



GPU Teaching Kit

Accelerated Computing



Module 9 – Parallel Computation Patterns (Reduction)

Lecture 9.2 - A Basic Reduction Kernel

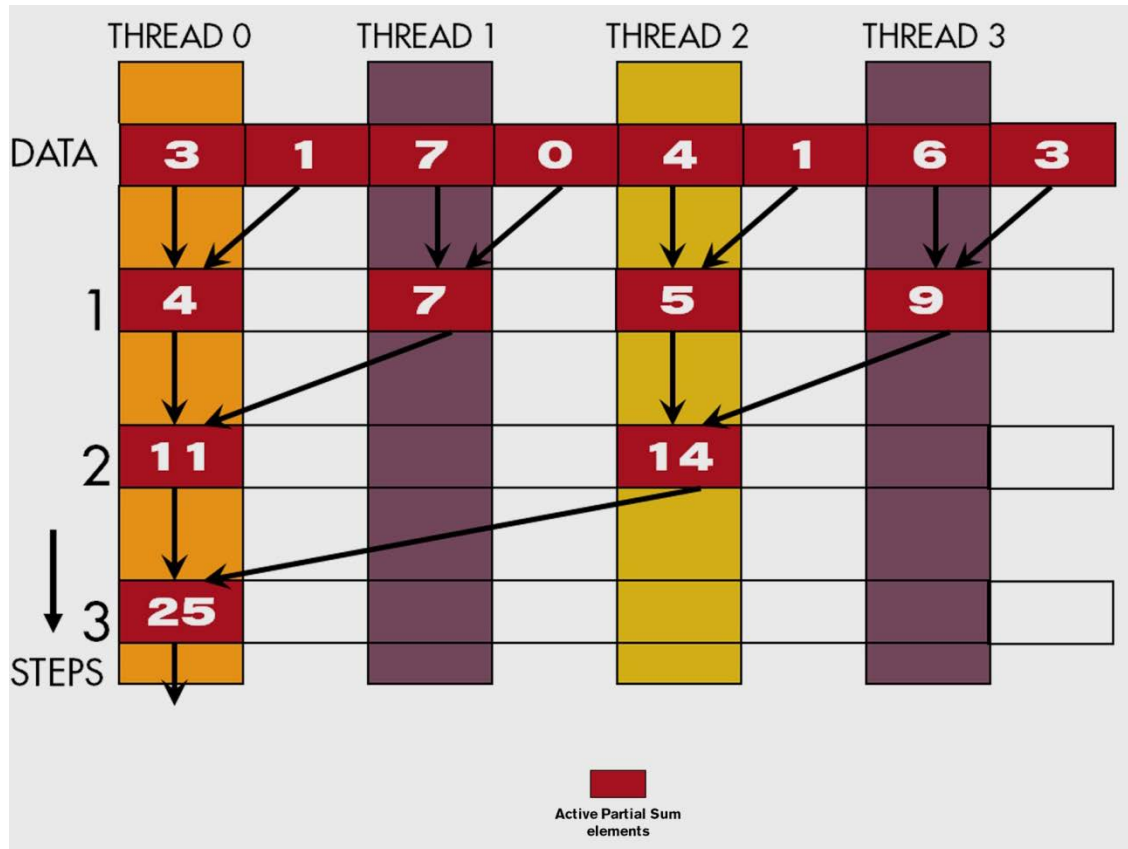
Objective

- To learn to write a basic reduction kernel
 - Thread to data mapping
 - Turning off threads
 - Control divergence

Parallel Sum Reduction

- Parallel implementation
 - Each thread adds two values in each step
 - Recursively halve # of threads
 - Takes $\log(n)$ steps for n elements, requires $n/2$ threads
- Assume an in-place reduction using shared memory
 - The original vector is in device global memory
 - The shared memory is used to hold a partial sum vector
 - Initially, the partial sum vector is simply the original vector
 - Each step brings the partial sum vector closer to the sum
 - The final sum will be in element 0 of the partial sum vector
 - Reduces global memory traffic due to reading and writing partial sum values
 - Thread block size limits n to be less than or equal to 2,048

A Parallel Sum Reduction Example



A Naive Thread to Data Mapping

- Each thread is responsible for an even-index location of the partial sum vector (location of responsibility)
- After each step, half of the threads are no longer needed
- One of the inputs is always from the location of responsibility
- In each step, one of the inputs comes from an increasing distance away

A Simple Thread Block Design

- Each thread block takes $2 \times \text{BlockDim.x}$ input elements
- Each thread loads 2 elements into shared memory

```
__shared__ float partialSum[2*BLOCK_SIZE];  
  
unsigned int t = threadIdx.x;  
unsigned int start = 2*blockIdx.x*blockDim.x;  
partialSum[t] = input[start + t];  
partialSum[blockDim+t] = input[start + blockDim.x+t];
```

The Reduction Steps

```
for (unsigned int stride = 1;
     stride <= blockDim.x; stride *= 2)
{
    __syncthreads();
    if (t % stride == 0)
        partialSum[2*t] += partialSum[2*t+stride];
}
```

Why do we need `__syncthreads()`?

Barrier Synchronization

- `__syncthreads()` is needed to ensure that all elements of each version of partial sums have been generated before we proceed to the next step

Back to the Global Picture

- At the end of the kernel, Thread 0 in each block writes the sum of the thread block in `partialSum[0]` into a vector indexed by the `blockIdx.x`
- There can be a large number of such sums if the original vector is very large
 - The host code may iterate and launch another kernel
- If there are only a small number of sums, the host can simply transfer the data back and add them together
- Alternatively, Thread 0 of each block could use atomic operations to accumulate into a global sum variable.



GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/).