ECE408 Spring 2020

Applied Parallel Programming

Lecture 16
Parallel Computation Patterns –
Parallel Scan (Prefix Sum)

1

1

# Objective

- to learn parallel scan (prefix sum) algorithms based on reductions and reverse reductions

- to learn the concept of double buffering

- to understand tradeoffs between work efficiency and latency

- to learn how to develop hierarchical algorithms (across multiple kernels)

2

2

# Scan Includes all Partial Results

Reductions are a simplified form of scans.

In scan / parallel prefix,
- we need all of the partial sums
- (or whatever the operator might be).

3

3

# (Inclusive) Scan (Prefix-Sum) Definition

**Definition:** *The scan operation takes a binary associative operator* $\oplus$*, and an array of n elements*
$$[x_0, x_1, ..., x_{n-1}],$$

*and returns the prefix-sum array*

$$[x_0, (x_0 \oplus x_1), ..., (x_0 \oplus x_1 \oplus ... \oplus x_{n-1})].$$

**Example:** If $\oplus$ is addition, the scan operation
on the array        [3  1  7  0  4  1  6  3],
returns        [3  4 11 11 15 16 22 25].

4

4

1

## Example: Sharing a Big Sandwich

You order a 100-inch sandwich to feed 10 people,
and you know how much each person wants in inches:
  [3  5  2  7  28  4  3  0  8  1].
  **How do you cut the bread quickly?**
  **How much of the sandwich is left over?**
Method 1: sequentially!
  Cut 3 inches, then cut 5 inches, then …
Method 2: **calculate cutting offsets with prefix-sum**
  [3, 8, 10, 17, 45, 49, 52, 52, 60, 61] (39 inches left)

5

---

## Typical Applications of Scan

A simple and useful parallel building block.

Convert sequential recurrences
```
for(j=1;j<n;j++)
    out[j] = out[j-1] + f(j);
```
into parallel:
```
forall(j) { temp[j] = f(j) };
scan(out, temp);
```

6

---

## Typical Applications of Scan

- Useful for many parallel algorithms:
  - radix sort
  - quicksort
  - String comparison
  - Lexical analysis
  - Stream compaction
  - Polynomial evaluation
  - Solving recurrences
  - Tree operations
  - Histograms
  - Etc.

7

---

## Other Applications

- Assigning camp slots
- Assigning farmer market space
- Allocating memory to parallel threads
- Allocating memory buffer to communication channels
- …

8

---

## An Inclusive Sequential Scan

Given a sequence $[x_0, x_1, x_2, ... ]$

Calculate output $[y_0, y_1, y_2, ... ]$

Such that

$$y_0 = x_0$$
$$y_1 = x_0 + x_1$$
$$y_2 = x_0 + x_1 + x_2$$

...

*Using a recursive definition*

$$y_i = y_{i-1} + x_i$$

9

---

## An Sequential C Implementation

```
y[0] = x[0];
for (i = 1; i < Max_i; i++) y[i] = y [i-1] + x[i];
```

Computationally efficient:

N additions needed for N elements - O(N)!

10

---

## A Naïve Inclusive Parallel Scan

• Assign one thread to calculate each y element

• Have every thread to add up all x elements needed for the y element

$$y_0 = x_0$$
$$y_1 = x_0 + x_1$$
$$y_2 = x_0 + x_1 + x_2$$

"Parallel programming is easy as long as you do not care about performance."

11

---

## Parallel Inclusive Scan using Reduction Trees

Calculate each output element as the reduction of all previous elements

• Some reduction partial sums will be shared among the calculation of output elements

• Based on hardware added design by Peter Kogge and Harold Stone at IBM in the 1970s – Kogge-Stone Trees

• Goal: low latency

12

---

9

10

11

12

## Slide 13

# A Kogge-Stone Parallel Scan Algorithm

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

1. Load input from global memory into shared memory array T

Each thread loads one value from the input (global memory) array into shared memory array T.

13

13

## Slide 14

# A Kogge-Stone Parallel Scan Algorithm

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |

Thread 5

1. (previous slide)
2. Assuming n is a power of 2. Iterate log(n) times, stride from 1 to n/2. Threads *stride* to *n-1 active:* add pairs of elements that are s*tride* elements apart.

