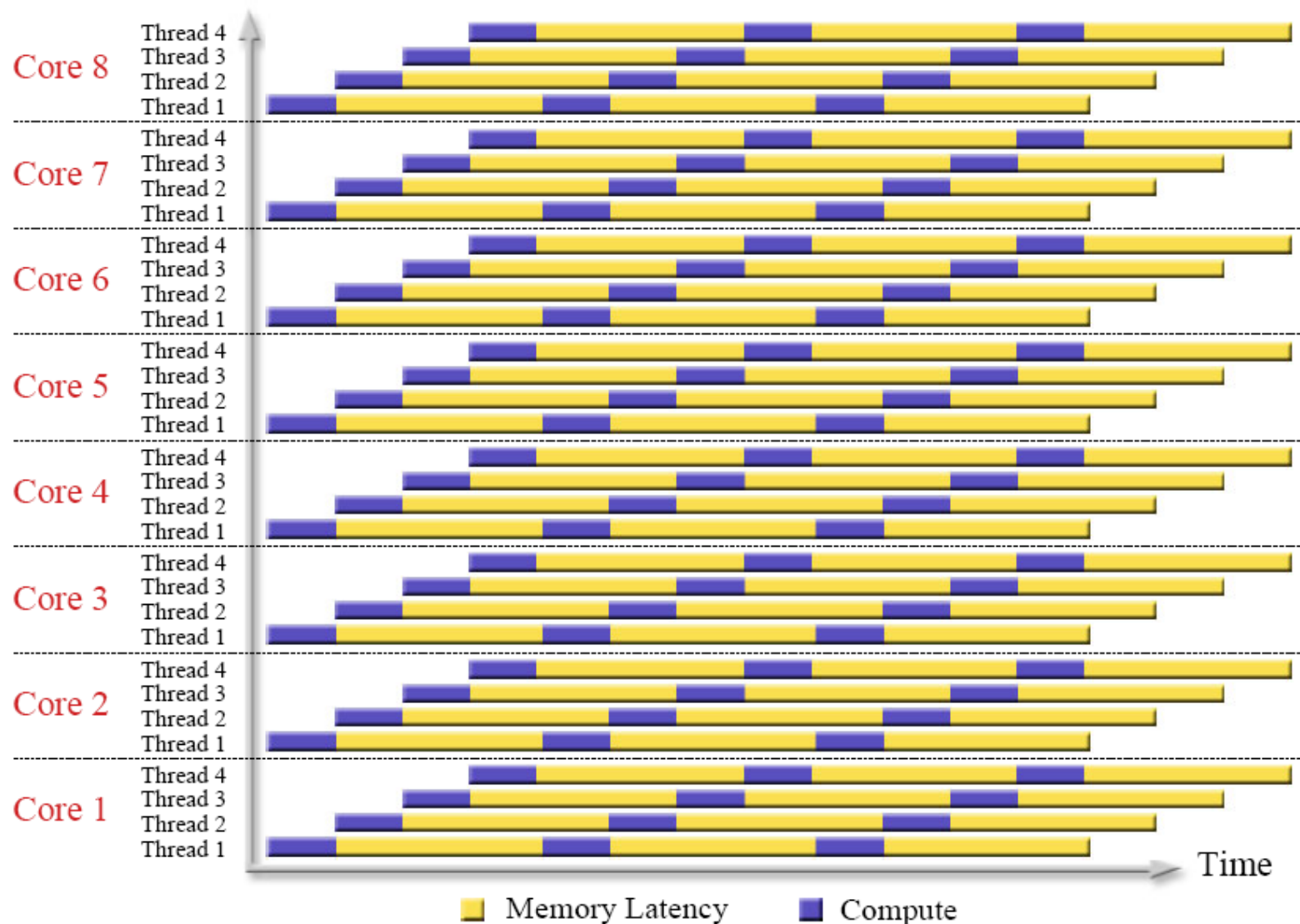
❑ Niagara

- ❑ Tolerar (soportar) la latencia de memoria mediante hilos concurrentes

