

Doble Grado en  
Informática y  
Matemáticas

# Arquitectura de Computadores

## Tema 5. Arquitecturas de Propósito Específico



*ugr*

Universidad  
de Granada

**ETSIIT**  
Escuela Técnica Superior  
de Ingenierías Informática  
y de Telecomunicación



Doble Grado en  
Informática y  
Matemáticas

# Arquitectura de Computadores

## Tema 5. Lección 14. Arquitecturas de Propósito Específico



ugr

Universidad  
de Granada

ETSIIT  
Escuela Técnica Superior  
de Ingenierías Informática  
y de Telecomunicación



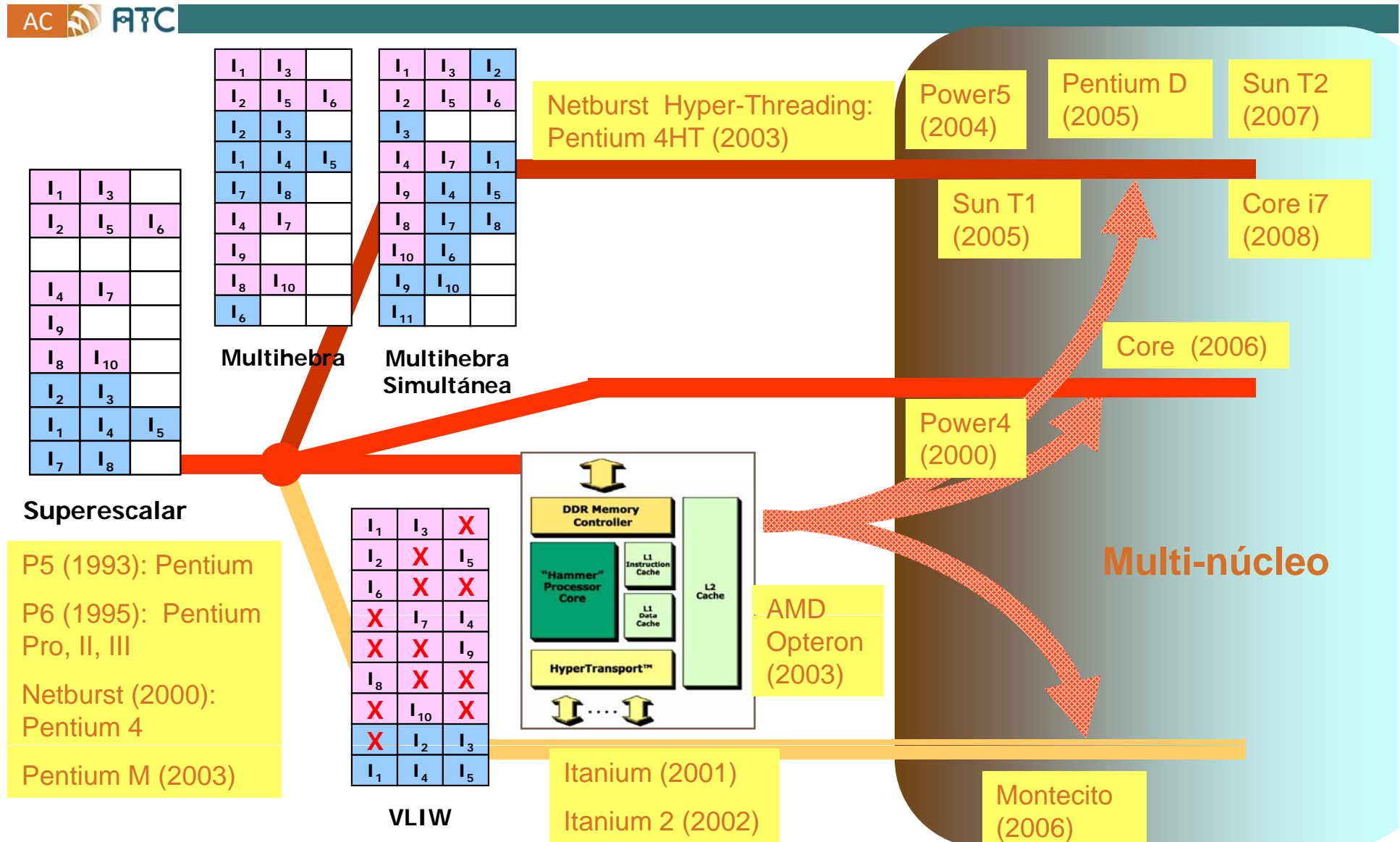
# Bibliografía

1. J. ORTEGA, M. ANGUITA y A. PRIETO. *Arquitectura de Computadores*, Thomson, 2005. (1.2.3, 4.7, 5.4.3)
2. V. AGARWAL et al.: *Clock rate versus IPC: The End of the Road for Conventional Microarchitectures*. ACM ISCA, pp.248-259, 2000.
3. IEEE Computer, Septiembre, 1997. (Microprocesadores en CI con más de 1000 Millones de Transistores)
4. IEEE Computer, Noviembre, 2005. (*Power-aware computing*)
5. M.J. FLYNN, P. HUNG: *Microprocessor design issues: thoughts on the road ahead*. IEEE Micro, pp. 16-31. Mayo-Junio, 2005.
6. S. AKHTER, J. ROBERTS. *Multi-core Programming. Increasing performance through software Multi-threading*, Intel Press, 2006.
7. ORTEGA, J.: *Entre la profecía de Moore y la ley de Amdahl*. Lección inaugural ETSIIT. Noviembre, 2008.  
[http://atc.ugr.es/~jesus/asignaturas/aci/docs/conferencia\\_etsiit\\_nov08.pdf](http://atc.ugr.es/~jesus/asignaturas/aci/docs/conferencia_etsiit_nov08.pdf)

