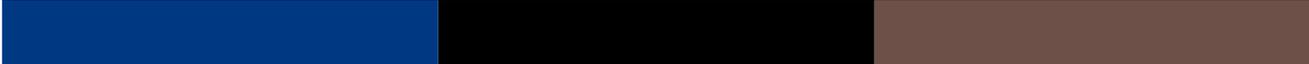




Institut  
Mines-Télécom

# Parallel Reduction in CUDA

A decorative horizontal bar at the bottom of the slide, composed of three segments: a blue segment on the left, a black segment in the middle, and a brown segment on the right.

Elisabeth Brunet  
Based on Mark Harris (Nvidia) talk



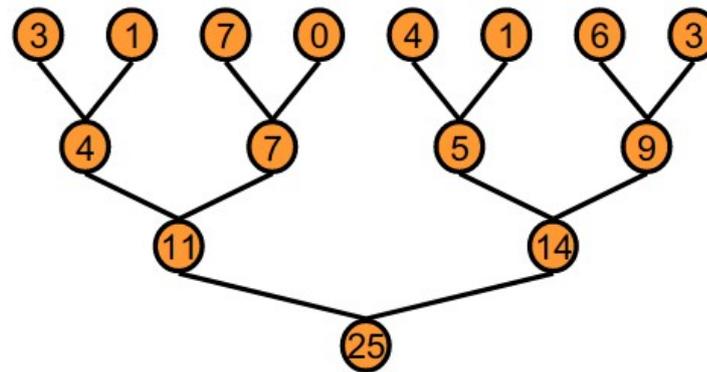
# Reduction : Naive algorithms

- Reduction by an only one thread
- An atomic operation in global memory
  - Huge synchronization

# Parallel Reduction



- **Tree-based approach used within each thread block**



- **Need to be able to use multiple thread blocks**
  - To process very large arrays
  - To keep all multiprocessors on the GPU busy
  - Each thread block reduces a portion of the array
- **But how do we communicate partial results between thread blocks?**



Problem :  
how sharing intermediate results

# Problem: Global Synchronization

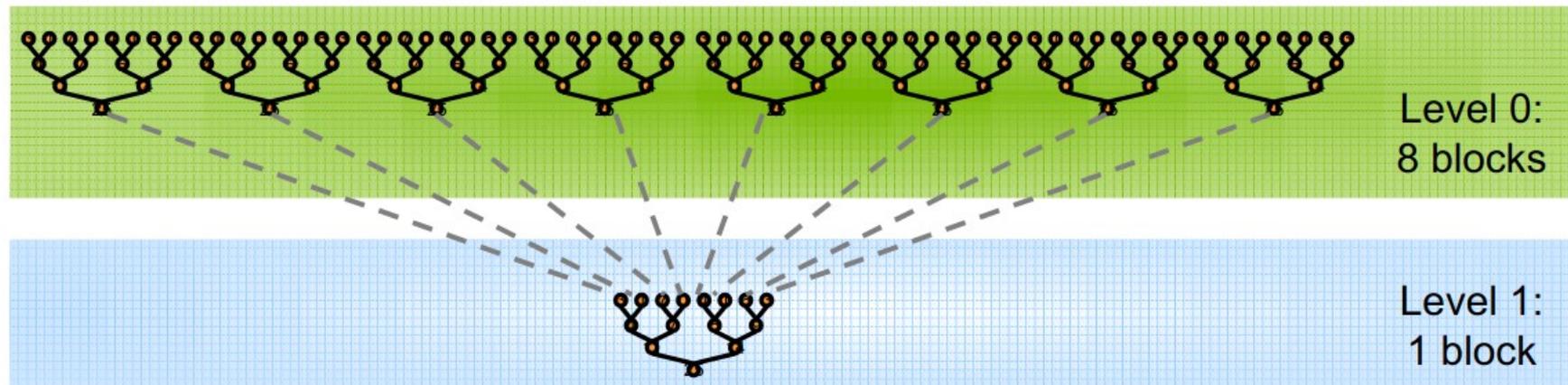


- If we could synchronize across all thread blocks, could easily reduce very large arrays, right?
  - Global sync after each block produces its result
  - Once all blocks reach sync, continue recursively
- But CUDA has no global synchronization. Why?
  - Expensive to build in hardware for GPUs with high processor count
  - Would force programmer to run fewer blocks (no more than  $\#$  multiprocessors \*  $\#$  resident blocks / multiprocessor) to avoid deadlock, which may reduce overall efficiency
- Solution: decompose into multiple kernels
  - Kernel launch serves as a global synchronization point
  - Kernel launch has negligible HW overhead, low SW overhead

