

컴퓨터구조03

2018년 3월 19일 월요일 오후 12:35



컴퓨터구조

03

GPU

- GPU are accelerators that supplement a CPU, so they don't need to be able to perform all the tasks of a CPU
 - GPU problem sizes: hundreds of megabytes to gigabytes.
 - GPUs rely on hardware multi-threading to hide the latency of CPUs // multilevel caches to overcome the long latency to memory. 디자인은 각 코어가 주어진 thread를 구동하기 가능하고 또한 프로그램을 지원합니다.
 - GPU memory bandwidth > latency.
 - GPU는 각 프로세서가 주어진 thread를 구동하기 가능하고 또한 프로그램을 지원합니다.

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.

o WP

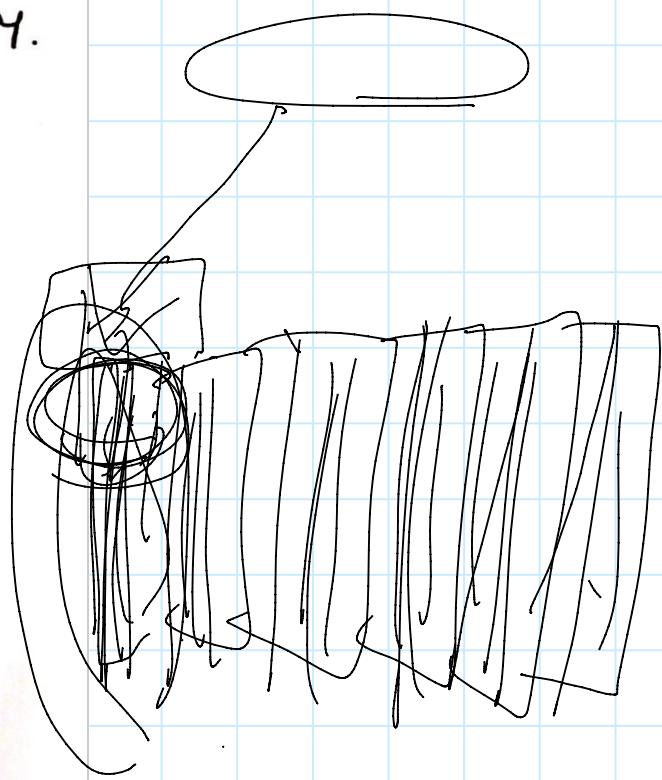
from memory

memory.

ju.

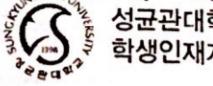
415

ds



Local memory: Shared by SIMD Lanes within multithreaded SIMD processor. But this memory is not shared between SIMD processors.

GPU memory: off-chip DRAM shared by the whole GPU and all threads.



Scanned by CamScanner

1. Tim application을 위한 대처策 (large ~~do~~ caches) X
smaller streaming caches and rely on extensive multi-threading of threads of SIMD instructions. to hide the long latency to DRAM, O

Working set of hundreds of megabytes 일정수준의 메모리 last level cache에 포함된다.

- Given the use of hardware multithreading to hide DRAM latency. the chip area used for cache in symmetric 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.

multithread

read blocks
학교
개발원

Scanner



Scanned by CamScanner

Feature	Multicore with SIMD	Multimedia extensions
SIMD processors	4 to 8	4 to 16
SIMD lanes/processor	2 to 4	4 to 16
Multithreading hardware support for SIMD threads	2 to 4	4 to 16
Largest cache size	8 MiB	16 MiB
Size of memory address	64-bit	64-bit
Size of main memory	8 GiB to 256 GiB	4 GiB to 128 GiB
Memory protection at level of page	Yes	Yes
Demand paging	Yes	Yes
Cache coherent	Yes	Yes

FIGURE 6.11 Similarities and differences between multicore with Multimedia extensions and recent GPUs.



Scanner

GPU
8 to 16
8 to 16
16 to 32
0.75 MiB
64-bit
GiB to 6 GiB
Yes
No
No

SIMD

Type	More descriptive name	Closest old term outside of GPUs	Official CUDA/NVIDIA GPU term	Book definition
Program abstractions	Vectorizable Loop	Vectorizable Loop	Grid	A vectorizable loop, executed on the GPU, made up of one or more Thread Blocks (bodies of vectorized loop) that can execute in parallel.
	Body of Vectorized Loop	Body of a (Strip-Mined) Vectorized Loop	Thread Block	A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more SIMD instructions. They can communicate via Local Memory.
	Sequence of SIMD Lane Operations	One iteration of a Scalar Loop	CUDA Thread	A vertical cut of a thread of SIMD instructions corresponding to one element executed by a SIMD Lane. Result is stored depending on predicate register.
Machine object	A Thread of SIMD Instructions	Thread of Vector Instructions	Warp	A traditional thread, but it contains just SIMD instructions that are executed on a multithreaded SIMD Processor. Results stored depending on per-element mask.
	SIMD Instruction	Vector Instruction	PTX Instruction	A single SIMD instruction executed across all Lanes.
Processing hardware	Multithreaded SIMD Processor	(Multithreaded) Vector Processor	Streaming Multiprocessor	A multithreaded SIMD Processor executes multiple threads of SIMD instructions, independent of other SIMD Processors.
	Thread Block Scheduler	Scalar Processor	Giga Thread Engine	Assigns multiple Thread Blocks (bodies of vectorized loop) to multithreaded SIMD Processors.
	SIMD Thread Scheduler	Thread scheduler in a Multithreaded CPU	Warp Scheduler	Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution.
	SIMD Lane	Vector lane	Thread Processor	A SIMD Lane executes the operations in a sequence of SIMD instructions on a single element. Results stored depending on mask.
Memory hardware	GPU Memory	Main Memory	Global Memory	DRAM memory accessible by all multithreaded SIMD Processors in a GPU.
	Local Memory	Local Memory	Shared Memory	Fast local SRAM for one multithreaded SIMD Processor, unavailable to other SIMD Processors.
	SIMD Lane Registers	Vector Lane Registers	Thread Processor Registers	Registers in a single SIMD Lane allocated to a full thread block (body of vectorized loop).

FIGURE 6.12 Quick guide to GPU terms. We use the first column for hardware terms. For convenience, we cluster these 12 terms. From top to bottom: Program Abstractions, Machine Objects, Processing Hardware, and Memory Hardware.

made
of
el.

ded
threads
ate via

ons
y one
mask

MD
rethreaded
g on a

SIMD

of

hreads
to
MD

thread
Results

aded

MD
essors.

l across
o).

ur groups
Hardware,