



360.252 - COMPUTATIONAL SCIENCE ON MANY-CORE ARCHITECTURES

WS 2020 - EXERCISE 5

Christian GOLLMANN, 01435044

Last update: November 22, 2020

Contents

1	Inclusive and Exclusive Scan 1	1
2	Inclusive and Exclusive Scan 2	6
3	Inclusive and Exclusive Scan 3	7
4	Inclusive and Exclusive Scan 4	8

1 Inclusive and Exclusive Scan 1

Let's start with `scan_kernel_1`. In my sketch I show how the situation looks like for a grid with 4 blocks and a blocksize of 6 threads per block. The for loop in line 24 goes over the values of `X` which belong to the respective block, in my case block 2. At the end of the for loop, we can write the 6 temporary scan results into the result vector `Y`. The current offset gets saved in `block_offset` and gets added to the values of the next iteration. Please refer to the sketch for further details. At the end of `scan_kernel_1`, every block holds its exclusively scanned respective values and the carry which is needed in the next steps.

In `scan_kernel_2` the respective carries get summed up so they can be added to the already exclusively scanned values in `scan_kernel_3`. Since a picture says more than thousand words, please also refer to the sketch in figure 2 for further details.

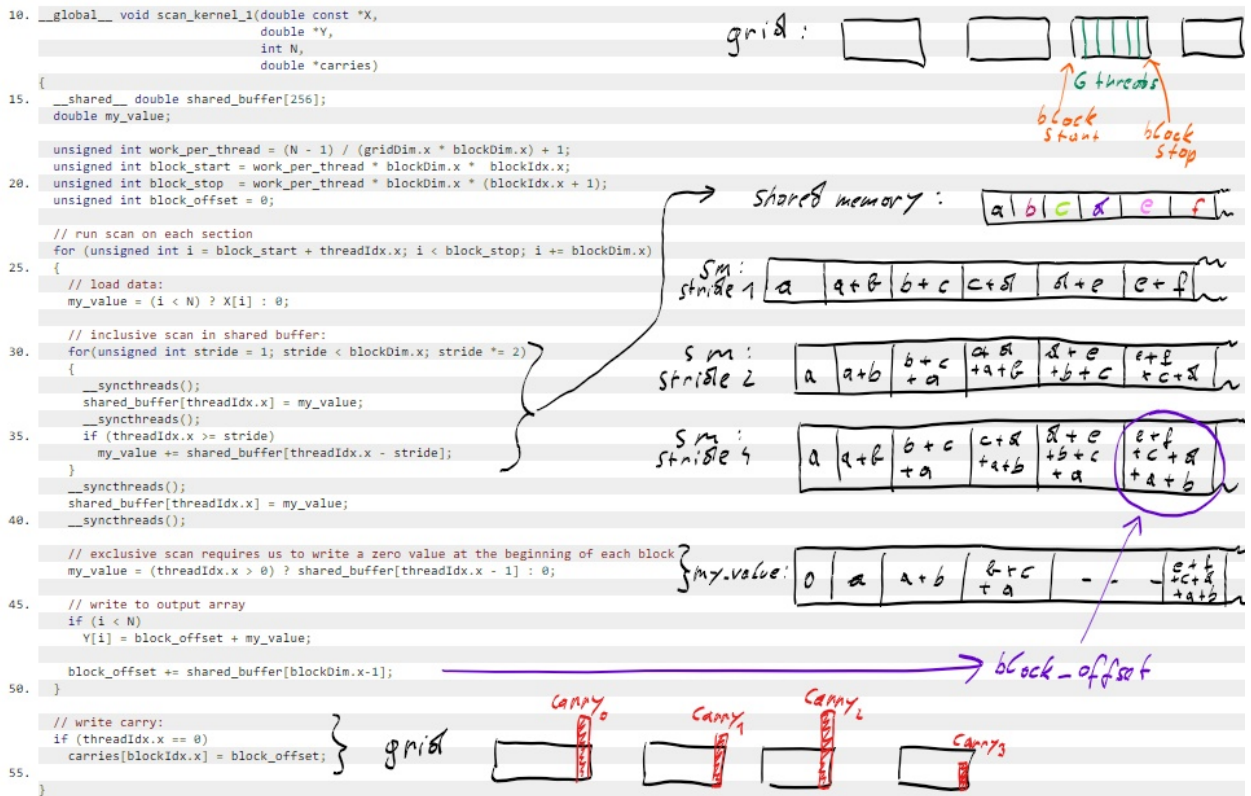


Figure 1: `scan_kernel_1`

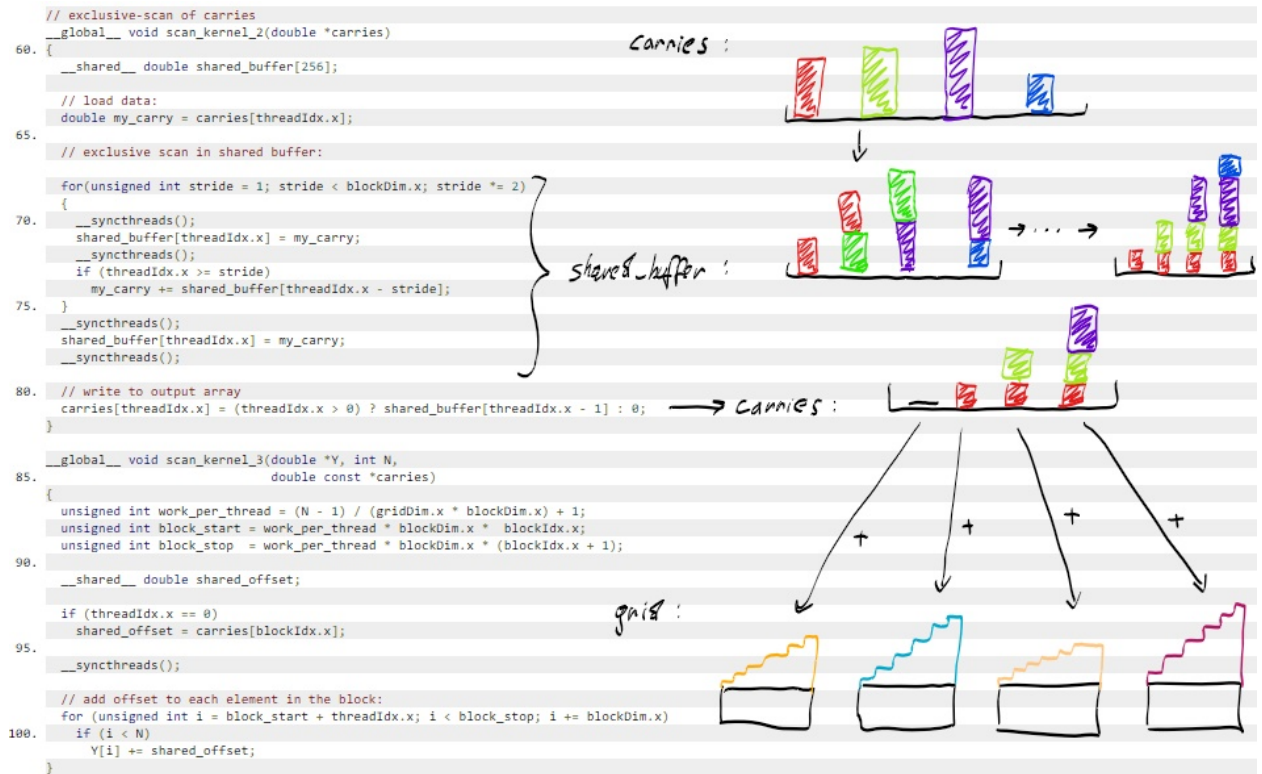


Figure 2: scan_kernel_2 and scan_kernel_3

For sake of completeness, I would also like to add the whole provided code here.

Listing 1: Provided Code for exclusive scan

```

1  #include "poisson2d.hpp"
2  #include "timer.hpp"
3  #include <algorithm>
4  #include <iostream>
5  #include <stdio.h>
6
7
8
9
10 __global__ void scan_kernel_1(double const *X,
11                             double *Y,
12                             int N,
13                             double *carries)
14 {
15     __shared__ double shared_buffer[256];
16     double my_value;
17
18     unsigned int work_per_thread = (N - 1) / (gridDim.x * blockDim.x) + 1;
19     unsigned int block_start = work_per_thread * blockDim.x * blockIdx.x;
20     unsigned int block_stop = work_per_thread * blockDim.x * (blockIdx.x + 1);
21     unsigned int block_offset = 0;
22
23     // run scan on each section
24     for (unsigned int i = block_start + threadIdx.x; i < block_stop; i += blockDim.x)
25     {

