# Parallel reduction

In this exercice a parallel reduction kernel will be implemented. Write a kernel performing the sum operation.

Starting point:

```
%% file parallel_reduction.cu

#include <stdio.h>

void cpu_sum(int *x, int n)
{
    int result = 0;
    for(unsigned int i=0; i < n; ++i) {
        result += x[i];
    }
    printf("CPU Sum is %d \n", result);
}

__global__ void gpu_sum(int *x)
{
    int tid = blockIdx.x * blockDim.x + threadIdx.x;

    // write your code here
    // tip: use `__syncthreads()` to synchronize the threads
}

int main()
{
    int h[] = {10, 1, 8, -1, 0, -2, 3, 5, -2, -3, 2, 7, 0, 11, 0, 2};

    int size = sizeof(h);
    int count = size/sizeof(int);

    int* d;
    cudaMalloc(&d, size);
    cudaMemcpy(d, h, size, cudaMemcpyHostToDevice);

    gpu_sum <<<1, count >>>(d);

    int result;
    cudaMemcpy(&result, d, sizeof(int), cudaMemcpyDeviceToHost);
    printf("GPU Sum is %d \n", result);

    //cpu_sum(h, count);
    cudaFree(d);
    return 0;
}
```
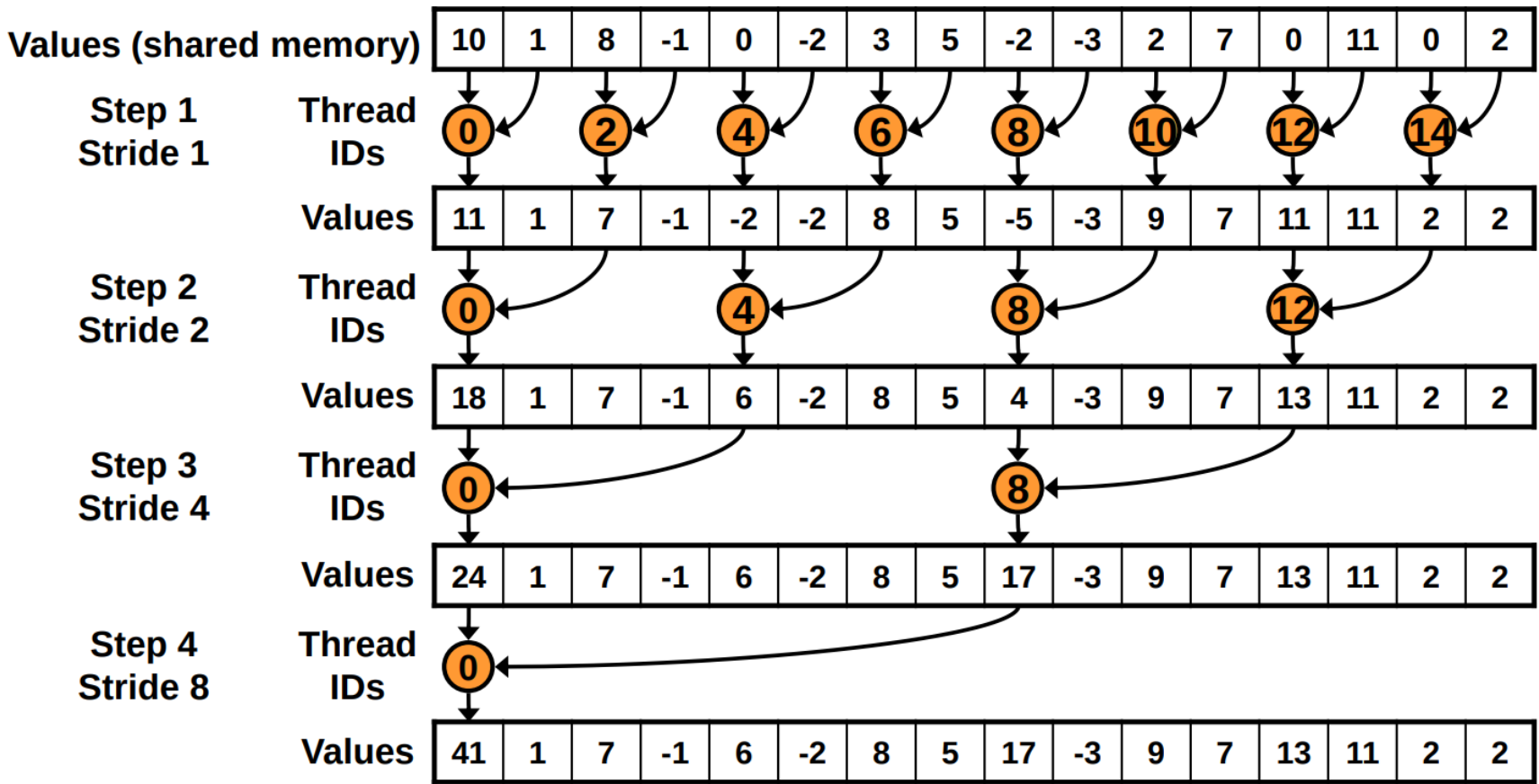
```
!nvidia-smi
```

```
%%bash

CUDA_SUFF=35
nvcc -gencode arch=compute_${CUDA_SUFF},code=sm_${CUDA_SUFF} ./parallel_reduction.cu -o parallel_reduction
./parallel_reduction
```

The algorithm can be implemented in two ways:

Naive memory access (interleaved addresing):

| Values (shared memory) | 10 | 1 | 8 | -1 | 0 | -2 | 3 | 5 | -2 | -3 | 2 | 7 | 0 | 11 | 0 | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Step 1 Stride 1 — Thread IDs | 0 | | 2 | | 4 | | 6 | | 8 | | 10 | | 12 | | 14 | |
| Values | 11 | 1 | 7 | -1 | -2 | -2 | 8 | 5 | -5 | -3 | 9 | 7 | 11 | 11 | 2 | 2 |
| Step 2 Stride 2 — Thread IDs | 0 | | | | 4 | | | | 8 | | | | 12 | | | |
| Values | 18 | 1 | 7 | -1 | 6 | -2 | 8 | 5 | 4 | -3 | 9 | 7 | 13 | 11 | 2 | 2 |
| Step 3 Stride 4 — Thread IDs | 0 | | | | | | | | 8 | | | | | | | |
| Values | 24 | 1 | 7 | -1 | 6 | -2 | 8 | 5 | 17 | -3 | 9 | 7 | 13 | 11 | 2 | 2 |
| Step 4 Stride 8 — Thread IDs | 0 | | | | | | | | | | | | | | | |
| Values | 41 | 1 | 7 | -1 | 6 | -2 | 8 | 5 | 17 | -3 | 9 | 7 | 13 | 11 | 2 | 2 |

Optimised memory access (sequantial addresing):

| 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 |