# Computational Science on Many-Core Architectures Exercise 5

## Example 1 Inclusive and Exclusive Scan (4 Points)

**a)**

```c
__global__ void scan_kernel_1(double const *X,
                              double *Y,
                              int N,
                              double *carries)
{
  __shared__ double shared_buffer[256];
  double my_value;

  unsigned int work_per_thread = (N - 1) / (gridDim.x * blockDim.x) + 1;
  unsigned int block_start = work_per_thread * blockDim.x *  blockIdx.x;
  unsigned int block_stop  = work_per_thread * blockDim.x * (blockIdx.x + 1);
  unsigned int block_offset = 0;

  // run scan on each section
  for (unsigned int i = block_start + threadIdx.x; i < block_stop; i += blockDim.x)
  {
    // load data:
    my_value = (i < N) ? X[i] : 0;

    // inclusive scan in shared buffer:
    for(unsigned int stride = 1; stride < blockDim.x; stride *= 2)
    {
      __syncthreads();
      shared_buffer[threadIdx.x] = my_value;
      __syncthreads();
      if (threadIdx.x >= stride)
        my_value += shared_buffer[threadIdx.x - stride];
    }
    __syncthreads();
    shared_buffer[threadIdx.x] = my_value;
    __syncthreads();

    // exclusive scan requires us to write a zero value at the beginning of each block
    my_value = (threadIdx.x > 0) ? shared_buffer[threadIdx.x - 1] : 0;

    // write to output array
    if (i < N)
      Y[i] = block_offset + my_value;

    block_offset += shared_buffer[blockDim.x-1];
  }

  // write carry:
  if (threadIdx.x == 0)
    carries[blockIdx.x] = block_offset;
```
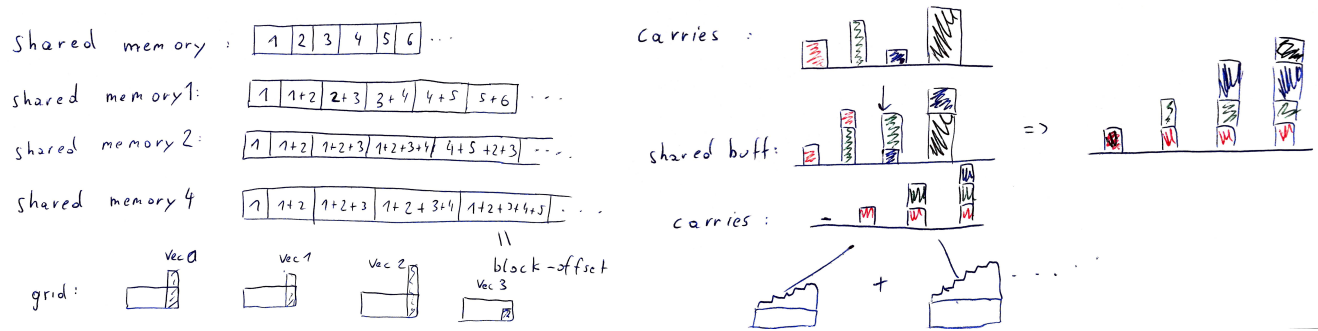
For the first kernel "kernel_1" and simulate it with 4 blocks and 6 threads per block.
Begin with line 24: the for loop iterates over the values X which belong to the block.
At the end of the for loop → write the temporary result of the scan into a vector Y and the offset stored in block_offset. At the end of the kernel every block contains it scanned value and the vector for the next step.

```c
// exclusive-scan of carries
__global__ void scan_kernel_2(double *carries)
{
  __shared__ double shared_buffer[256];

  // load data:
  double my_carry = carries[threadIdx.x];

  // exclusive scan in shared buffer:

  for(unsigned int stride = 1; stride < blockDim.x; stride *= 2)
  {
    __syncthreads();
    shared_buffer[threadIdx.x] = my_carry;
    __syncthreads();
    if (threadIdx.x >= stride)
      my_carry += shared_buffer[threadIdx.x - stride];
  }
  __syncthreads();
  shared_buffer[threadIdx.x] = my_carry;
  __syncthreads();

  // write to output array
  carries[threadIdx.x] = (threadIdx.x > 0) ? shared_buffer[threadIdx.x - 1] : 0;
}

__global__ void scan_kernel_3(double *Y, int N,
                              double const *carries)
{
  unsigned int work_per_thread = (N - 1) / (gridDim.x * blockDim.x) + 1;
  unsigned int block_start = work_per_thread * blockDim.x *  blockIdx.x;
  unsigned int block_stop  = work_per_thread * blockDim.x * (blockIdx.x + 1);

  __shared__ double shared_offset;

  if (threadIdx.x == 0)
    shared_offset = carries[blockIdx.x];

  __syncthreads();

  // add offset to each element in the block:
  for (unsigned int i = block_start + threadIdx.x; i < block_stop; i += blockDim.x)
    if (i < N)
      Y[i] += shared_offset;
}
```