```

```

26 // load data:
27 my_value = (i < N) ? X[i] : 0;
28
29 // inclusive scan in shared buffer:
30 for(unsigned int stride = 1; stride < blockDim.x; stride *= 2)
31 {
32     __syncthreads();
33     shared_buffer[threadIdx.x] = my_value;
34     __syncthreads();
35     if (threadIdx.x >= stride)
36         my_value += shared_buffer[threadIdx.x - stride];
37 }
38 __syncthreads();
39 shared_buffer[threadIdx.x] = my_value;
40 __syncthreads();
41
42 // exclusive scan requires us to write a zero value at the beginning of each block
43 my_value = (threadIdx.x > 0) ? shared_buffer[threadIdx.x - 1] : 0;
44
45 // write to output array
46 if (i < N)
47     Y[i] = block_offset + my_value;
48
49     block_offset += shared_buffer[blockDim.x-1];
50 }
51
52 // write carry:
53 if (threadIdx.x == 0)
54     carries[blockIdx.x] = block_offset;
55
56 }
57
58 // exclusive-scan of carries
59 __global__ void scan_kernel_2(double *carries)
60 {
61     __shared__ double shared_buffer[256];
62
63     // load data:
64     double my_carry = carries[threadIdx.x];
65
66     // exclusive scan in shared buffer:
67
68     for(unsigned int stride = 1; stride < blockDim.x; stride *= 2)
69     {
70         __syncthreads();
71         shared_buffer[threadIdx.x] = my_carry;
72         __syncthreads();
73         if (threadIdx.x >= stride)
74             my_carry += shared_buffer[threadIdx.x - stride];
75     }
76     __syncthreads();
77     shared_buffer[threadIdx.x] = my_carry;
78     __syncthreads();
79
80     // write to output array
81     carries[threadIdx.x] = (threadIdx.x > 0) ? shared_buffer[threadIdx.x - 1] : 0;
82 }

```

```

83
84 --global-- void scan_kernel_3(double *Y, int N,
85                               double const *carries)
86 {
87     unsigned int work_per_thread = (N - 1) / (gridDim.x * blockDim.x) + 1;
88     unsigned int block_start = work_per_thread * blockDim.x * blockIdx.x;
89     unsigned int block_stop = work_per_thread * blockDim.x * (blockIdx.x + 1);
90
91     --shared-- double shared_offset;
92
93     if (threadIdx.x == 0)
94         shared_offset = carries[blockIdx.x];
95
96     __syncthreads();
97
98     // add offset to each element in the block:
99     for (unsigned int i = block_start + threadIdx.x; i < block_stop; i += blockDim.x)
100         if (i < N)
101             Y[i] += shared_offset;
102 }
103
104
105
106
107 void exclusive_scan(double const * input,
108                    double * output, int N)
109 {
110     int num_blocks = 256;
111     int threads_per_block = 256;
112
113     double *carries;
114     cudaMalloc(&carries, sizeof(double) * num_blocks);
115
116     // First step: Scan within each thread group and write carries
117     scan_kernel_1<<<num_blocks, threads_per_block>>>(input, output, N, carries);
118
119     // Second step: Compute offset for
120     each thread group (exclusive scan for each thread group)
121     scan_kernel_2<<<1, num_blocks>>>(carries);
122
123     // Third step: Offset each thread group accordingly
124     scan_kernel_3<<<num_blocks, threads_per_block>>>(output, N, carries);
125
126     cudaFree(carries);
127 }
128
129
130
131
132
133
134 int main() {
135
136     int N = 200;
137
138     //
139     // Allocate host arrays for reference

```

```

140 //
141 double *x = (double *)malloc(sizeof(double) * N);
142 double *y = (double *)malloc(sizeof(double) * N);
143 double *z = (double *)malloc(sizeof(double) * N);
144 std::fill(x, x + N, 1);
145
146 // reference calculation:
147 y[0] = 0;
148 for (std::size_t i=1; i<N; ++i) y[i] = y[i-1] + x[i-1];
149
150 //
151 // Allocate CUDA-arrays
152 //
153 double *cuda_x, *cuda_y;
154 cudaMalloc(&cuda_x, sizeof(double) * N);
155 cudaMalloc(&cuda_y, sizeof(double) * N);
156 cudaMemcpy(cuda_x, x, sizeof(double) * N, cudaMemcpyHostToDevice);
157
158
159 // Perform the exclusive scan and obtain results
160 exclusive_scan(cuda_x, cuda_y, N);
161 cudaMemcpy(z, cuda_y, sizeof(double) * N, cudaMemcpyDeviceToHost);
162
163 //
164 // Print first few entries for reference
165 //
166 std::cout << "CPU y: ";
167 for (int i=0; i<10; ++i) std::cout << y[i] << " ";
168 std::cout << " ... ";
169 for (int i=N-10; i<N; ++i) std::cout << y[i] << " ";
170 std::cout << std::endl;
171
172 std::cout << "GPU y: ";
173 for (int i=0; i<10; ++i) std::cout << z[i] << " ";
174 std::cout << " ... ";
175 for (int i=N-10; i<N; ++i) std::cout << z[i] << " ";
176 std::cout << std::endl;
177
178 //
179 // Clean up:
180 //
181 free(x);
182 free(y);
183 free(z);
184 cudaFree(cuda_x);
185 cudaFree(cuda_y);
186 return EXIT_SUCCESS;
187 }

