



GPU Teaching Kit

Accelerated Computing



Module 10.4 – Parallel Computation Patterns (scan)

More on Parallel Scan

Objective

- To learn more about parallel scan
 - Analysis of the work efficient kernel
 - Exclusive scan
 - Handling very large input vectors

Work Analysis of the Work Efficient Kernel

- The work efficient kernel executes $\log(n)$ parallel iterations in the reduction step
 - The iterations do $n/2, n/4, \dots, 1$ adds
 - Total adds: $(n-1) \rightarrow O(n)$ work
- It executes $\log(n)-1$ parallel iterations in the post-reduction reverse step
 - The iterations do $2-1, 4-1, \dots, n/2-1$ adds
 - Total adds: $(n-2) - (\log(n)-1) \rightarrow O(n)$ work
- Both phases perform up to no more than $2x(n-1)$ adds
- The total number of adds is no more than twice of that done in the efficient sequential algorithm
 - The benefit of parallelism can easily overcome the 2X work when there is sufficient hardware

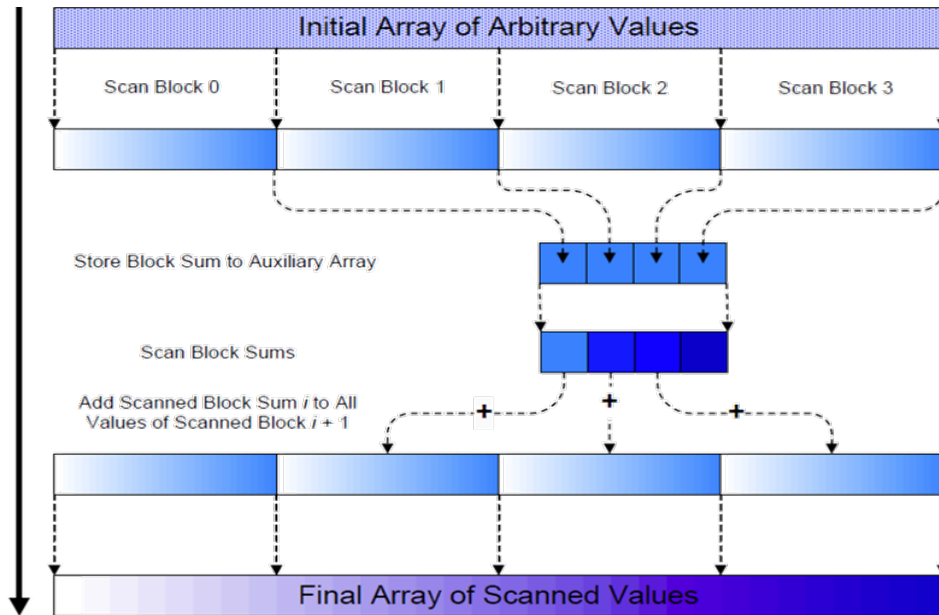
Some Tradeoffs

- The work efficient scan kernel is normally more desirable
 - Better Energy efficiency
 - Less execution resource requirement
- However, the work inefficient kernel could be better for absolute performance due to its single-phase nature (forward phase only)
 - There is sufficient execution resource

Handling Large Input Vectors

- Build on the work efficient scan kernel
- Have each section of $2 \times \text{blockDim.x}$ elements assigned to a block
 - Perform parallel scan on each section
- Have each block write the sum of its section into a `Sum[]` array indexed by `blockIdx.x`
- Run the scan kernel on the `Sum[]` array
- Add the scanned `Sum[]` array values to all the elements of corresponding sections
- Adaptation of work inefficient kernel is similar.

Overall Flow of Complete Scan



Exclusive Scan Definition

Definition: The exclusive scan operation takes a binary associative operator \oplus , and an array of n elements

$$[x_0, x_1, \dots, x_{n-1}]$$

and returns the array

$$[0, x_0, (x_0 \oplus x_1), \dots, (x_0 \oplus x_1 \oplus \dots \oplus x_{n-2})].$$

Example: If \oplus is addition, then the exclusive scan operation on the array `[3 1 7 0 4 1 6 3]`, would return `[0 3 4 11 11 15 16 22]`.

Why Use Exclusive Scan?

- To find the beginning address of allocated buffers
- Inclusive and exclusive scans can be easily derived from each other; it is a matter of convenience

[3 1 7 0 4 1 6 3]

Exclusive [0 3 4 11 11 15 16 22]

Inclusive [3 4 11 11 15 16 22 25]

A Simple Exclusive Scan Kernel

- Adapt an inclusive, work inefficient scan kernel
- Block 0:
 - Thread 0 loads 0 into $XY[0]$
 - Other threads load $X[\text{threadIdx.x}-1]$ into $XY[\text{threadIdx.x}]$
- All other blocks:
 - All thread load $X[\text{blockIdx.x}*\text{blockDim.x}+\text{threadIdx.x}-1]$ into $XY[\text{threadIdx.x}]$
- Similar adaption for work efficient scan kernel but ensure that each thread loads two elements
 - Only one zero should be loaded
 - All elements should be shifted to the right by only one position

Read the Harris article (Parallel Prefix Sum with CUDA) for a more intellectually interesting approach to exclusive scan kernel implementation.



GPU Teaching Kit



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/).