

## HPC

### \* Warp, SIMD Hardware (Thread Execution Efficiency)

- 1) Cuda me threads block me hote hai.
- 2) Har block ko 32 thread ke groups me divide kiya jaata hai.
- 3) 32 thread group is called warp.

Warp = Scheduled unit of SM (Streaming Multiprocessor)

- 4) Every warp has 32 threads.
- 5) Not a part of CUDA programming model, hardware implementation detail hai.
- 6) Future warps GPU's me warp size change ho skta hai.
- 7) Warp ke andar threads SIMD model me execute hote hai.

Simple Instruction,  
Multiple data

### \* Multidimensional

#### Thread Blocks

- 1) 2D/3D blocks ko pehle 1D linear order me convert ho kiya jaata hai.
- 2) Order is :-  
 $X \text{ dimension} \rightarrow Y \text{ dimension} \rightarrow Z \text{ dimension}$
- 3) Warps are groups of linearized threads.

### \* Warp as Scheduling Unit

- 1) CPU ka SM (Streaming Multiprocessor) warp level par schedule karta hai.
- 2) SM at a time executes warps not threads solo

## \* Warp Rules and Synchronization

- 1.) Warps ke andar ya warp ke beech se dependence par depend mt karo.
- 2.) Agar threads ke beech dependency ho  $\rightarrow$  sync threads () zaroori hai.
- 3.) " is also compulsory when shared memory is used.

## \* SM = SIMD Processor

- 1.) SM Ke andar :- Have registers + shared memory
  - Control unit common hota hai.
  - Multiple ALU's
- 2.) Benefits :-
  - Instruction fetch and decode ek baar hota hai.
  - Control overhead kam hota hai.

## \* Control Divergence (Performance Impact)

- 1.) When same warp ke threads different execution paths follow krte hai its known as control divergence.
- 2.) Control divergence slows why :-
  - Warp can run 1 instruction at a time
  - GPU puchhe if path execute krega
  - fir else ke saath execute krega.

Threads go us path me nahi h, idle rehte hai  
if (threadId.x < 16)

— / — / —

- 3) Reasons of Divergence → if-else based on threadIdx

  - Loop jisme → Some threads less iteration

`if (threadIdx.x > 2)      if (blockIdx.x > 2)`

- Example of divergence
  - Dift decision inside warp
  - Divergence occurs (bad)
  - Control Divergence
    - Boundary checks are imp.
    - Impact is negligible in large datasets.
  - Example of non divergence
  - Block level decision
    - All threads of the warp follows same path.

## \* DRAM Bandwidth (Memory Access Performance)

- 1.) Importance of Memory Bandwidth

  - GPU → Computes very fast  
    → But if data is slow GPU will be idle

That's why memory BW = first order performance factor.

## 2) DRAM Internals

- Transistor cell = capacitor
  - To Capacitor me charge detect karne slow hota hai
  - Core array speed < interface speed

Eg: DDR3 → core speed = 1/8 interface speed.

## \* DRAM Bursting Concept

- 1) DRAM stands for Dynamic random access memory.
- 2) DRAM ek address dekh kota hai.
- 3) Poche internal buffer me load hota hai.
- 4) Data burst mode me transfer hota hai.

Sequential access = full benefit

Random access = data waste

## \* DRAM Banks

- 1) Multiple banks = parallel row access
- 2) Bank switching se latency hide hote hai.

## \* GPU Memory Channels

- 1) Eg :-
  - RTX 6000,
  - Required bw = 672 GB/s
  - Single channel  $\approx$  56 GB/s
  - That's why 12 channels are used

## \* Memory Coalescing

- Def :- 1) Jab warp ke saare threads :-

- Same instruction
- Adjacent memory locations access karao.

Coalesced Access GPU sirf ek DRAM request karta hai.

