# ✳ Intro to GPU-

- age of parallel computing → multi tasking phones, tablets, laptops, etc.
- Provide novel, fast & rich - experience to user, developer need to understand & use various parallel plat form.
- CPU performance increasing no. of cores, transistor & other features &
  always has been a relatable way of enhancing performance but ~only way~
  reliable

- **GPU's-**
  - specialized electronic circuit. Helps alter memory rapidly to acc img,
    given as output to display device.
  - GPUs are in - Mobile phones [Adreno/ Mali], PCs (AMD/ NVIDIA)
    consoles [ PS4/5, XBOX]
  - GPUs are great at img processing & controlling computer
    graphics
  - They are highly parallel.

- **CUDA** — • Compute Unified Device Architecture
  - launched by NVIDIA in 2007. as an programming
    ~language~ interface providing parallel computation using GPUs.
  - 

- **CUDA architecture -**
  - Before CUDA archi. vertex & pixel shaders were used for parallel computing
  - pixel shader is a GPU component that can operate per pixel.
  - vertex shader is also - a GPU component like pixel & it's
    assembly language specific, used for geometrical operations.
  - • an extension of C++ features to support GPU acceration.
  - components → • parallel compute engines inside NVIDIA GPUs
    - OS kernel level support for hardware.
    - • user mode driver, • device level API for developers
    - • PTX instruction set architecture for parallel computing

• Bioinformatics.
• Fast Video Transcoding

• Appl<sup>n</sup> of CUDA ⟶ •Deep Learning •Computing • Computer Vision.
• CUDA → general purpose computing • [ • Data Science • Ultra sound imaging
• 3Device Level APIs → i) OpenCL 2) Direct X 3) CUDA Driver API.
• Language Level — 3) i) Fortan ii) C++ iii) C.

\* Processing flow of Cuda-C.
    1) Host code Exceeution.
    2) Kernel Launch
    3) Kernel Exceeution on GPU.
    4) Data Access & Processing
    5) Synchronization (optinal).
    6) Results Transfer (eptional).
    7) CleanUp—

\* CUDA Kernel program —    \* CUDA add<sup>n</sup> of two vectors

```
#include <iostream>                              cuda

__global__ void kernel (void) {                  _global_ void vectorAdd
}                                                (float *a, float b, float *c, intn)
int main () {                                    { int i = blockIdx + blockDim.x
    kernel <<< 1, 1 >>> ();                           + threadIdx.x;
    printf ( "Hello World! \n");                      if (i<n)
    return 0;                                          c[i] = a[i] + b[i];
}                                                }
```

• global keyword → helps call kernel function using host,
• global helps run the kernel fun<sup>n</sup> on GPU.
• use NVCC compiler to complet this code instead of gcc.
• kernel <<< 1, 1 >>> () is a CUDA specific syntax specifies call to
• specifies a call to device code.    device code.
• parameters inside "<<< >>>" are called execution configuration.

**\* Device-**
- Refers to GPU (Graphics Processing Unit) in CUDA programming - executes device code & performs parallel computations on data.
- Managed by CUDA runtime system for task scheduling & memory mgm.

**\* Host-**
- represents Central Processing Unit) in CUDA programming.
- Executes host code responsible for managing device operations, memory
- Data transfer & control execution flow of CUDA programs & interact with GPU.

**\* Device Code -**
- code written specifically to be executed on GPU in CUDA program.
- Utilizes CUDA libraries & functions for parallel processing tasks-
- can be invoked parallel processing tasks from host code to perform computations on GPU cores.

**\* Kernel -**
- function in CUDA programming designed to be executed in GPU.
- functions are identified by name & configuration, no of threads pa block & blocks per grid.
- Each thread executes kernel code independently processing. diff portions of data simultaneously to active parallelism.

**\* warps** — set of 32 concurrent threads in a block.

**\* Global Memory -**
- Largest memory space available to CUDA programs.
- Resides on device & is typically used to store data needs to multiple threads.

**\* Shared Memory -**
- Shared memory is fast, on-chip memory space shared by all threads within single block.
- Much faster than global memory but limited capacity.
- Shared memory is used for data that needs to be accessed frequently & efficiently within a block.

**\* Constant Memory -**
- Is also located on device & is read only for all threads.
- It is cached & provides fast access to data accessed uniformly by all threads within a block.