- Active threads: *stride* to *n*-1 (*n - stride* active threads)
- Thread *j* adds elements T[*j*] and T[*j-stride*] and writes result into element T[*j*]
- Each iteration requires two syncthreads
  - make sure that input is in place
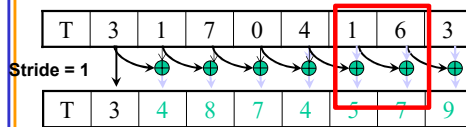  - make sure that all input elements have been used

Iteration #1
Stride = 1

14

14

## Slide 15

# A Kogge-Stone Parallel Scan Algorithm

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |

Thread 6

1. (previous slide)
2. Assuming n is a power of 2. Iterate log(n) times, stride from 1 to n/2. Threads *stride* to *n-1 active:* add pairs of elements that are s*tride* elements apart.

- Active threads: *stride* to *n*-1 (*n - stride* active threads)
- Thread *j* adds elements T[*j*] and T[*j-stride*] and writes result into element T[*j*]
- Each iteration requires two syncthreads
  - syncthreads(); // make sure that input is in place
  - float temp = T[*j*] + T[*j-stride*];
  - syncthreads(); // make sure that previous output has been consumed
  - T[*j*] = temp;

Iteration #1
Stride = 1
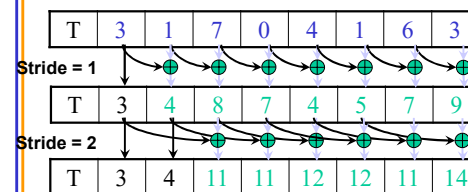
15

15

## Slide 16

# A Kogge-Stone Parallel Scan Algorithm

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |

Stride = 2

| T | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |

Iteration #2
Stride = 2

16

16

4

## A Kogge-Stone Parallel Scan Algorithm

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |
|---|---|---|---|---|---|---|---|---|

**Stride = 1**

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |
|---|---|---|---|---|---|---|---|---|

**Stride = 2**

| T | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |
|---|---|---|----|----|----|----|----|----|

**Stride = 4**

| T | 3 | 4 | 11 | 11 | 15 | 16 | 22 | 25 |
|---|---|---|----|----|----|----|----|----|

1. ...
2. ...
3. Write output from shared memory to device memory

Iteration #3
Stride = 4

17

---

## Sharing Computation in Kogge-Stone

| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |
|---|---|---|---|---|---|---|---|---|

**Stride = 1**

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |
|---|---|---|---|---|---|---|---|---|

**Stride = 2**

| T | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |
|---|---|---|----|----|----|----|----|----|

**Stride = 4**

| T | 3 | 4 | 11 | 11 | 15 | 16 | 22 | 25 |
|---|---|---|----|----|----|----|----|----|

Iteration #3
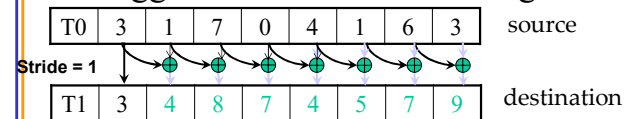Stride = 4

18

---

## Double Buffering

- Use two copies of data T0 and T1
- Start by using T0 as input and T1 as output
- Switch input/output roles after each iteration
  - Iteration 0: T0 as input and T1 as output
  - Iteration 1: T1 as input and T0 and output
  - Iteration 2: T0 as input and T1 as output
- This is typically implemented with two pointers, source and destination that swap their contents from one iteration to the next
- This eliminates the need for the second __syncthreads() call

19

---

## A Double-Buffered Kogge-Stone Parallel Scan Algorithm

| T0 | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |
|----|---|---|---|---|---|---|---|---|

source

**Stride = 1**

| T1 | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |
|----|---|---|---|---|---|---|---|---|

destination

- source = &T0[0]; destination = &T1[0];
- Each iteration requires only one syncthreads()
  - syncthreads(); // make sure that input is in place
  - float destination[*j*] = source[*j*] + source[*j-stride*];
  - temp = destination; destination = source; source = temp;
- After the loop, write destination contents to global memory

Iteration #1
Stride = 1

20

---

17

18

19

20

5

## A Kogge-Stone Parallel Scan Algorithm



| T0 | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T1 | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |    source

Stride = 2

| T0 | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |    destination

Iteration #2
Stride = 2

21

---

## Sharing Computation in Kogge-Stone



| T0 | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T1 | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |

