# Introduction to Parallel and Distributed Processing

## CUDA Reduce and Prefix

Jaroslaw 'Jaric' Zola

http://www.jzola.org/
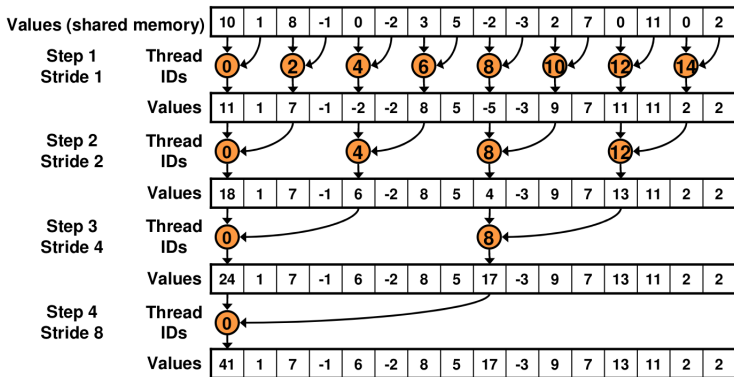
# Classic Problem

- Given vector $X = (x_1, \ldots, x_n)$ compute $\displaystyle\sum_{i=1}^{n} x_i$

- Smells like a simple 1D problem...

# Direct Approach

```
1   __global__ void reduce0(int* gin, int* gout) {
2     extern __shared__ int sdata[];
3
4     int tid = threadIdx.x;
5     int i = blockIdx.x * blockDim.x + tid;
6
7     sdata[tid] = gin[i];
8     __syncthreads();
9
10    for (int s = 1; s < blockDim.x; s *= 2) {
11      if (tid % (2 * s) == 0) sdata[tid] += sdata[tid + s];
12      __syncthreads();
13    }
14
15    if (tid == 0) gout[blockIdx.x] = sdata[0];
16  } // reduce0
```
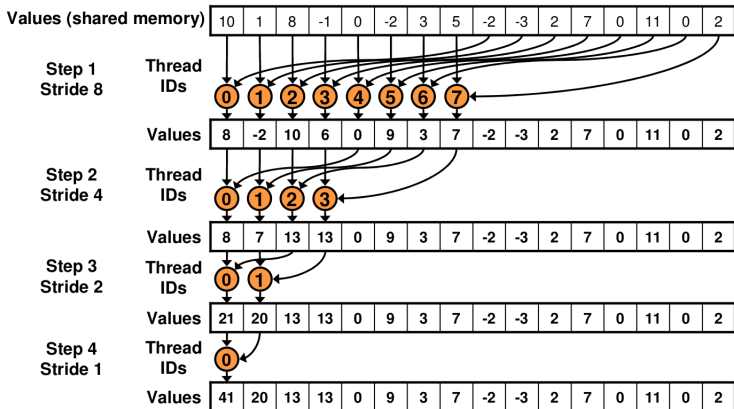
# Direct Approach – Problem

Divergent branching:

# Better Approach

Sequential addressing:

| Values (shared memory) | 10 | 1 | 8 | -1 | 0 | -2 | 3 | 5 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Step 1 / Stride 8** — Thread IDs: 0 1 2 3 4 5 6 7

| Values | 8 | -2 | 10 | 6 | 0 | 9 | 3 | 7 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Step 2 / Stride 4** — Thread IDs: 0 1 2 3

| Values | 8 | 7 | 13 | 13 | 0 | 9 | 3 | 7 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Step 3 / Stride 2** — Thread IDs: 0 1

| Values | 21 | 20 | 13 | 13 | 0 | 9 | 3 | 7 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Step 4 / Stride 1** — Thread IDs: 0

| Values | 41 | 20 | 13 | 13 | 0 | 9 | 3 | 7 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

## Sequential Addressing

```
1  __global__ void reduce1(int* gin, int* gout) {
2    extern __shared__ int sdata[];
3
4    int tid = threadIdx.x;
5    int i = blockIdx.x * blockDim.x + tid;
6
7    sdata[tid] = gin[i];
8    __syncthreads();
9
10   for (int s = blockDim.x >> 1; s > 0; s >>= 1) {
11     if (tid < s) sdata[tid] += sdata[tid + s];
12     __syncthreads();
13   }
14
15   if (tid == 0) gout[blockIdx.x] = sdata[0];
16 } // reduce1
```

# Even Better Sequential Addressing

```
1  // we halve the number of blocks
2  __global__ void reduce2(int* gin, int* gout) {
3    extern __shared__ int sdata[];
4
5    int tid = threadIdx.x;
6    int i = blockIdx.x * (blockDim.x * 2) + threadIdx.x;
7
8    sdata[tid] = gin[i] + gin[i + blockDim.x];
9    __syncthreads();
10
11   for (int s = blockDim.x >> 1; s > 0; s >>= 1) {
12     if (tid < s) sdata[tid] += sdata[tid + s];
13     __syncthreads();
14   }
15
16   if (tid == 0) gout[blockIdx.x] = sdata[0];
17 } // reduce2
```
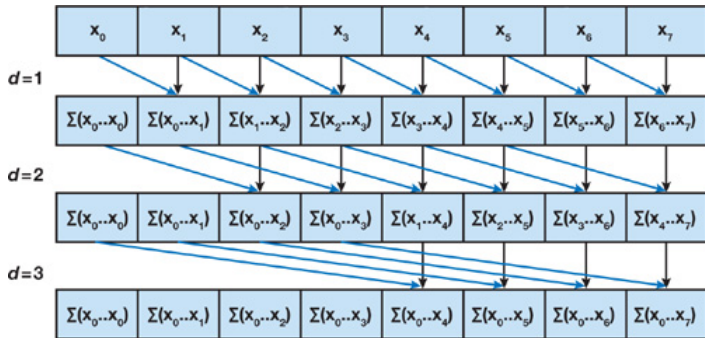
# What About Parallel Prefix?