# Solution: Kernel Decomposition



- Avoid global sync by decomposing computation into multiple kernel invocations



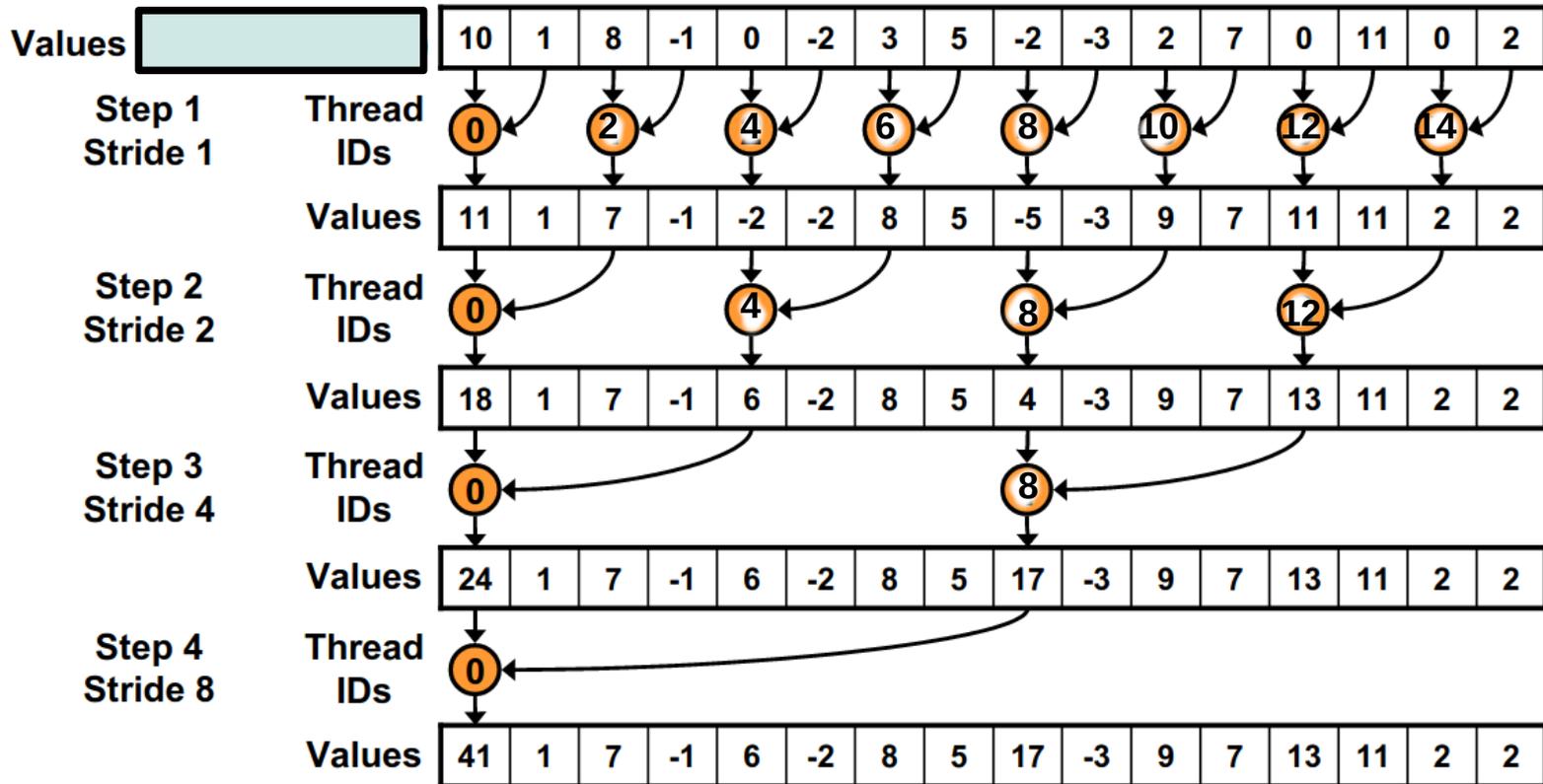
- In the case of reductions, code for all levels is the same
  - Recursive kernel invocation