Stride = 2

| T0 | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |    source

Stride = 4

| T1 | 3 | 4 | 11 | 11 | 15 | 16 | 22 | 25 |    destination

- Each iteration requires only one syncthreads()
  - syncthreads(); // make sure that input is in place
  - float destination[$j$] = source[$j$] + source[$j$-$stride$];
  - temp = destination; destination = source; source = temp;
- After the loop, write destination contents to global memory

Iteration #3
Stride = 4

22

---

## Work Efficiency Analysis
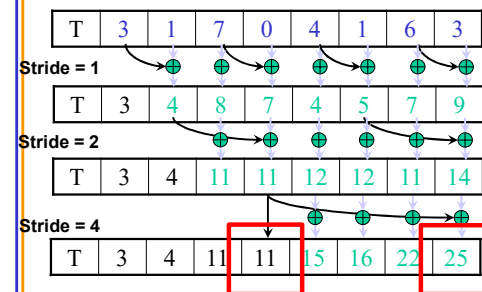
- A Kogge-Stone scan kernel executes log(n) parallel iterations
  - The steps do (n-1), (n-2), (n-4),..(n- n/2) add operations each
  - Total # of add operations: n * log(n)  - (n-1) → O(n*log(n)) work

- This scan algorithm is not very work efficient
  - Sequential scan algorithm does *n* adds
  - A factor of log(n) hurts: 20x for 1,000,000 elements!
  - Typically used within each block, where n ≤ 1,024

- A parallel algorithm can be slow when execution resources are saturated due to low work efficiency

23

---

## A Kogge-Stone Parallel Scan Algorithm



| T | 3 | 1 | 7 | 0 | 4 | 1 | 6 | 3 |

Stride = 1

| T | 3 | 4 | 8 | 7 | 4 | 5 | 7 | 9 |

Stride = 2

| T | 3 | 4 | 11 | 11 | 12 | 12 | 11 | 14 |

Stride = 4

| T | 3 | 4 | 11 | 11 | 15 | 16 | 22 | 25 |

24

---

## Improving Efficiency

- A common parallel algorithm pattern:

  *Balanced Trees*
  - Build a balanced binary tree on the input data and sweep it to and from the root
  - Tree is not an actual data structure, but a concept to determine what each thread does at each step

- For scan:
  - Traverse down from leaves to root building partial sums at internal nodes in the tree
    - Root holds sum of all leaves
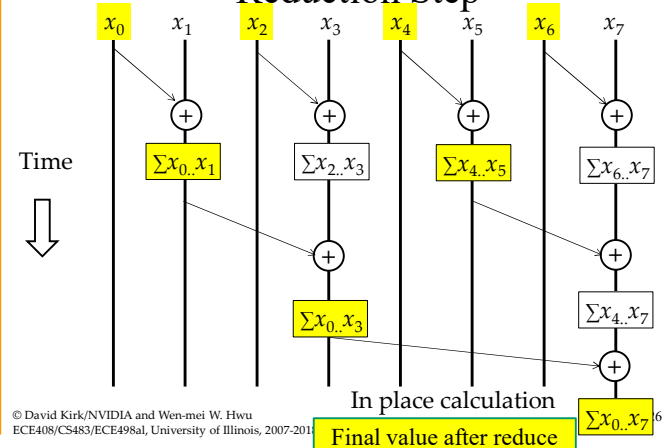  - Traverse back up the tree building the scan from the partial sums

25

---

## Brent-Kung Parallel Scan - Reduction Step



$x_0$ $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$

Time

$\sum x_{0..}x_1$ $\sum x_{2..}x_3$ $\sum x_{4..}x_5$ $\sum x_{6..}x_7$

$\sum x_{0..}x_3$ $\sum x_{4..}x_7$

In place calculation

$\sum x_{0..}x_7$

Final value after reduce

26

---

## Inclusive Post-Scan Step

$x_0$ $\sum x_{0..}x_1$ $x_2$ $\sum x_{0..}x_3$ $x_4$ $\sum x_{4..}x_5$ $x_6$ $\sum x_{0..}x_7$

$\sum x_{0..}x_5$

Move (add) a critical value to a central location where it is needed

27

---

## Inclusive Post Scan Step

$x_0$ $\sum x_{0..}x_1$ $x_2$ $\sum x_{0..}x_3$ $x_4$ $\sum x_{4..}x_5$ $x_6$ $\sum x_{0..}x_7$

