# 6.S894 **Accelerated Computing**

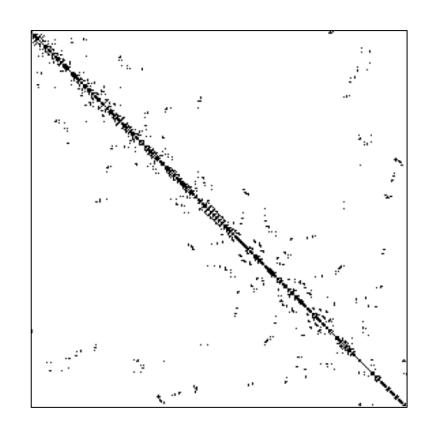
Lecture 7: Data-Parallel

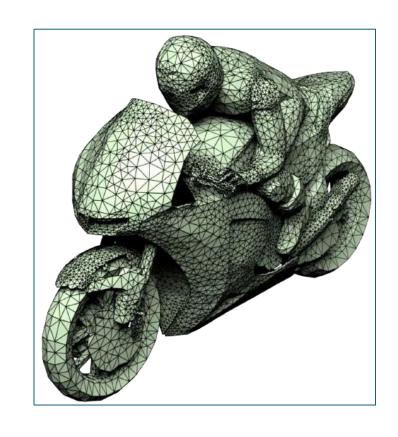
**Primitives** 

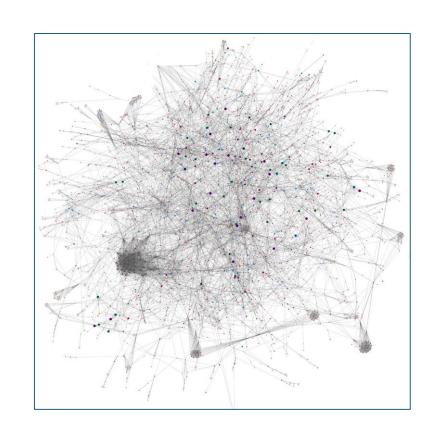
Ahmed Mahmoud



# My work: irregular computation







Sparse Matrices

Meshes

Graphs

## **Outline:**

- Thinking in "data-parallel patterns"
- Overview of data-parallel primitives
  - Map
  - Stencil
  - Reduce
  - Scan
  - Compaction
  - Merge

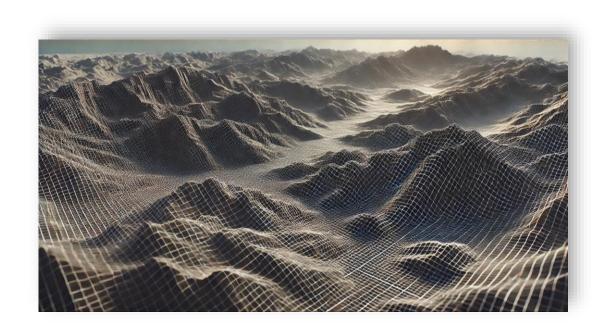
Main idea: high-performance parallel implementations of data-parallel primitives exist, allowing programs written with these primitives to run efficiently on the GPU.

# How bumpy is a surface?

#### Steps:

- 1. For each node, compute the difference between the node's height and the average height of its neighbors, then square that difference
- 2. Sum up all those differences
  - Don't sum all the differences that are zero

(*Note: This is a fake application!*)





```
result = 0
for all nodes:
   average =
       0.25 * (height[x - 1, y] +
               height[x + 1, y ] +
               height[x , y - 1] +
               height[x , y + 1] )
   diff[x, y] = (height[x, y] - average)^2
for all nodes where diff !=0:
   result += diff
return result
```



result += diff

return result

```
result = 0
for all nodes:
   average =
       0.25 * (height[x - 1, y] +
              height[x + 1, y ] +
               height[x , y - 1] +
              height[x , y + 1])
   diff[x, y] = (height[x, y] - average)^2
for all nodes where diff !=0:
```



## Pattern #1: Map

#### Given:

- Sequence of data elements A
- Unary function f(x)

map(A,f) = applies f(x) to all  $a_i \in A$ 

```
result = 0
```

```
for all nodes:
   average =
       0.25 * (height[x - 1, y] +
               height[x + 1, y ] +
               height[x , y - 1] +
               height[x , y + 1])
   diff[x, y] = (height[x, y] - average)^2
for all nodes where diff !=0:
   result += diff
return result
```

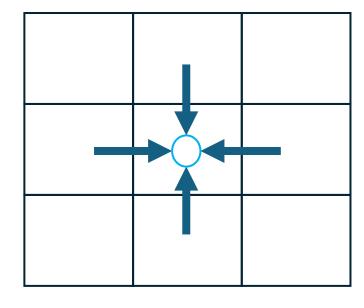


Given an input sequence a and a range I, stencil produces an output sequence p such that

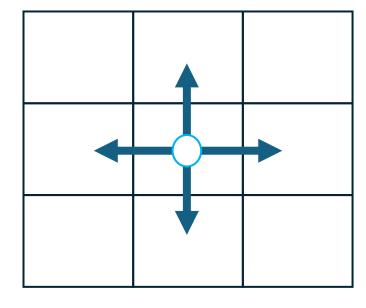
$$p = a[I]$$

Given an input sequence a and a range I, *stencil* produces an output sequence p such that

$$p = a[I]$$

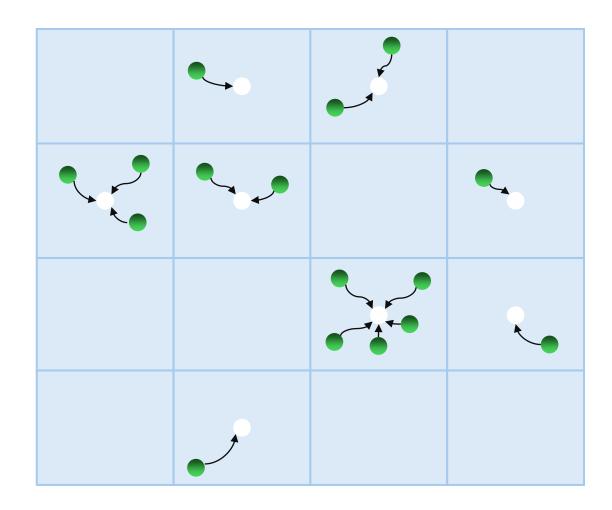


Gather



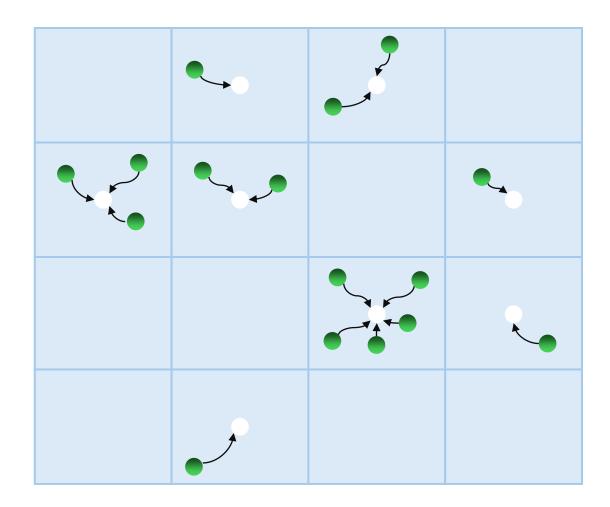
Scatter

Particles-to-grid (P2G)



