



**FaMAF | GPGPU
Computing Group**

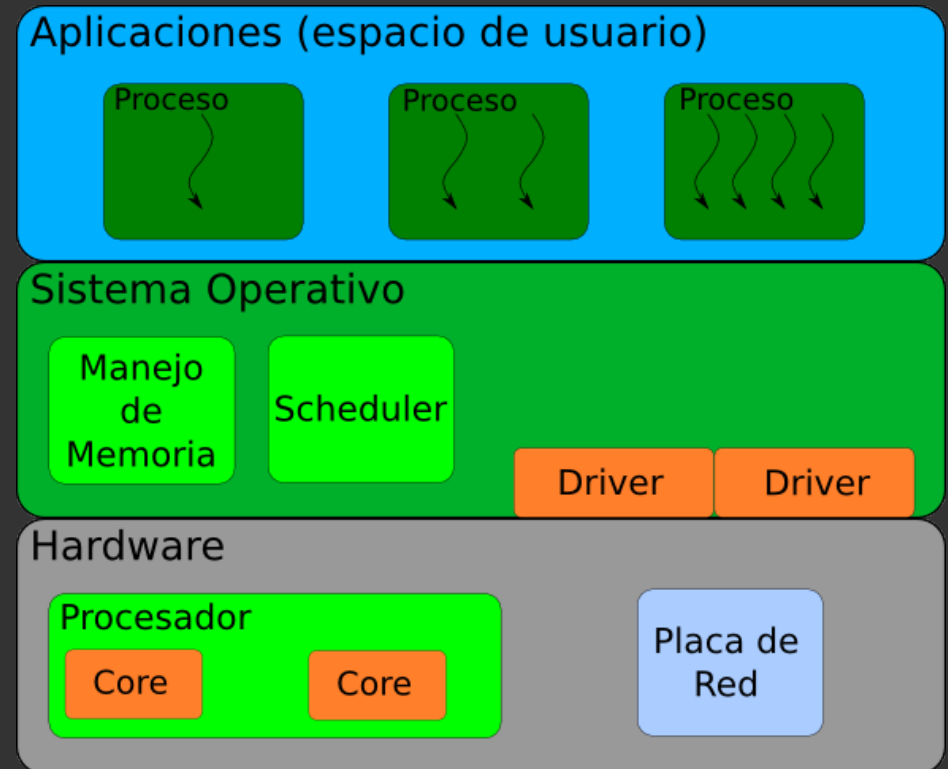


Honrarás tus recursos y tu hardware

Matías Bellone

PCs

- Muy estructurada
- Muy versátil
- Estructura base



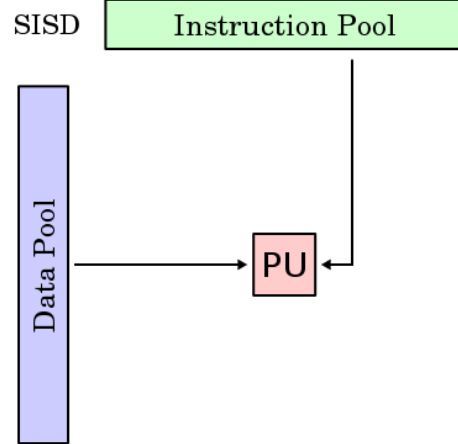
Clusters

- Muchas PCs trabajando en conjunto
- Organización particular
 - Sistema de archivos
 - Memoria
 - Comunicación
- La arquitectura determina un modelo de paralelismo ideal (o necesario)

