

GPU

1 GPU are accelerators that supplement a CPU, so they do not need be able to perform all the tasks of a CPU

1 GPU problems sizes: hundreds of megabytes to gigabytes.

1 GPUs rely on hardware multi-threading to hide the latency to memory
CPU's " multilevel caches to overcome the long latency to memory.
미리미리 데이터를 가져오는 동안 그걸 다 안걸었음은 수백개의 threads를 돌림.

1 GPU memory bandwidth > latency.

1 GPU는 각 프로세서가 수백의 thread를 동시하게 돌릴 수 있는 구조로 되어 있어 프로세서 수 자체가 많다.

GPU hardware has two levels of hardware schedulers

1. The thread Block Scheduler that assigns blocks of threads to multi-threaded SIMD processors, and
2. the SIMD Thread Scheduler within a SIMD processor, which schedules when SIMD threads should run.

Local memory! Shared by SIMD Lanes within multithreaded SIMD processor. but this memory is not shared between multithreading SIMD processors.

GPU memory! off chip DRAM shared by the whole GPU and all thread blocks



성균관대학교
학생인재개발원

1. SIMD application은 대량 데이터를 처리 (large data caches X)

Smaller streaming caches and rely on extensive.

multi-threading of threads of SIMD instructions. to hide the long latency to DRAM. 0

• working set이 hundreds of megabytes 인 경우
대용량 last level cache에 안맞는다.

• Given the use of hardware multi-threading to hide DRAM latency. the chip area used for cache in system processors is spent instead on computing resources and on the large number of registers to hold the state of the many threads of SIMD instructions.