# Contenidos

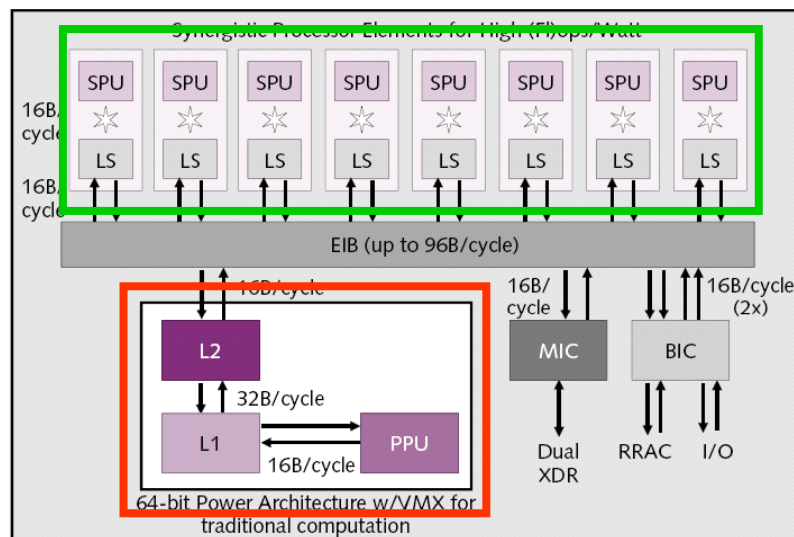
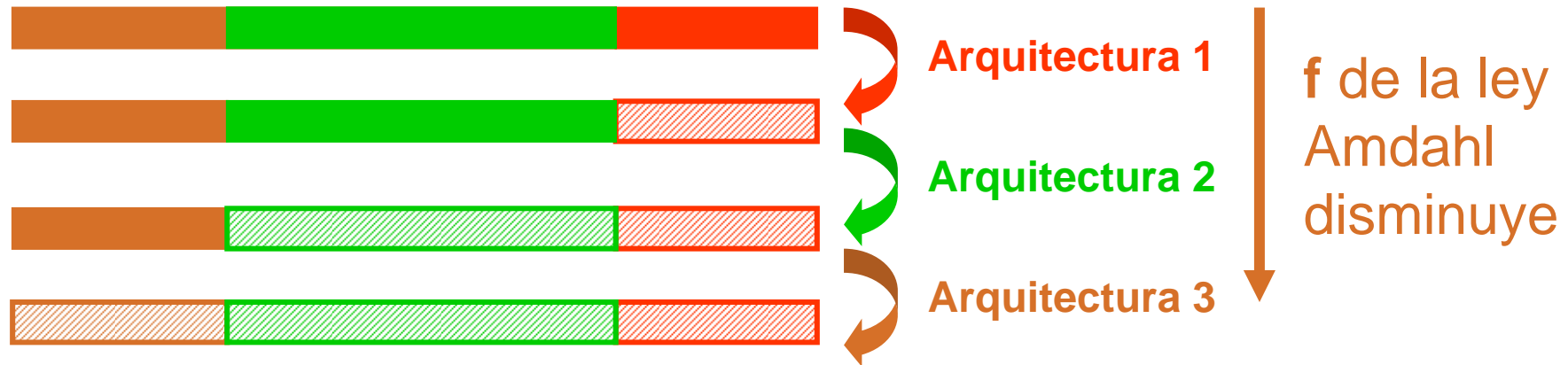
- Tendencias en el diseño de microarquitecturas
- Unidades de Procesamiento de Gráficos (GPU)
- Procesadores de Red