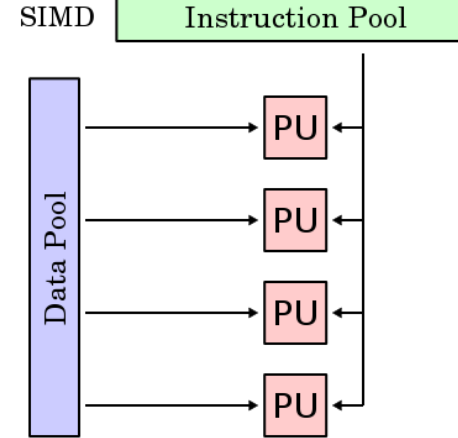
Flynn's Taxonomy

Single Instruction

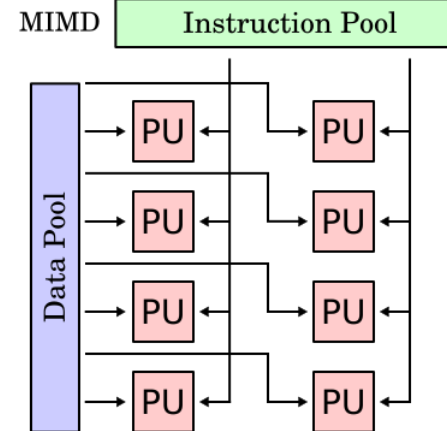
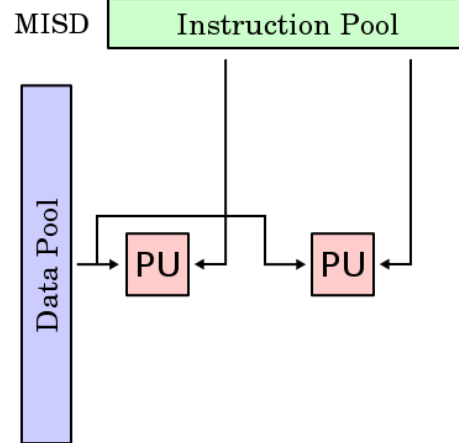
Single Data



