# Parallel reduction

- Code example: vector_reduction_small

$$S = \sum_n a_n$$

A: |_____|
   0                      n-1

**V1** A:

$s_0$ $s_1$ $s_2$ $s_3$ $s_4$ $s_5$ $s_6$ $s_7$
$d_0$ $d_1$ $d_2$ $d_3$ $d_4$ $d_5$ $d_6$ $d_7$

loaded elements from global to shared memory

sync

```
stride = 1
if (tid/2 * stride == 0)
    A[tid] += A[tid+stride]
```

```
stride = 2 * stride = 2
if (tid/2 * stride == 0)
    A[tid] += A[tid+stride]
```

sync

```
stride = 2 * stride = 4
if (tid/2 * stride == 0)
    A[tid] += A[tid+stride]
```

- Tree style reduction
- Computational complexity: $O(\log n)$

A[0]

Performance issue:
- Branch divergence within warp

---

**V2  Tree style reduction**

$s_0$ $s_1$ $s_2$ $s_3$ $s_4$ $s_5$ $s_6$ $s_7$
$d_0$ $d_1$ $d_2$ $d_3$ $d_4$ $d_5$ $d_6$ $d_7$

```
stride = 8/2 = 4
if (tid < stride)
    A[tid] += A[tid+stride]
```

sync

```
stride = stride/2 = 4/2 = 2
if (tid < stride)
    A[tid] += A[tid+stride]
```

sync

Computational complexity: $O(\log n)$

```
stride = stride/2 = 2/2 = 1
if (tid < stride)
    A[tid] += A[tid+stride]
```

A[0]

Version V2 shows

reduced branch divergence:



No branch divergence till we hit one warp (32 elements)

Branch diverge ← 32 threads

Reducing large arrays:

A:

Shared memory

TB0   TB1   TBk-1

Stride = blockDim.x * gridDim.x

n-1

0

ΣS
global memory.

Shared variable must be protected via a mutex or use atomic operations supported by NVIDIA's ISA.

· # of elements must be power of 2 (since it is tree style)

· What if not power of 2?
  · Pad out shared memory area with zeros to nearest power of 2

① Each thread calculates a partial sum and writes result to shared memory

② Reduce partial sums to a single value within each thread block

③ Designate one thread within a thread block to accumulate the reduced value to a shared variable living in GPU global memory

```
if (threadIdx.x == 0) {        (or)   if (threadIdx.x == 0)
                                            atomicAdd(S)
    atomicCAS(mutex);
    // accumulate into S
    mutex = 0;
}
```

# Parallel scan

- Code example: scan

$\oplus$ : scan operator
ex: addition

$I$ : identity element

## Scan:

input : $[a_0 \ a_1 \ \dots \ a_{n-1}]$

## Inclusive Scan:

output : $[a_0, \ a_0 \oplus a_1 \ \ a_0 \oplus a_1 \oplus a_2, \ \dots \ a_0 \oplus a_1 \oplus \cdots \oplus a_{n-1}]$

## Exclusive scan:

output : $[I \ \ a_0 \ \ a_0 \oplus a_1 \ \ a_0 \oplus a_1 \oplus a_2 \dots a_0 \oplus a_1 \oplus \cdots \oplus a_{n-2}]$

