# GPU Computing

# 规约算法

Hui Liu

Email: hui.sc.liu@gmail.com

# Reduction Steps

- Replace the divergent branching code:
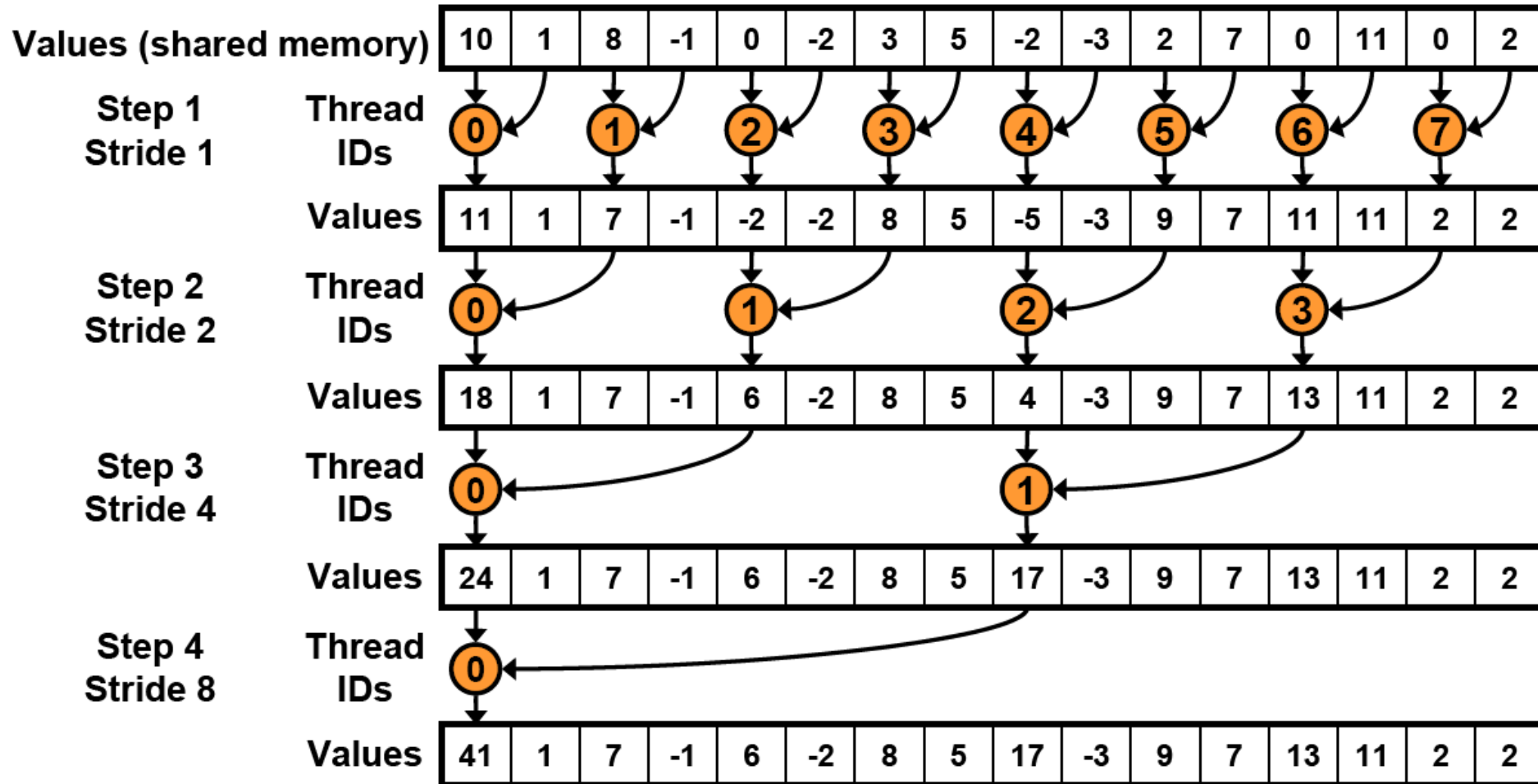
```
// do reduction in shared mem
for (unsigned int s=1; s < blockDim.x; s *= 2) {
        if (tid % (2*s) == 0) {
                    sdata[tid] += sdata[tid + s];
        }
        __syncthreads();
}
```

