

INF351 – Computación de Alto Desempeño

Microarquitecturas y Características

PROF. ÁLVARO SALINAS

Microarquitecturas

Las microarquitecturas corresponden a las distintas organizaciones de un dispositivo respecto a diseños y tecnologías.

En las GPU de NVIDIA capaces de utilizar CUDA se distinguen las siguientes según su orden de lanzamiento:

- Tesla
- Fermi
- Kepler
- Maxwell
- Pascal
- Volta
- Turing

Tesla

Lanzamiento: 2006

Compute capability: 1.x

CUDA cores per SM: 8

Principales cambios respecto a su predecesor:

- Soporte para CUDA
- Unified shader model
- Streaming multiprocessors (SM)

Fermi

Lanzamiento: 2009

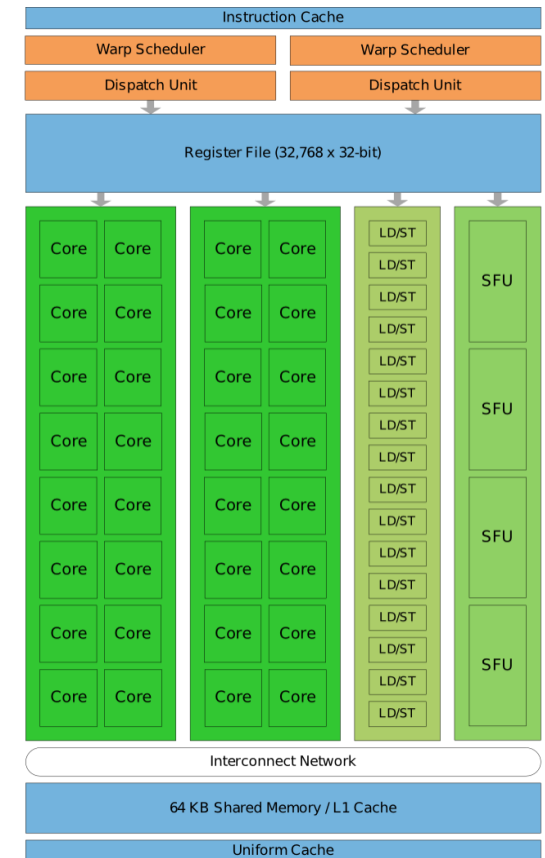
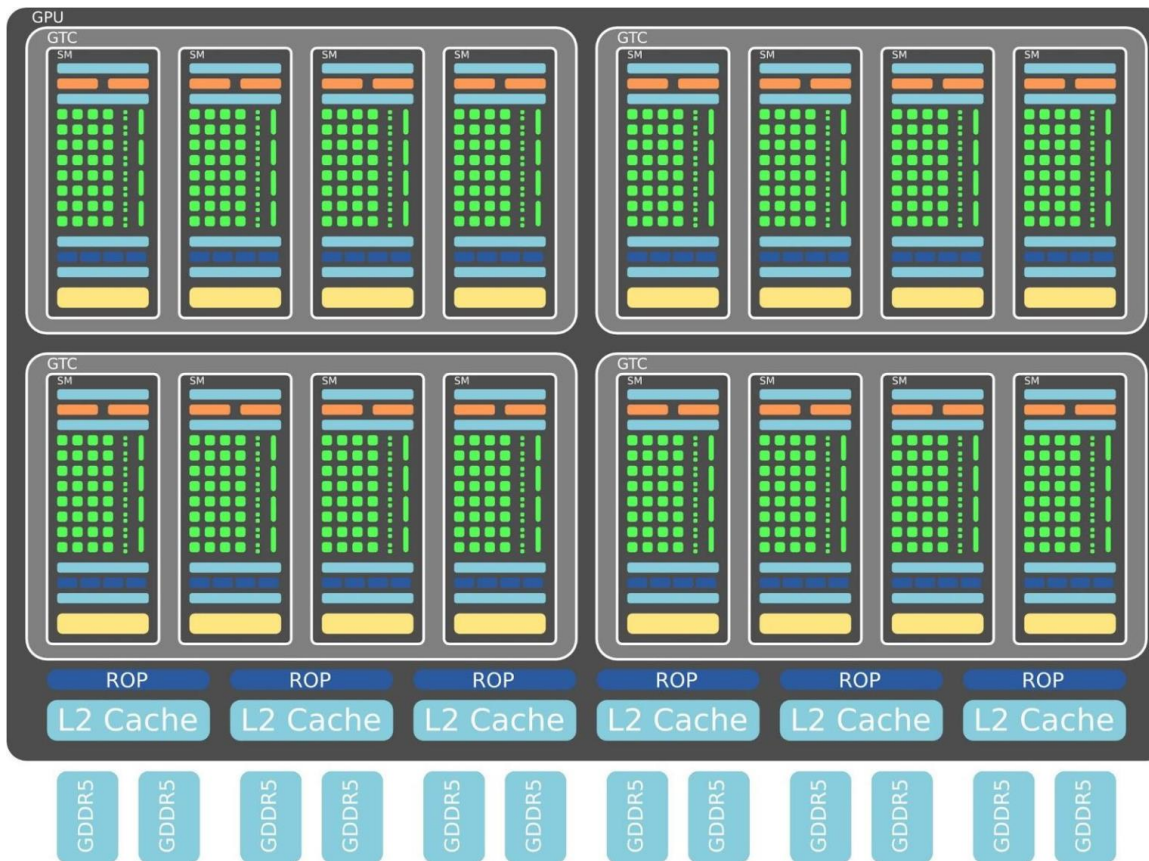
Compute capability: 2.x

