# HPC ASSIGNMENT 2 – CUDA PARALLELIZATION

GROUP 3 – STEFANO PIGOZZI, CATERINA GAZZOTTI AND FABIO ZANICHELLI

# Some changes from last time…

- We noticed the algorithm wasn't allocating memory correctly when computing **A$^T$AX = Y**!
  - It did not account for non-square matrices, segfaulting if `NX` ≠ `NY`!
  - We fixed all these problems before proceeding to the optimization

- We noticed that Jeston Nano cannot perform atomic adds with double variables
  - So, execution is performed using float variables

```
/* Variable declaration/allocation. */
POLYBENCH_2D_ARRAY_DECL(A, DATA_TYPE, NX, NY, nx, ny);
POLYBENCH_1D_ARRAY_DECL(x, DATA_TYPE, NY, ny);
POLYBENCH_1D_ARRAY_DECL(y, DATA_TYPE, NY, ny);
POLYBENCH_1D_ARRAY_DECL(tmp, DATA_TYPE, NX, nx);
```

# kernel_atax optimization with CUDA

- Each thread adds Y vector to A matrix

- First lines of the function find the correct index's number of the current thread

- The atomicAdd at the end of the function doesn't work with doubles

- If HPC_USE_CUDA macro is not defined then sequential code is executed

```
283  #ifdef HPC_USE_CUDA
284  __global__ static void kernel_atax_cuda(DATA_TYPE* A, DATA_TYPE* X, DATA_TYPE* Y)
285  {
286      // Find out how many threads there are
287      unsigned int threads = gridDim.x * blockDim.x;
288
289      // Find how many iterations should be performed by each thread
290      unsigned int perThread = NX / threads + 1;
291
292      // Find the index of the current thread, even if threads span multiple blocks
293      unsigned int blockThreadIdx = blockIdx.x * blockDim.x + threadIdx.x;
294
295      // Have each thread perform the previously determined number of iterations
296      for(int stride = 0; stride < perThread; stride++)
297      {
298          // Iterate over x; y is not parallelized
299          unsigned int x = threads * stride + blockThreadIdx;
300
301          // Prevent the thread from accessing unallocated memory
302          if(x < NX)
303          {
304              // The same tmp as earlier
305              DATA_TYPE tmp = 0;
306
307              for (unsigned int y = 0; y < NX; y++)
308              {
309                  tmp += A[a_index(x, y)] * X[y];
310              }
311
312              for (unsigned int y = 0; y < NX; y++)
313              {
314                  // THIS DOES NOT WORK ON THE NANO, AS IT IS TOO OLD TO SUPPORT ATOMIC ADDITION WITH DOUBLES!
315                  // If you want to use the Nano, swap this for something else, or change atax.hu to use float instead of double
316                  atomicAdd(&Y[x], A[a_index(x, y)] * tmp);
317              }
318          }
319      }
320  }
321  #endif
```

# `init_array` optimization with CUDA

```
94   #ifdef HPC_USE_CUDA
95   __device__ static void init_array_cuda_x(DATA_TYPE* X, unsigned int threads)
96   {
97       // Find how many iterations should be performed by each thread
98       unsigned int perThread = NY / threads + 1;
99
100      // Find the index of the current thread, even if threads span multiple blocks
101      int blockThreadIdx = blockIdx.x * blockDim.x + threadIdx.x;
102
103      // Have each thread perform the previously determined number of iterations
104      for(int stride = 0; stride < perThread; stride++)
105      {
106          // Find the index of the current iteration
107          // This is equal to `y` of the init_array function
108          unsigned int iterationIdx = threads * stride + blockThreadIdx;
109
110          // Prevent the thread from accessing unallocated memory
111          if(iterationIdx < NY)
112          {
113              // Set the array element
114              X[iterationIdx] = iterationIdx * M_PI;
115          }
116      }
117  }
118  #endif
```

```
156  #ifdef HPC_USE_CUDA
157  __device__ static void init_array_cuda_a(DATA_TYPE* A, unsigned int threads)
158  {
159      // Find how many elements should be written in total
160      unsigned int elements = NX * NY;
161
162      // Find how many iterations should be performed by each thread
163      unsigned int perThread = elements / threads + 1;
164
165      // Find the index of the current thread, even if threads span multiple blocks
166      int blockThreadIdx = blockIdx.x * blockDim.x + threadIdx.x;
167
168      // Have each thread perform the previously determined number of iterations
169      for(int stride = 0; stride < perThread; stride++)
170      {
171          // Find the index of the current iteration
172          // This is equal to `y` of the init_array function
173          unsigned int iterationIdx = threads * stride + blockThreadIdx;
174
175          // Determine current x and y
176          unsigned int y = iterationIdx % NY;
177          unsigned int x = iterationIdx / NY;
178
179          // Prevent the thread from accessing unallocated memory
180          if(iterationIdx < elements)
181          {
182              // Set the array element
183              A[iterationIdx] = (DATA_TYPE)(x * (y + 1)) / NX;
184          }
185      }
186  }
187  #endif
```

This function is similar to the initialization of y

# What's the speedup?

- The optimized version takes 0,278 seconds to execute all the program with large dataset

- So, the speedup is $\frac{old\ time}{new\ time} = \frac{0,746}{0,278} = 2,68$

- Where do we achieve this speedup?

# Profiling

- It is better to optimize the `kernel_atax` part or the initialization?
  - The `kernel_atax` part, with large dataset, takes 0,35s with sequential code
  - The optimized version takes 0,26s, so the speedup is 1,35

  - The initialization part takes 0,396s to execute with sequential code
  - The optimized version takes 0,0178s with a speedup of 21,9

- The initialization part generates much more speedup, but also `kernel_atax` generates some speedup

$\rightarrow$ **So, the best choice is to optimize both parts**

# Experiments on other datasets

| | Mini dataset | Small dataset | Standard dataset | Large dataset | Extralarge dataset |
|---|---|---|---|---|---|
| **Sequential times** | $1,27 * 10^{-5}s$ | 0,00344s | 0,188s | 0,746s | 1,68s |
| **Optimized times** | $1,61 * 10^{-3}s$ | 0,0112s | 0,0647s | 0,278s | 0,665s |
| **Speedup** | **0,0079** | **0,31** | **2,90** | **2,68** | **2,52** |

*Speedups written in red are slowdowns

# OpenMP **vs** CUDA

| | Speedup with OpenMP* | Speedup with CUDA |
|---|---|---|
| **Mini dataset** | **0,476** | 0,0079 |
| **Small dataset** | **3,32** | 0,31 |
| **Standard dataset** | **3,36** | 2,90 |
| **Large dataset** | **3,37** | 2,61 |
| **Extralarge dataset** | **3,47** | 2,53 |

*OpenMP programs are executed with float variables too.
For mini dataset there is no best option than sequential code; OpenMP has only less slowdown

Thanks for the attention!