Problem : divergency



# Parallel Reduction: Interleaved Addressing



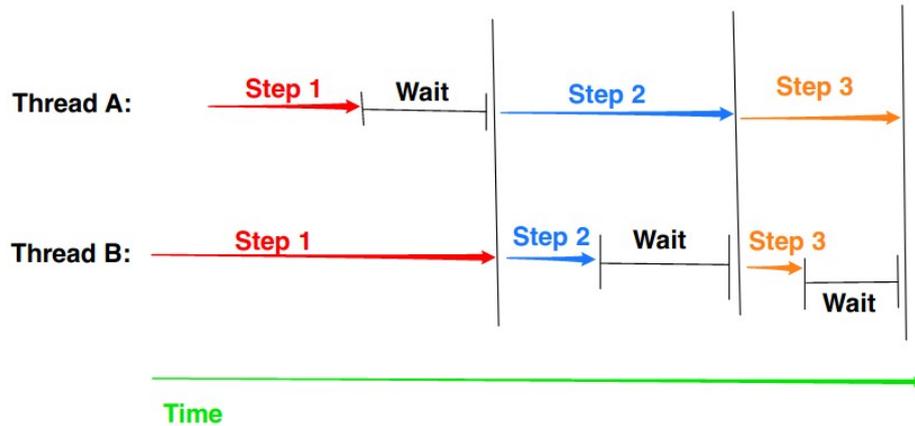
```

for (unsigned int s=1; s < blockDim.x; s *= 2) {
    if (tid % (2*s) == 0) {
        sdata[tid] += sdata[tid + s];
    }
    __syncthreads();
}
    
```

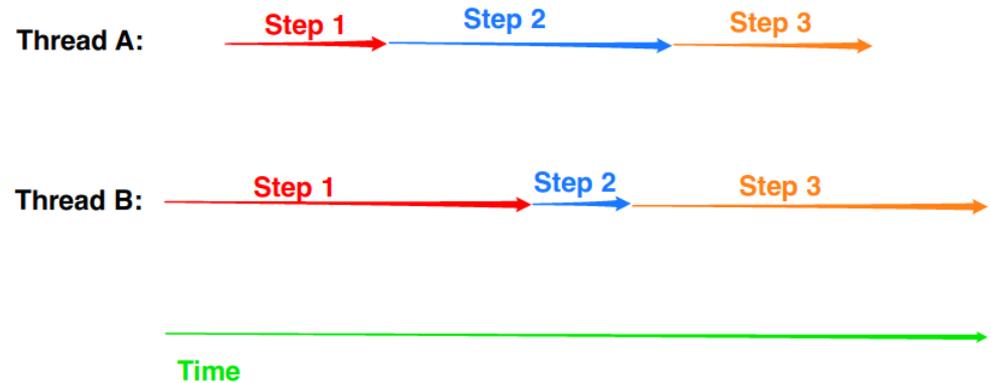
**Problem: highly divergent branching results in very poor performance!**

# Divergency

## Threads in the same warp:



## Threads in different warps:



# Let all work !

Just replace divergent branch in inner loop:

```
for (unsigned int s=1; s < blockDim.x; s *= 2) {  
    if (tid % (2*s) == 0) {  
        sdata[tid] += sdata[tid + s];  
    }  
    __syncthreads();  
}
```