- Again, 1D problem PREFIX_SUM
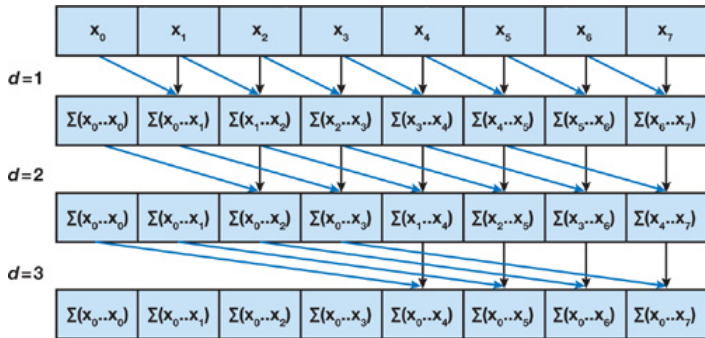
**Input:** $[x_0, x_1, \ldots, x_{n-1}]$
1: **for** $d = 1 \ldots \log(n)$ **do**
2:   **for** $j = 1 \ldots n$ **pardo**
3:     **if** $j \geq 2^d$ **then**
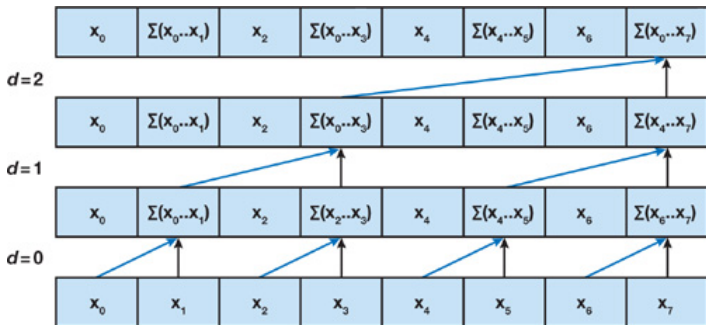4:       $x[j] = x[j - 2^{d-1}] + x[j]$

# What About Parallel Prefix?

# What About Parallel Prefix?



$d=1$

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

$\Sigma(x_0..x_0)$ $\Sigma(x_0..x_1)$ $\Sigma(x_1..x_2)$ $\Sigma(x_2..x_3)$ $\Sigma(x_3..x_4)$ $\Sigma(x_4..x_5)$ $\Sigma(x_5..x_6)$ $\Sigma(x_6..x_7)$

$d=2$

$\Sigma(x_0..x_0)$ $\Sigma(x_0..x_1)$ $\Sigma(x_0..x_2)$ $\Sigma(x_0..x_3)$ $\Sigma(x_1..x_4)$ $\Sigma(x_2..x_5)$ $\Sigma(x_3..x_6)$ $\Sigma(x_4..x_7)$

$d=3$

$\Sigma(x_0..x_0)$ $\Sigma(x_0..x_1)$ $\Sigma(x_0..x_2)$ $\Sigma(x_0..x_3)$ $\Sigma(x_0..x_4)$ $\Sigma(x_0..x_5)$ $\Sigma(x_0..x_6)$ $\Sigma(x_0..x_7)$
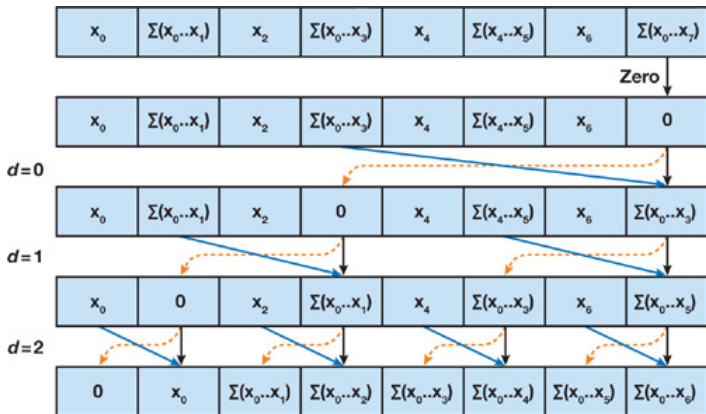
Problem: total work $O(n \log(n))$

# Going Work Efficient

- We resolve to our old friend :-)
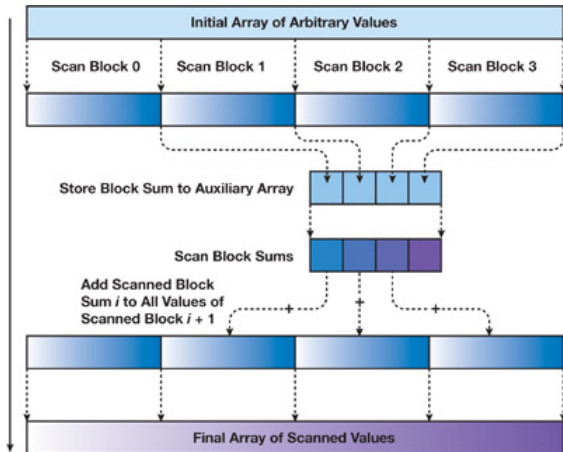- Decompose problem into two phases: up-sweep (reduction) and down-sweep (update)

# Going Work Efficient

- We resolve to our old friend :-)
- Decompose problem into two phases: up-sweep (reduction) and down-sweep (update)

# Working on Large Arrays

# For Fun

- Modify the prefix code to work with large arrays.