# Evolución de las microarquitecturas



# Arquitecturas Heterogéneas de Propósito Específico

# T1



## Procesador CELL (IBM, Sony, Toshiba)

## Inicialmente para la Playstation 3

## El Roadrunner de IBM lo utiliza

## Tecnología SOI de 90 nm / 4.6 GHz / 80W

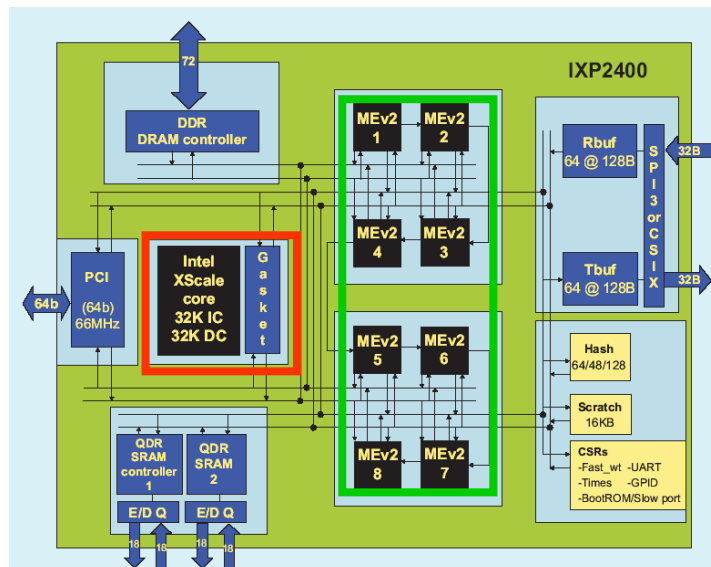
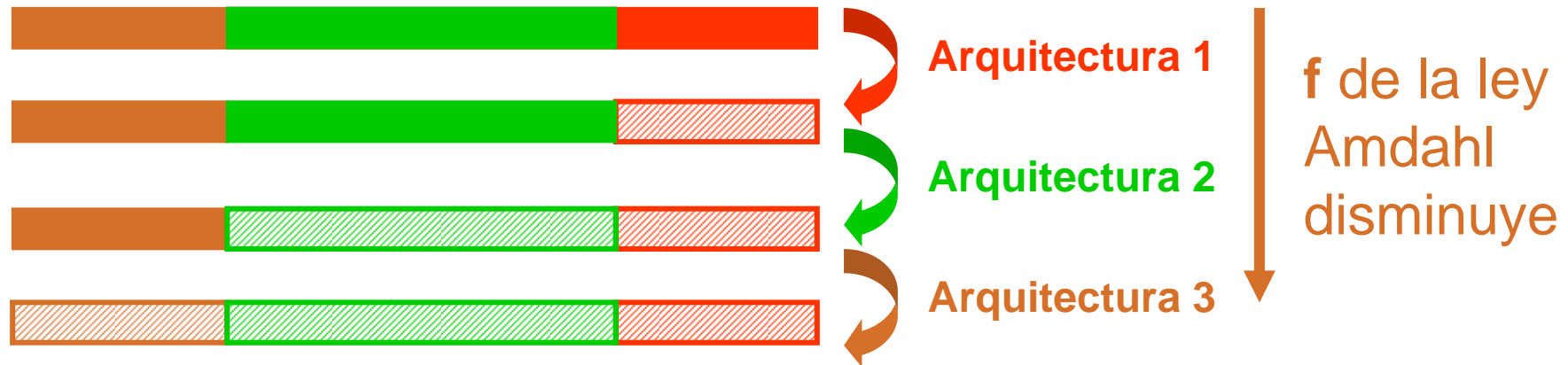
Velocidad pico: 256 GFlops (8 proc SPE a 4GHz)

Area: 235 mm<sup>2</sup>      Transistores: 235 Millones

# Arquitecturas Heterogéneas de Propósito Específico: Procesadores de Red



T1



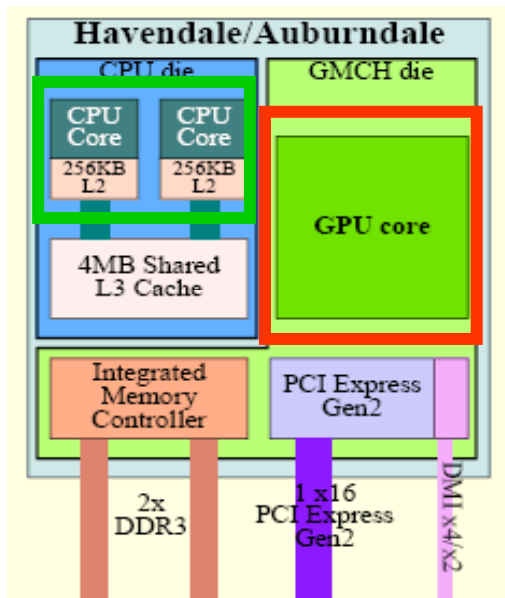
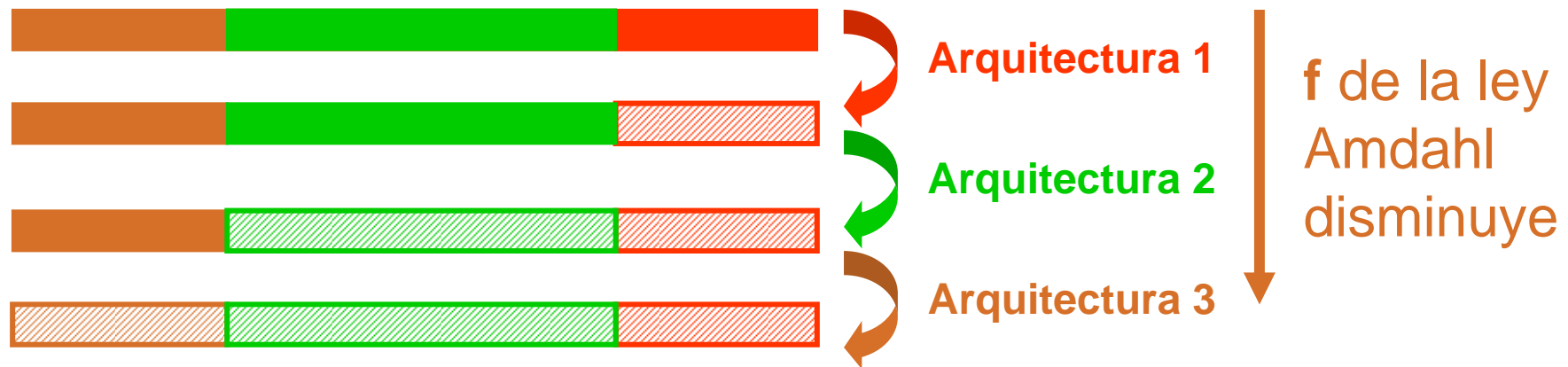
## Procesadores de Red

Arquitecturas específicas para el procesamiento de paquetes (aceleran las interfaces de red, *externalizan* servicios como la detección de intrusiones, errores, etc.)