Multiple Data



Multiple Instruction

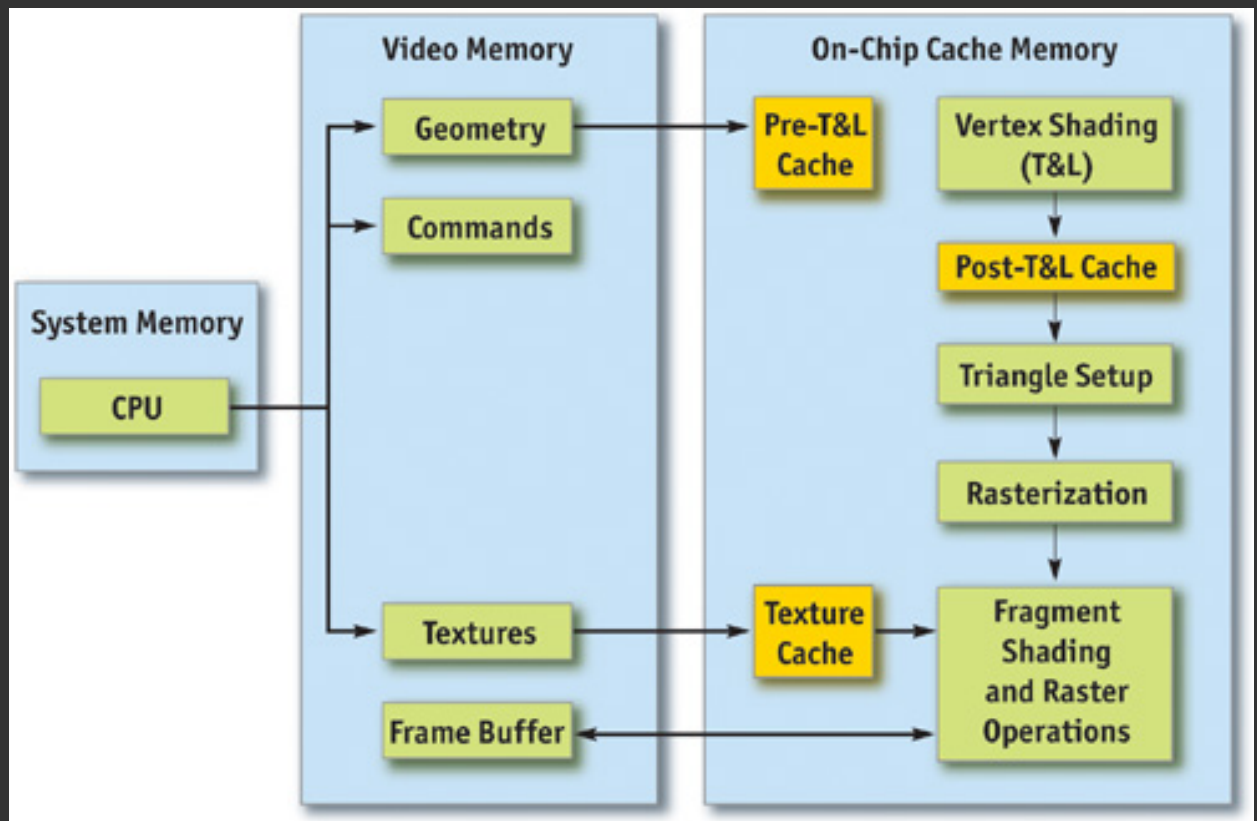


Co-procesadores

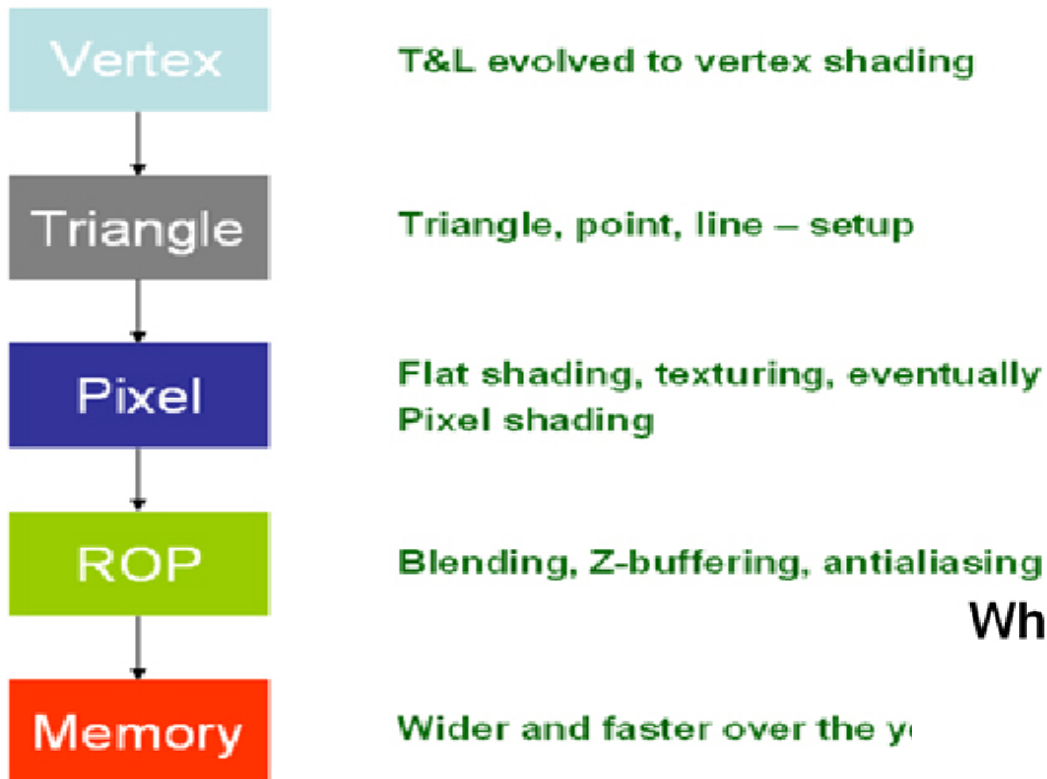
- Ayudar al procesador con una tarea
- Más popular: co-procesador matemático
 - Para XTs, 386s y 486s
- Después vinieron otros
 - Sonido
 - Video