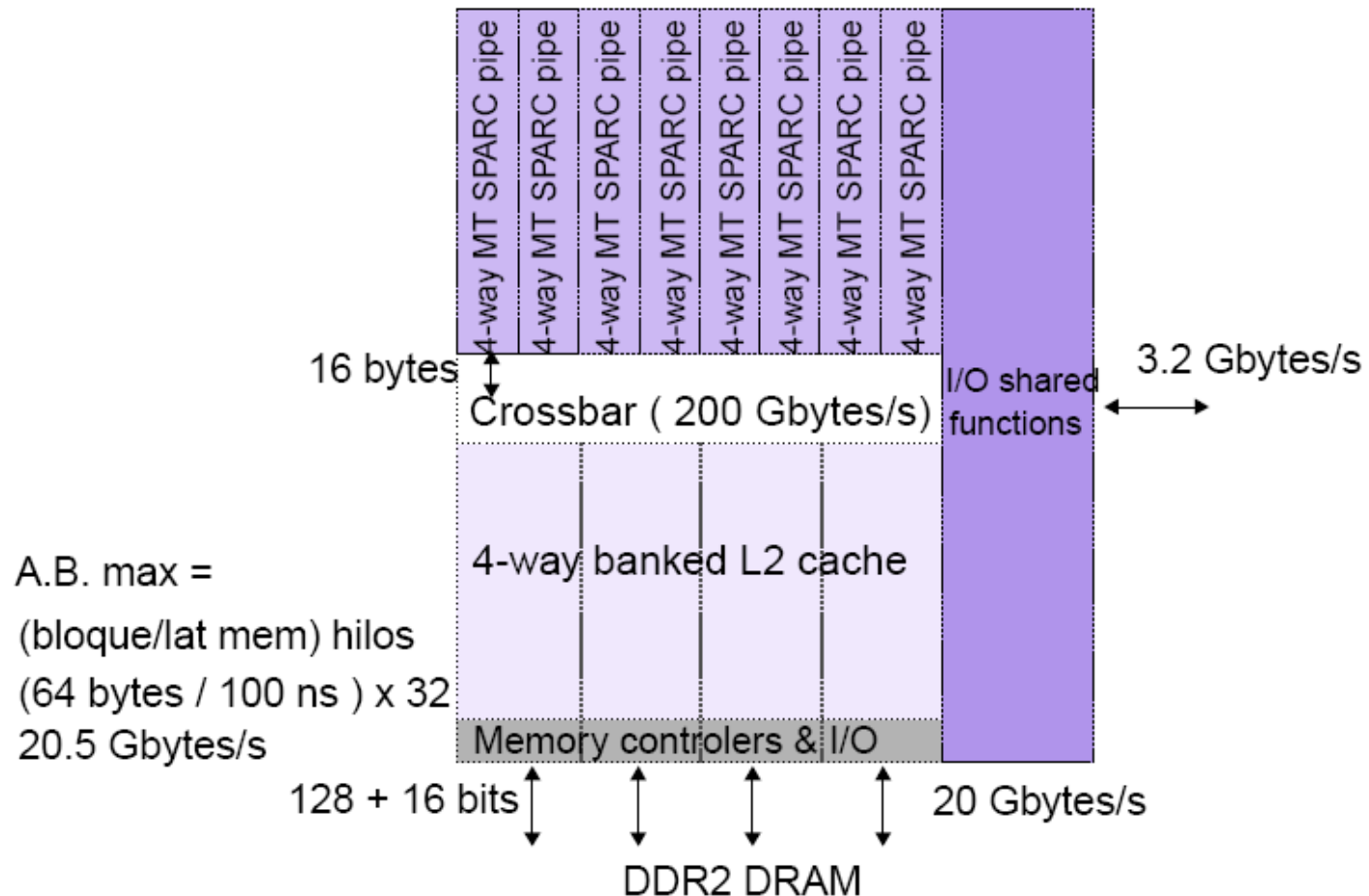


- ✓ Incrementa la utilización del procesador
- ✓ Es necesario un gran ancho de banda
 - 4 accesos concurrentes a memoria