**IXP2800** (1.4 GHz) dispone de 17 núcleos (16 MicroEngines + 1 procesador Xscale)

# Arquitecturas Heterogéneas de Propósito Específico: GPU

T1



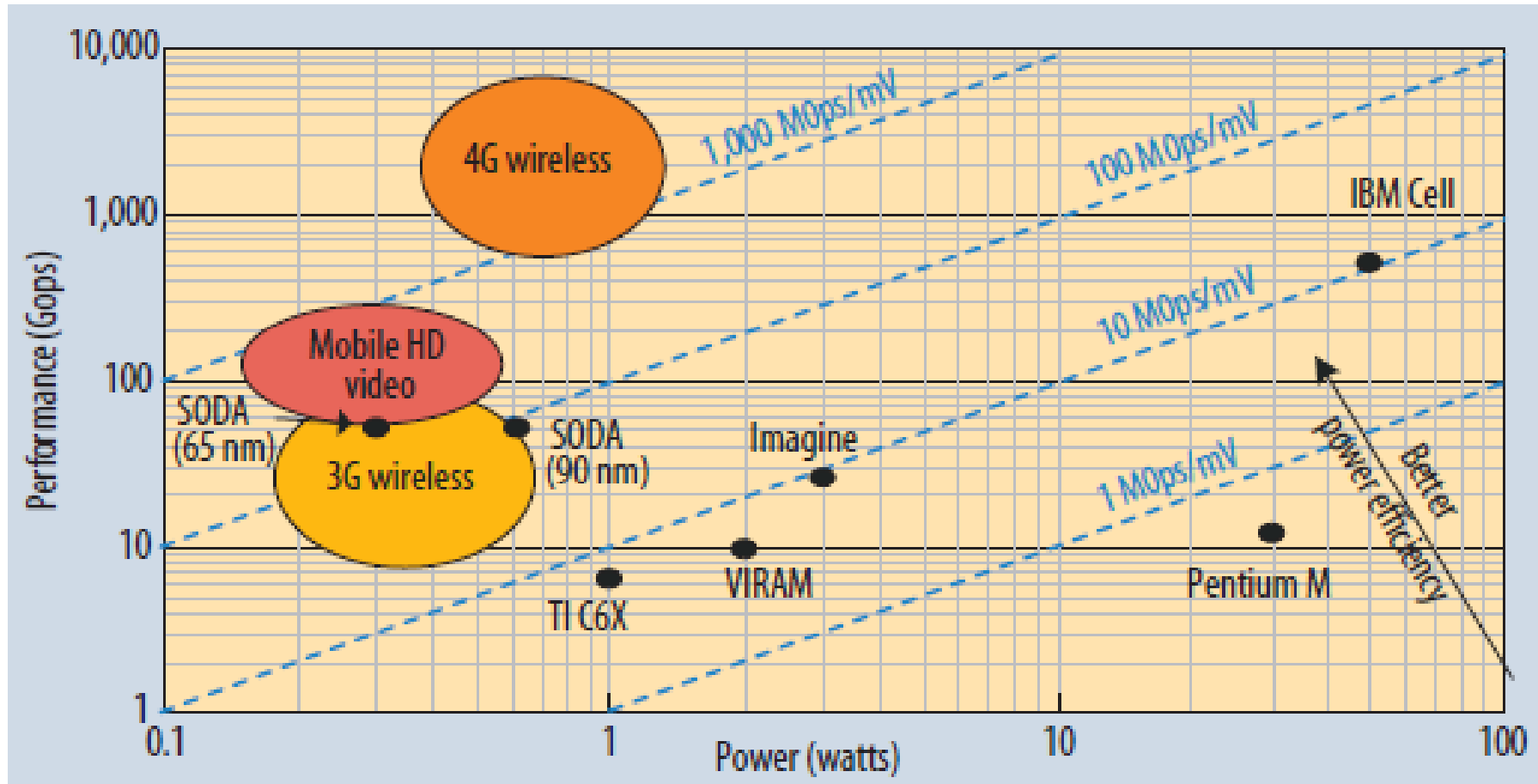
**Microarquitectura multi-núcleo Havendale de Intel para 2009.**

Tecnología de 45 nm / 3 GHz / 75 -93 W

Incluye 2 procesadores SMT de dos hebras y una unidad de procesamiento de gráficos (GPU)