Aceleradoras de Video

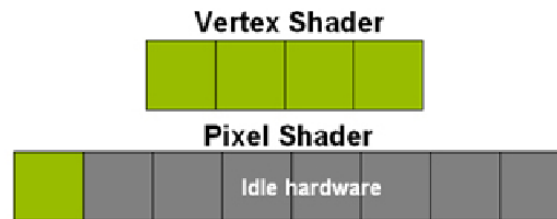
- Co-procesadores vectoriales
- APIs de Video
 - OpenGL
 - DirectX



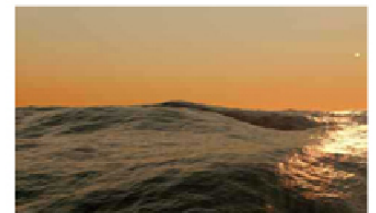
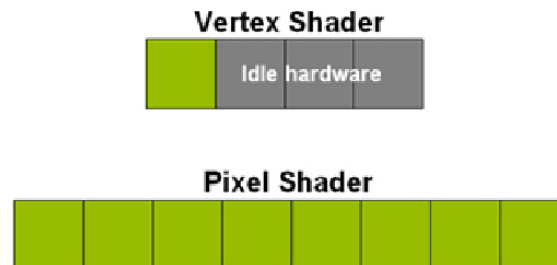
GPUs



Why unify?



Heavy Geometry
Workload Perf = 4

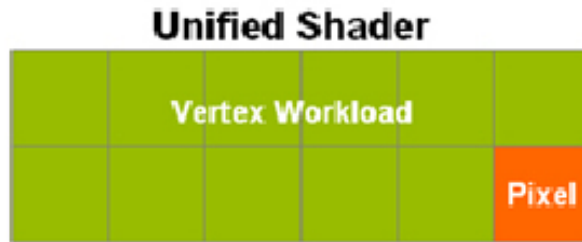
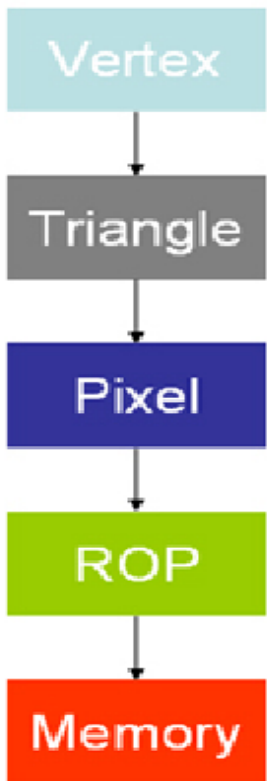


Heavy Pixel
Workload Perf = 8

GPUs

T&L evolved to vertex shading

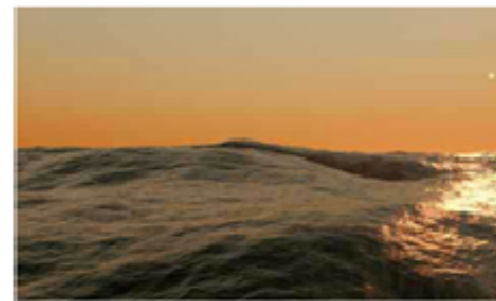
Why unify?