❑ Niagara: Múltiples cores-múltiples thread



❑ Niagara UltraSparc T1

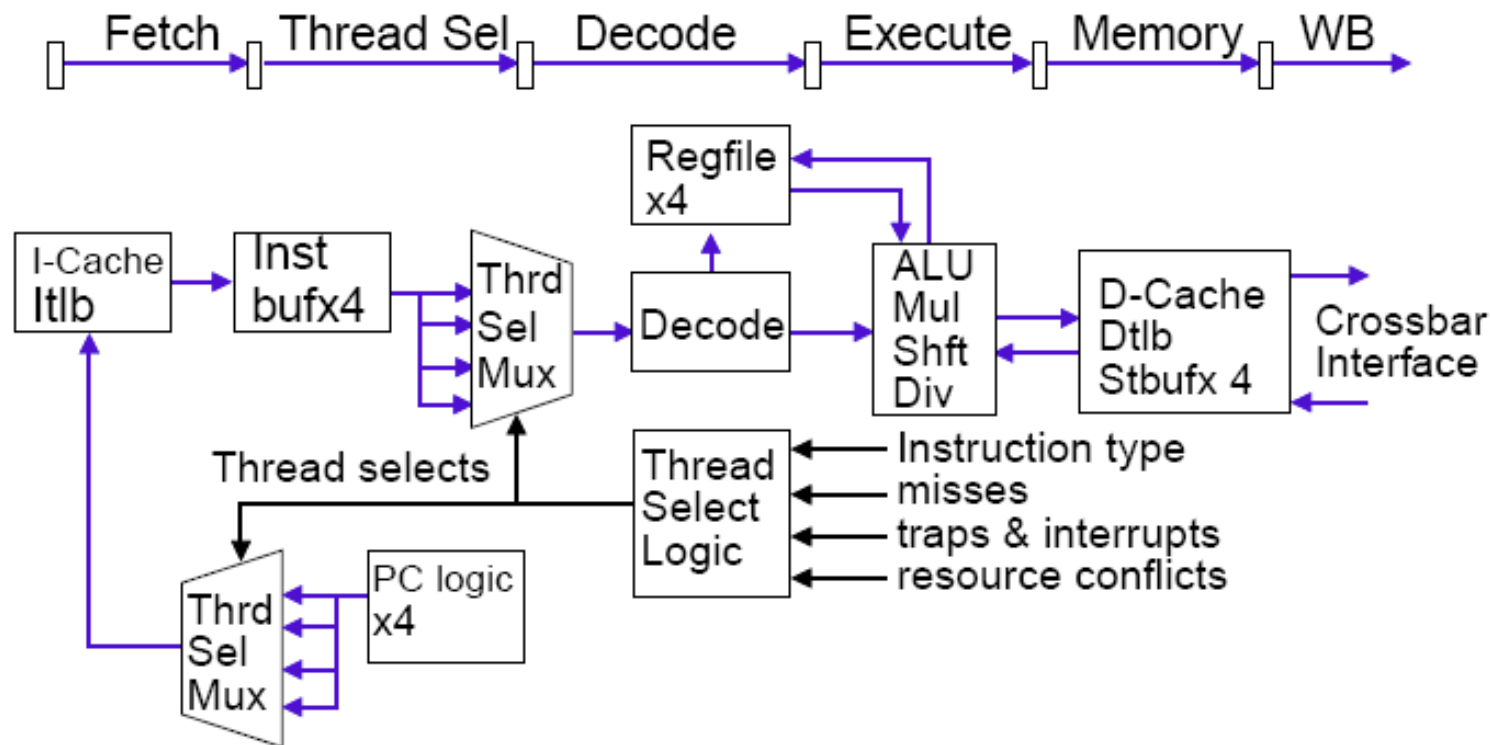


1Ghz, 1 instrucción por ciclo, 4 thread por core, 60W

I-L1 16Kb(4W), D-L1 8Kb (4W), escritura directa, L2, 3Mb(12W)

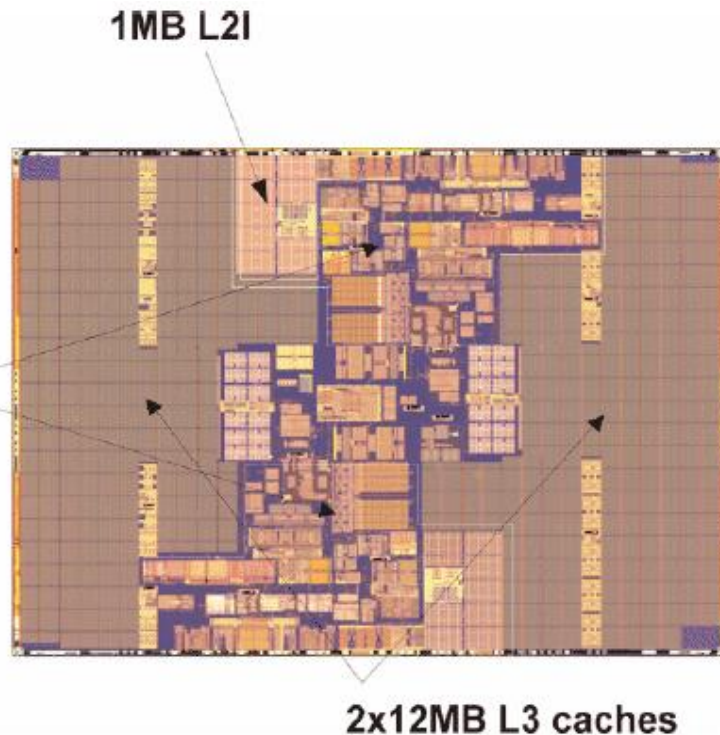
Crossbar no bloqueante, No predictor de saltos, no FP (1 por chip)

❑ Niagara UltraSparcT1



- 6 etapas
- 4 thread independientes
- algunos recursos x4: Banco de registros, contador de programa, store buffer, buffer de instrucciones
- el control de selección de hilo determina el hilo en cada ciclo de reloj
 - ✓ cada ciclo elige un hilo

