

PMC Lecture 19

Gianmaria Romano

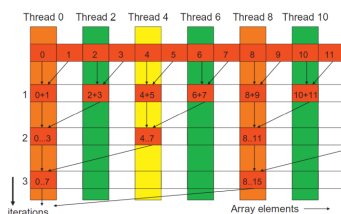
17 December 2025

Chapter 6

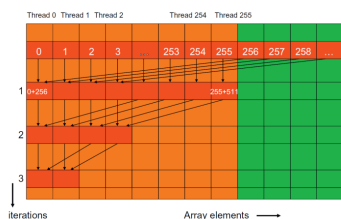
Compute Unified Device Architecture

6.1 Reduction operations in CUDA

While atomic operations can be used to implement reductions in CUDA, they often suffer from poor performance because, at each step, only half of the threads actually performs computations, resulting in a utilization drop.



For this reason, it is convenient to structure reductions so that threads within the same warp will always execute the same task, although only a part of the used threads may actually contribute to computations as the reduction progresses, therefore improving utilization.



Example: The following code provides an efficient implementation of a reduction operation on an array:

```

1 #include <stdio.h>
2 #include <cuda.h>
3
4 __shared__ float partialSum[SIZE];
5 partialSum[threadIdx.x] = x[blockIdx.x * blockDim.x + threadIdx.x];
6
7 unsigned int t = threadIdx.x;
8 for (unsigned int stride = blockDim.x / 2; stride >= 1; stride = stride>>1) {
9     // N.B.: The loop divides the stride by 2 at each iteration.
10    __syncthreads();
11    if (t < stride) {
12        partialSum[t] += partialSum[t + stride];
13    }
14 }

```

Note that, if the number of elements in the array is larger than the number of threads within a block, it is possible to introduce some additional kernels and work on partial results that will be merged at the end.

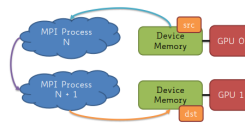
6.2 GPU data transfers

When dealing with distributed GPU applications, data transfers between GPUs are often handled using message passing interface.

In particular, if the message passing interface is not GPU-aware, data must be explicitly copied from the sender device to the host memory before making the desired call.

On the receiving side, data is transferred from the host to the receiver device.

- Source MPI process:
 - `cudaMemcpy(tmp, src, cudaMemcpyDeviceToHost)`
 - `MPI_Send(tmp, ..., N-1, ...)` // send tmp to N+1
- Destination MPI process:
 - `MPI_Recv(tmp, ..., N, ...)` // recv tmp from N+1
 - `cudaMemcpy(dst, tmp, cudaMemcpyHostToDevice)`



On the other hand, if the message passing interface is GPU aware, it can directly access device memory, allowing to handle communications without having to involve the host memory during data transfers.

