

## Fermi architecture.

### • The G80 Architecture.

- ① the first GPU to support C
- ② the first GPU to replace the separate vertex and pixel pipelines with a single unified processor that executed vertex, geometry, pixel, and computing programs.
- ③ the first GPU to utilize a scalar thread processor eliminating the need for programmers to manually manage vector registers.

Scalar thread processor



Vector processor

- ④ Introduced the single instruction multiple thread (SIMT) execution model.

- ⑤ introduced shared memory and barrier synchronization for inter-thread communication.

## • Fermi

- ① Improve double precision performance.
- ② ECC support (memory error control)
- ③ True cache hierarchy (parallel algorithms)  
Fermi GPU의 shared memory는 48KB로  
768x96x128로 구성됨
- ④ More shared memory (SM shared memory  
중 16KB는 application data)
- ⑤ Faster context switching, Atomic operations.

## • Third Generation Streaming Multiprocessor

- 32 CUDA cores per SM
- Dual Warp scheduler  
두 개의 warps를 동시에 파이프라인을  
하루 dispatch 함

- 64KB of RAM은 shared memory  
at L1 cache에 저장됨

(! shared memory + L1 cache = 64KB)

## • Second Generation Parallel Thread Execution ISA

- Memory access instructions to support transition to 64 bit addressing
- Improved performance through Prediction.

## • Improved Memory Subsystem

- NVIDIA Parallel. data cache hierarchy with configurable L1 and Unified L2 caches
- Greatly improved atomic memory operation performance

(read → no diff → write)

## • Nvidia Eight Thread Engine

- faster context switching
- concurrent kernel execution
- Out of order thread block execution

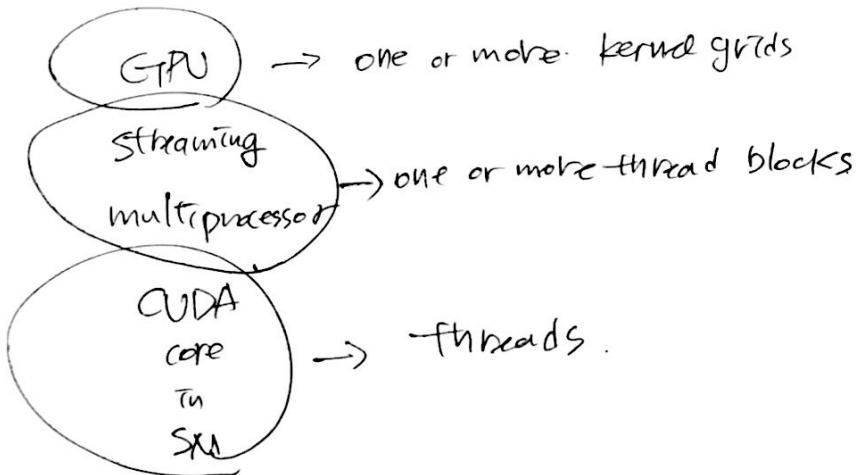


# Hardware Execution.

CUDA's hierarchy of threads



hierarchy of processors on the GPU.



SM : 32 threads == Warp.

## An Overview of The Fermi Architecture

$$512 \times 2^9 \text{ CUDA cores} \\ = 16 \times 2^4 \text{ SM} \times 32 \times 2^5 \text{ CUDA cores per SM}$$

- 64 bit memory partitions
- 384 bit memory interface.  
 $\Rightarrow$  6GB of GDDR5 DRAM Max.
- Giga Thread global scheduler distributes thread blocks to SM thread schedulers.

## 512 High Performance. CUDA cores

- 1 SM 32 CUDA cores
- CUDA core ALU + FPU pipelined
- IEEE 754-2008 floating point standard
- fused multiply-add (FMA) instruction.
- FMA improves over a multiply-add (MAD) instruction for both single and double precision arithmetic  $\rightarrow$  precision 30%
- 32 bit ALU + 64 bit instruction  $\rightarrow$  64 bit
- 64 bit extended precision operation. 240

## 16 Load/Store Units.

- source and destination addresses to be calculated for sixteen threads per clock, 160

- cache or DRAM

## 4 Special Function Units (SFUs)

- execute transcendental instructions ex)  $\sin$ ,  $\cos$ ,  $\exp$ ,  $\log$ ,  $\sqrt{x}$ ,  $\text{reciprocal}$ ,  $\text{square root}$ , ...
- one instruction per thread per clock  
warp 64 clocks 0.5!
- The SFU pipeline is decoupled from the dispatch unit, allowing the dispatch unit to issue to other execution units while the SFU is occupied



## Dual Warp Scheduler

- 1 SM (32 parallel threads) = 1 warp
- 1 SM 2 Warp scheduler  
2 Instruction dispatch units.

→ 두개의 warps가 동시에 dispatch

### 1. dual warp scheduler

① warps 2개 dispatch

② 각 warp마다 한 Instruction만

16 cores, 16 load/store units, or 4 SFU를  
만들기 가능

③ warps 간에 dispatch → scheduler가

Instruction stream 안에 dependencies  
check 필요

④ 여러 Instruction은 dual issued

2 integer 2 floating point, mixed  
integer, floating point, load/store, SFU...  
등으로 dispatch 가능

⑤ Double precision instructions do not  
support dual dispatch with any  
other operation.

## 64KB Configurable Shared Memory and L1 Cache

- on chip Shared memory.
- thread block 간에  
reuse of on-chip data ↑  
off-chip traffic ↓

• 64KB on chip memory

= 48KB shared memory  
+

16KB of L1 cache

or

16KB shared memory.

+  
48KB L1 cache.