# Sistemas Empotrados: Prestaciones y Consumo



Smartphones

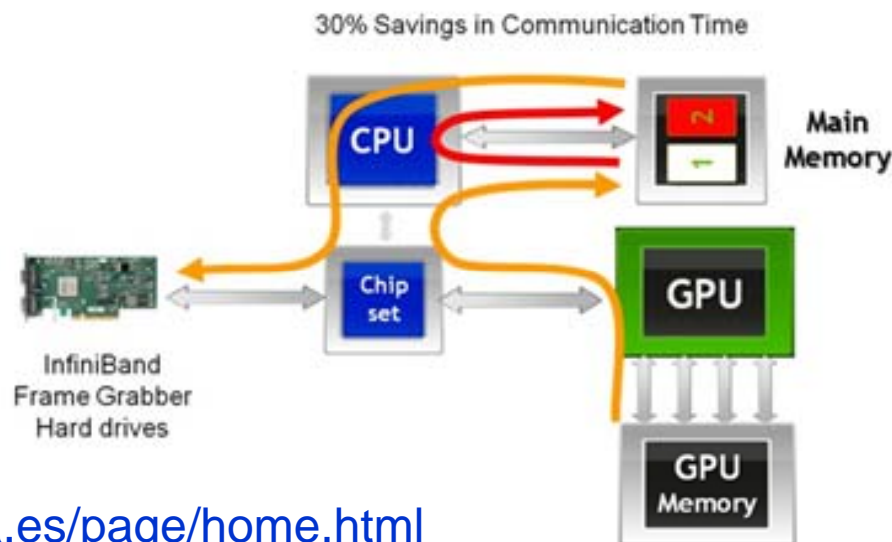
# Contenidos

- Tendencias en el diseño de microarquitecturas
- Unidades de Procesamiento de Gráficos (GPU)
- Procesadores de Red

# Unidades de Procesamiento de Gráficos (Graphics Processing Units, GPU)

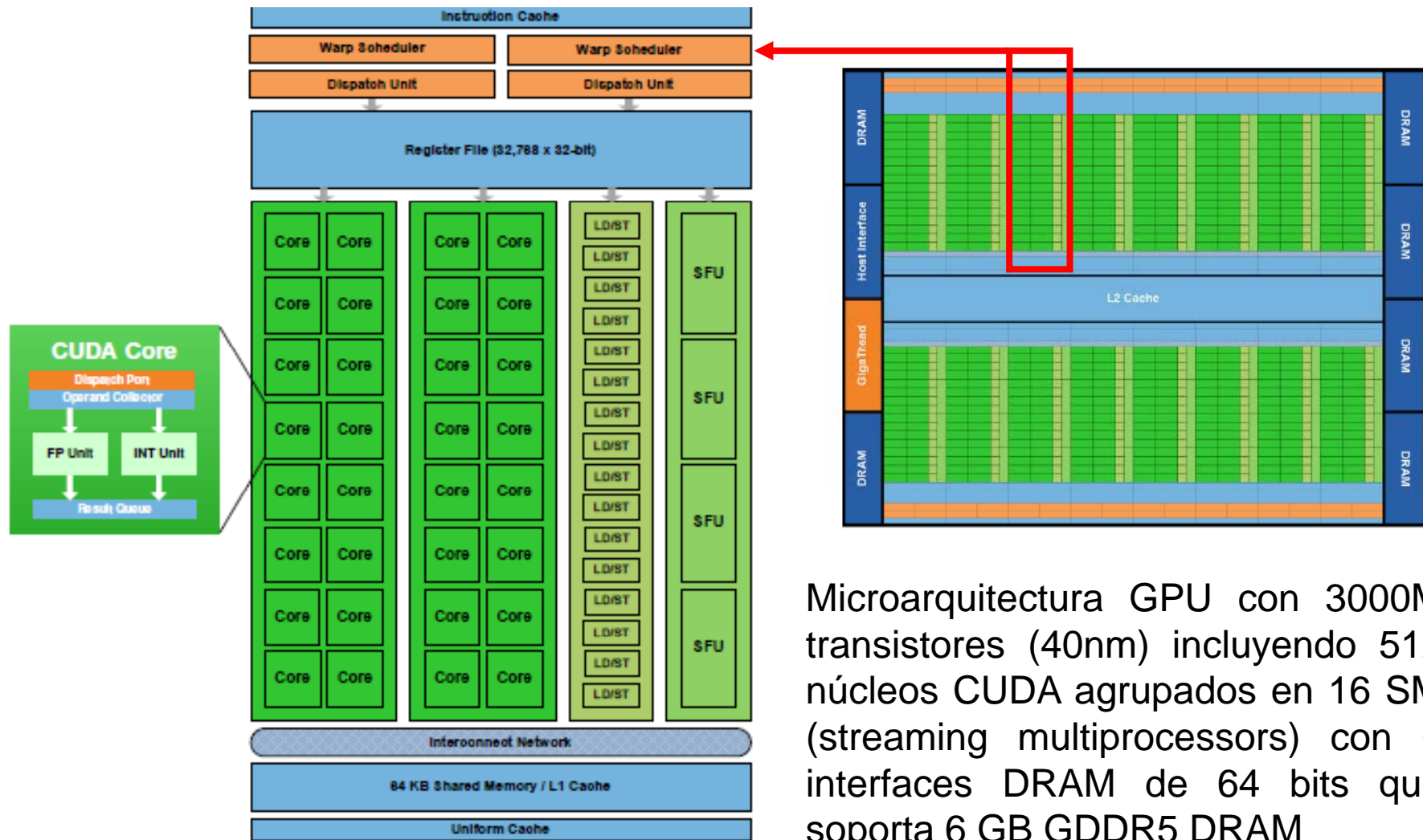


- Procesadores especializados en la manipulación y alteración de la memoria para acelerar la manipulación de imágenes
- Dado el nivel de paralelismo que implementan (para aprovechar el paralelismo de datos) se están utilizando también para acelerar otras aplicaciones (GPGPU, General Purpose GPU)
- [http://www.nvidia.es/object/fermi\\_architecture\\_es.html](http://www.nvidia.es/object/fermi_architecture_es.html)