## Example:   Scan operator = +

Input :  $[3 \ 1 \ 7 \ 0 \ 4 \ 1 \ 6 \ 3]$

IScan :  $[3 \ 4 \ 11 \ 11 \ 15 \ 16 \ 22 \ 25]$

EScan :  $[0 \ 3 \ 4 \ 11 \ 11 \ 15 \ 16 \ 22]$
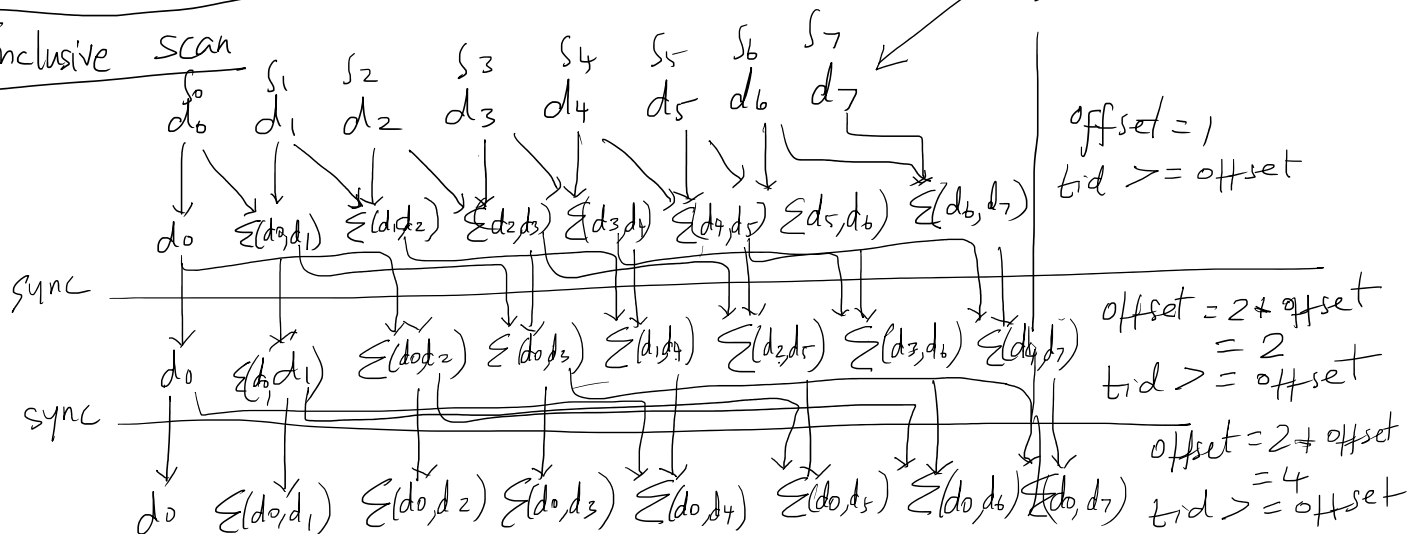
## Serial code:

EScan :

```
scan (input, output, n) {
    output[0] = 0;
    for (i=1; i<n; i++)
        output[i] = output[i-1] + input[i];
}
```

## Parallel prefix scan:

### Inclusive scan

values are in shared memory



offset = 1
tid >= offset

offset = 2 + offset
= 2
tid >= offset

offset = 2 + offset
= 4
tid >= offset

Use ping-pong buffers to store values between steps to prevent read after write hazards

# Recap: atomics on the GPU

Thursday, February 25, 2021     8:14 AM

- Use *atomicCAS()* to achieve mutual exclusion to a critical section  of code
    - Do not use *atomicCAS()* as a synchronization mechanism between threads belonging to the same thread warp --> may lead to deadlocks
    - Use *atomicCAS()* as synchronization mechanism between thread blocks
        - Ideally: designated one thread in each thread block, say threadidx.x = 0, to be the writer to the critical section:

            ```
            If (threadIdx.x == 0) {
                    Lock()
                    /* Critical section code */
                    Release()
            }
            ```
- For operations such as addition, subtraction, etc., it is preferrable to use the atomic variants of these operations for better performance:
    - *atomicAdd(), atomicSub(), …*
    - See: https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#atomic-functions