Pattern  $A[\text{base} + \text{threadIdx.x}] \rightarrow$  Fast, BW efficient

## 2.) Uncoalesced Access

- Threads scattered addresses access karte hain
- Multiple DRAM requests
- Slow execution

## 3.) Matrix Multiplication Case

- B Matrix access  $\rightarrow$  coalesced
- A Matrix  $\rightarrow$  uncoalesced

Solution  $\rightarrow$  Tiling, Shared Memory, Coherency

## \* Histogram

### 1.) Data to bins me count karna is histogram

- Example  $\rightarrow$  Alphabet frequency
- $\rightarrow$  Img Intensity count

### 2.) Parallel histogram challenge

- Multiple threads
- Same bin update
- Memory Access Inefficient

### 3.) Sectioned Partitioning (Bad)

- Each thread continuous chunk
- Adjacent threads  $\rightarrow$  far memory
- No coalescing

### 4.) Interleaved Partitioning (Good)

- Threads alternate data access
- Coalesced memory, Better bandwidth (BW)

A bin is a memory location that stores the count of data elements falling into a specific range or category.

## \* Data Races (DR)

- 1) DR are :- When ~~multiple threads~~ Read-Modify-Write  
  - Multiple threads
  - Same memory location
  - Read-Modify-Write
  - without synchronization

→ Output is unpredictable.

$n = n + 1$ ,  
2 → threads

result kabhi 1

kabhi 2

- 2) Real life eg :- Bank cash counting  
 Ticket booking, customer numbering system

## \* Atomic Operations (AO)

- 1) AO are :-  
  - Single hardware instruction
  - Read + Modify + Write together
  - No interference allowed
- 2) CUDA Atomic functions  
  - atomicAdd, atomicMin / Max
  - atomicSub, atomicCAS
- 3) Histogram with Atomics  
  - Threads input ko stride pattern me process karte hai.
  - Histogram update atomicAdd se hota hai
  - Data race eliminate ho jata hai.

## \* Kernel and SPMD parallelism

- 1) Kernel = GPU par run hove wala fxn.  
CPU se launch hota h, but parallel threads GPU par execute krta hai.
- 2) SPMD = Single Program Multiple Data, means same code (kernel) but different data for every thread

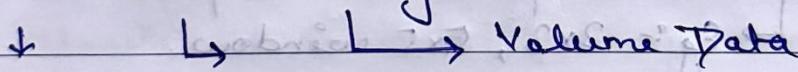
Eg :- Img ke hove pixel ke liye same kernel  
Hare thread ko ek pixel process krta hai

Imp: Hare threads ka unique ID hota hai  
Same kernel sab threads ke liye sun hota hai  
GPU massive parallelism achieve krta hai

- CPU follows SPMD model where the same kernel program is executed by multiple threads on diff. data elements.

## \* Multidimensional Kernels

- 1) 1D vs 2D vs 3D grids.



Vector      Image /

Matrix width -  $n = \text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x}$   
 $y = \text{threadIdx.y} + \text{blockIdx.y} * \text{blockDim.y}$       Dim.x

MD Kernels are used to naturally map threads to 2D or 3D like images and matrices

I use hare thread apna pixel (x,y) find krta hai

- 2) Why Imp :- Img are naturally 2D

Matrix multiplication also uses 2D mapping

## \* Color to Grayscale (Img Processing)

1) RGB img  $\rightarrow$  Hm pixel = (R, G, B)  
 $\rightarrow$  Value & range = 0-255

2) Grayscale img  $\rightarrow$  Surf intensity value  
Color info. removed ho gati hai

### Formula

$$\text{Gray} = 0.21R + 0.71G + 0.07B \quad (\text{Green is max})$$

bcz human eye

### CUDA Approach

(is more sensitive to green)

$\hookrightarrow$  Hm thread = 1 pixel

RGB thread  $\rightarrow$  grayscale write

RGB to grayscale conversion is done using a weighted sum to preserve perceived brightness

## \* Blur Kernel (Image Blurring)

1) Pixel ke neighbours ke avg se replace kernel is image blurring, img smooth ho gati hai

2) Blur box  $\rightarrow$  Center pixel ke ase pass ka square  
 $\hookrightarrow$  Eg:- 3x3, 5x5 window.

3) CUDA working  $\rightarrow$  Hm thread = 1 pixel

$\rightarrow$  Neighbour pixels read

$\rightarrow$  Avg calculate, Output img me store

if (row < height, & 8 col < width)

Image blurring uses a 2D kernel where each thread computes the avg. of neighbouring pixels

## \* Thread Scheduling and Transparent Scalability

### 1) Transparent Scalability

- Programmer ko GPU size ka tension nahi
- Block can be executed in any order
- Hardware automatically schedule karta hai

### 2) SM (Streaming Multiprocessor)

- Thread block granularity me SM par assign hote h,EK SM me multiple blocks ho skte hai

### 3) Warps (32 threads)

- Scheduling unit
- Warp ke sab threads same instruction execute karte h (SIMD)

### 4) zero Overhead Scheduling

- Jabi 1 warp wait karta hai (memory), dura warp execute ho jata hai, latency hide ho jati hai.