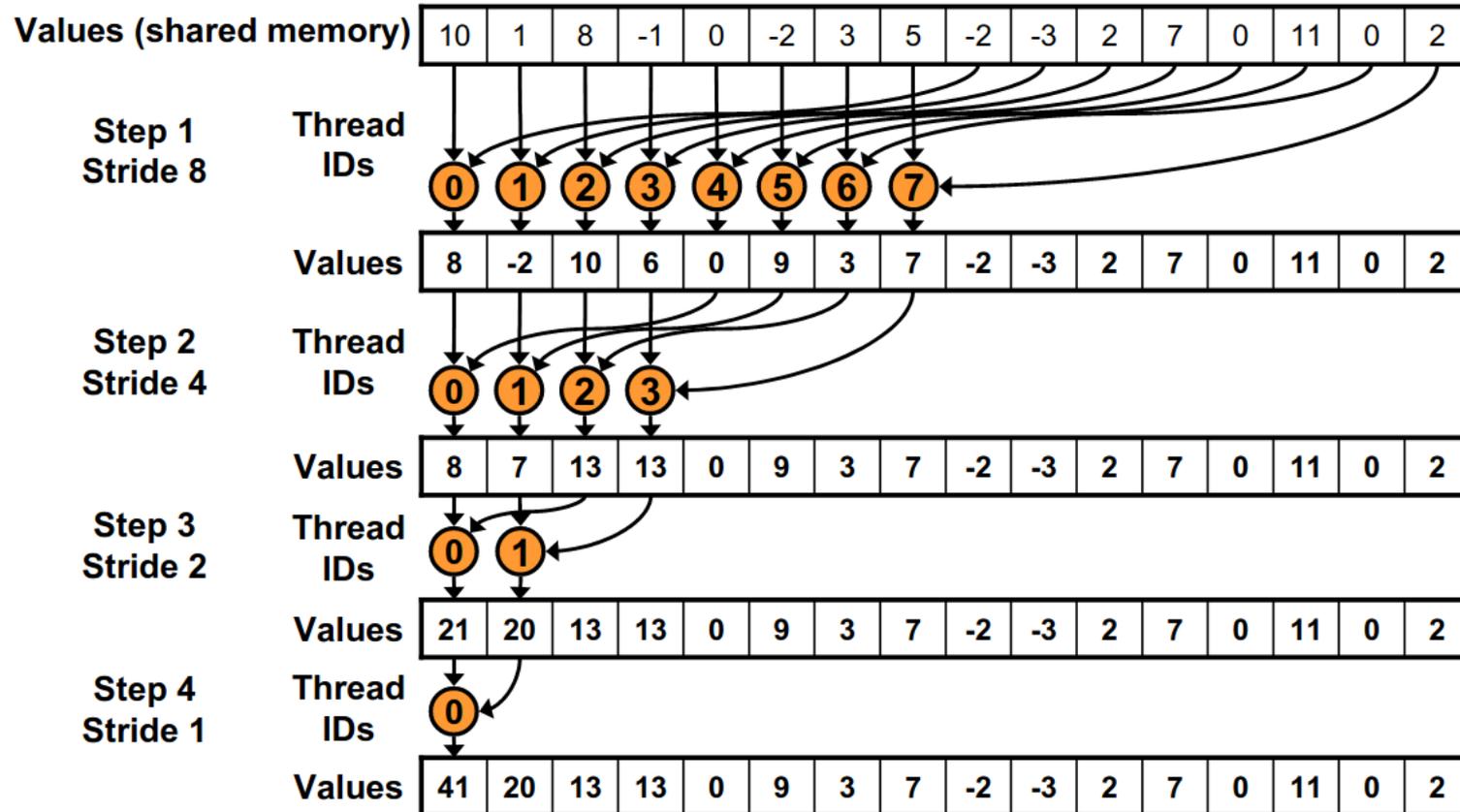
With strided index and non-divergent branch:

```
for (unsigned int s=1; s < blockDim.x; s *= 2) {  
    int index = 2 * s * tid;  
  
    if (index < blockDim.x) {  
        sdata[index] += sdata[index + s];  
    }  
    __syncthreads();  
}
```



17

# Parallel Reduction: Sequential Addressing





# Problem : Load unbalancing



# Idle Threads



**Problem:**

```
for (unsigned int s=blockDim.x/2; s>0; s>>=1) {  
    if (tid < s) {  
        sdata[tid] += sdata[tid + s];  
    }  
    __syncthreads();  
}
```

**Half of the threads are idle on first loop iteration!**

**This is wasteful...**



# To go further

- First add during global load
- Unroll last warp
- Completely unrolled
- Multiple adds per thread
  
- Use intermediate memory → next week.