□ Itanium2 9000



Itanium2 9000- Montecito

1720 Mtrans, 90 nm, 595 mm², 104 W
1,6Ghz, 2 cores y 2 threads/core
Cache L2 separada

- 1 MByte L2 Icache
- 256KB L2 Dcache

L3 mas grande

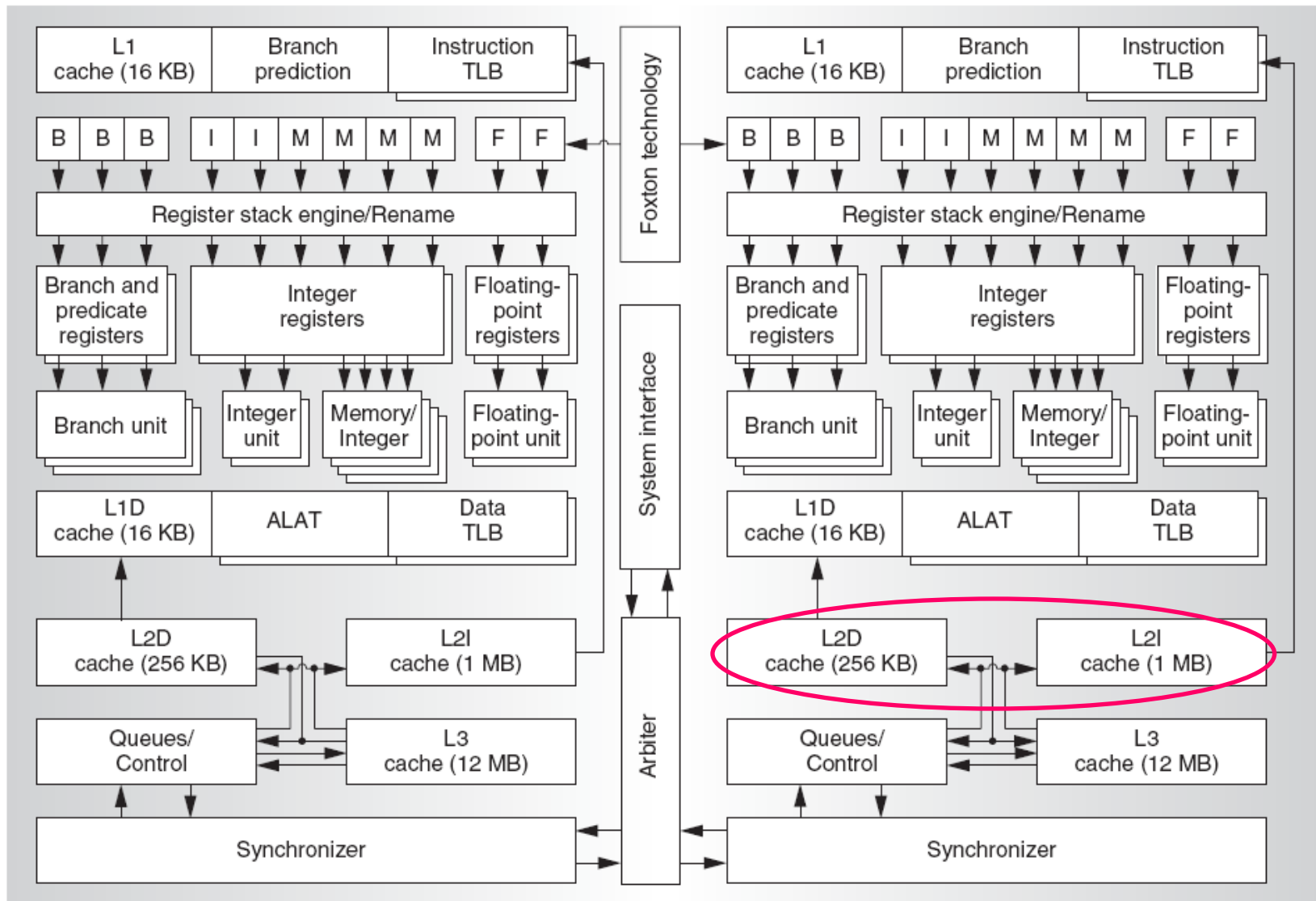
- L3 crece a 12 Mbyte por core (24MB)
- Mantiene latencia Itanium® 2