⇒ CUDA provides transparent scalability by allowing blocks to execute in any order across SM's.

## \* CUDA Memories

### Memory

- Shared
- Register
- Global
- Constant

### Scope

- Block
- Thread

### Grid

- Slot
- Cached

### Speed

- Very fast
- fastest

1) Registers → per thread, automatic variables, fast  
Shared Memory but limited

↳ per block, global se kaafi fast,  
stored b/w threads.

Global Memory in between shared memory

↳ sabse slow, high latency, sab  
threads access kar skte hai

Efficient use of shared memory is critical  
to improve CUDA Kernel performance.

### \* Tiled Parallel Algorithms

1) Problem → Global memory is slow  
↳ Many operation me global access → low performance

2) Tiling Idea → Data ko small tiles me break karo  
Tiles ko shared memory me load karo.  
Baar baar use karo.

Real life analogy → Carpooling = tiling  
↳ Zeyada log kam gaadi = kam traffic

#### 3) Steps of tiling

- Tile identity kro
- Global → shared load
- - sync threads ()
- Computation
- Next tile

⇒ Tiling reduces global memory access by  
reusing data from fast on chip shared memory

## \* Tiled Matrix Multiplication (TMM)

### 1) Basic MM

- One thread = 1 output element
- Global memory access is very much

### 2) Tiled approach

- Matrix  $M \times N$  block width  $\times$  block width tiles me
- $M$  and  $N$  tile shared memory divide me load
- Phase by phase computation

### 3) Phases

- Phase 0 :- Tile 0 load + compute
- Phase 1 :- Tile 1 load + compute
- ... repeat

Tiled matrix multiplication improves arithmetic intensity by inc. computation per memory access

## \* TMM Kernel

### 1) Shared Memory Usage

- shared - float ds - M (TILE) (TILE);
- shared - float ds - N (TILE) (TILE);

### 2) Kernel flow :- Thread index calculate

- Tile loaded into shared memory
- syncthreads ()
- Partial sum calculation
- Next phase
- Final result write

size  
Proper tile selection  
balances shared  
memory usage and  
thread occupancy.

## \* Course Info and Overview

1) Aim of course :-

- Heterogeneous parallel computing systems ke program karna :-

→ High performance

Energy efficiency

Maintainability

Scalability (future hardware ke liye)

Portability (diff vendors ke devices)

EK hi program GPU's/ CPU's par bhi acha

chle bina

dobara likhe

## \* Heterogeneous Parallel Computing

Q:-  
Kya  
hotah.

Job system me diff types ke processors milte kaam kerte hain :-

- CPU, GPU, DSP, Specialised hardware
- ↓  
↳ Global Graphics

Central Processing Unit

DSP - Digital Signal Processor

### 2) CPU vs GPU

CPU

- Latency Oriented
- Powerful Cores
- Large Cache
- Fast response time
- Best for sequential code

Eg:- OS tasks  
control logic

GPU

- Throughput Oriented
- Bolt zyada cores
- Small cache
- Slow
- Thousands of threads

Eg:- Img processing  
Matrix Multiplication

- 3) Why use both :- CPU  $\rightarrow$  Sequential part  
 GPU  $\rightarrow$  Parallel part

Sequential code : CPU 10x faster

Parallel " " : GPU "

#### 4) Applications of Heterogeneous Computing

- Scientific simulation
- Medical Imaging
- Video processing
- Machine Learning
- Financial Analysis

#### \* Portability and Scalability

- 1) Software cost is imp :- Software lines faster grow

- Har naya hardware      Koi ekai h than hardware
- ane par code rewrite is costly.