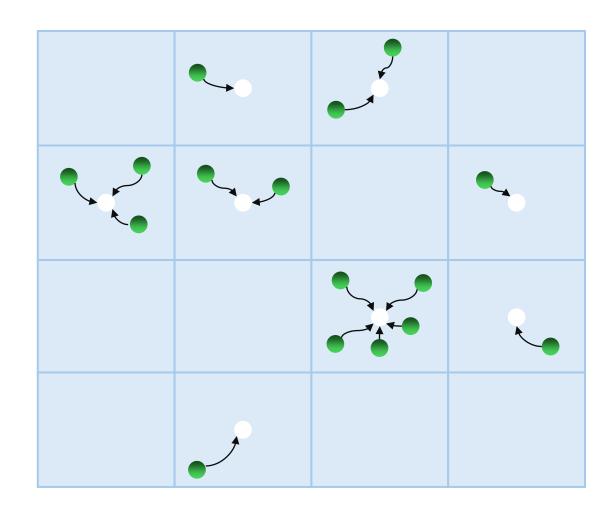
#### Particles-to-grid (P2G)

1. **Gather:** Each cell collects information from the particles within it



#### Particles-to-grid (P2G)

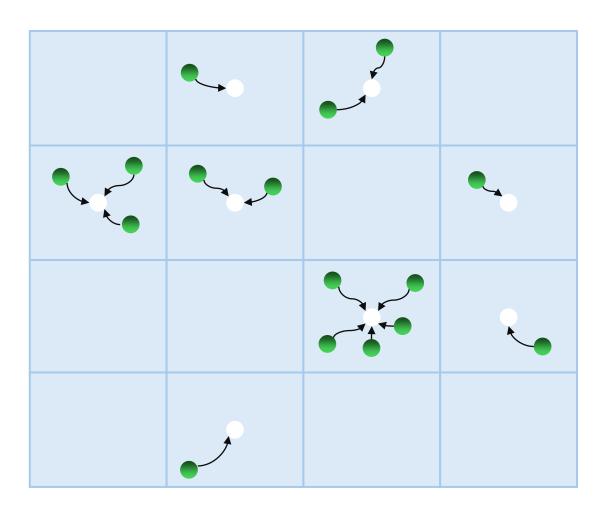
- 1. **Gather:** Each cell collects information from the particles within it
  - Needs to collect particles in each cell
  - Different number of particles per cell
  - Many empty cells
  - Workload imbalance



#### Particles-to-grid (P2G)

- 1. **Gather:** Each cell collects information from the particles within it
  - Needs to collect particles in each cell
  - Different number of particles per cell
  - Many empty cells
  - Workload imbalance

- 2. **Scatter:** Each particle distributes information to the cell it belongs to
  - Needs synchronization (e.g., atomics)



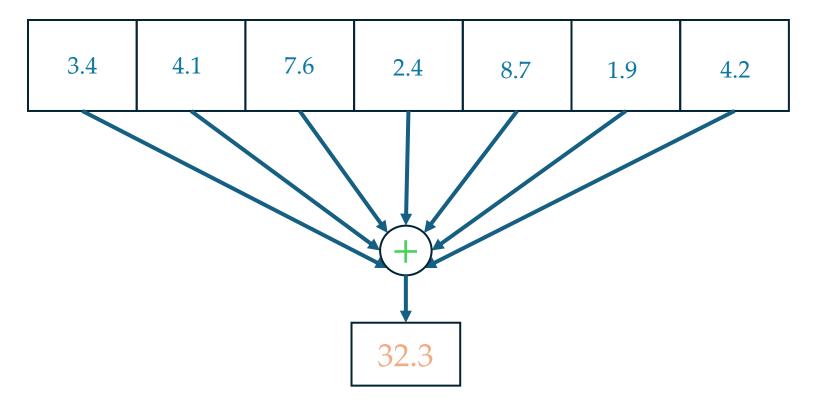
return result

```
result = 0
for all nodes:
   average =
      0.25 * (height[x - 1, y] +
             height[x + 1, y ] +
             height[x , y - 1] +
             diff[x, y] = (height[x, y] - average)^2
for all nodes where diff !=0:
   result += diff
```



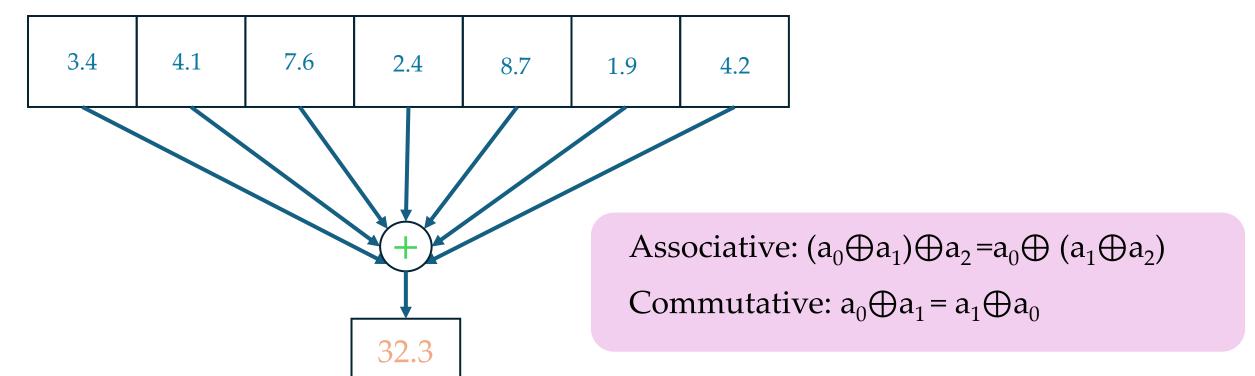
Given a sequence of input  $S = [a_0, a_1, ..., a_{n-1}]$  and a binary associative and commutative operator  $\oplus$  (e.g., +, \*, min, max)

Reduce(
$$\oplus$$
, S) =  $a_0 \oplus a_1 \oplus ... \oplus a_{n-1}$ 



Given a sequence of input  $S = [a_0, a_1, ..., a_{n-1}]$  and a binary associative and commutative operator  $\oplus$  (e.g., +, \*, min, max)

Reduce(
$$\oplus$$
, S) =  $a_0 \oplus a_1 \oplus ... \oplus a_{n-1}$ 





The total amount of work done over by all processors



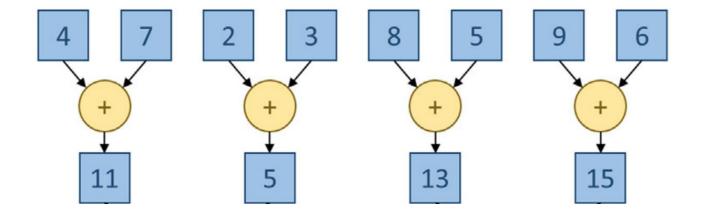
#### Work efficiency:

The total amount of work done over by all processors

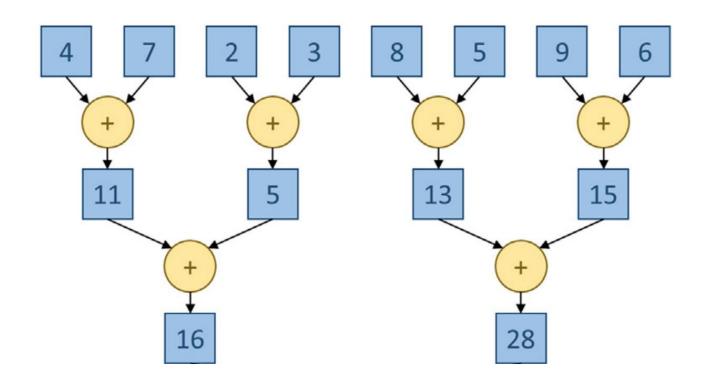


The number of steps it takes to do that work

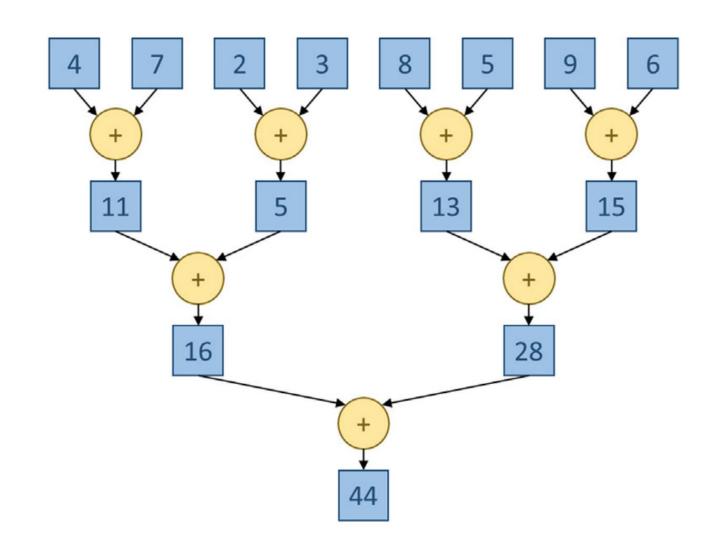
- Parallel Reduction
  - Add two halves of domain together repeatedly