<http://www.nvidia.es/page/home.html>

# Arquitectura GPU Fermi (Nvidia)



Microarquitectura GPU con 3000M transistores (40nm) incluyendo 512 núcleos CUDA agrupados en 16 SM (streaming multiprocessors) con 6 interfaces DRAM de 64 bits que soporta 6 GB GDDR5 DRAM

# Arquitectura GPU Fermi (Nvidia)

GPU	G80	GT200	Fermi
Transistors	681 million	1.4 billion	3.0 billion
CUDA Cores	128	240	512
Double Precision Floating Point Capability	None	30 FMA ops / clock	256 FMA ops /clock
Single Precision Floating Point Capability	128 MAD ops/clock	240 MAD ops / clock	512 FMA ops /clock
Special Function Units (SFUs) / SM	2	2	4
Warp schedulers (per SM)	1	1	2
Shared Memory (per SM)	16 KB	16 KB	Configurable 48 KB or 16 KB
L1 Cache (per SM)	None	None	Configurable 16 KB or 48 KB
L2 Cache	None	None	768 KB
ECC Memory Support	No	No	Yes
Concurrent Kernels	No	No	Up to 16
Load/Store Address Width	32-bit	32-bit	64-bit

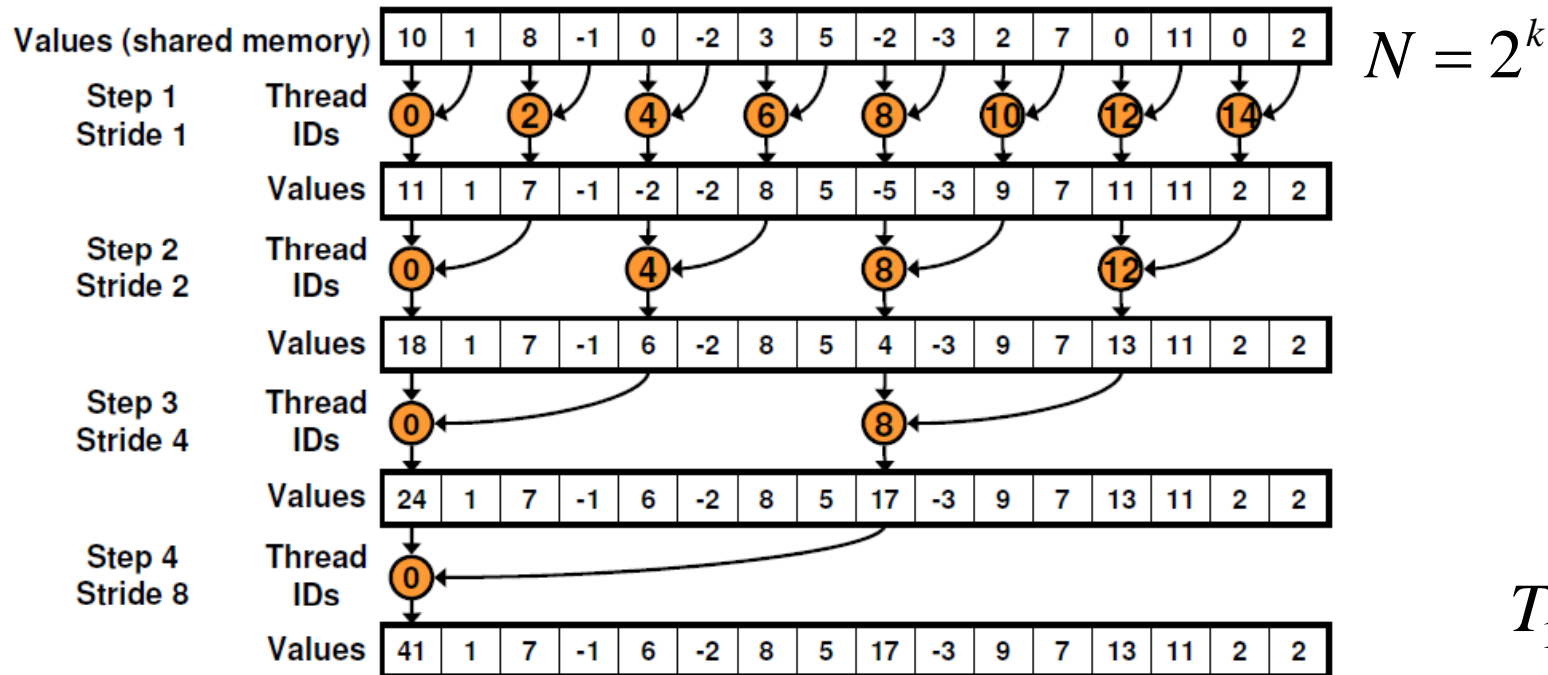
## Multiply-Add (MAD):