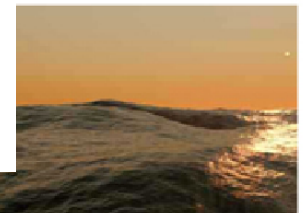
Heavy Geometry
Workload Perf = 12



Heavy Geometry
Workload Perf = 4



Heavy Pixel
Workload Perf = 12

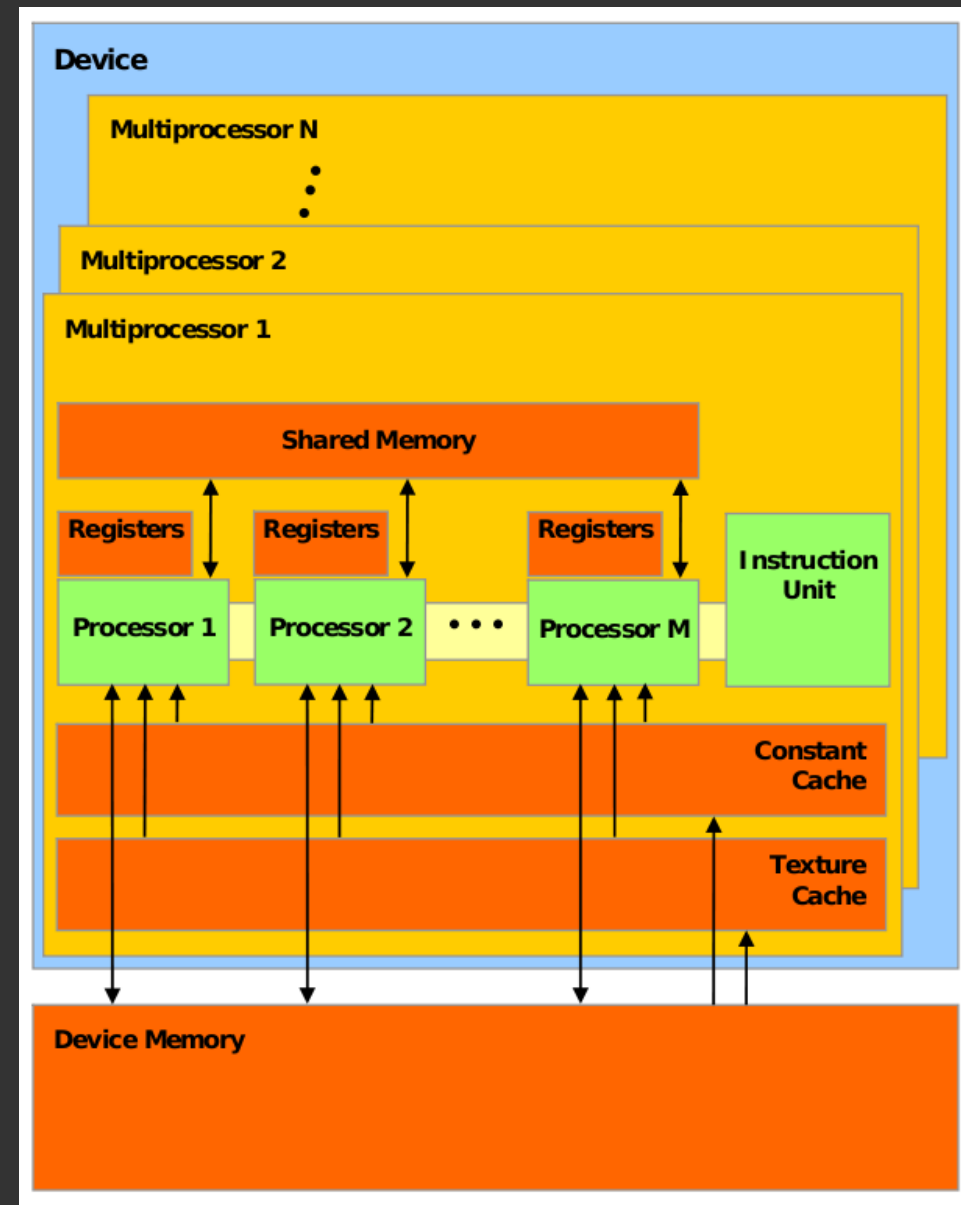


Heavy Pixel
Workload Perf = 8



