

GPU Computations

Ilia Nechaev

10.01.2025

Overview

- 1 GPU vs CPU
- 2 Intro to GPU architectures
- 3 Massive Parallelism
- 4 OpenCL
 - Terminology
 - Approach
- 5 Writing Kernels
- 6 Coalesced Memory Access

GPU vs CPU: Main Differences

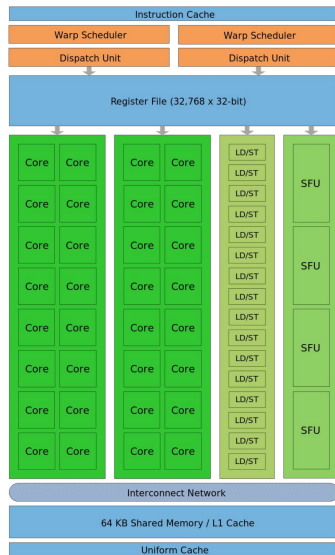
- **CPU (Central Processing Unit):**

- Optimized for low latency and complex tasks
- Few powerful cores (usually 4–16)
- Large cache memory per core
- Good at serial (sequential) processing

- **GPU (Graphics Processing Unit):**

- Optimized for high throughput
- Hundreds or thousands of simpler cores
- Highly parallel architecture
- Single instruction pointer per warp/wavefront
- Cache per warp/wavefront, but without synchronisation between them

GPU Architectures



Wavefront and Other Vendor Terminology

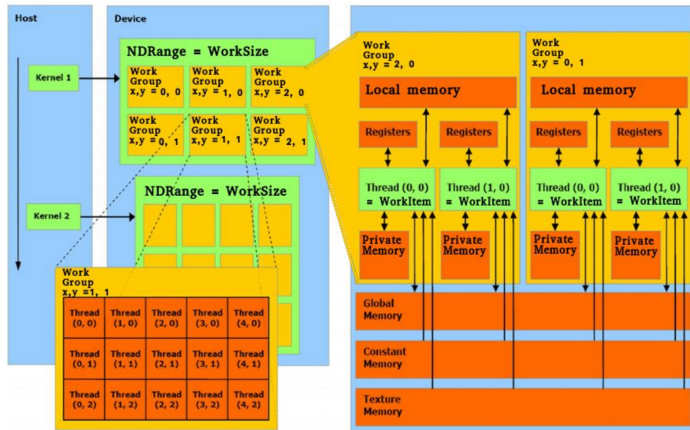
- **Wavefront (AMD)**: The set of work-items executed in lockstep (often 32 or 64 threads)
- **Warp (NVIDIA)**: Similar concept to Wavefront; typically 32 threads
- **Wave (Intel)**: May use different names or sizes depending on the architecture
- Conceptually, these terms all refer to a hardware-level group of threads that share instruction decoding

Why Massive Parallelism on GPUs?

- GPUs can execute thousands of threads simultaneously
- Hardware schedulers distribute threads among available cores
- Fine-grained parallelism allows splitting large tasks into many small subtasks
- Achieves significant speedups for vectorized and matrix-based operations

- **Kernel:** A function executed on a device (GPU) across a set of work-items
- **Work-Item:** A single execution instance of a kernel
- **Work-Group:** A group of work-items that execute together on a compute unit
- **Global Size:** The total number of work-items
- **Local Size:** The number of work-items in a single work-group

OpenCL approach




```

1  #define VALUES_PER_WORK_ITEM 32
2  #define WORKGROUP_SIZE 32
3
4  __kernel void atomic_sum(__global const int *arr,
5                          __global unsigned int *sum,
6                          unsigned int n)
7  {
8      unsigned int id = get_global_id(0);
9      if (id < n)
10     {
11         atomic_add(sum, arr[id]);
12     }
13 }
14
15 __kernel void loop_sum(__global const int *arr,
16                       __global unsigned int *sum,
17                       unsigned int n)
18 {
19     const unsigned int idx = get_global_id(0);
20     unsigned int res = 0;
21     for (int i = idx * VALUES_PER_WORK_ITEM; i < (idx + 1) * VALUES_PER_WORK_ITEM; ++i)
22     {
23         if (i < n)
24         {
25             res += arr[i];
26         }
27     }
28
29     atomic_add(sum, res);
30 }

```

What just happened?

Important syntax

- **__kernel**: Specifies a kernel function
- **__global**: Denotes a pointer to global memory
- **__local**: Denotes a pointer to local memory (cache)
- **barrier(int flag)**: Synchronizes work-items within a work-group
 - **CLK_LOCAL_MEM_FENCE**: Synchronizes local memory
 - **CLK_GLOBAL_MEM_FENCE**: Synchronizes global memory
- **get_global_id(int dim)**: Returns the global ID of the current work-item
- **get_local_id(int dim)**: Returns the local ID of the current work-item
- **get_group_id(int dim)**: Returns the group ID of the current work-item

Coalesced Memory Access

When different work items access consecutive memory locations, the GPU can combine these requests into a single memory transaction. This is called **coalesced memory access**.

Example: Sum of vector elements (Coalesced vs Non-coalesced Access)

```
1  __kernel void loop_sum
2  (
3      __global const int *arr
4      , __global unsigned int *sum
5      , unsigned int n
6  )
7  {
8      const unsigned int idx = get_global_id(0);
9      unsigned int res = 0;
10     for (int i = idx * VALUES_PER_WORK_ITEM; i <
11         (idx + 1) * VALUES_PER_WORK_ITEM; ++i)
12     {
13         if (i < n)
14         {
15             res += arr[i];
16         }
17     }
18     atomic_add(sum, res);
19 }
```

```
1
2  __kernel void loop_coalesced_sum
3  (
4      __global const int *arr
5      , __global unsigned int *sum
6      , unsigned int n
7  )
8  {
9      const unsigned int lid = get_local_id(0);
10     const unsigned int wid = get_group_id(0);
11     const unsigned int grs = get_local_size(0);
12
13     unsigned int res = 0;
14     for (int i = 0; i < VALUES_PER_WORK_ITEM; ++
15         i)
16     {
17         int idx = wid * grs * VALUES_PER_WORK_ITEM
18             + i * grs + lid;
19         if (idx < n)
20         {
21             res += arr[idx];
22         }
23     }
24     atomic_add(sum, res);
25 }
```