- With strided index and non-divergent branch
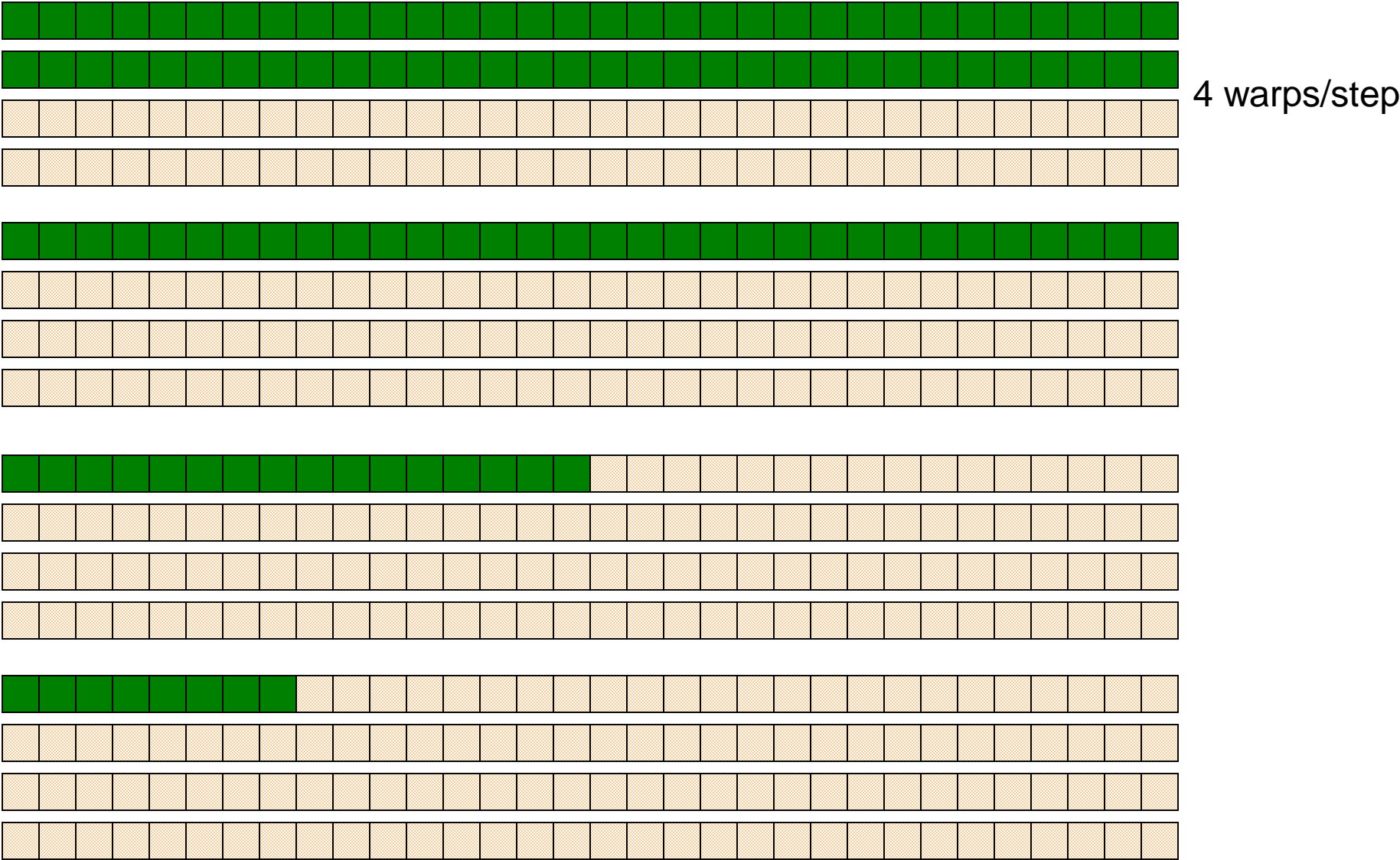
```
// do reduction in shared mem
for (unsigned int s=1; s < blockDim.x; s *= 2) {
        int index  = 2 * s * tid;

        if (index < blockDim.x / s) {
                    sdata[index] += sdata[index + s];
        }
        __syncthreads();
}
```
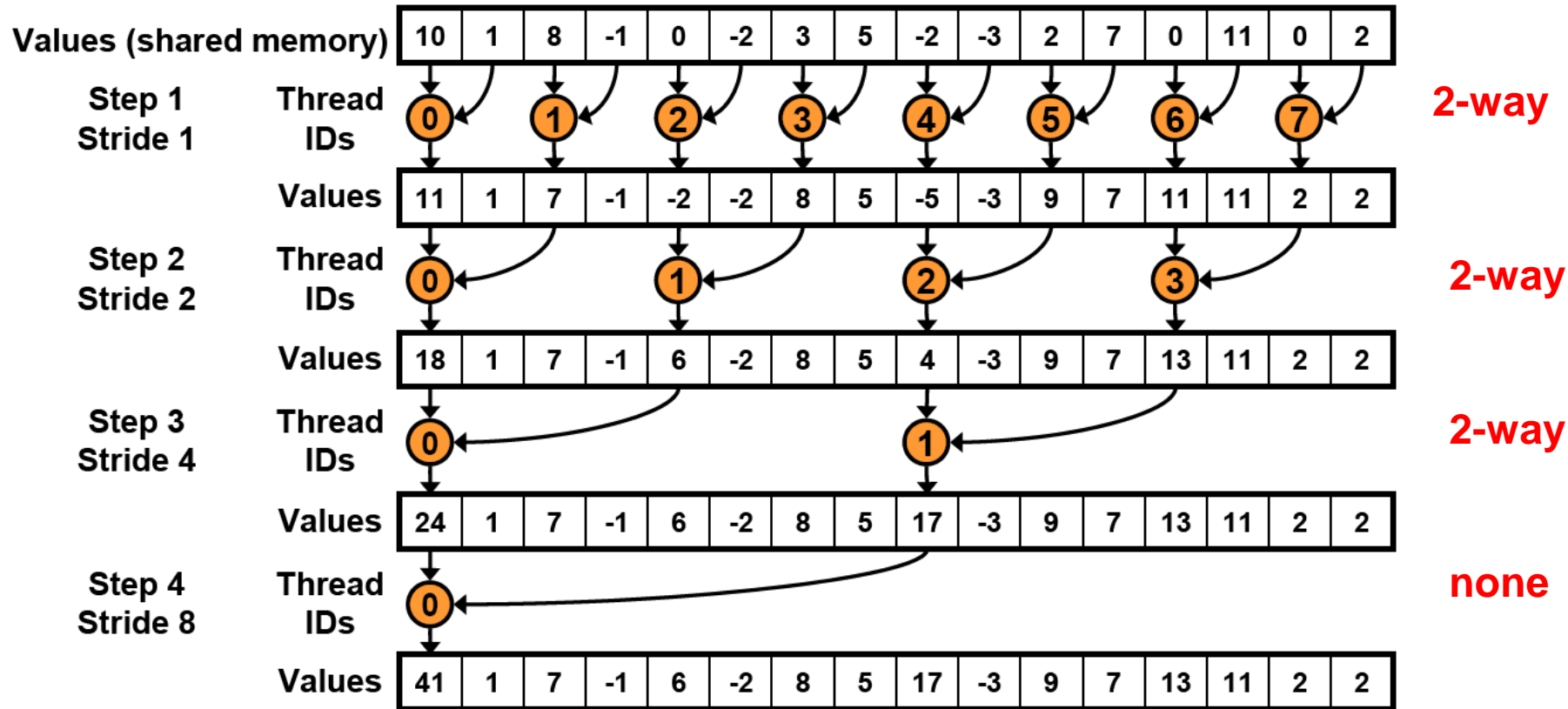
4 warps/step

# Performance for 4M element reduction

| | Time ($2^{22}$ ints) | Bandwidth | Step Speedup | Cumulative Speedup |
|---|---|---|---|---|
| **Kernel 1:** interleaved addressing with divergent branching | **8.054 ms** | **2.083 GB/s** | | |
| **Kernel 2:** interleaved addressing non-divergent branching | **3.456 ms** | **4.854 GB/s** | **2.33x** | **2.33x** |

**2-way bank conflicts at every step**
**Recall there are more than 16 threads**
**To see the conflicts see what happens with 128 threads**

# THANK YOU