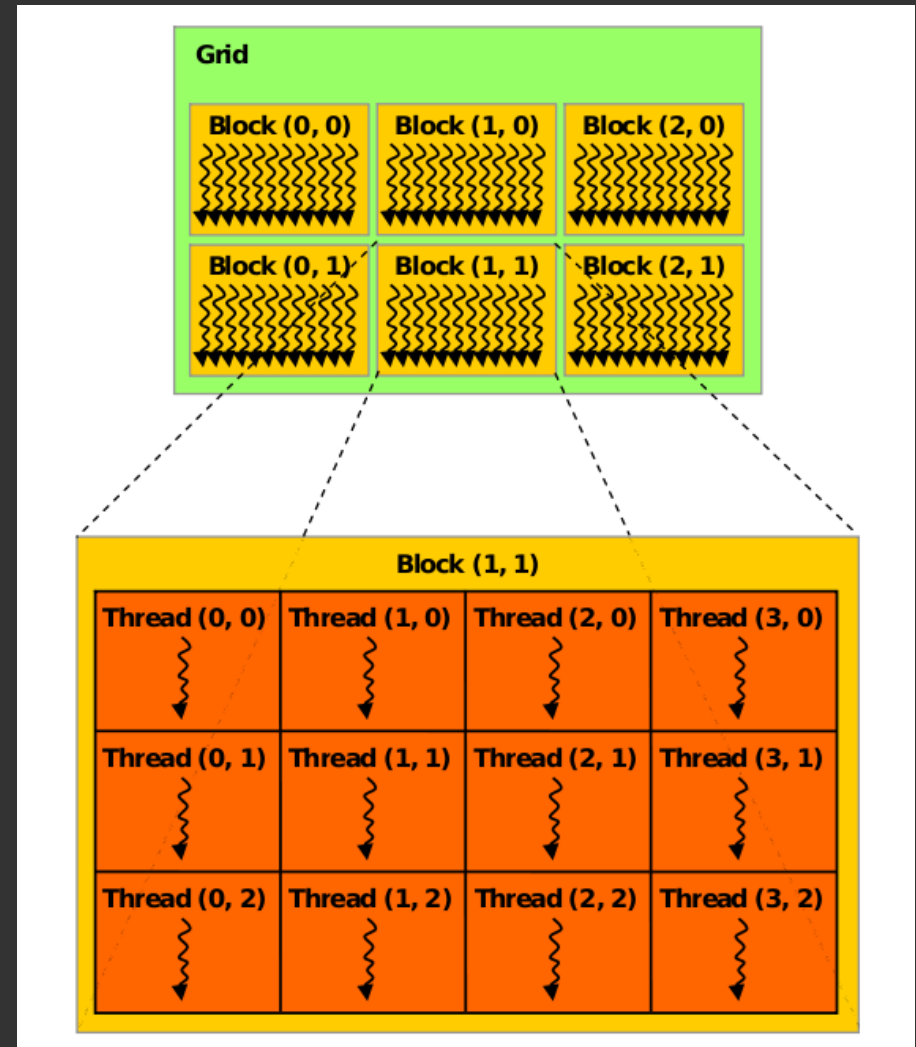
SIMT

- Multi-Procesadores
 - Muchas FPUs
 - Una sola unidad de instrucción
 - Pocas DPUs
- Memoria
 - Memoria general
 - Memoria específica
- Scheduler de 0 overhead (en hardware)