## Fused Multiply-Add (FMA)



# Ejemplo SIMD: Reducción I



$$T_1 = N \times t_{op}$$

$$T_N = \log_2(N) \times (t_{op} + t_{com})$$

$$S_N = \frac{N}{\log_2(N)} \times \frac{1}{1 + \frac{t_{com}}{t_{op}}} \approx \frac{N}{\log_2(N)} (t_{com} \ll t_{op})$$

$$E_N = \frac{1}{\log_2(N)} (t_{com} \ll t_{op})$$

# Contenidos

- Tendencias en el diseño de microarquitecturas
- Unidades de Procesamiento de Gráficos (GPU)
- Procesadores de Red

# Aumento de la velocidad de los enlaces

Tecnología	Ancho de banda Gb/s	Paquetes pequeños (64B) (Kp/s)	Paquetes grandes (1518B) (Kp/s)	Tiempo por paquete pequeño (microseg.)	Tiempo por paquete grande (microseg)
10Base-T	0.010	19.5	0.8	51.2	1214.40
100Base-T	0.100	195.3	8.2	5.12	121.44
OC-3	0.156	303.8	12.8	3.29	78.09
OC-12	0.622	1214.8	51.2	0.82	19.52
1000Base-T	1.000	1953.1	82.3	0.51	12.14
OC-48	2.488	4860.0	204.9	0.21	4.88
OC-192	9.953	19440.0	819.6	0.05	1.22
OC-768	39.813	77760.0	3278.4	0.01	0.31



# Aumento de la velocidad de los enlaces

## 10Base-T:

19500 paquetes/s (51.2  $\mu$ s entre paquetes de 64 bytes)

Procesamiento de un paquete: 5000 – 10000 instrucciones

**100 MIPS – 200 MIPS (Paquete cada 102400 ciclos a 2GHz)**

OC-3

0.156

303.8

12.8

3.29

78.09

## OC-192:

19440 Kpaquetes/s (0.05  $\mu$ s entre paquetes de 64 bytes)

Procesamiento de un paquete: 5000 – 10000 instrucciones

**100 GIPS – 200 GIPS (!!)** (Paquete cada 100 ciclos a 2GHz)

OC-768

39.813

77760.0

3278.4

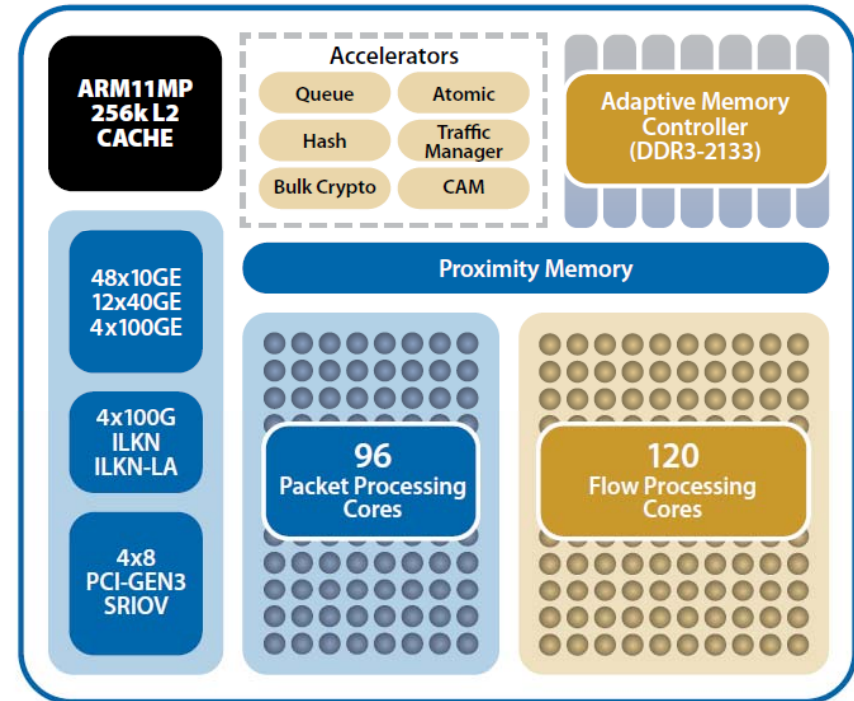
0.01

0.31

# Procesador NFP-6XXX (Netronome)



<http://netronome.com/>



- Performance
  - Over 256 Gbps wire speed programmable stateful flow and packet processing
  - 384 million packets per second
  - 307 billion instructions per second

## 200 Gbps Programmable Flow Processors

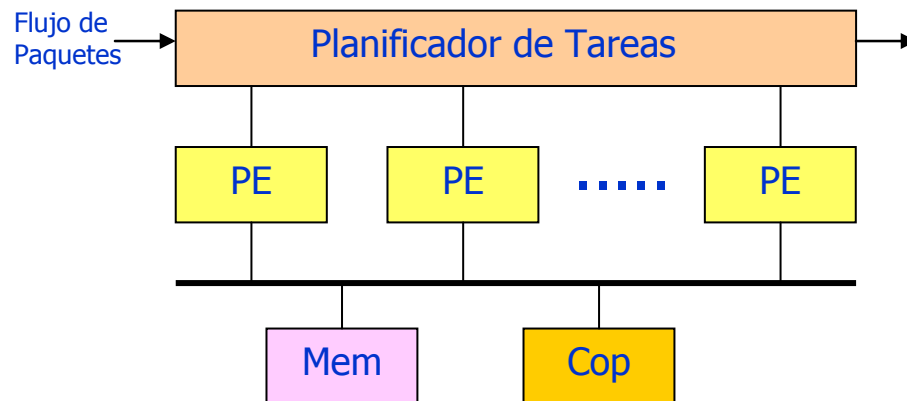
# Procesadores de Red II

AC



ATC

## Modelo de Procesador de Red Paralelo



### Ejemplo:

Línea a **10 Gbps**; Paquetes de **50 Bytes**

Llega un paquete cada **40 ns** (se tienen **25 Mpps**)

