

Assignment #3

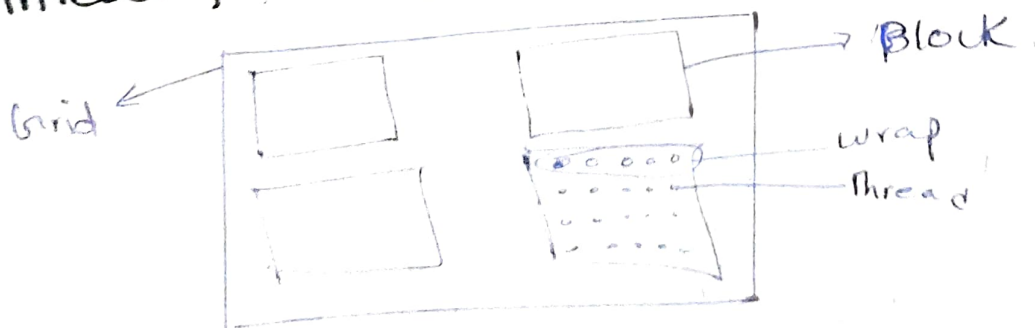
Abdullah Latif
SP23BCS-005

(a)

In GPU computing, these terms define the hierarchy of execution.

1. **Thread**: The smallest unit of execution, running a single sequence of instruction. Each thread has its own Registers and a unique ID to calculate which piece of Data it should work on.
2. **Block**: A group of threads that cooperate and share data through shared memory.
- 3) **Grid**: A collection all thread blocks launched in single kernel execution.
- 4) **Wrap**: A group of 32 consecutive threads in a block executed simultaneously by a streaming multiprocessor.

Thread, block, Grid, Wrap.



Thread Organization in a Block

Thread in a block can be organized in 1D, 2D or 3D, accessed via ThreadIdx

(b)

Memory Hierarchy in GPUs

Registers. Fast, on-chip, private to each thread

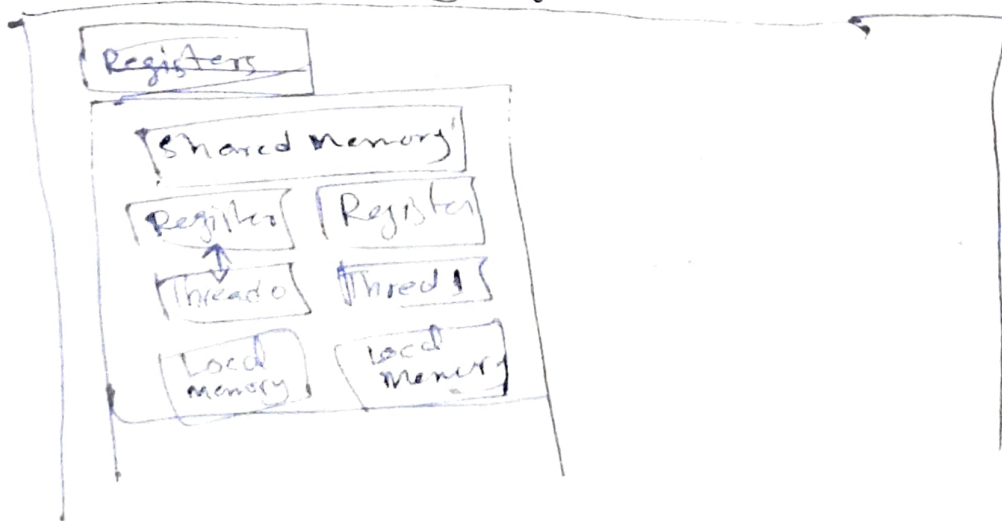
Shared memory: Very fast on-chip, shared by threads in a block.

Global Memory: Largest capacity, slow, accessible by all threads.

Constant / Texture Memory

Read only, cached accessible by all threads.

Access: Thread use registers for local variables, shared memory for block wide communication and global memory for large datasets.



GPU

Block 0

shared Memory

Registers

Registers

Thread 0

Local Mem

Thread 1

Local mem

Block 1

shared Mem

Reg

Reg

T₂

T₄

L Mem

L Mem

Global memory
Constant Memory
Texture memory

CPU

wrap switching and Latency Hiding (C)

Fast wrap switching is enabled by having dedicated hardware and zero-overhead content switching.

Content-switching

GPU maintain multiple wraps on each SM.

- * when one wrap stall.
- → The SM instantly switch to another wrap.

Why Fast.

→ state already stored in H.

No costly O.

Why effective.

- keep ALUs busy
- Hides long mem latency.
- Increase throughput.

(d)

Branch Instruction.

* GPU use SIMT

* If thread wrap take diff branch

* Execution become serialized

* This is called ~~the~~ wrap divergence

Performance Impact.

* Reduce parallel efficiency

* Best to keep wrap thread same path. following.

(e)

No of wraps in 16×16 .

T per block = 256

W size = 32

$$256 / 32 = 8 \text{ wraps}$$

8 wraps.

(f)

Num of Block on a SM with 1536Ts.

T per block = 256

Max T = 1536

$$1536 / 256 = 6 \text{ Blocks.}$$

(9)

Maximi Occupancy

$$8 \times 8 \text{ block} = 8 \times 64 = 512 \quad 33\%$$

$$16 \times 16 \text{ Block} = 6 \times 256 = 1536 \quad 100\%$$

$$24 \times 24 \text{ Block} = 2 \times \del{25} 576 = 1152 \quad 75\%$$

conclusion.

The 16×16 would be most efficient.