# 规约算法

Hui Liu

Email: hui.sc.liu@gmail.com

- Get maximum GPU performance



- Two components:
  - Compute Bandwidth: GFLOPs
  - Memory Bandwidth: GB/s

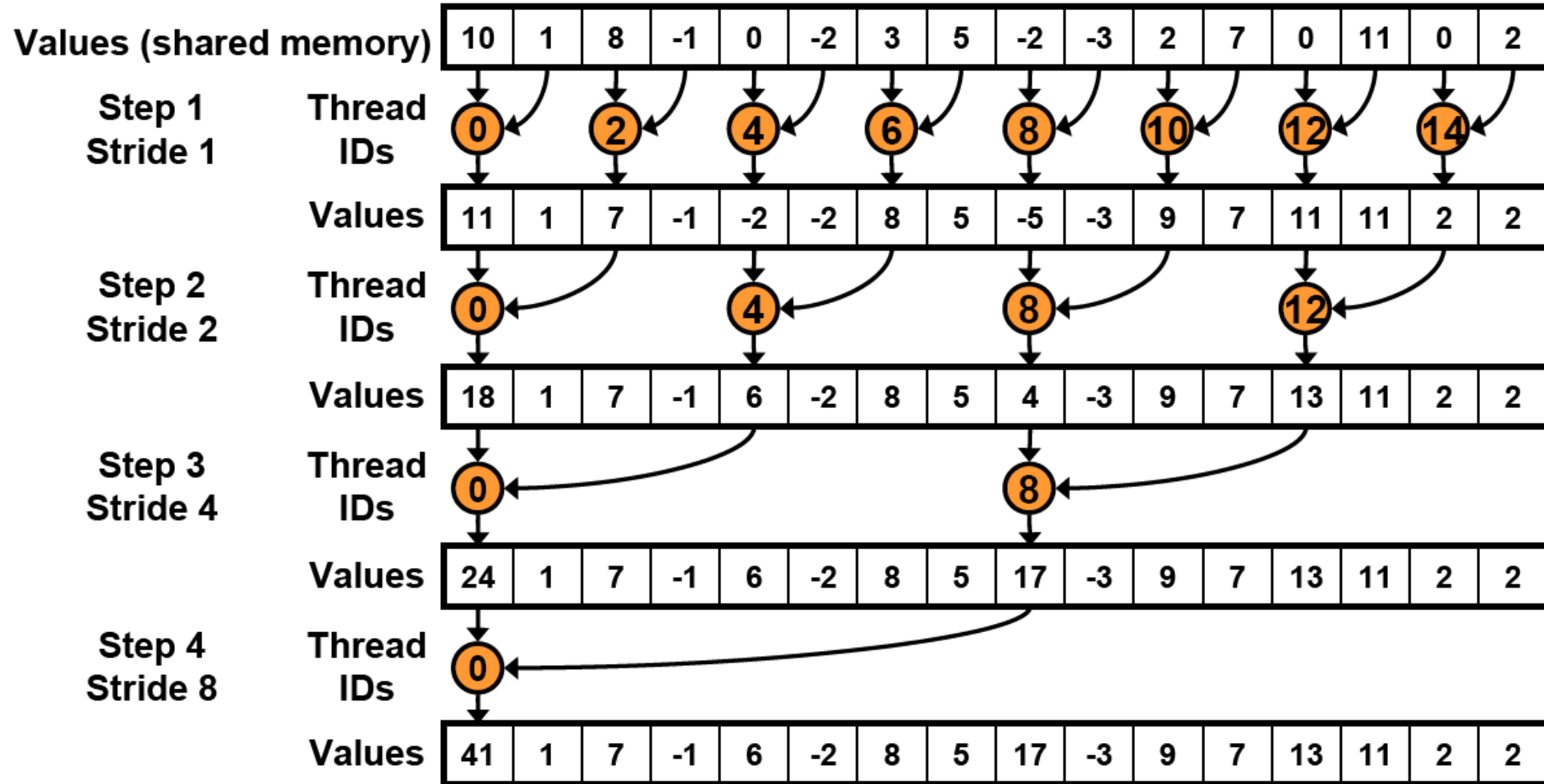- Reductions typically have low arithmetic intensity
  - FLOPs/element loaded from memory
- So, bandwidth will be the limiter

- For the **ASUS ENGTX280 OC**
  - 512-bit interface, 1.14GHz DDR3
  - 512 / 8 x 1.14 x 2 = **145.92 GB/s**

- Load data:
  - Each thread loads one element from global memory to shared memory
- Actual Reduction: Proceed in log*N* steps
  - A thread reduces two elements
    - The first two elements by the first thread
    - The next two by the next thread
    - And so, on
  - At the end of each step:
    - Deactivate half of the threads
  - Terminate: when one thread left
- Write back to global memory

```
__global__ void reduce0(int *g_idata, int *g_odata) {
  extern __shared__ int sdata[];

  // each thread loads one element from global to shared mem
  unsigned int tid = threadIdx.x;
  unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
  sdata[tid] = g_idata[i];
  __syncthreads();

  // do reduction in shared mem
  for (unsigned int s=1; s < blockDim.x; s *= 2) { // step = s x 2
      if (tid % (2*s) == 0) { // only threadIDs divisible by the step participate
              sdata[tid] += sdata[tid + s];
      }
      __syncthreads();
  }

  // write result for this block to global mem
  if (tid == 0) g_odata[blockIdx.x] = sdata[0];
}
```

| | Time ($2^{22}$ ints) | Bandwidth |
|---|---|---|
| **Kernel 1:** interleaved addressing with divergent branching | **8.054 ms** | **2.083 GB/s** |