- Parallel Reduction
  - Add two halves of domain together repeatedly

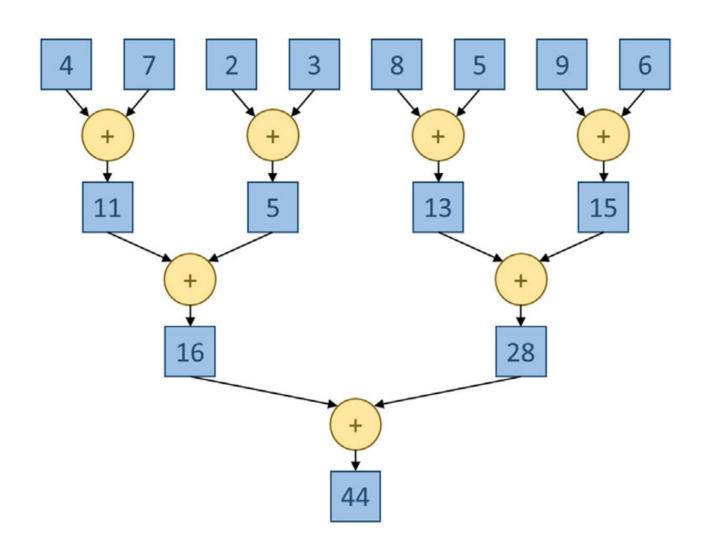


- Parallel Reduction
  - Add two halves of domain together repeatedly



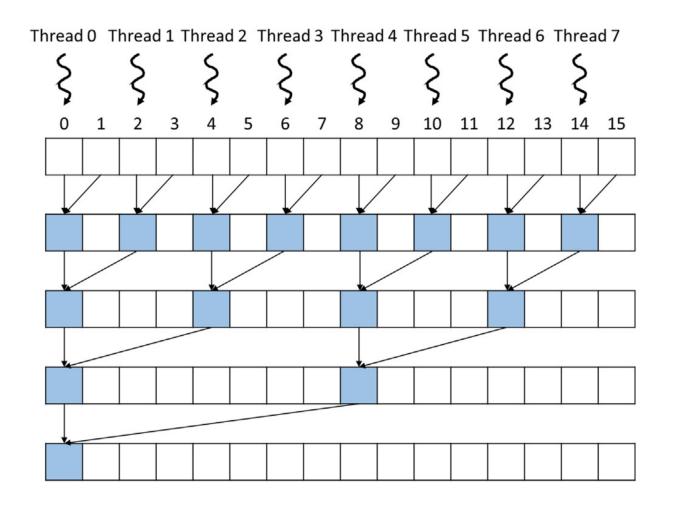
- Parallel Reduction
  - Add two halves of domain together repeatedly

O(log<sub>2</sub>N) Steps O(N) Work

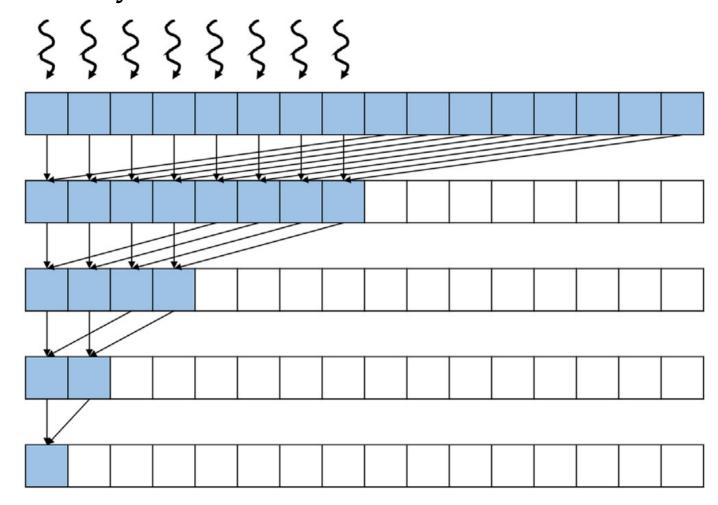


- Parallel Reduction
  - log(N) parallel steps, each step S does N/2S independent ops
    - Step complexity is O(log<sub>2</sub> N)
  - Performs N/2+N/4+...+1=N-1 operations
    - Work complexity is O(N)
    - It is work-efficient, i.e., does not perform more operations than a sequential algorithm