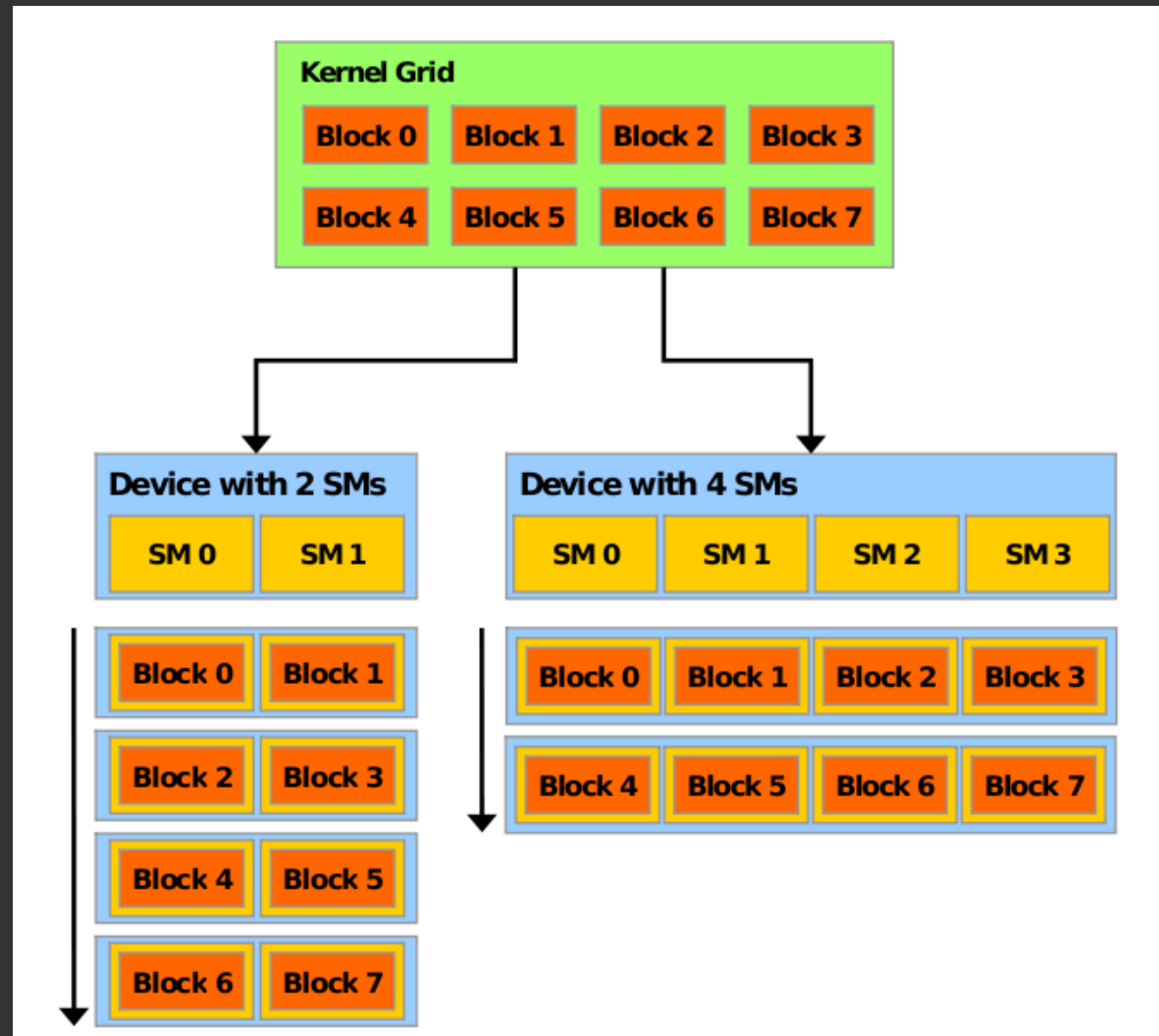
Grillas de Threads

- Organización de threads en bloques
- Organización de bloques en una grilla
- Facilita el mapeo de threads a datos
- Límites máximos impuestos por la arquitectura



