# Exercise 2: Solutions

Lecture, Summer 2019

Prof. Dr. Matthias S. Müller

Dr. Christian Terboven

Simon Schwitanski

Julian Miller

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

# Hardware Accelerators

## K-means Clustering

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

# Task 1

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

# Task 1.1 Preparations

- **This exercise can be done on the RWTH cluster environment CLAIX**
- **Access via lecture project lect0034**
  - → Register for an account "Hochleistungsrechnen RWTH Aachen" at the selfservice (http://www.rwth-aachen.de/selfservice)
  - → Sent your TIM ID to contact@hpc.rwth-aachen.de in case you are no member yet
- **Any problems with the cluster environment?**

- **Implement the proposed k-means algorithm as denoted in the code by TODO: task 1.2.**

```c
void k_means(int niters, point_t *points, point_t *centroids,
             int *assignment, point_t *result, int n, int k) {
  for (int iter = 0; iter < niters; ++iter) {
    // determine nearest centroids
    for (int i = 0; i < n; ++i) {
      double optimal_dist = DBL_MAX;      // Calculate Euclidean
distance to each centroid and
          determine the closest mean
      for (int j = 0; j < k; ++j) {
        double dist = (points[i].x - centroids[j].x) *
                      (points[i].x - centroids[j].x) +
                      (points[i].y - centroids[j].y) *
                      (points[i].y - centroids[j].y);
        if (dist < optimal_dist) {
          optimal_dist = dist;
          assignment[i] = j;
    } } }
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

```
// Calculate new positions of centroids
int count[k];
double sum_x[k];
double sum_y[k];
for (j = 0; j < k; ++j) {
  count[j] = 0;
  sum_x[j] = 0.0;
  sum_y[j] = 0.0;
}
for (i = 0; i < n; ++i) {
  count[assignment[i]]++;
  sum_x[assignment[i]] += points[i].x;
  sum_y[assignment[i]] += points[i].y;
}
for (j = 0; j < k; ++j) {
  if (count[j] != 0.0) {
    centroids[j].x = sum_x[j] / count[j];
    centroids[j].y = sum_y[j] / count[j];
} } } }
```

$ make run-small
Executing k-means clustering with
20 iterations, 1000 points, and 5
centroids...
Time Elapsed: 0.000235 s

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Verify that the results obtained in task 1.2 are correct (template: TODO: task 1.3).**

  → Write initial position of centroids to result before execution the algorithm:

  ```
  for (int i = 0; i < k; ++i) {
    result[i].x = centroids[i].x;
    result[i].y = centroids[i].y;
  }
  ```

  → Store results after each iteration:
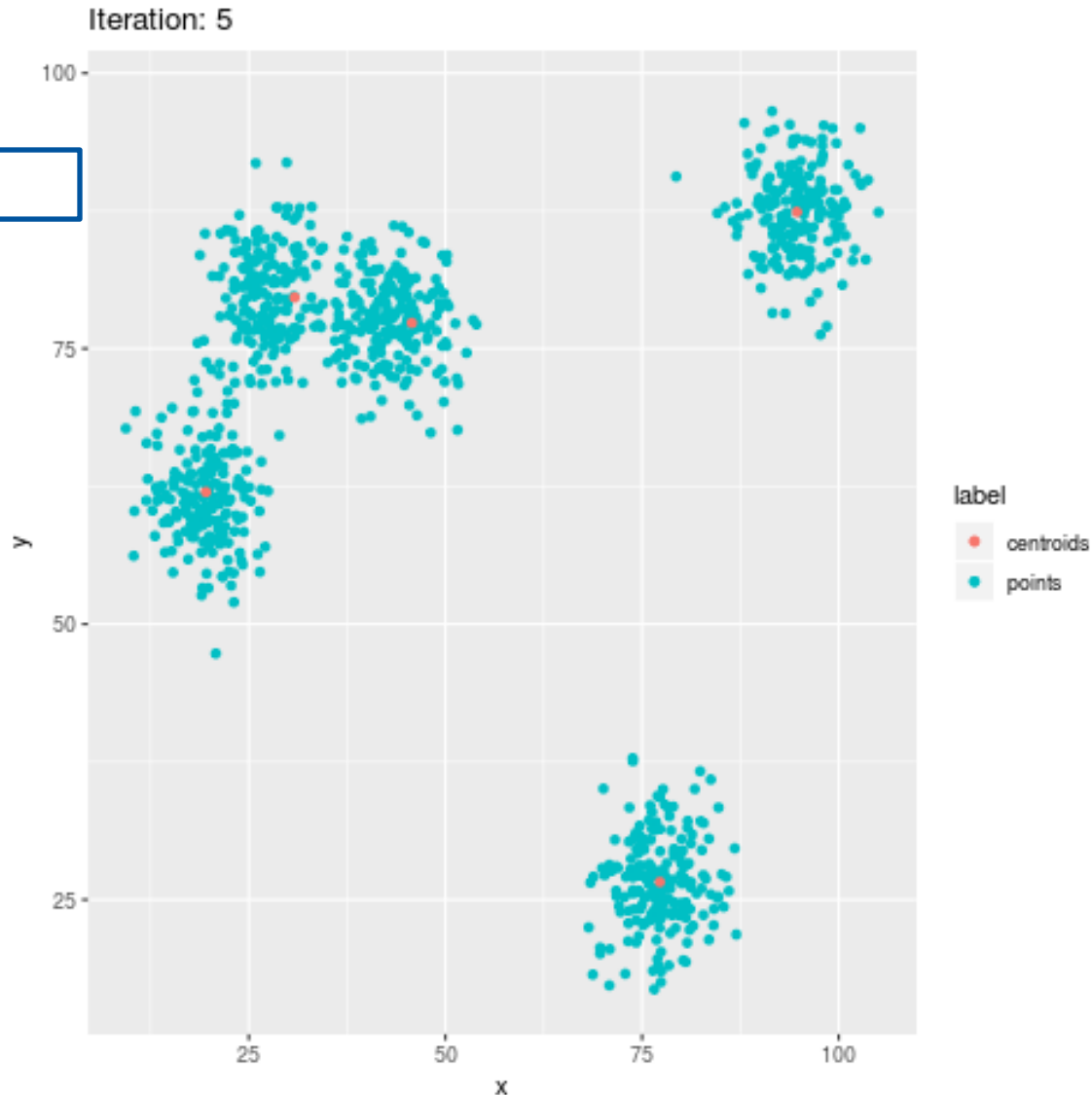
  ```
  for (int j = 0; j < k; ++j) {
    if (count[j] != 0) {
      centroids[j].x = sum_x[j] / count[j];
      centroids[j].y = sum_y[j] / count[j];
    }
    result[(iter + 1) * k + j].x = centroids[j].x;
    result[(iter + 1) * k + j].y = centroids[j].y;
  } }
  ```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

$ make vis-small

**Parallel and Data-c**
Chair for High Performance Computing, RWTH Aachen University

- **Optimize the serial implementation.**

- **Remember AoS vs. SoA**

```
struct point_t {
  double x;
  double y;
};
```

```
struct point_aos_t {
  double *x;
  double *y;
};
```

| Address | 0 | 8 | 16 | 24 | 32 | 40 |
|---------|------|------|------|------|------|------|
| AoS | x[0] | y[0] | x[1] | y[1] | x[2] | y[2] |
| SoA | x[0] | x[1] | x[2] | y[0] | y[1] | y[2] |

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Change all accesses from *[i].x to *->x[i]**

```
$ make run-small
Executing k-means clustering with
20 iterations, 1000 points, and 5
centroids...
Time Elapsed (AoS): 0.000235 s
Time Elapsed (SoA): 0.000235 s
```

```c
int main(int argc, const char* argv[]) {
  ...
  // Initialize points and centroids in SoA format
  point_soa_t points_soa;
  points_soa.x = (double*) malloc(n*sizeof(double));
  points_soa.y = (double*) malloc(n*sizeof(double));
  for (int i = 0; i < (niters + 1) * k; ++i) {
    result[i].x = -1.0;
    result[i].y = -1.0;
  }
  for (int i = 0; i < k; ++i) {
    result->y[i] = centroids->x[i];
    result->y[i] = centroids->y[i];
  }
  ...
}
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Why is there no performance difference?**

    → Optimal coalescing: all data read is requested

- **Would the following data structure change things?**

    → Yes, to compute distance only x and y are needed

    → Degree of coalescing $= \dfrac{\text{\#bytes requested}}{\text{\#bytes read}} = \dfrac{16 \text{ bytes}}{24 \text{ bytes}} = \dfrac{2}{3}$

```
struct point_t {
    double x;
    double y;
    int assignment;
};
```

| Address | 0 | 4 | 8 | 12 | 16 | 20 | 24 | 28 | 32 | 36 | 40 | 44 | 48 | 52 | 56 | 60 | 64 | 68 |
|---------|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| AoS | x[0] | | y[0] | | a[0] | | x[1] | | y[1] | | a[1] | | x[2] | | y[2] | | a[2] | |
| SoA | x[0] | | x[1] | | x[2] | | y[0] | | y[1] | | y[2] | | a[0] | a[1] | a[2] | | | |