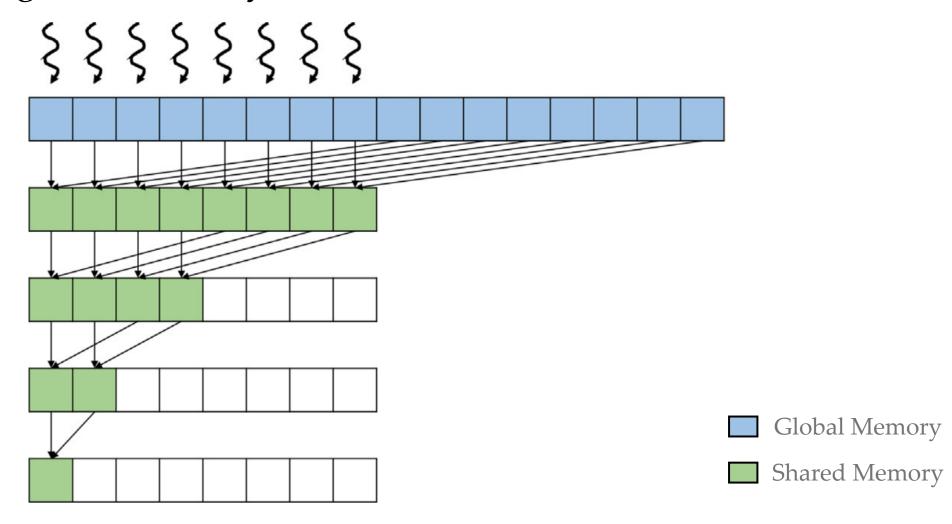
Parallel Reduction Kernel



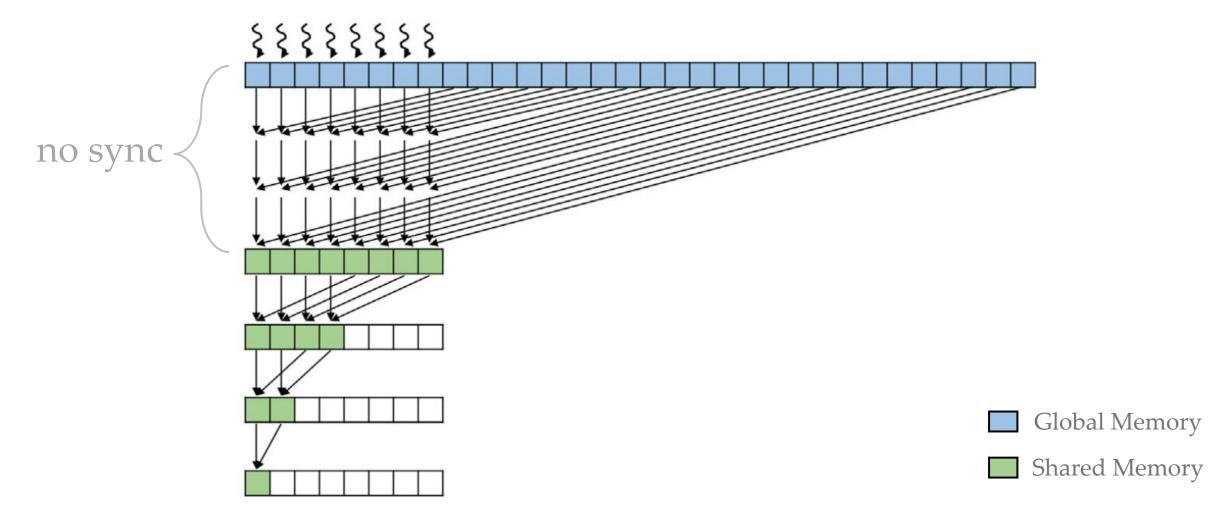
• Improve memory accesses



• Minimize global memory accesses



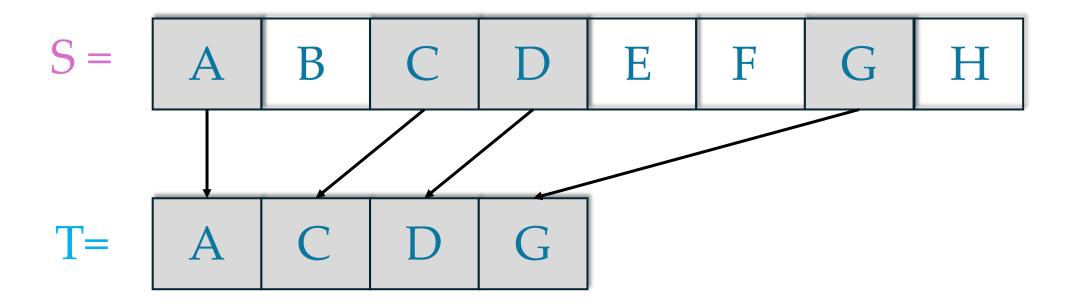
• Thread coarsening to reduce overhead

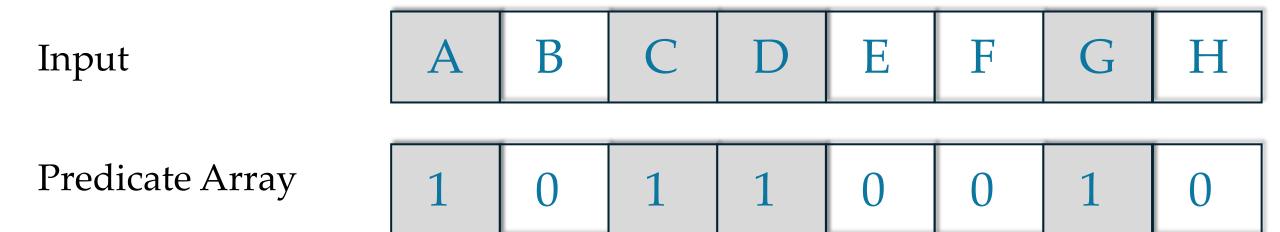


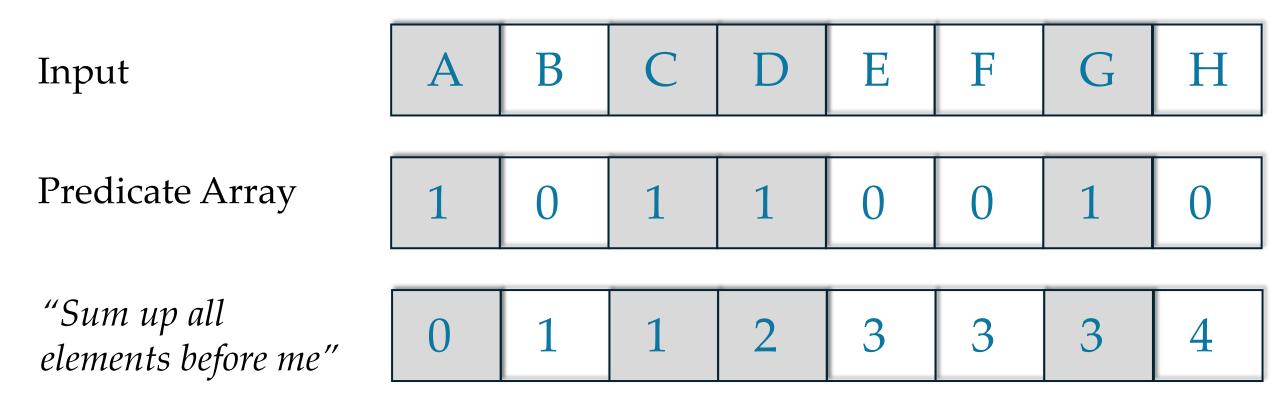
```
result = 0
for all quads:
   average =
      0.25 * (height[x - 1, y] +
             height[x + 1, y ] +
             height[x , y - 1] +
             diff[x, y] = (height[x, y] - average)^2
for all nodes where diff !=0:
   result += diff
return result
```

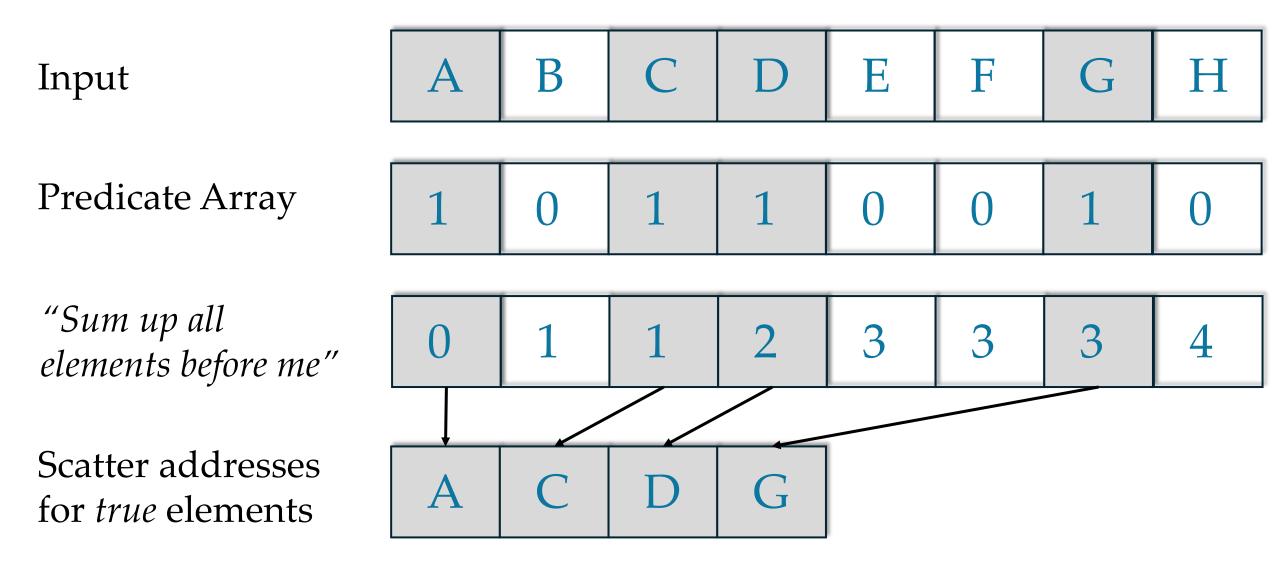


Given an input sequence S and a **predicate** P, output another sequence T that only *contains* elements that satisfy the predicate









"Where do I write my output?"

# "Where do I write my output?"

- The answer is:

"That depends on how much the other threads need to write!"

# "Where do I write my output?"

- The answer is:

"That depends on how much the other threads need to write!"