**b)**

Listing 1: kernel for inclusive_scan)

```
1  __global__ void Inclusive(double *Y, int N, const double *X)
2  {
3      for (int i = blockDim.x * blockIdx.x + threadIdx.x; i < N-1;
4      i += gridDim.x * blockDim.x)
5      {
6          Y[i] = Y[i+1];
7      }
8      if (blockDim.x * blockIdx.x + threadIdx.x == 0)
9          // First step: Scan within each thread group and write carries
10     scan_kernel_1<<<num_blocks, threads_per_block>>>(input, output, N, carries);
11
12     // Second step: Compute offset for each thread group
13     (exclusive scan for each thread group)
14     scan_kernel_2<<<1, num_blocks>>>(carries);
15
16     // Third step: Offset each thread group accordingly
17     scan_kernel_3<<<num_blocks, threads_per_block>>>(output, N, carries);
18
19     // Make inclusive
20     makeInclusive<<<num_blocks, threads_per_block>>>(output, N, input);
21
22     cudaFree(carries);
23     {
24         Y[N-1] += X[N-1];
25     }
26 }
27
28 void exclusive_scan(double const * input, double* output, int N)
29 {
30     int num_blocks = 512;
31     int threads_per_block = 512;
32
33     double *carries;
34     cudaMalloc(&carries, sizeof(double) * num_blocks);
35
36     // First step: Scan within each thread group and write carries
```

```
37        scan_kernel_1<<<num_blocks, threads_per_block>>>(input, output, N, carries);
38
39     // Second step: Compute offset for each thread group
40     (exclusive scan for each thread group)
41     scan_kernel_2<<<1, num_blocks>>>(carries);
42
43     // Third step: Offset each thread group accordingly
44     scan_kernel_3<<<num_blocks, threads_per_block>>>(output, N, carries);
45
46     // Make inclusive
47     makeInclusive<<<num_blocks, threads_per_block>>>(output, N, input);
48
49     cudaFree(carries);
50 }
```

**c)**

Only nessesary to remove the folloing code snipped:
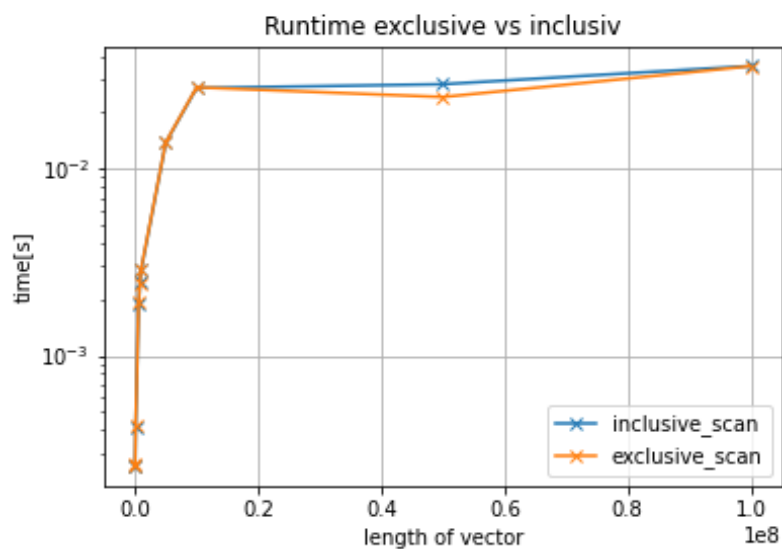
Listing 2: kernel for inclusive_scan)

```
1 // exclusive scan requires us to write a zero value
2 at the beginning of each block
3 my_value = (threadIdx.x > 0) ? shared_buffer[threadIdx.x - 1] : 0;
```

**d)**



There is bassicaly no differences.