# Task 2

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

# Preparations

- **Frontend for development and short test**

  → `$ ssh login18-g-1.hpc.itc.rwth-aachen.de`

  → Load NVIDIA compiler

    → `$ module load cuda/100`

  → List loaded modules

    → `$ module list`

- **Backend for performance measurements**

  → `$ sbatch kmeans_gpu.sh`

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Identify the hotspot of the algorithm of the optimized and checked serial version (task 1.4) and evaluate which parts are suitable for the GPU. Investigate which steps of the algorithm can be parallelized and how it can be achieved. Reason about dependencies.**

  → Profiler shows 100% time spent in `kmeans()`

  → 3 nested loops:

  → `for (iter = 0; iter < NO_ITER; ++iter) {`

  → Outer loop over iterations not parallelizable

  → `for (i = 0; i < n; ++i) {`

  → 2nd loop parallelizable (no data dependencies)

  → `for (j = 0; j < k; ++j) {`

  → 3rd loop parallelizable with a reduction on `optimal_dist`

- **Model the execution time $t_{GPU}$ of the hotspot on the V100 GPU based on the performance model introduced in the lecture. What limits the execution time of the hotspot?**

  → No data dependencies in 2nd loop

  → Many computations to offload

- **n*k operations per iteration with**

  → Simplification: leave out update to `optimal_dist`

  → 7 DP operations

  → 4 READS from main memory (down to 0 reads if `points` and `centroids`

  can be cached)

```
for (int iter = 0; iter < niters; ++iter) {
   for (int i = 0; i < n; ++i) {
      double optimal_dist = DBL_MAX;
      for (int j = 0; j < k; ++j) {
         double dist = (points[i].x - centroids[j].x) *
                       (points[i].x - centroids[j].x) +
                       (points[i].y - centroids[j].y) *
                       (points[i].y - centroids[j].y);
         if (dist < optimal_dist) {
            optimal_dist = dist;
            assignment[i] = j;
   } } }
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

→ Count, sum_x and sum_y

can be (for reasonable large k)

cached

→ 2n READ of points and

n READ of assignment

→ 4k WRITE to centroids

and result

→ 2n DP operations (additions on

sum_x and sum_y) and 2k DP

operations (divisions for cal-

culation of centroids)

```c
int count[k];
double sum_x[k];
double sum_y[k];
for (j = 0; j < k; ++j) {
  count[j] = 0;
  sum_x[j] = 0.0; sum_y[j] = 0.0;
}
for (i = 0; i < n; ++i) {
  count[assignment[i]]++;
  sum_x[assignment[i]] += points[i].x;
  sum_y[assignment[i]] += points[i].y;
}
for (int j = 0; j < k; ++j) {
  if (count[j] != 0) {
    centroids[j].x = sum_x[j] / count[j];
    centroids[j].y = sum_y[j] / count[j];
  }
  result[(iter + 1) * k + j].x =
      centroids[j].x;
  result[(iter + 1) * k + j].y =
      centroids[j].y;
} }
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Total kernel：**

  → n*k*7 + 2k +2n DP operations

  → n*k*4 + 2n + 4k DP READs/ WRITEs and n Integer READs

  → Dominated by main loop over n and k

  → Operational intensity of main loop: 7 Flops/32 byte

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

## Assess: Performance Modeling of K-Means

■ **Performance Model**

→ $t_{kernel} = \max(t_{compute}, t_{memory})$

→ $t_{compute} = \dfrac{arithmetic\ operations\ [Flop]}{P_{max}}$

→ No. DP operations: no_iter * n * k * 7 Flop

→ $P_{max}$ = 5120 cores * 1.3 GHz = 6656 GFlop/s

→ $t_{memory} = \dfrac{data\ transfers(LOAD, STORE)\ [words]}{b_s}$

→ Data transfers: no_iter * n * k * 4 Words

→ $b_s$ measured with BabelStream benchmark: 865 GB/s

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

■ **Performance Model**

→ $t_{GPU} = t_{H2D} + t_{kernel} + t_{D2H}$

→ $t_{data} = t_{H2D} + t_{D2H}$

→H2D: points (2n doubles), centroids (2k doubles), assigments (n integers),

result (2k (no_iter+1) doubles)

→D2H: result (2k (no_iter+1) doubles)