Isliye :-  
 ↗ Scalability

↗ Portability

- 2) Scalability :- Same app efficiently run kare

- New CPU generations
- Zyada cores par
- " threads par

- 3) Achieving Scalability  $\rightarrow$  Fine grained parallelism

H/w details pe depend  $\leftrightarrow$  Dynamic thread scheduling

na karna

- 4) Portability :- Same program efficiently run kare

- CPU, GPU par
- X86, ARM par
- SIMD, threading par
- Shared and distributed memory systems par

## \* CUDA C vs Thrust vs Libraries

- 1.) Libraries → Easy, fast  
Compiler Directives → Portable, simple  
Programming Lang. (CUDA C) → Max performance
- 2.) Libraries → Ready made GPU made code  
Minimal changes required  
Eg: → Thrust, CUBLAS  
Tbb pattern already available ho → libraries best
- 3.) Thrust → STL kinda interface  
Easy vector operations, less code, less control
- 4.) Compiler Directives (OpenACC)
  - #pragma based
  - Compiler parallelism handle karta hai
  - Performance compiler pe depend karta hai
- 5.) CUDA C → Most powerful, full control, max performance

CUDA C offers max. flexibility and performance at the cost of programming complexity.

## \* Memory Allocation and Data transfer

- 1.) CUDA execution model
  - Host (CPU) → serial code
  - Device (GPU) → parallel Kernel

## 2) Explicit Memory Management Flow

- Host memory allocate
- Device " "
- Data copy Host  $\rightarrow$  Device
- Kernel launch
- Result copy Device  $\rightarrow$  Host
- Device memory free

Explicit memory management requires manual allocation and data transfer b/w host and device.

## \* Threads and Kernel functions

### 1) CUDA thread model

- Kernel  $\rightarrow$  grid of threads
- Sab threads same code run krte hai (SPMD)

### 2) Threads as Virtual Processor

- Har thread ke pass :-

→ Registers, program counter, control logic

Thread = Virtual Von-Neumann processor

### 3) Thread Indexing

$$i = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x};$$

- Har thread apna data element calculate karta h.

### 4) Thread Blocks

- Threads ko blocks me divide kya jaata hai
- Block ke andar : $\rightarrow$  Shared memory, synchronization
- Blocks ke beech no communication possible

⇒ Thread blocks enable scalable cooperation among threads using shared memory and synchronization.

## \* Unified Memory (UM)

- 1) UM is a single memory space.
  - CPU aur GPU deno access kar skte hai
  - Data migration automatically hota hai
- 2) Benefits → Code simple, less bugs, easy learning
- 3) Advanced UM
  - Compute capability 6.x ke baad :-
  - Page faulting, on demand migration, memory oversubscription possible.
  - ⇒ Unified memory simplifies programming by providing a single address space managed by hardware and software.

## \* FORMULAS

- 1) Warp size = 32 threads / warp
- 2) Warps per block =  $\frac{\text{Threads per block}}{\text{Warp size}}$
- 3) Blocks =  $\frac{\text{Total elements}}{\text{Threads per block}}$
- 4) Threads in grid = Block  $\times$  Threads per block
- 5) Idle threads = Threads in grid - Useful threads
- 6) Block threads =  $\frac{\text{Max threads per SM}}{\text{Threads per block}}$

— / —

$$7) \frac{\text{Blocks registers}}{\text{Registers per thread}} = \frac{\text{Total registers per SM}}{\text{Registers per thread} \times \text{Thread per block}}$$

$$8) \frac{\text{Blocks shared memory}}{\text{Shared memory per block}} = \frac{\text{Shared memory per SM}}{\text{" " " block}}$$

$$9) \text{Blocks per SM} = \min(6, 7, 8, \text{HW block limit})$$

$$10) \text{Active warps} = \frac{\text{Blocks per SM}}{\text{Threads per block}} \times \frac{\text{Threads per block}}{\text{Warp size}}$$

$$11) \text{Active threads} = \frac{\text{Active warps}}{\text{per SM}} \times \frac{\text{Warp size}}{\text{per SM}}$$

$$12) \text{Max warps per SM} = \frac{\text{Max threads per SM}}{\text{Warp size}}$$

$$13) \text{Occupancy} = \frac{\text{Active threads per SM}}{\text{Max " " "}}$$

threads ki jagh warps . new formula.

$$14) \text{Unutilized warps} = \frac{\text{Max warps per SM}}{\text{Threads}} - \frac{\text{Active warps per SM}}{\text{Threads}}$$

$$\text{Threads} = \frac{\text{Max threads per SM}}{\text{Threads}}$$

$$15) \text{Required warps} = 1 + \left( \frac{L}{M} \right) \quad L = \text{Memory latency}$$

M = useful compute cycles b/w memory access