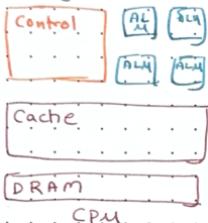


## CUDA Introduction

CPU designed to minimize latency.



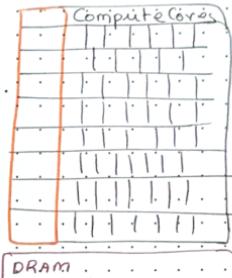
### CPU

- majority of silicon is dedicated to:
  - Advanced Control logic
  - a large cache

→ Cores are called Cuda Cores

→ To Access Compute cores they are accessed by Kernels which are Basic Building Blocks to access the massive Compute cores (CUDA cores)

### GPU



These massive cores can compute basic operation in parallel and faster.

## Kernels

Function that runs on GPU.

## Threads

Kernels execute as a set of parallel threads.



### CUDA

⇒ Heterogeneous Host + device

⇒ C with extensions

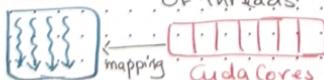
allow heterogeneous programming on Nvidia GPU's

⇒ The transfr. Bl/w host to device & vice versa are costly operations.

## Concept of threads:

- kernels execute as a set of parallel threads.
- CUDA is designed to execute 1000's or even more threads.
- Threads executions are in SIMD / SIMD T.
- single instruction multiple thread
- Threads are similar to data parallel tasks.
- Each thread one a subset of data.
- Threads execute independently.

\* **Blocks** They have 1D or 2D or 3D structure of threads:



\* when the kernel is launched Block gets mapped to Correspondingly set of CUDA cores.

\* Threads do not execute at the same rate  
→ This gives rise to various performance bottlenecks and can be optimized further.

Organization & hierarchy of threads:

\* Thread      \* Blocks      \* Grid

\* Thread lowest level in hierarchy  
Execution of kernel on a single piece of data

\* mapping      Each thread gets mapped to one core in GPU

\* Thread "CUDA Core" → GridDim =  $3 \times 2 = 6$  blocks

\* Blocks are grouped into Grids

\* Each kernel launch creates a single Grid

\* Grids have dimensions 1D or 2D or 3D of blocks

Grid 0

