# PMPH - Assignment 3

Mikkel Willén

1. november 2022

# Indhold

## Task 1

The first inner loop isn't parallel, since it needs to access the previously calculated element, to calculate the new element.

The second inner loop isn't parallel, since it also needs to access previously calculated elements, to calculate new elements.

The outer loop isn't parallel, since array A is overwriten in each iteration of the loop.

Since each element of array A is changed with each iteration, we can privatize array A.

```
1   float A[N, 2*M];
2
3   // parallel
4   for (int i = 0; i < N; i++) {
5       A[i, 0] = N;
6   }
7
8   // parallel
9   for (int i = 0; i < N; i++) {
10      // sequential
11      for (int k = 1; k < 2 * M; k++) {
12          A[i, k] = sqrt(A[i, k - 1] * i * k);
13      }
14  }
15
16  // sequential
17  for (int i = 0; i < N; i++) {
18      // sequential
19      for (int j = 0; j < M; j++) {
20          B[i + 1, j + 1] = B[i, j] * A[i, 2 * j    ];
21          C[i,     j + 1] = C[i, j] * A[i, 2 * j + 1];
22      }
23  }
```

Computing the direction vector for the two last loops, we get:

$$
\begin{aligned}
&S_1 \to S_1 : (i+1, j+1) &&= (i, j) \\
&i_1 < i_2 \quad \& \quad j_1 < j_2 \\
&S_2 \to S_2 : (i, j+1) &&= (i, j) \\
&i_1 = i_2 \quad \& \quad j_1 < j_2 \\
&S_1 \to S_1 : [<, <] \\
&S_2 \to S_2 : [=, <]
\end{aligned}
$$

With this we see, that we safely can make a loop interchange, in the last loop nest, since we get an $<$ in the outer most loop and we get the following code, with maximum parallelism:

```
1   float A[N, 2*M];
2
3   // parallel
4   for (int i = 0; i < N; i++) {
5       A[i, 0] = N;
6   }
7
```

```
8    // parallel
9    for (int i = 0; i < N; i++) {
10       // sequential
11       for (int k = 1; k < 2 * M; k++) {
12           A[i, k] = sqrt(A[i, k - 1] * i * k);
13       }
14   }
15
16   // parallel
17   for (int j = 0; j < N; j++) {
18       // sequential
19       for (int i = 0; i < M; i++) {
20           B[i + 1, j + 1] = B[i, j] * A[i, 2 * j    ];
21           C[i,     j + 1] = C[i, j] * A[i, 2 * j + 1];
22       }
23   }
```

## Task 2

The outer loop is not parallel, since the accumulator is changed with each iterations of the inner loop, and reset with each iterations of the outer loop.

This could be fixed by privatizing the accum value and defining tmpA in the inner loop.

```
1    float A[N,64];
2    float B[N,64];
3    float accum[N]
4    for (int i = 0; i < N; i++) { // outer loop
5        accum[i] = 0;
6        for (int j = 0; j < 64; j++) { // inner loop
7            float tmpA = A[i, j];
8            accum[i] = sqrt(accum[i]) + tmpA * tmpA; // (**)
9            B[i,j] = accum[i];
10       }
11   }
```

The inner loop isn't parallel, since each iteration depends on the previous calculation of the accumulator.

The following is semantically-equivalent futhark code to if line (**) is rewriten as accum = accum + tmpA * tmpA

```
1    scan (+) 0 (map(/a -> a*a) A)
```

## Task 3

```
1    transposeTiled<float, TILE>(d_A, d_Atr, num_thds, width);
2    transfProg<<< num_blocks, block >>>(d_Atr, d_Btr, num_thds);
3    transposeTiled<float, TILE>(d_Btr, d_B, width, num_thds);
```

```
1    __global__ void
2    transfProg(float* Atr, float* Btr, unsigned int N) {
3        const unsigned int lid = threadIdx.x;
4        const unsigned int gid = blockIdx.x * blockDim.x + lid;
```

```
5        if (gid < N) {
6            float accum = 0;
7            for (int j = 0; j < 64; j++) {
8                float tmpA = Atr[gid + j * N];
9                accum = sqrt(accum) + tmpA * tmpA;
10               Btr[gid + j * N] = accum;
11           }
12       }
13   }
```

Since the number of rows and the number of columns are switched, when we compose A and B, we have to iterate through elements of the array indexed by the threadId plus the blockId times the size of the block plus the number of columns of the transposed arrays times the number of the current iteration of the a sequential loop, from 0 to the number of rows in the tranposed arrays. This is done in the above program.

Testing the program, we get following speeds:

```
Naive Memcpy GPU Kernel runs in: 994 microsecs, GB/sec: 540.11


Original Program runs on GPU in: 33463 microsecs, GB/sec: 16.043718
GPU PROGRAM   VALID!
Coalesced Program with manifested transposition runs on GPU in: 3135 microsecs,
GB/sec: 171.250702
GPU PROGRAM   VALID!
```

We get a speedup of around 10.68 times, though it is still around 3.15 times slower than naive memcpy.

## Task 4