Colas y Control

- Mas L3 y L2 victim buffers
- Mejor control de las colas

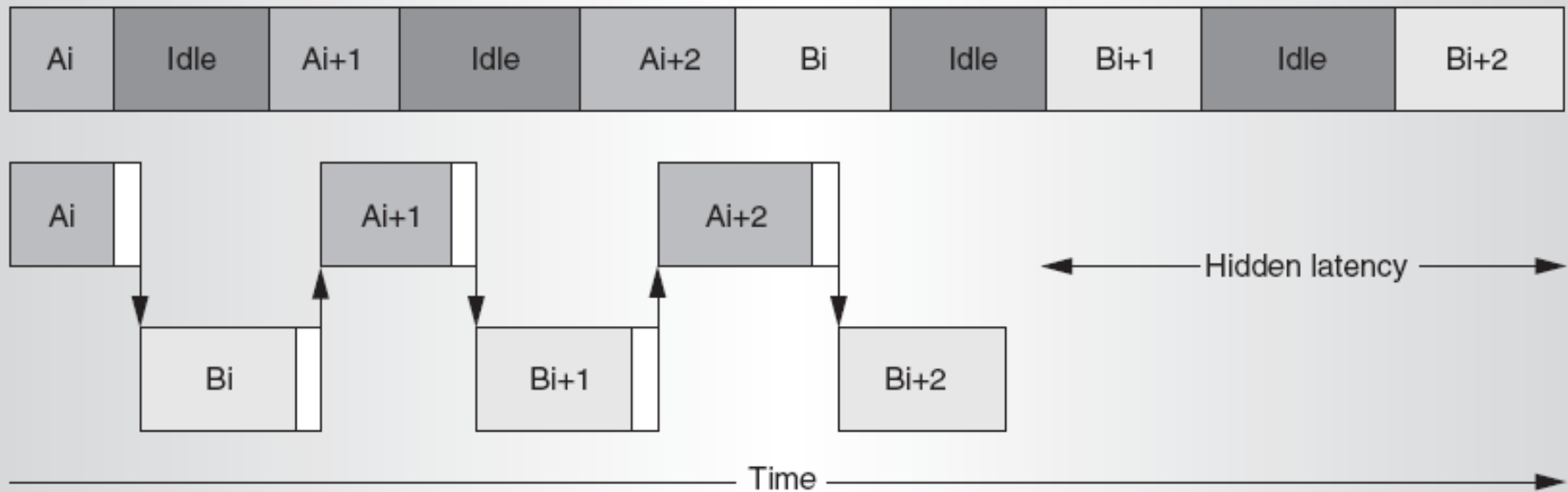
10,66 GBps AB con memoria

Itanium2 9000

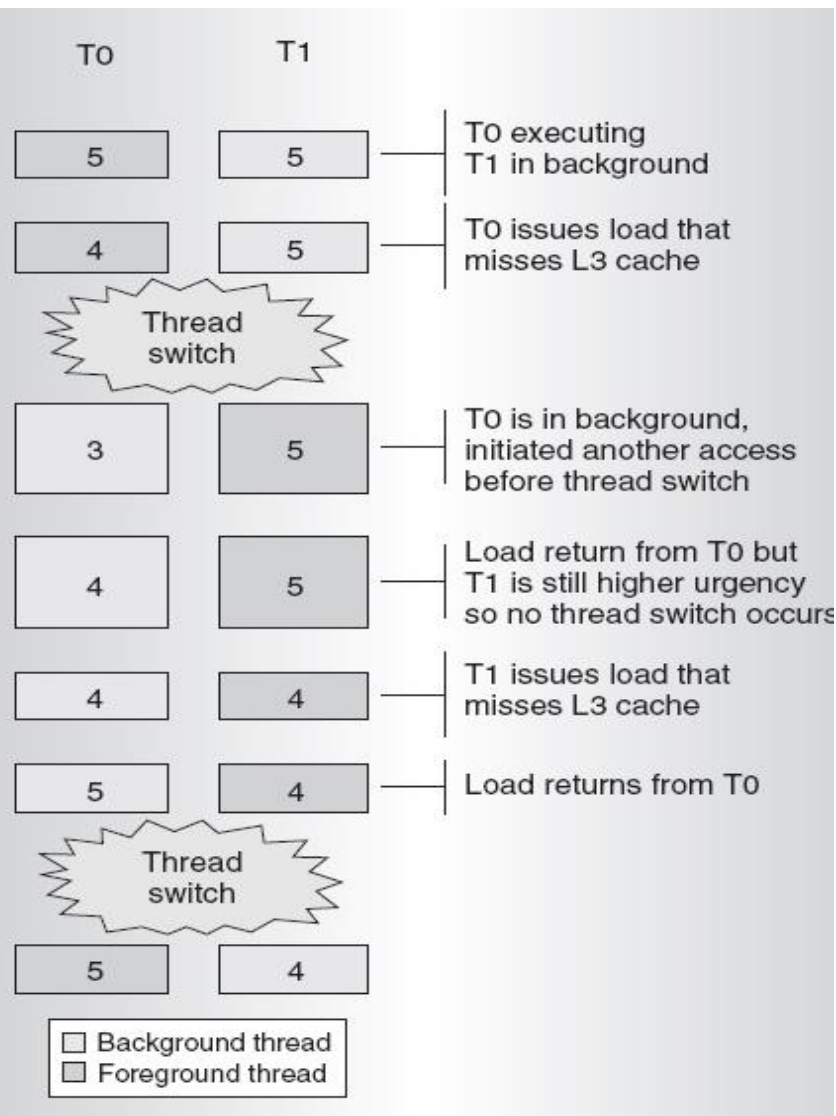


□ Itanium2 9000 Multi-Threading

- Dinámicamente asigna recursos en función del el uso efectivo a realizar.
 - Un evento de alta latencia determina la cesión de recursos por parte del thread activo.