## 2 Poisson equation (5 Points)

First I have to change all row_offsets, indices to an int datatype because there was an error with the matching datatype for the residual_norm function in "poisson2D.hpp".

**a)**

Listing 3: kernel for serching for zeros)

```cpp
__global__ void count_nonzero_entries(int* row_offsets, int N, int M) {
    for(int row = blockDim.x * blockIdx.x + threadIdx.x; row < N * M;
     row += gridDim.x * blockDim.x) {
        int nnz_for_this_node = 1;
        int i = row / N;
        int j = row % N;

        if(i > 0) nnz_for_this_node += 1;
        if(j > 0) nnz_for_this_node += 1;
        if(i < N-1) nnz_for_this_node += 1;
        if(j < M-1) nnz_for_this_node += 1;

        row_offsets[row] = nnz_for_this_node;
    }
}
```

**b)**

Again changed the double const *input to a int const * input.

Listing 4: given exclusive scan)

```cpp
void exclusive_scan(int const * input,
int        * output, int N)
{
    int num_blocks = 512;
    int threads_per_block = 512;

    int *carries;
    cudaMalloc(&carries, sizeof(int) * num_blocks);

    // First step: Scan within each thread group and write carries
    scan_kernel_1<<<num_blocks, threads_per_block>>>(input, output, N, carries);

    // Second step: Compute offset for each thread group (exclusive scan for each
    scan_kernel_2<<<1, num_blocks>>>(carries);

    // Third step: Offset each thread group accordingly
    scan_kernel_3<<<num_blocks, threads_per_block>>>(output, N, carries);

    cudaFree(carries);
}
```

**c)**

Listing 5: assemble the matrix A)

```
1  __global__ void assemble_A_GPU(double* values, int* columns, int* row_offsets, i
2  {
3      for(int row = blockDim.x * blockIdx.x + threadIdx.x;
4      row < N*M; row += gridDim.x * blockDim.x)
5      {
6          int i = row / N;
7          int j = row % N;
8          int counter = 0;
9
10         if ( i > 0)
11         {
12             values[(int)row_offsets[row] + counter] = -1;
13             columns[(int)row_offsets[row] + counter] = (i-1)*N+j;
14             counter++;
15         }
16
17         if ( j > 0)
18         {
19             values[(int)row_offsets[row] + counter] = -1;
20             columns[(int)row_offsets[row] + counter] = i*N+(j-1);
21             counter++;
22         }
23
24         values[(int)row_offsets[row] + counter] = 4;
25         columns[(int)row_offsets[row] + counter] = i*N+j;
26
27         counter++;
28
29         if ( j < M-1)
30         {
31             values[(int)row_offsets[row] + counter] = -1;
32             columns[(int)row_offsets[row] + counter] = i*N+(j+1);
33             counter++;
34         }
35         if ( i < N-1)
36         {
37             values[(int)row_offsets[row] + counter] = -1;
38             columns[(int)row_offsets[row] + counter] = (i+1)*N+j;
39             counter++;
40         }
41     }
42 }
```

## d)

Just adapt the GD program that you gave us and changed the names from the row_offsets.
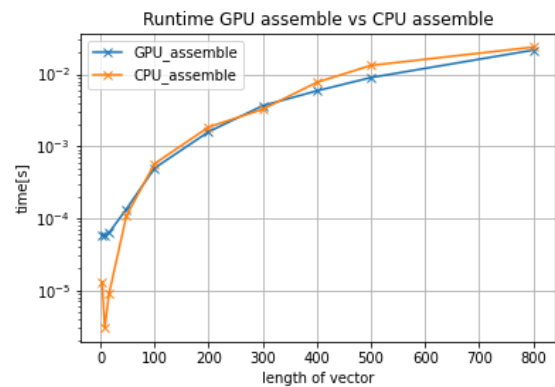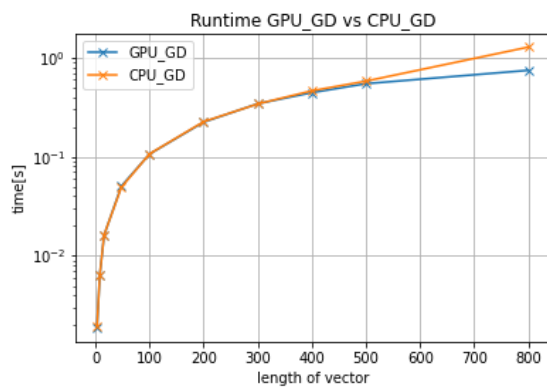Code is in the appendix in the mail.