- Scan is an efficient way to answer this question in parallel
- Examples: marching cubes, collision detection, sorting, building trees

Given a sequence of input  $S = [a_0, a_1, ..., a_{n-1}]$  and a binary associative operator  $\oplus$  with identity I

Scan(S, 
$$\oplus$$
, I) = [I,  $a_0$ ,  $a_0 \oplus a_1$ , ...,  $(a_0 \oplus a_1 ... \oplus a_{n-2})$ ]

When  $\oplus$  is addition, this corresponds to the **prefix sum** 

Input	3	1	7	0	4	1	6	3
Scan	0	3	4	11	11	15	16	22

Example for when  $\oplus$  is addition, i.e., prefix sum

Input	3	1	7	0	4	1	6	3
Exclusive Scan	0	3	4	11	11	15	16	22

Example for when  $\oplus$  is addition, i.e., prefix sum

Input	3	1	7	0	4	1	6	3
Exclusive Scan	0	3	4	11	11	15	16	22
<i>Inclusive</i> Scan	3	4	11	11	15	16	22	25

#### Sequential implementation

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

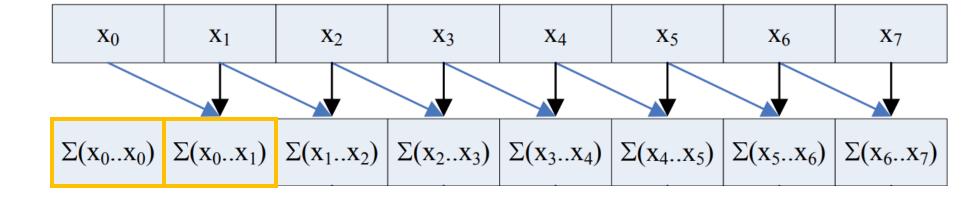
```
Input [x_0, x_1, x_2, ....]
Output [y_0, y_1, y_2, ....]
```

#### Such that

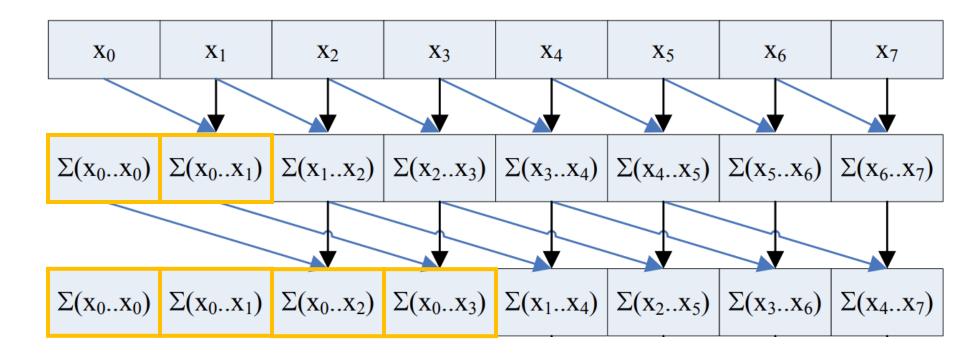
```
y_0 = x_0
y_1 = x_0 + x_1
y_2 = x_0 + x_1 + x_2
...
```

```
y[0] = x[0]
for (i=1; i<len; ++i)
    y[i] = y[i-1] + x[i]</pre>
```

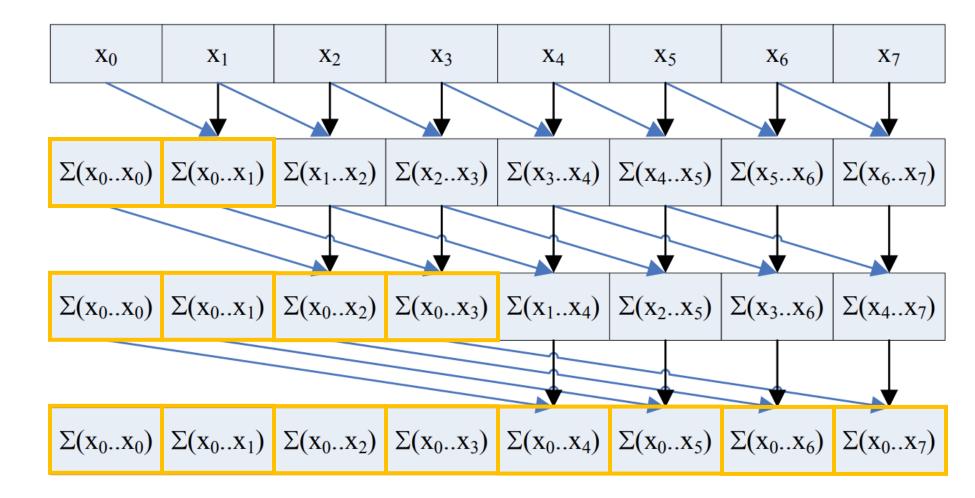
$\mathbf{x}_0$	$\mathbf{X}_1$	X2	X3	$\mathbf{X}_4$	X5	X <sub>6</sub>	<b>X</b> <sub>7</sub>	
110	1.1	112	123	114	1.5	110	11/	



iter 0, stride 2<sup>0</sup>



iter 1, stride 2<sup>1</sup>



iter 2, stride 2<sup>2</sup>

• Kogge-Stone algorithm

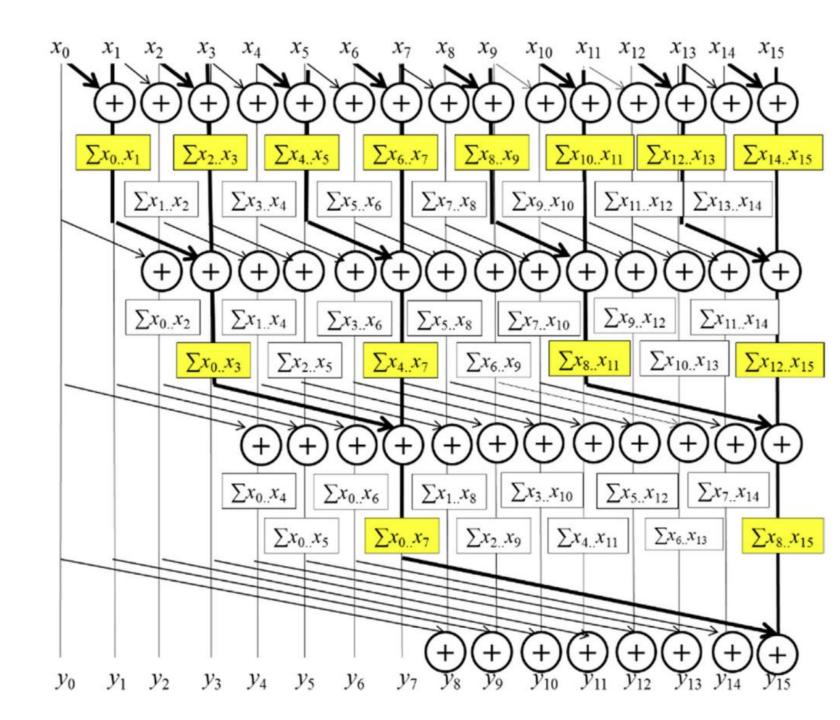
 $O(log_2N)$  Steps  $O(N log_2N)$  Work

 $\mathbf{X}_{\mathbf{0}}$  $\mathbf{X}_{1}$  $X_2$  $X_3$  $X_4$  $X_5$ X6 **X**7  $\Sigma(x_0..x_0) \mid \Sigma(x_0..x_1) \mid \Sigma(x_1..x_2) \mid \Sigma(x_2..x_3) \mid \Sigma(x_3..x_4) \mid \Sigma(x_4..x_5) \mid \Sigma(x_5..x_6) \mid \Sigma(x_6..x_7)$  $\Sigma(x_0..x_0)$   $\Sigma(x_0..x_1)$   $\Sigma(x_0..x_2)$   $\Sigma(x_0..x_3)$   $\Sigma(x_1..x_4)$   $\Sigma(x_2..x_5)$   $\Sigma(x_3..x_6)$   $\Sigma(x_4..x_7)$  $\Sigma(x_0..x_0)$   $\Sigma(x_0..x_1)$   $\Sigma(x_0..x_2)$   $\Sigma(x_0..x_3)$   $\Sigma(x_0..x_4)$   $\Sigma(x_0..x_5)$   $\Sigma(x_0..x_6)$   $\Sigma(x_0..x_7)$ 