Montecito Multi-Threading



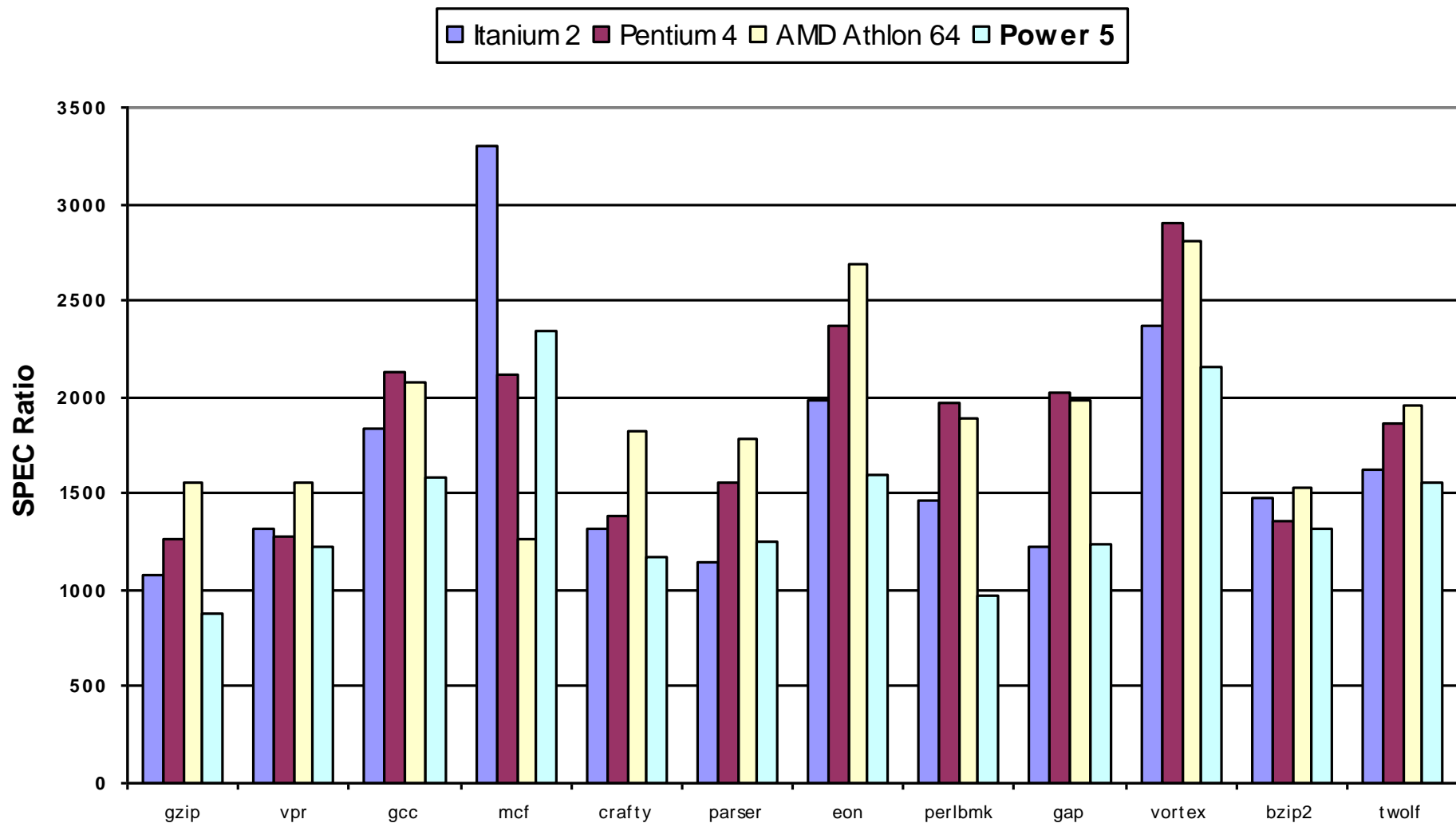
Conmutación de Threads

- Eventos de alta latencia producen un "stall" de ejecución
 - L3 misses o accesos a datos no "cacheable"
 - Agotar el "time slice" asignado a un thread
- Prioridad permite la conmutación de thread