## e)

Simply replace the conditions for the loop in the assembling kernel with this.
Code is in the appendix in the mail.

Listing 6: kernel for serching for zeros)

```
1  for(int row = 0; row < N*M; row++)
```

In total the differences are not that much. Maybe for higher dimensions of the matrix $A$ but I got an SEGMENTATION FAULT at a certain point of the benchmark $N > 800$. I tried to change int into a long int but then I got also some errors. I think the difference should be much higher for bigger dimensions of $A$.

### f) Bonus point

Very ineffective way to get the results for the plotting. The two for loops outputs the first line of the solution matrix.
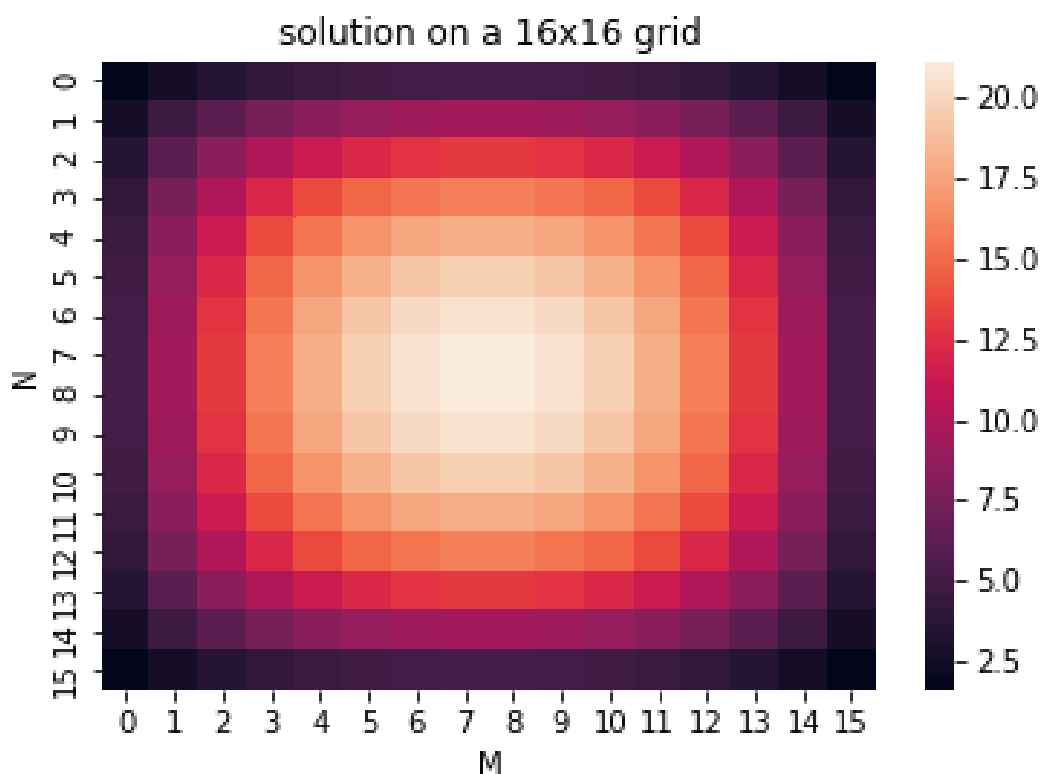
Listing 7: kernel for serching for zeros)

```
1   int ck = 0;
2   for (int i = 0; i < N/max; i++)
3   {
4       for (int j = 0; j < N/max; j++)
5       {
6           std::cout << solution[ck] << "," << std::endl;
7           ck = ck + 1;
8       }
9       std::cout << " " << std::endl;
10      std::cout << " " << std::endl;
11      std::cout << " " << std::endl;
12      std::cout << " " << std::endl;
13      std::cout << " " << std::endl;
14  }
```

I copied the output by hand in python and plot it with the function "heat_map". The dimension of the solution is $16 \times 16$ so in total 256 unknowns. The result looks nice. It is obviously the solution to the poisson equation.



The color bar on the side represents the amplitude of the solution. Where $N$ and $M$ are the cordinates of the solution matrix. In our case is the unite square given in the matrix assemble function.