Step efficient

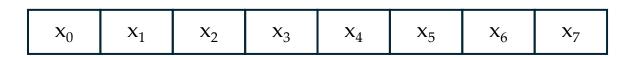
Not work efficient 🔀

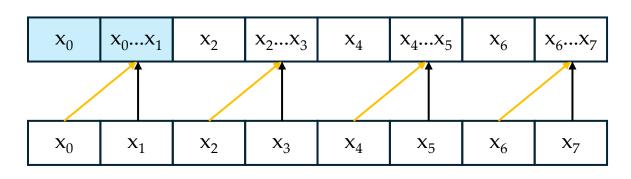
• Kogge-Stone algorithm



iter = 0  

$$(tid+1)\%2^{1} == 0$$
  
 $stride = 2^{0}$ 

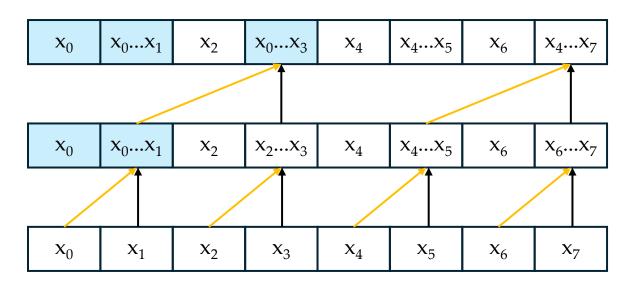




iter = 0  

$$(tid+1)\%2^{1} == 0$$
  
 $stride = 2^{0}$ 

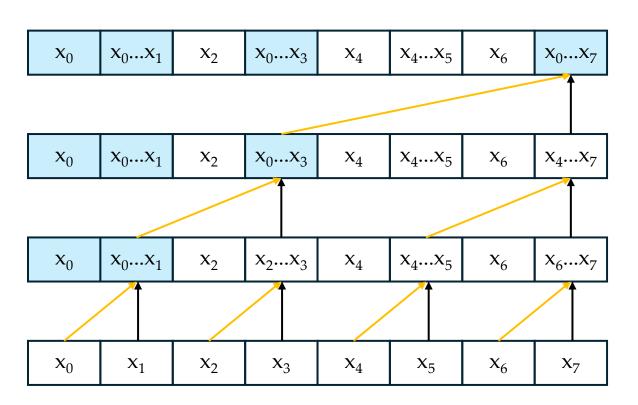
Up-sweep Phase



iter = 1  

$$(tid+1)\%2^2 == 0$$
  
 $stride = 2^1$ 

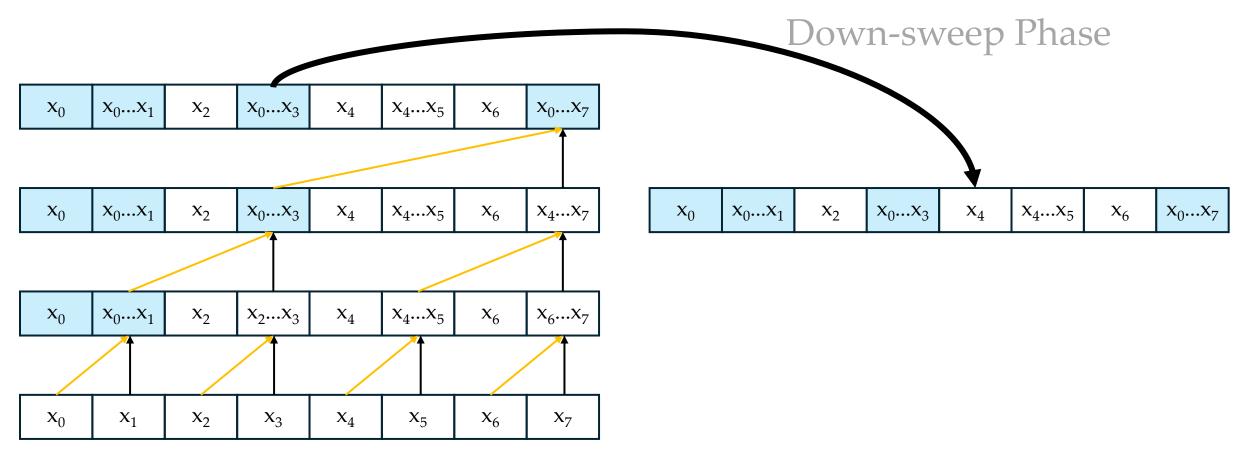
Up-sweep Phase



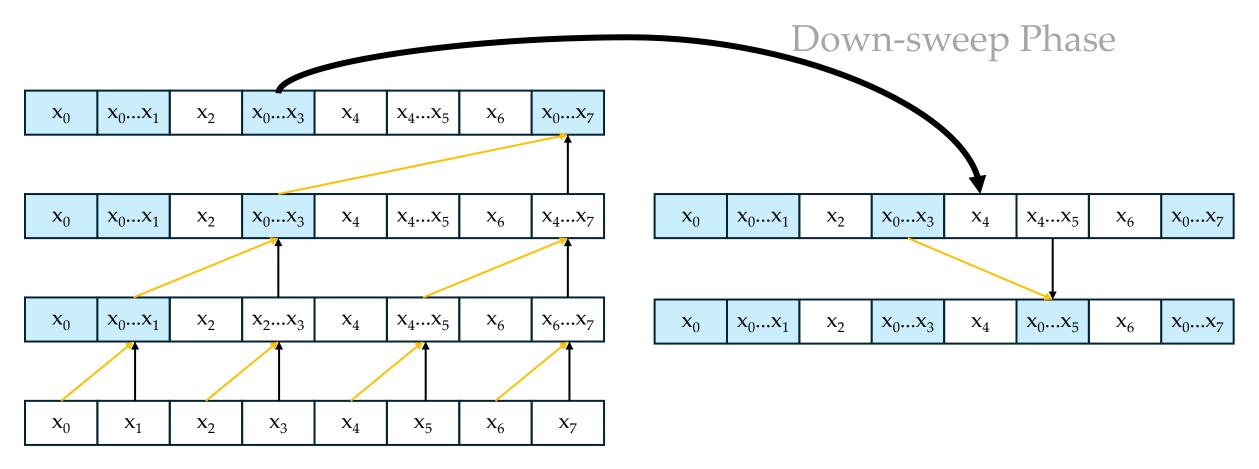
iter = 2  

$$(tid+1)\%2^3 == 0$$
  
 $stride = 2^2$ 

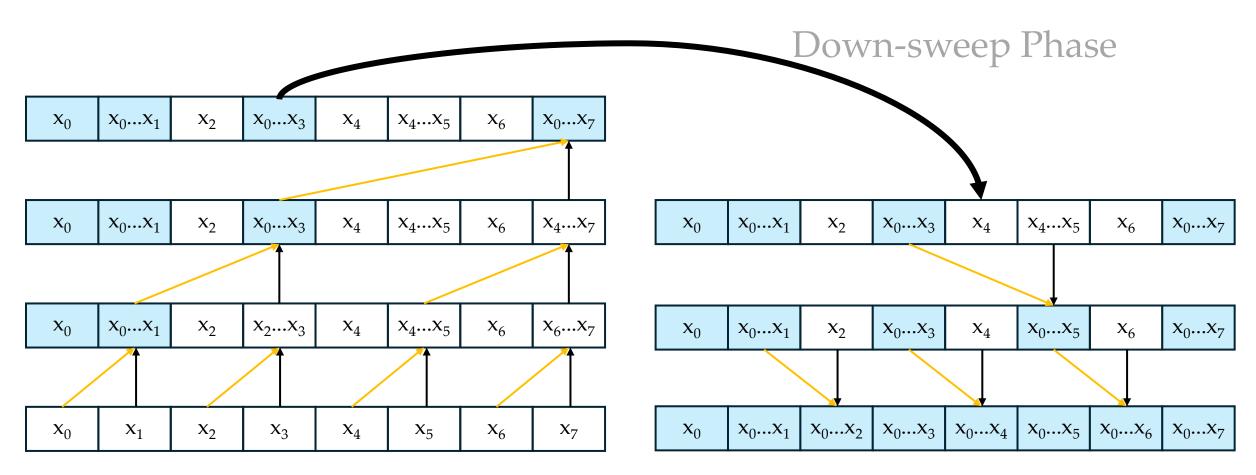
Up-sweep Phase