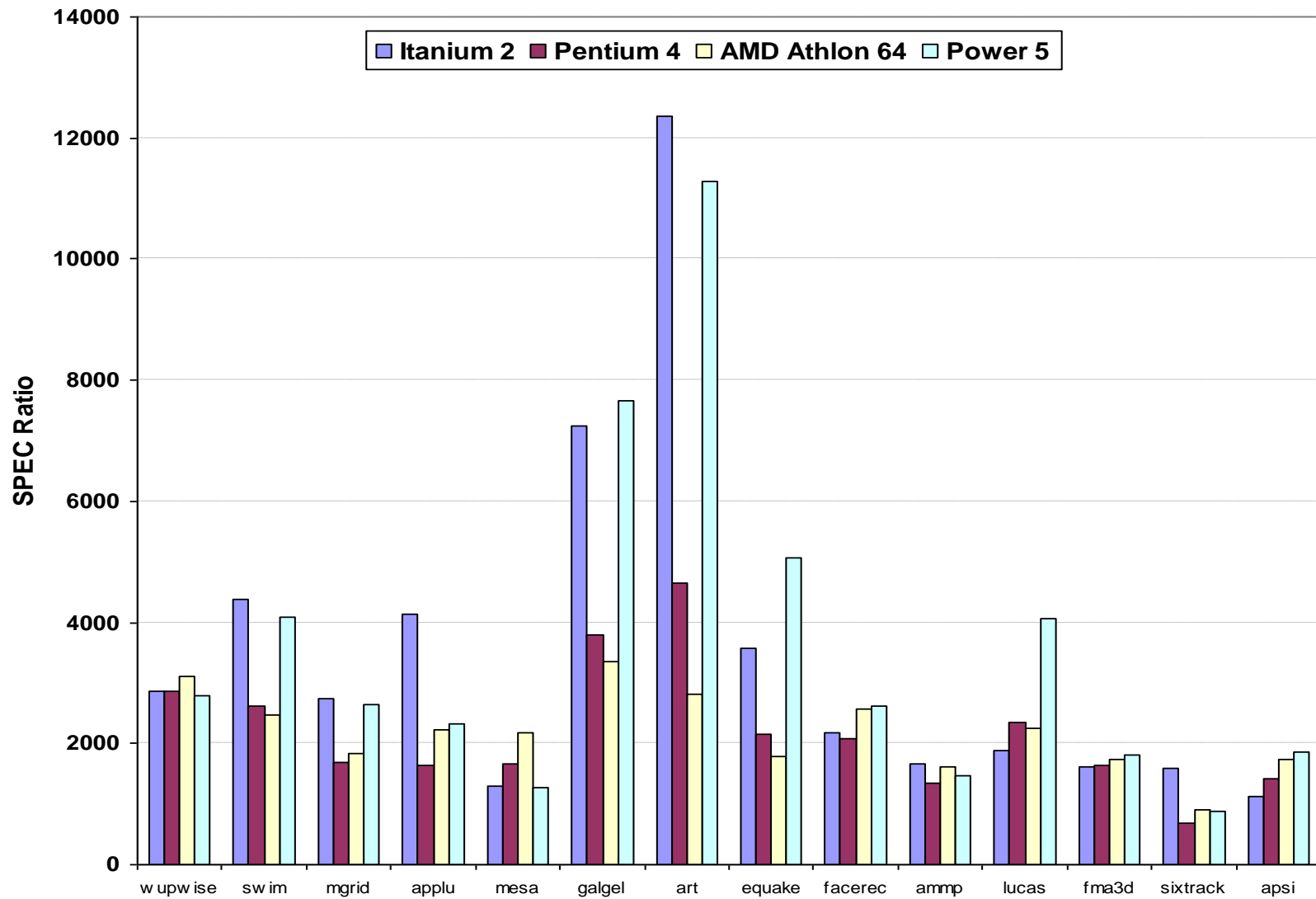
¿ Quien es mejor?

Procesador	Microarquitectura	Fetch / Issue / Execute	FU	Clock Rate (GHz)	Transis- tores Die size	Power
Intel Pentium 4 Extreme	Especulativo con planificación dinámica; Pipe profundo; SMT	3/3/4	7 int. 1 FP	3.8	125 M 122 mm ²	115 W
AMD Athlon 64 FX-57	Especulativo con planificación dinámica.	3/3/4	6 int. 3 FP	2.8	114 M 115 mm ²	104 W
IBM Power5 (1 CPU only)	Especulativo con planificación dinámica; SMT 2 CPU cores/chip	8/4/8	6 int. 2 FP	1.9	200 M 300 mm ² (est.)	80W (est.)
Intel Itanium 2	Planificación estática VLIW-style	6/5/11	9 int. 2 FP	1.6	592 M 423 mm ²	130 W