$\sum x_{0..}x_5$

$\sum x_{0..}x_2$ $\sum x_{0..}x_4$ $\sum x_{0..}x_6$

28

7

## Putting it Together (Data View)

$x_0$ $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$ $x_8$ $x_9$ $x_{10}$ $x_{11}$ $x_{12}$ $x_{13}$ $x_{14}$ $x_{15}$



© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498al, University of Illinois, 2007-2018

29

---

## Reduction Step Kernel Code

```
// float T[2*BLOCK_SIZE] is in shared memory
// for previous slide, BLOCK_SIZE is 8
int stride = 1;
while(stride < 2*BLOCK_SIZE)
  {
    __syncthreads();
    int index = (threadIdx.x+1)*stride*2 - 1;
    if(index < 2*BLOCK_SIZE && (index-stride) >= 0)
      T[index] += T[index-stride];
    stride = stride*2;
  }
```

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498al, University of Illinois, 2007-2018

30

30

---

## Reduction Step Kernel Code

```
// float T[2*BLOCK_SIZE] is in shared memory

int stride = 1;
while(stride < 2*BLOCK_SIZE)
{
    int index = (threadIdx.x+1)*stride*2 - 1;
    if(index < 2*BLOCK_SIZE && (index-stride) >=0)
      T[index] += T[index-stride];
    stride = stride*2;     // For previous example,
                           // threadIdx.x+1  = 1, 2, 3, 4, 5, 6, 7,8
                           // stride = 1, index = 1, 3, 5, 7, 9, 11, 13, 15
    __syncthreads();
  }
```
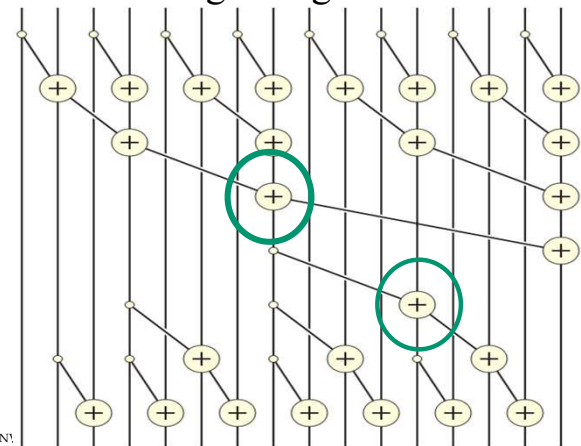
© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498al, University of Illinois, 2007-2018

31

31

---

## Putting it Together



© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498al, University of Illinois, 2007-2018

32

8

## Post Scan Step

```
int stride = BLOCK_SIZE/2;
while(stride > 0)
  {
      __syncthreads();
      int index = (threadIdx.x+1)*stride*2 - 1;
      if((index+stride) < 2*BLOCK_SIZE)
      {
          T[index+stride] += T[index];
      }
      stride = stride / 2;      // for the previous example,
                                // BLOCK_SIZE is 8
                                // stride will go 4, 2, 1
                                // for the first iteration, the active thread
                                // will be thread 0, with index = 7 and
  }                             // index+stride =  11
```
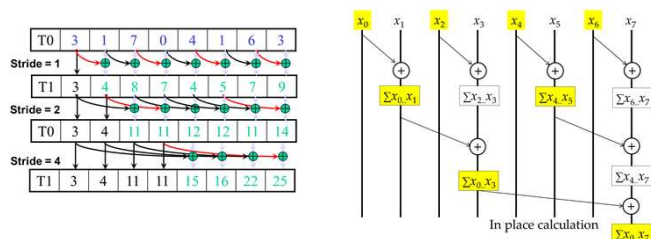
33

---

## Work Analysis

- The parallel Inclusive Scan executes $2* \log(n)$ parallel iterations
  - $\log(n)$ in reduction and $\log(n)$ in post scan
  - The iterations do $n/2, n/4,..1, (2-1), ...., (n/4-1), (n/2-1)$ useful adds
  - In our example, $n = 16$, the number of useful adds is $16/2 + 16/4 + 16/8 + 16/16 + (16/8-1) + (16/4-1) + (16/2-1)$
  - Total adds: $(n-1) + (n-2) - (\log(n) -1) = 2*(n-1) - \log(n)\rightarrow$ O(n) work
- 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