Up-sweep Phase

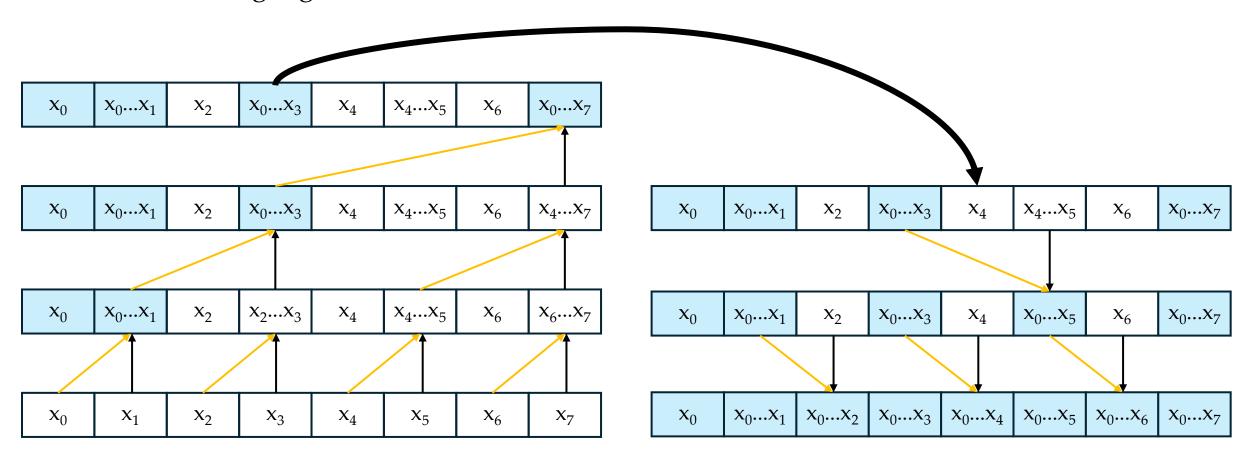


Up-sweep Phase



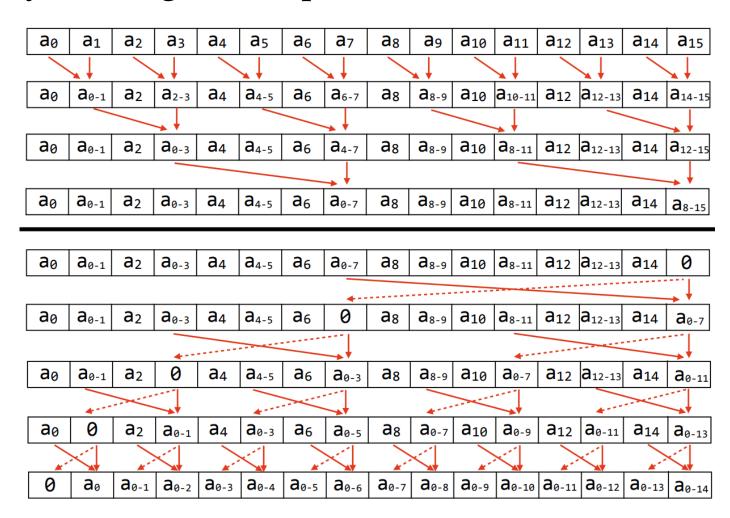
Up-sweep Phase

• Brent-Kung algorithm

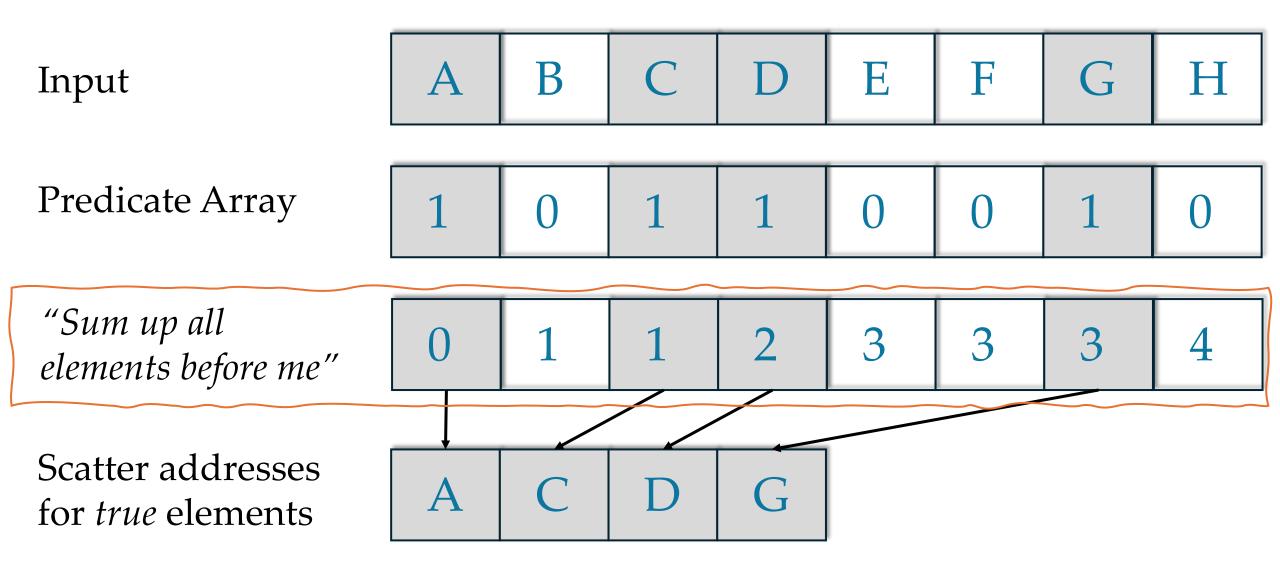


• Improve work efficiency by reusing of computation results

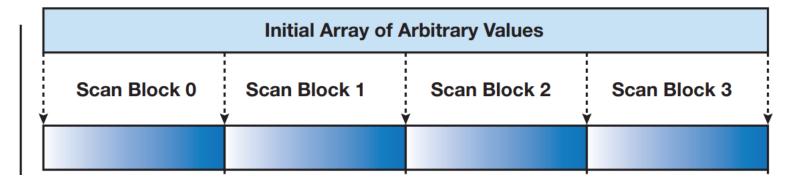
O(log<sub>2</sub>N) Steps O(N) Work



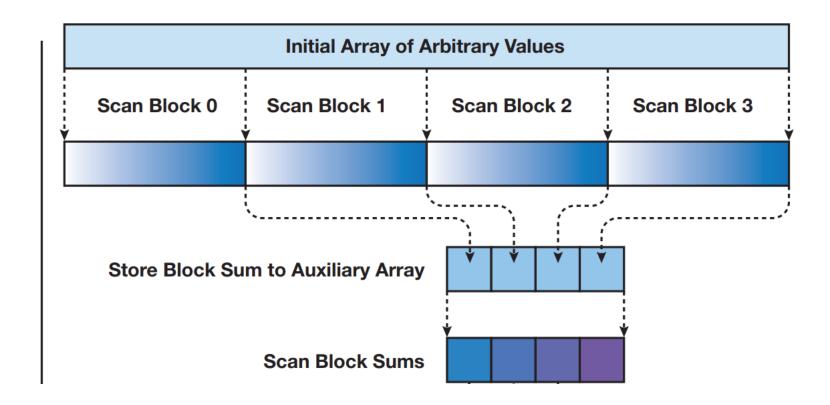
# Back to Stream Compaction



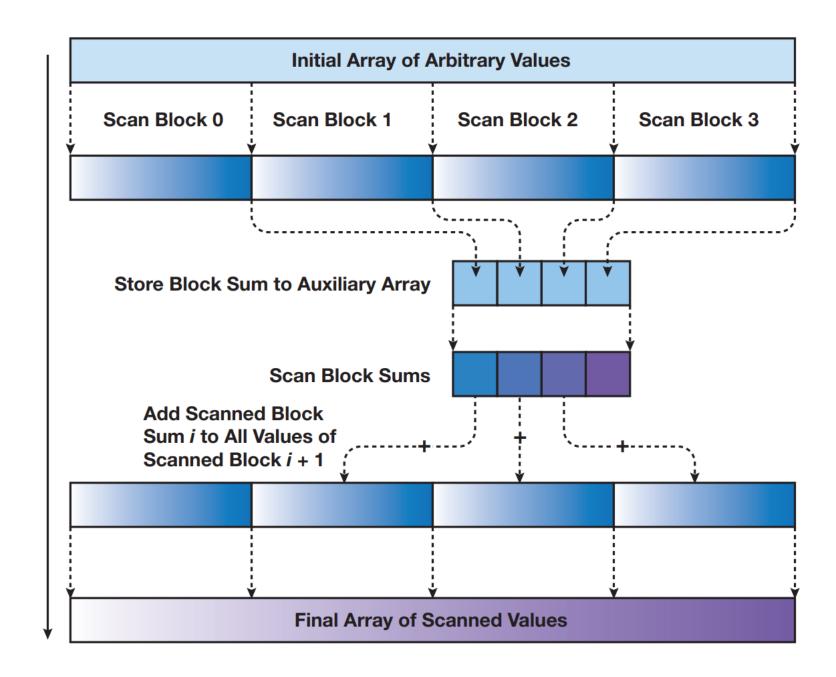
Scan across blocks



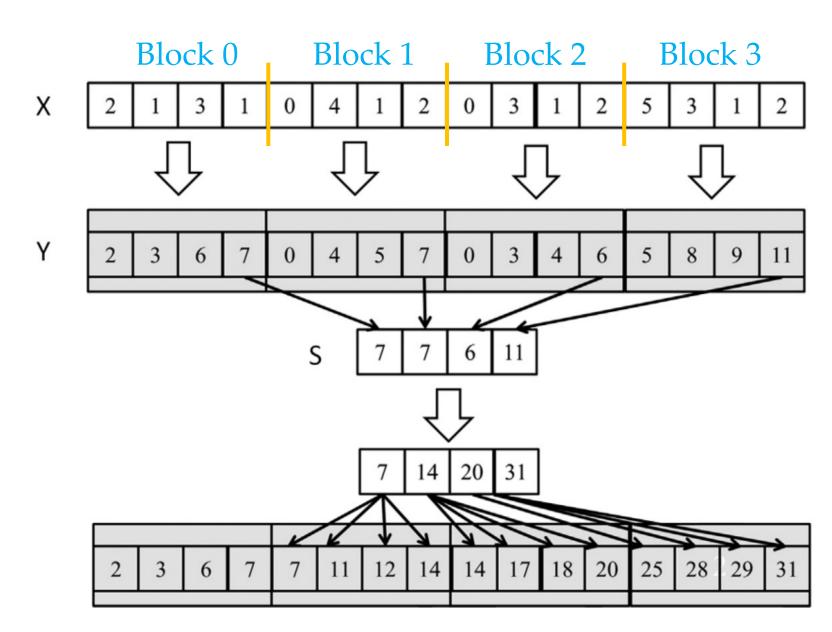
Scan across blocks



Scan across blocks

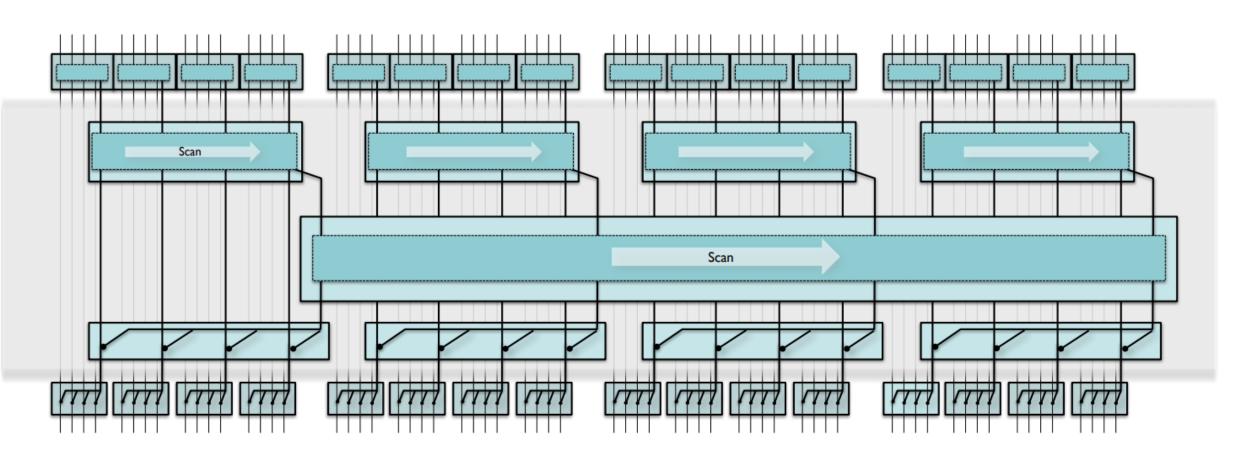


Scan across blocks