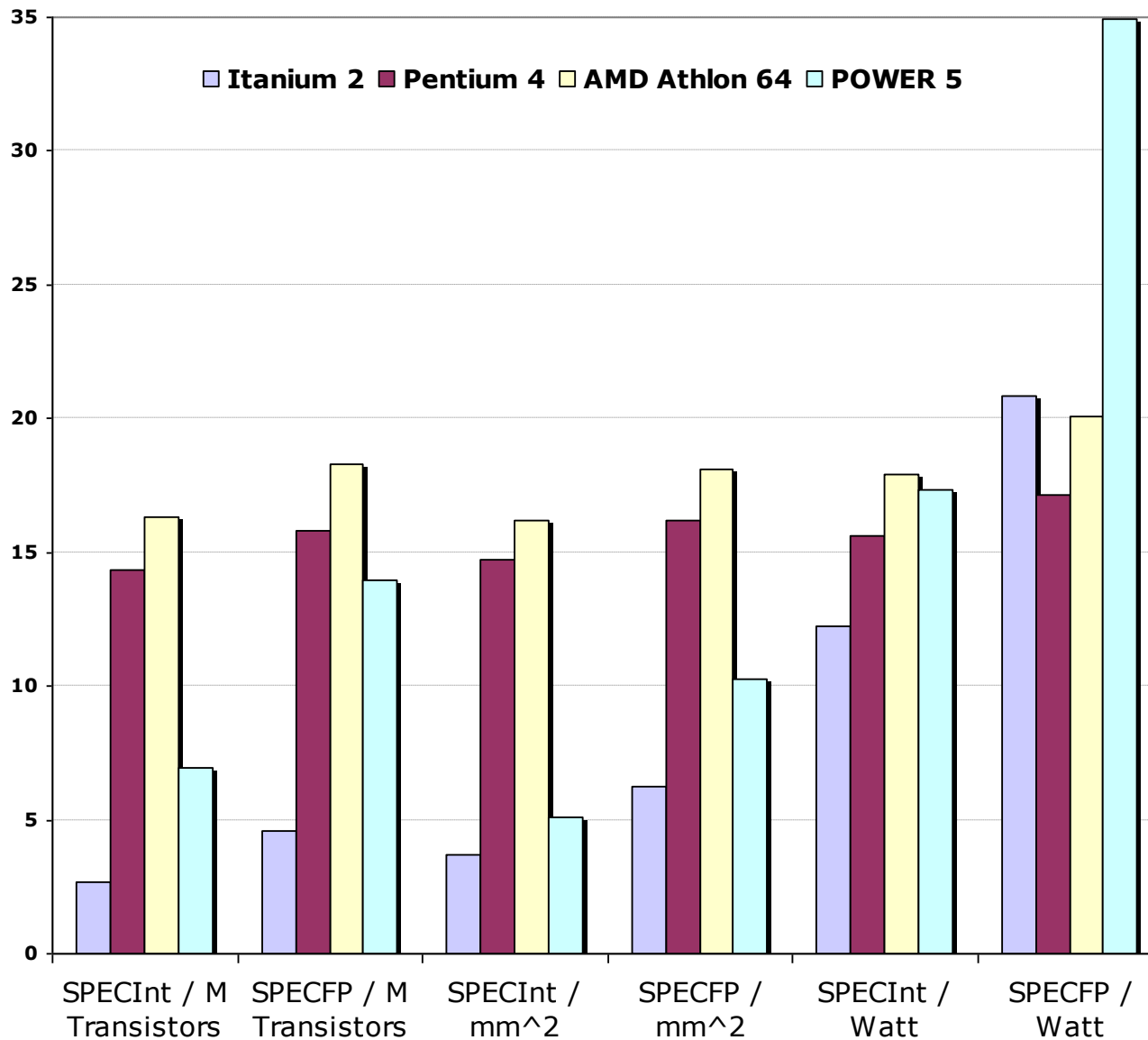
■ SPECint2000



□ SPECfp2000



□ Rendimiento normalizado: Eficiencia



	I t a n i u m 2	P e n t i u m 4	A t h l o n	P o w e r 5
Rank	2	4	n	5
Int/Trans	4	2	1	3
FP/Trans	4	2	1	3
Int/area	4	2	1	3
FP/area	4	2	1	3
Int/Watt	4	3	1	2
FP/Watt	2	4	3	1

- ❑ No hay un claro ganador en todos los aspectos
- ❑ El AMD Athlon gana en SPECInt seguido por el Pentium4, Itanium 2, y Power5
- ❑ Itanium 2 y Power5, tienen similares rendimientos en SPECFP, dominan claramente al Athlon y Pentium 4
- ❑ Itanium 2 es el procesador menos **eficiente** en todas las medidas menos en SPECFP/Watt.
- ❑ Athlon y Pentium 4 usan bien los transistores y el área en términos de eficacia
- ❑ IBM Power5 es el mas eficaz en el uso de la energía sobre los SPECFP y algo menos sobre SPECINT

- ❑ Doblar en ancho de emisión (issue rates) sobre los valores actuales 3-6 instrucciones por ciclo, a digamos 6 a 12 instrucciones requiere en el procesador
 - o de 3 a 4 accesos a cache de datos por ciclo,
 - o Predecir-resolver de 2 a 3 saltos por ciclo,
 - o Renombrar y acceder a mas de 20 registros por ciclo,
 - o Buscar de la cache de instrucciones de 12 a 24 instrucciones por ciclo.
- ❑ La complejidad de implementar estas capacidades implica al menos sacrificar la duración del ciclo e incrementa de forma muy importante el consumo.

- ❑ La mayoría de las técnicas que incrementan rendimiento incrementan también el consumo.
- ❑ Una técnica es *eficiente en energía* si incrementa más rendimiento que el consumo.
- ❑ Todas las técnicas de emisión múltiple son poco eficientes desde el punto de vista de la energía.

- ❑ La arquitectura Itanium **no** representa un paso adelante en el incremento el ILP, eliminado los problemas de complejidad y consumo.
- ❑ En lugar de seguir explotando el ILP, los diseñadores se han focalizado sobre multiprocesadores en un chip (CMP, multicores,...)
- ❑ En el 2000, IBM abrió el campo con el 1º multiprocesador en un chip, el Power4, que contenía 2 procesadores Power3 y una cache L2 compartida. A partir de este punto todos los demás fabricantes han seguido el mismo camino. (Intel, AMD, Sun, Fujitsu,...).
- ❑ El balance entre ILP y TLP a nivel de chip no es todavía claro, y probablemente será muy dependiente del tipo de aplicaciones a que se dedique el sistema.