Suponiendo que se tienen **16 PE**

Cada PE puede procesar un paquete durante un tiempo igual a **16x40 ns = 640 ns** y tiene que tener un rendimiento (throughput) de **25/16 = 1.56 Mpps**

Los elementos de proceso (**PE**) tienen arquitecturas RISC con algunas instrucciones especiales de manipulación de bits adecuadas para procesar paquetes y no incluyen aritmética compleja

Tienen caches de instrucciones y de datos pequeñas (la mayor parte de los datos no se reutilizan)

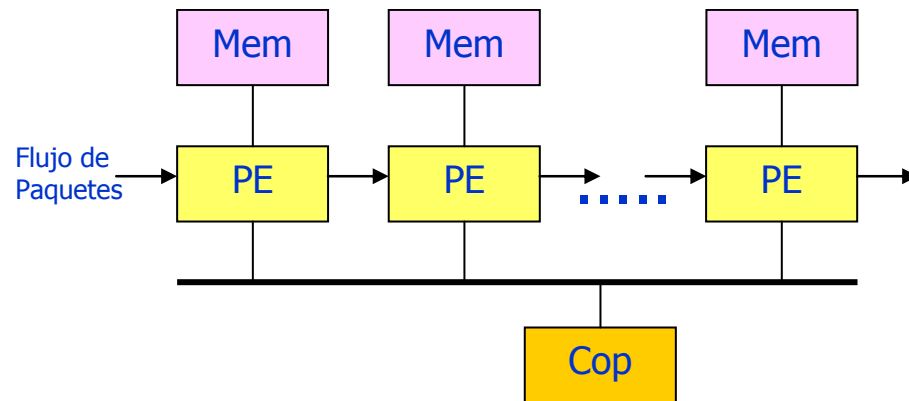
El **Planificador de Tareas** se encarga de asignar paquetes a los PE a medida que llegan y es responsable de preservar la secuencia de paquetes. Puede implementarse de forma que sea programable o no.

Los **coprocesadores** cooperan con los PE para acelerar funciones de comunicación comunes que son computacionalmente costosas (clasificación de paquetes, funciones criptográficas, o de búsqueda en tablas). Son compartidos por varios PE, que acceden a ellos a través de mensajes, memoria compartida o instrucciones especiales

En un PE las fuentes de latencia principales son los accesos a memoria y a los coprocesadores. El uso de **multihebra** contribuye a ocultar latencias de una hebra solapándolas con el procesamiento de otras hebras.

# Procesadores de Red III

## Modelo de Procesador de Red Segmentado



El procesamiento del paquete se divide en múltiples etapas, cada una de las cuales se encarga de una tarea de procesamiento necesaria para realizar la función de procesamiento.

Cada PE puede estar optimizado para realizar una tarea concreta.

De esta forma, el paquete va pasando por los distintos PE.

Los coprocesadores cooperan con los PE para acelerar algunas funciones de comunicación comunes que son computacionalmente costosas (clasificación de paquetes, funciones criptográficas, o de búsqueda en tablas).

### Ejemplo:

Línea a **10 Gbps**; Paquetes de **50 Bytes**

Llega un paquete cada **40 ns** (se tienen **25 Mpps**)

Suponiendo que se tienen **16 PE** en el cauce

Cada PE puede procesar un paquete durante un tiempo igual a **40 ns** y tiene que tener un rendimiento (throughput) de **25 Mpps** (en el modelo paralelo, los PE tienen unos requisitos de throughput menores) pero el tiempo de procesamiento del paquete puede ser mayor **16x40 ns = 640 ns**

# Para ampliar .....

## ➤ Páginas Web:

- <http://www.semichips.org/home.cfm> . SIA (Semiconductor Industry Association). *"The technology roadmap for semiconductors"*.
- <http://www.cs.wisc.edu/arch/www> . Página de Arquitectura de Computadores (*WWW Computer Architecture Page*)

## ➤ Artículos de Revistas (para comprobar si las predicciones se cumplen):

- IEEE Micro, Julio-Agosto, 2000. (Microprocesadores para el siglo XXI)
- IEEE Computer: "Multicore and Many-Core Architectures". Diciembre, 2009.
- Wittenbrink, C.M.; et al.: "Fermi GF100 GPU Architecture". IEEE Micro, pp.50-59. March/April, 2011.
- Ran Giladi: "Network Processors: Architecture, Programming, and Implementation". Morgan Kaufmann, 2008.

## ➤ Algunas revistas para mantenerse informado:

IEEE Computer, IEEE Micro, IEEE Spectrum, Commun. ACM,....