Scan : +

Input : $\begin{bmatrix} 3 & 1 & 7 & 0 & 4 & 1 & 6 & 3 \end{bmatrix}$

iScan : $\begin{bmatrix} 3 & 4 & 11 & 11 & 15 & 16 & 22 & 25 \end{bmatrix}$

eScan : $\begin{bmatrix} 0 & 3 & 4 & 11 & 11 & 15 & 16 & 22 \end{bmatrix}$
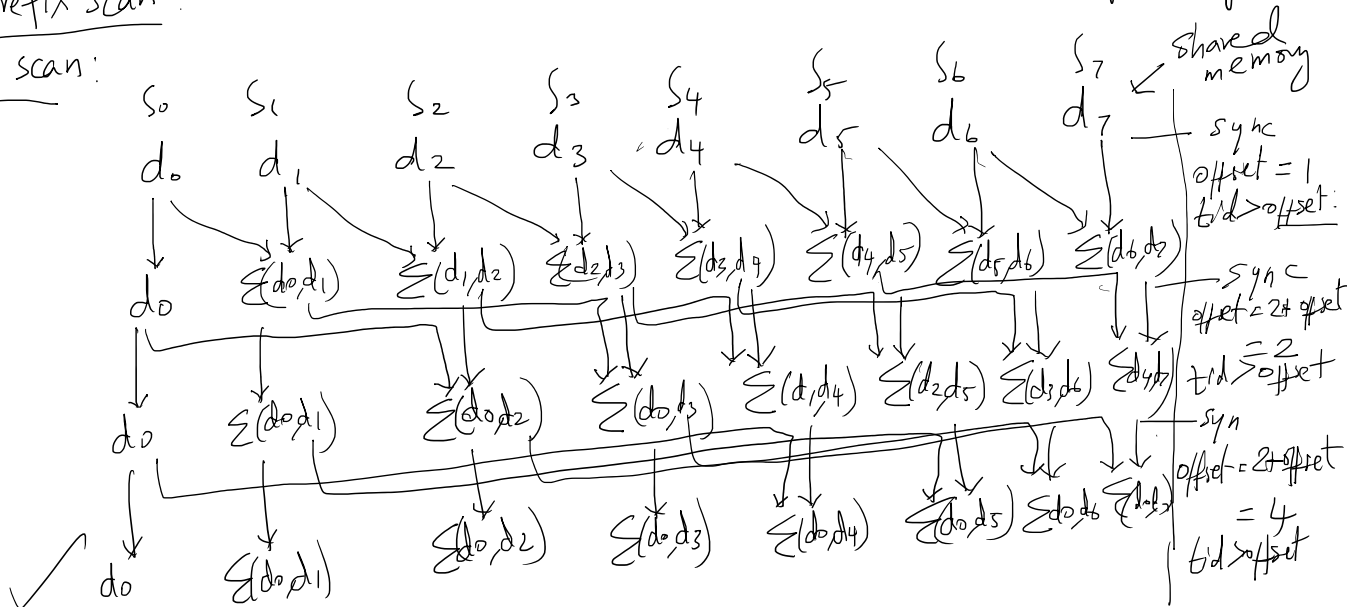
eScan:

output[0] = 0

for (i = 1; i < n; i++)

   output[i] = output[i-1] + input[i-1];
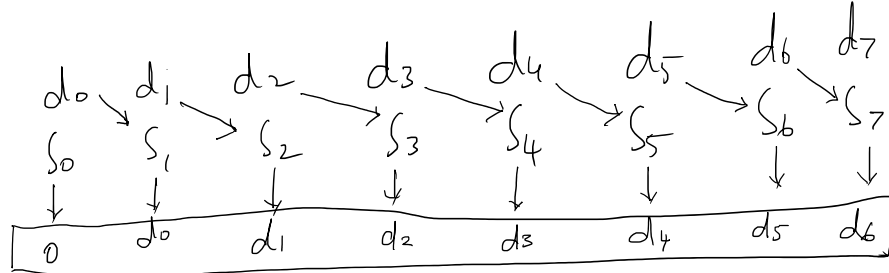
   ↑ loop carried dependency

## Parallel prefix scan

### Inclusive scan:
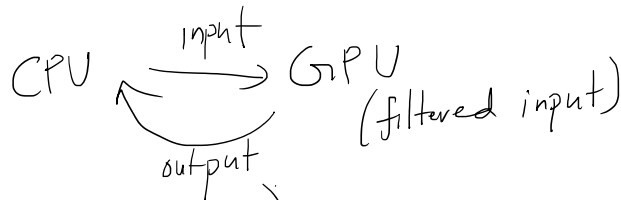


$S_0$  $S_1$  $S_2$  $S_3$  $S_4$  $S_5$  $S_6$  $S_7$ ← shared memory

$d_0$  $d_1$  $d_2$  $d_3$  $d_4$  $d_5$  $d_6$  $d_7$ — sync
offset = 1
tid > offset