```

2 Inclusive and Exclusive Scan 2

In order to make the scan exclusive, using what was already provided, I simply shift the end result.

Listing 2: Provided Code for exclusive scan

```
1  __global__ void makeInclusive(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          Y[i] = Y[i+1];
6      }
7      if (blockDim.x * blockIdx.x + threadIdx.x == 0)
8          Y[N-1] += X[N-1];
9  }
10
11
12 void inclusive_scan(double const * input ,
13                    double          * output , int N)
14 {
15     int num_blocks = 256;
16     int threads_per_block = 256;
17
18     double *carries;
19     cudaMalloc(&carries , sizeof(double) * num_blocks);
20
21     // First step: Scan within each thread group and write carries
22     scan_kernel_1<<<(num_blocks , threads_per_block>>>)(input , output , N, carries);
23
24     // Second step: Compute offset for each
25     thread group (exclusive scan for each thread group)
26     scan_kernel_2<<<1, num_blocks>>>(carries);
27
28     // Third step: Offset each thread group accordingly
29     scan_kernel_3<<<(num_blocks , threads_per_block>>>)(output , N, carries);
30
31     // Make inclusive
32     makeInclusive<<<(num_blocks , threads_per_block>>>)(output , N, input);
33
34     cudaFree(carries);
35 }
```

3 Inclusive and Exclusive Scan 3

In order to make the scan exclusive, modifying the existing code, it is enough to remove line 43 from the provided code.

4 Inclusive and Exclusive Scan 4

The different implementations perform pretty similar what was no big surprise I think because they are almost the same. Just for higher values of N , the inclusive Scan that reuses the existing code and therefore has an additional kernel, is slower. This can be explained by the fact that there is one more kernel to be executed.

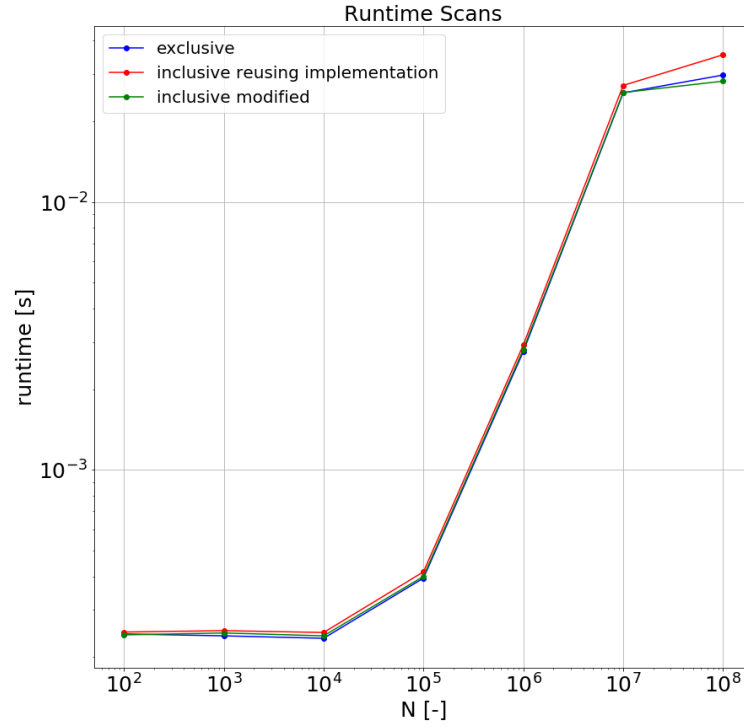


Figure 3: runtimes for different scan implementations