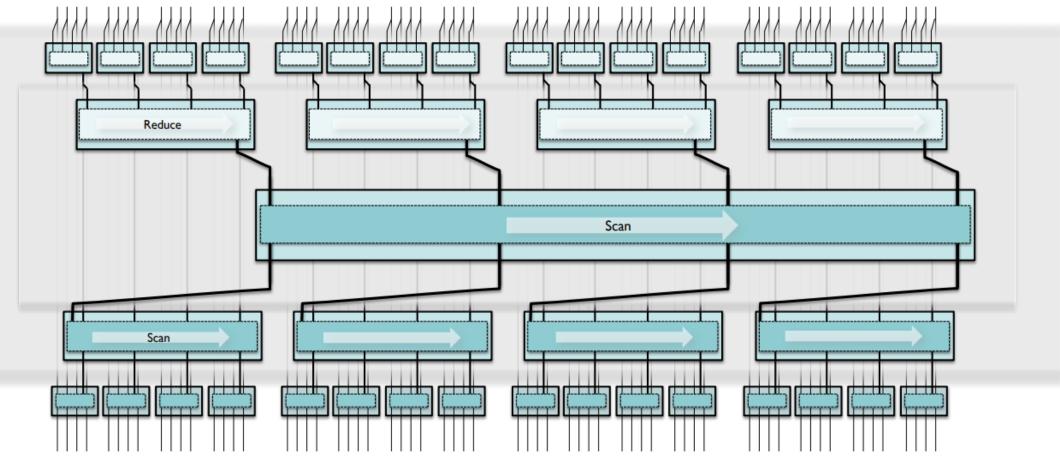
#### Scan across blocks

• Scan-and-add strategy (4n read/write)



#### Scan across blocks

• Reduce-then-scan strategy (3n read/write)



# Questions?!

#### **Credits:**

#### This lecture is primarily derived from:

- John Owens's course on Modern Parallel Computing (EEC 289Q, UC Davis, Winter 2018)
- Kayvon Fatahalian's course on Parallel Computing (CS149, Stanford, Fall 2023)
- Programming Massively Parallel Processors A Hands-on Approach book, 4<sup>th</sup> edition by Wen-mei W. Hwu, David B. Kirk, and Izzat El Hajj, 2023

#### Scan across blocks

• Decoupled look-back