$d_0$  $\Sigma(d_0,d_1)$  $\Sigma(d_1,d_2)$  $\Sigma(d_2,d_3)$  $\Sigma(d_3,d_4)$  $\Sigma(d_4,d_5)$  $\Sigma(d_5,d_6)$  $\Sigma(d_6,d_7)$ — sync
offset = 2+offset = 2
tid ≥ offset

$d_0$  $\Sigma(d_0,d_1)$  $\Sigma(d_0,d_2)$  $\Sigma(d_0,d_3)$  $\Sigma(d_1,d_4)$  $\Sigma(d_2,d_5)$  $\Sigma(d_3,d_6)$  $\Sigma(d_4,d_7)$ — syn
offset = 2+offset = 4
tid > offset

$d_0$  $\Sigma(d_0,d_1)$  $\Sigma(d_0,d_2)$  $\Sigma(d_0,d_3)$  $\Sigma(d_0,d_4)$  $\Sigma(d_0,d_5)$  $\Sigma(d_0,d_6)$  $\Sigma(d_0,d_7)$  ✓

### Exclusive Scan:

· Use same technique as above

· Threads load data into shared memory as ——



$d_0$  $d_1$  $d_2$  $d_3$  $d_4$  $d_5$  $d_6$  $d_7$

$S_0$  $S_1$  $S_2$  $S_3$  $S_4$  $S_5$  $S_6$  $S_7$

| 0 | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|---|---|---|---|---|---|---|---|

# Stream compaction

CPU $\xrightarrow{\text{input}}$ GPU (filtered input)

output

filter:
only +ve values
$> 0$
$\leq 0$ ✗

$$\text{size (output)} \leq \text{size (input)}$$

Input :    1  -1   3   4   -6   5   -8   10
Output :   1   3   4   5   10

|  | $S_0$ | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ | $S_6$ | $S_7$ |
|---|---|---|---|---|---|---|---|---|
| input | ①  | -1  | ③  | ④  | -6  | ⑤  | -8  | ⑩  |
| flag | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| escan(flag) | 0 | 1 | 1 | 2 | 3 | 3 | 4 | 4 |

↖ locations in the output array for each thread.

Output : | 1  3  4  5  10 |
          0  1  2  3  4

#elements = 5

Counting sort :     $D : \{d_0, d_1, \dots d_{n-1}\}$     $d_i \in$ range

↑ integers          ↑ for example, $[0, 255]$

- Non-comparison based
- Complexity : $O(n)$

input :     8   5   1   3   7   8   6   5   3   8

bin :

| | $S_0$ | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ | $S_6$ | $S_7$ | $S_8$ | create histogram |
|---|---|---|---|---|---|---|---|---|---|---|
| | bin0 | bin1 | bin2 | bin3 | bin4 | bin5 | bin6 | bin7 | bin8 | |
| | 0 | 1 | 0 | 2 | 0 | 2 | 1 | 1 | 3 | |

escan(bin)     0   0   1   1   3   3   5   6   7  ← generate starting address for each thread in sorted array

output

| 1 | 3 | 3 | 5 | 5 | 6 | 7 | 8 | 8 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

# Histogram generation

$a_i$: integer (in some range)

- Code example: histogram_generation

A: | $TB_0$ | $TB_1$ | $\cdots$ | $TB_{k-1}$ | |

freq: | Local hist for $TB_0$ | Local hist for $TB_1$ | local hist for $TB_{k-1}$ | } stored in shared memory

bins          bins          bins

freq | Shared histogram between thread blocks | } stored in GPU global memory

bins

① Thread in thread block strides across elements in A and bins each in a local histogram shared by threads in that thread block

② Use atomics when incrementing bins

③ Selected threads within each thread block accumulate local histogram bin values into the shared histogram in GPU global memory

④ Use atomics when accumulating bin values in the shared histogram