**\* Thread Hierarchy**
1) Grids → • grid is highest level of organizations in CUDA & multiple blocks.
   - represents overall computational workload & is executed on GPU device.
2) Blocks → • subdivided into blocks, each containing threads.
   - executed independently & scheduled on any multiprocessor on GPU.
3) Threads → • smallest unit of execution in CUDA & organized into blocks.
   - threads within a block can cooperate & synchronize with each other using shared memory & barriers.

- Block dimensions - • refers to no. of no. of threads per block.
  - specified as a 3-dimensional array of integers.
  - each block can contain up to a max. no. of threads (depend on GPU).
  - Dimension is crucial for organizing threads.

**\* Grid Dimension -**

- Refers to no. of blocks in grid.
- specified as a 3d array of integers.
- grid dimensions determines overall size of computation &
  how blocks are organized for execution on GPU.
- Multiple blocks withing grid can execute concurrently on
  diff streaming multiprocessors (SMs) on GPU.

CUDA kernel for adding two vectors element -wise-

```
_global __void vectorAddition (float a, float b, float
result, int size) {
int index = blockIdx.x × blockDim.x + threadIdx.x;
    if (index < size) {
        result [index] = a [index] + b [index];
    }
}
```
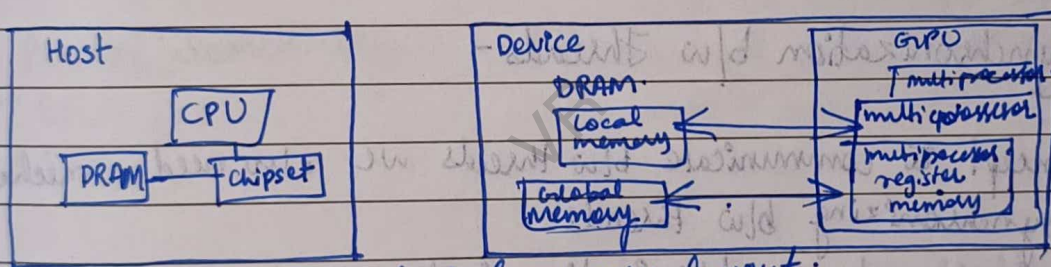
**\* Kernel ~~Execution~~ in CUDA -**

- CUDA → kernel is a function that is executed on GPU.
- written using CUDA c/c++ syntax & is responsible for performing
  parallel platforms computations.
- kernes are launched from CPU & executed by multiple
  threads on GPU.
- executes same kernel code with diff data, enabling massive
  parallelism.
- invoke CUDA kernel from cpu.
  - kernel is launched, CUDA runtime system, manages allocation
    of resources on GPU & schedules execution of threads.
  - CUDA kernel, while GPU executes kernel in parallel.

✳ CUDA kernel & Handling errors :-
- CUDA as extension of C consists of host code (program control) & device code (GPU) ·combined in a single C Program·
- CUDA is C for parallel processor·
- you can write a program for one thread & instantiate it on many parallel threads·

✳ CUDA memory model / GPU memory·
① CPU & GPU have seprate memory spaces·
② Device pointers point to GPU memory·
③ Host pointer point to GPU memory·
④ Host CPU code Manages device GPU Memory·



physical memory layout·

⑤ does following tasks —
  i) Allocate or free money·
  ii) copy data to & from device·
  iii) Applies to global device memory (DRAM)·
⑥ local memory resldures in device DRAM·
⑦ Host Can read & write global memory but not shared memory)·
⑧ CUDA, threads may access data from multiple (memory spaces during execution·
⑨ each thread has private local memory·
⑩ All threads have acces to same global memory.

**\* Manage communication bet^n threads -**

- possible & many times necessary for threads within same block to communicate with each other.
- Need to share data b/w threads to compute final result or outputs
- communication b/w threads → possible using shared memory.
- block of threads shares memory called shared memory.
- variables in shared memory treated different than typical variable variabled by CUDA C complier.
- creates a copy of variable for each block that you launch on
- threads with same block communicate data with each other using shared memory or global memory.

**\* Synchronization b/w Threads -**

- except to communicate b/w threads we also need a mechanism for synchronizing b/w threads.
- Main use of synchronization is to prevent RAW, WAR, WAW hazards during communication b/w mutiple threads.
  - RAW (Read After Write)
  - WAR (Write After Read)
  - WAW (Write After Write).

synchronize to commit all memory writers, reads &
③ CUDA synchronizes threads using function - syncthreads().
④ there are two types of synchronization in CUDA.
  1) Implicit Synchronization. 2) Explicit Synchronization.
⑤ Barrier a programmer can place the synchronization barrier explicitly to synchronize tasks.