Note: Block size = 128 for all experiments

**Caveat:** results are for a G80 processor

**Bandwidth calculation:**

Each block processes 128 elements and does:
    128 reads & 1 write
    We only care about global memory

**At each kernel/step:**
    N (element reads) + N / 128 (element writes)
Every kernel/step reduces input size by 128x
    next set N = N / 128
So for:
    N = 4194304
    **Accesses = 4227394**
Each access is four bytes

```
__global__ void reduce0 (int *g_idata, int *g_odata) {
    extern __shared__ int sdata[];

    // each thread loads one element from global to shared mem
    unsigned int tid = threadIdx.x;
    unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
    sdata[tid] = g_idata[i];
    __syncthreads();

    // 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();
    }

    // write result for this block to global mem
    if (tid == 0) g_odata[blockIdx.x] = sdata[0];
}
```
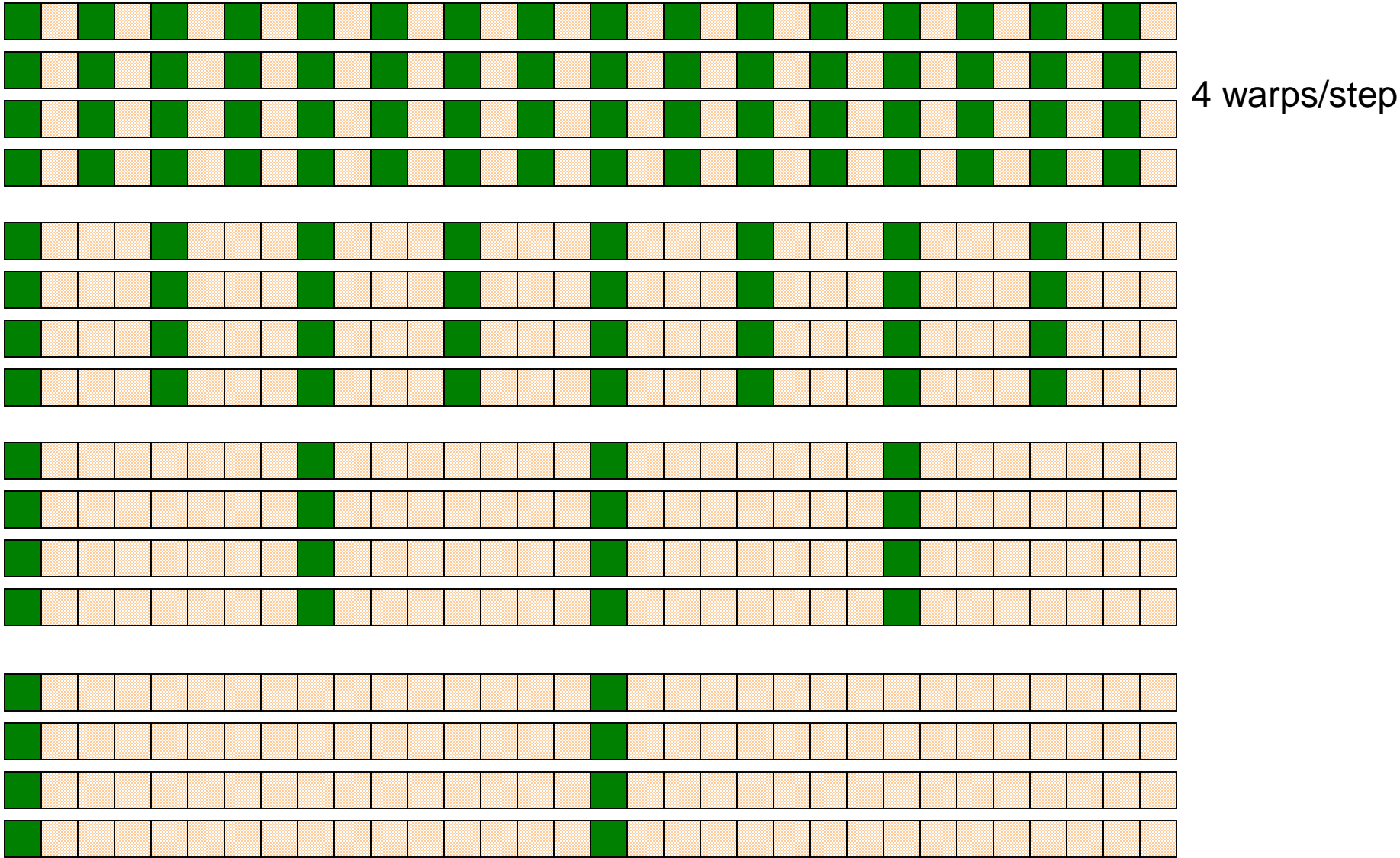
**Highly divergent code leads to very poor performance**

# Divergent Branching: Warp Control Flow



4 warps/step

# THANK YOU