34

---

## Kogge-Stone vs. Brent-Kung
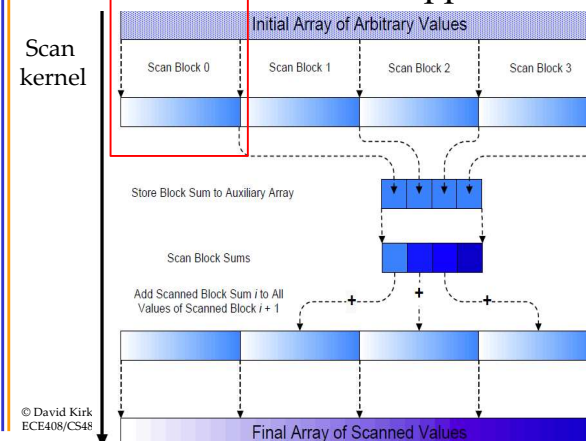


In place calculation

- Brent-Kung uses half the number of threads compared to Kogge-Stone
  - Each thread should load two elements into the shared memory
- Brent-Kung takes twice the number of steps compared to Kogge-Stone
  - Kogge-Stone is more popular for parallel scan with blocks in GPUs

35

---

## Overall Flow of Complete Scan
## A Hierarchical Approach

Scan kernel

36

9

## Using Global Memory Contents in CUDA

- Data in registers and shared memory of one thread block are not visible to other blocks
- To make data visible, the data has to be written into global memory
- However, any data written to the global memory are not visible until a memory fence. This is typically done by terminating the kernel execution
- Launch another kernel to continue the execution. The global memory writes done by the terminated kernels are visible to all tead blocks.

37

---

## Scan of Arbitrary Length Input

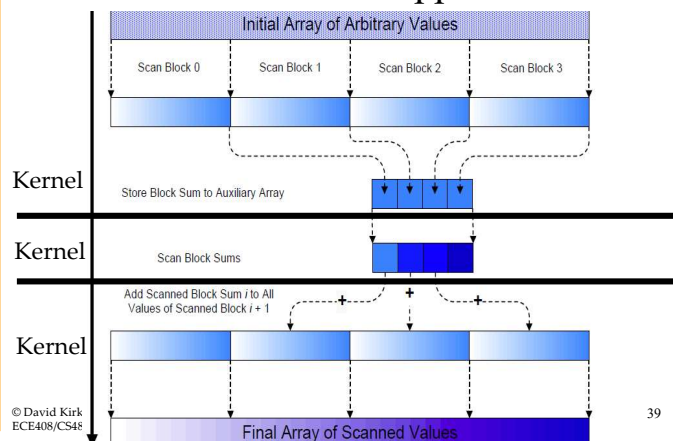- Build on the scan kernel that handles up to 2*blockDim.x elements from Brent-Kung
  - For Kogge-Stone, have each section of blockDim.x elements assigned to a block
- Have each block write the sum of its section into a Sum array using its blockIdx.x as index
- Run parallel scan on the Sum array
  - May need to break down Sum into multiple sections if it is too big for a block
- Add the scanned Sum array values to the elements of corresponding sections

38

---

## Overall Flow of Complete Scan
## A Hierarchical Approach



Initial Array of Arbitrary Values

Scan Block 0    Scan Block 1    Scan Block 2    Scan Block 3

Kernel — Store Block Sum to Auxiliary Array

Kernel — Scan Block Sums

Kernel — Add Scanned Block Sum *i* to All Values of Scanned Block *i* + 1

Final Array of Scanned Values

39

---

## (Exclusive) Scan Definition

**Definition:** *The exclusive* scan *operation takes a binary associative operator* $\oplus$, *and an array of n elements*
$$[x_0, x_1, \ldots, x_{n-1}]$$

*and returns the array*

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

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

40

10

## Why 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]

41

41

## A simple exclusive scan kernel

- Adapt an inclusive, Kogge-Stone scan kernel
  - Block 0:
    - Thread 0 loads 0 into (shared) XY[0]
    - Other threads load (global) X[threadIdx.x-1] into XY[threadIdx.x]
  - All other blocks:
    - All thread load X[blockIdx.x*blockDim.x+threadIdx.x-1] into XY[threadIdex.x]
- Similar adaption for Brent-Kung kernel but pay attention that each thread loads two elements
  - Only one zero should be loaded
  - All elements should be shifted by only one position

- Intellectual contribution vs. practical contribution

42

42

## ANY MORE QUESTIONS?

43

43