CUDA cores per SM: 32

Principales cambios respecto a su predecesor:

- Más CUDA cores
- GigaThread global scheduler
- Soporte para FP64

Fermi



Kepler

Lanzamiento: 2012

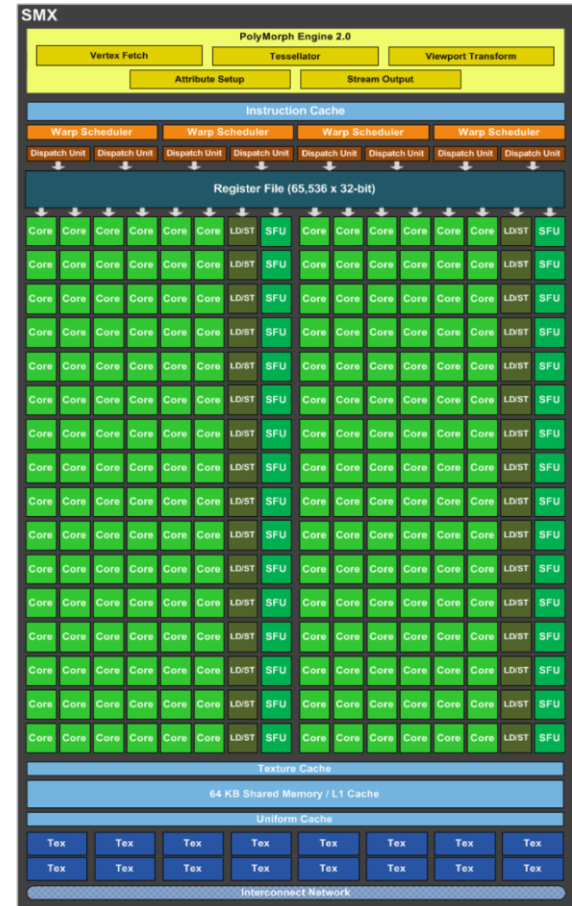
Compute capability: 3.x

CUDA cores per SM (SMX): 192

Principales cambios respecto a su predecesor:

- Más CUDA cores
- Warps de 32 subprocesos
- SMX son más eficientes en power consumption

Kepler



Kepler



Maxwell

Lanzamiento: 2014

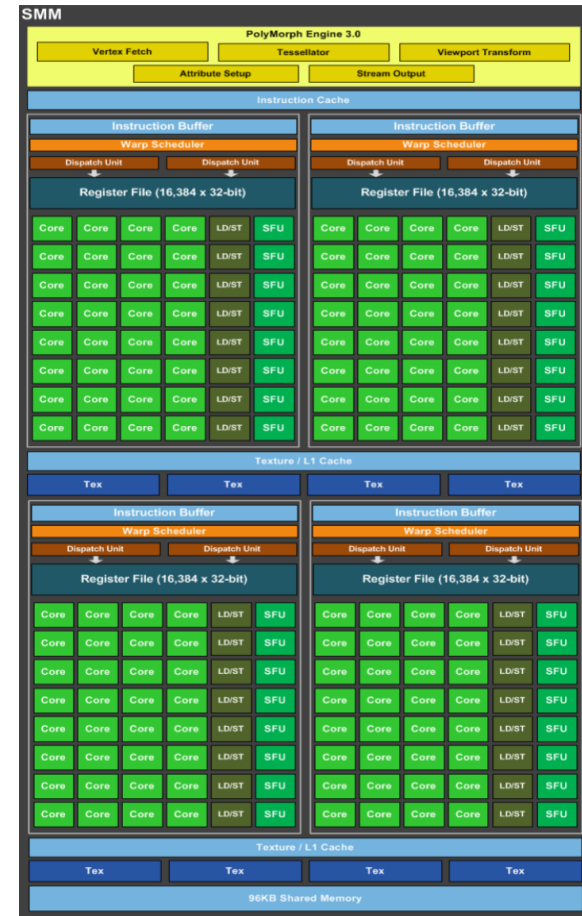
Compute capability: 5.x

CUDA cores per SM (SMM): 128

Principales cambios respecto a su predecesor:

- Gran aumento de memoria L2 cache
- SMM sacrifican un poco de performance por mucha eficiencia

Maxwell



Pascal

Lanzamiento: 2016

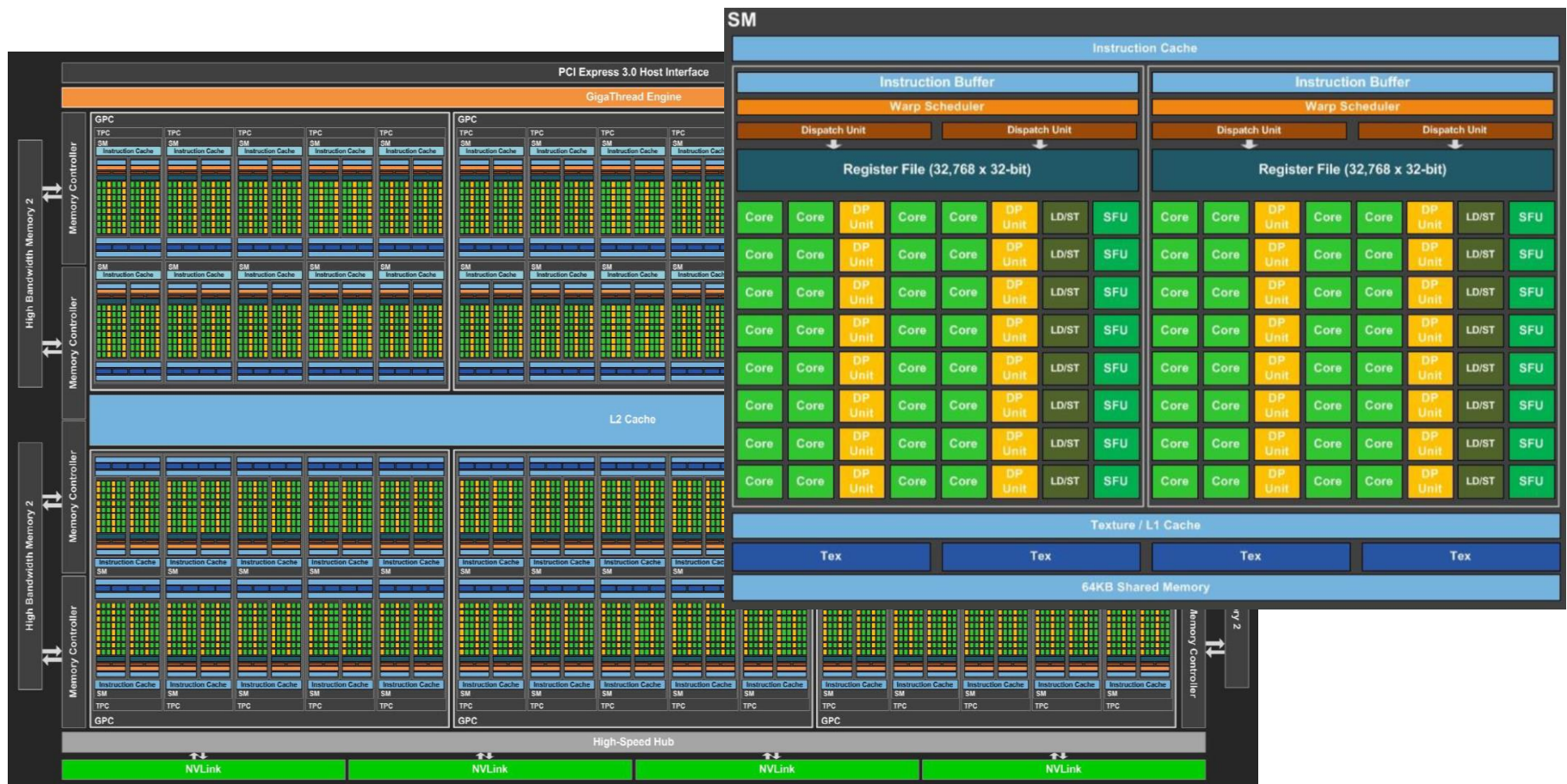
Compute capability: 6.x

CUDA cores per SM: 64

Principales cambios respecto a su predecesor:

- High Bandwith Memory 2
- Operaciones de FP16 se realizan el doble de rápido que operaciones de FP32, y éstas a su vez, el doble de rápido que las de FP64
- Más registros por CUDA core y más Shared Memory
- Dynamic scheduler

Pascal



Volta

Lanzamiento: 2017

Compute capability: 7.x

CUDA cores per SM: 64

Principales cambios respecto a su predecesor:

- Tensor cores
- NVLink 2.0 (solo Tesla)

Volta



Turing

Lanzamiento: 2018

Compute capability: 7.5

CUDA cores per SM: 64

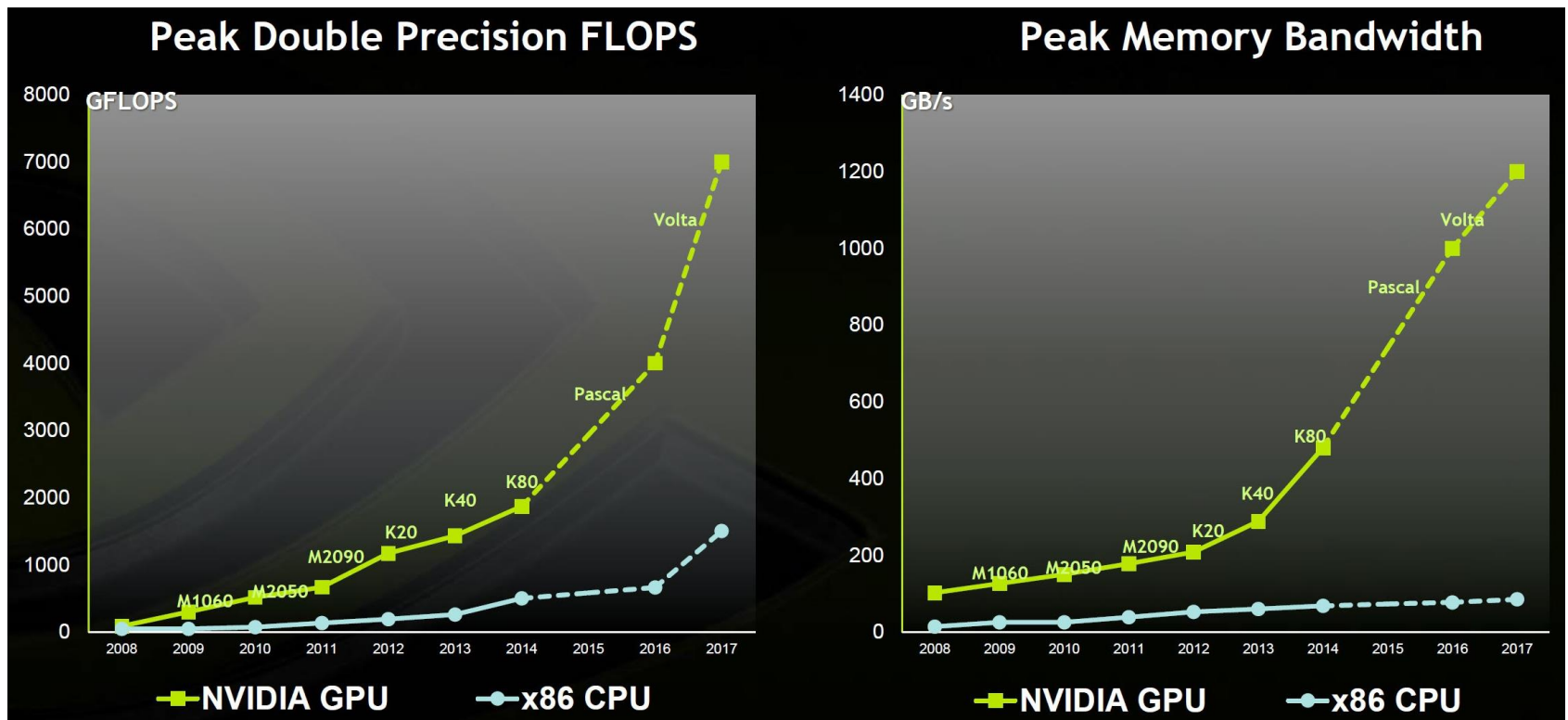
Principales cambios respecto a su predecesor:

- Ray-Tracing (RT) cores

Turing



Evolución



Diferencias según compute capability

Para saber cuáles son las características soportadas por cada compute capability, así también como las especificaciones técnicas de cada generación, consulte el siguiente link:

<https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#features-and-technical-specifications>

Para consultar la compute capability de su dispositivo, acceda al siguiente link:

https://en.wikipedia.org/wiki/CUDA#GPUs_supported