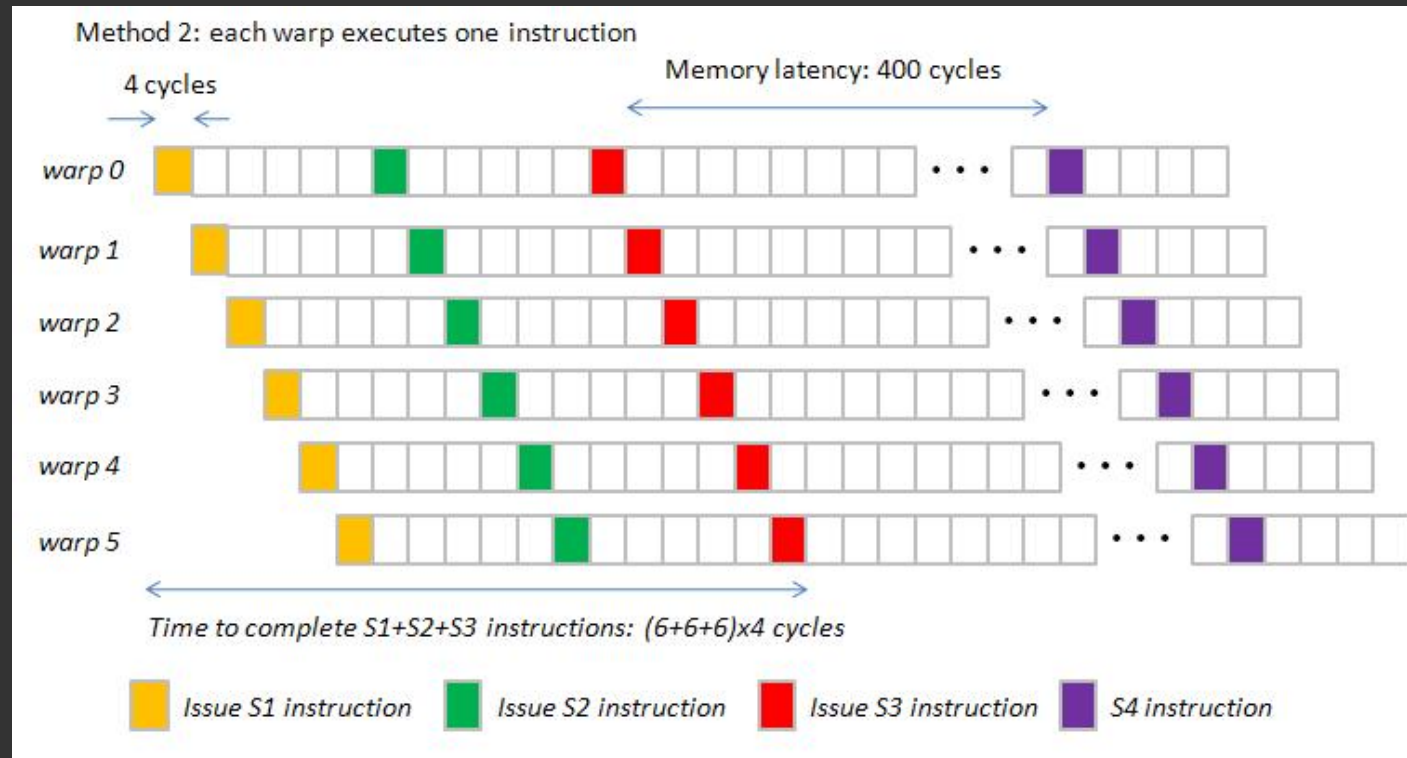
Bloques y MPs

- Bloque: unidad lógica de trabajo
- Se van a asignar tantos bloques como se pueda a un MP
- Thread: unidad mínima de ejecución
- Warp: unidad mínima de sincronización



Occupancy

- El acceso a memoria es caro
- Los threads son baratos
- Esconder latencia y aprovechar poder de procesamiento



Limitaciones

- Primitivas de sincronización
- Manejo de memoria
- No soporta todo C
 - nada de struct y union
 - nada de punteros a funciones (inlining)
 - sólo arrays uni-dimensionales
- Propietario (Hardware y Software)

Soluciones

- GPU como dumb-worker
- Es una tecnología en evolución
 - Mejora del lenguaje
 - Agregado de funcionalidad
 - Mejora en el hardware
- Existen soluciones equiparables
 - OpenCL
 - Ocelot