→$b_{PCI}$ measured: 12 GB/s

→$t_{data} = 2\alpha + \dfrac{2(n+k+k(no\_iter+1)) * 8\ bytes + n * 4\ bytes}{12\ GB/s} + \dfrac{2k(no\_iter+1) * 8\ bytes}{12\ GB/s}$

- **Offload the identified hotspot in task 2.1 to a GPU while taking care of the required data.**

```c
int main(int argc, const char* argv[]) {
  ...
  // Allocate memory for GPU
  point_t *d_points = 0; point_t *d_centroids = 0;
  int *d_assignments = 0; point_t *d_result = 0;
  cudaMalloc((void**)&d_points, N * sizeof(point_t));
  cudaMalloc((void**)&d_centroids, K * sizeof(point_t));
  cudaMalloc((void**)&d_assignments, N * sizeof(int));
  cudaMalloc((void**)&d_result, (niters + 1) * k * sizeof(point_t));
  // Copy data to GPU
  double runtime_all = get_time();
  cudaMemcpy(d_points, h_points, N * sizeof(point_t),
             cudaMemcpyHostToDevice);
  cudaMemcpy(d_centroids, h_centroids, K * sizeof(point_t),
             cudaMemcpyHostToDevice);
  cudaMemcpy(d_assignments, h_assignments, N * sizeof(int),
             cudaMemcpyHostToDevice);
  cudaMemcpy(d_result, result, (niters + 1) * k * sizeof(point_t),
             cudaMemcpyHostToDevice);
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

```
// Apply k-means algorithm
double runtime_kernel = get_time();
k_means<<<(n + THREADSPERBLOCK - 1)/THREADSPERBLOCK,
         THREADSPERBLOCK>>>(d_points, d_centroids, d_assignment
         d_result, n, k);
cudaDeviceSynchronize();
runtime_kernel = get_time() - runtime_kernel;

// Copy results back to host
cudaMemcpy(result, d_result, (niters + 1) * k * sizeof(point_t),
         cudaMemcpyDeviceToHost);
runtime_all = get_time() - runtime_all;

// Free memory
cudaFree(d_result);
cudaFree(d_assignment);
cudaFree(d_centroids);
cudaFree(d_points);
```

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Intuitive solution**

```
__global__ void k_means(point_t *points,
                        point_t *centroids,
                        int *assignment) {
int iter, j;
int tid = blockDim.x * blockIdx.x + threadIdx.x;
for (iter = 0; iter < NO_ITER; ++iter) {
  if (tid < n) {
    double optimal_dist = DBL_MAX;
    // Calculate Euclidean distance to each centroid and
       determine the closest mean
    for (j = 0; j < K; ++j) {
      double dist = (points[tid].x - centroids[j].x) *
                    (points[tid].x - centroids[j].x) +
                    (points[tid].y - centroids[j].y) *
                    (points[tid].y - centroids[j].y);
      if (dist < optimal_dist) {
        optimal_dist = dist;
        assignment[tid] = j;
  } } }
```

```
// Calculate new positions of centroids
if (tid < k) {
  int count = 0;
  double sum_x = 0.0;
  double sum_y = 0.0;
  for (int j = 0; j < n; ++j) {
    if (assignment[j] == tid) {
      sum_x += points[j].x;
      sum_y += points[j].y;
      count++;
    }
  }
  if (count != 0.0) {
    centroids[tid].x = sum_x / count;
    centroids[tid].y = sum_y / count;
  }
  result[(iter + 1) * k + tid].x = centroids[tid].x;
  result[(iter + 1) * k + tid].y = centroids[tid].y;
} } }
```

> $ make run-small
> Executing k-means clustering with 20 iterations, 1000 points, and 5 centroids...
> Time Elapsed (kernel): 0.000492 s
> Time Elapsed (total): 0.000554 s

- **Wrong results! Why?**

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Two separate kernels are required for synchronization**

```
__global__ void calc_distances(point_t *points, point_t *centroids,
                               int *assignment, int n, int k) {
  int tid = blockDim.x * blockIdx.x + threadIdx.x;
  if (tid < n) {
    double optimal_dist = DBL_MAX;
    for (int j = 0; j < k; ++j) {
      double dist = (points[tid].x - centroids[j].x) *
                    (points[tid].x - centroids[j].x) +
                    (points[tid].y - centroids[j].y) *
                    (points[tid].y - centroids[j].y);
      if (dist < optimal_dist) {
        optimal_dist = dist;
        assignment[tid] = j;
} } } }
```

```
__global__ void update_centroids(int iter, point_t *points,
                                 point_t *centroids, int *assignment,
                                 point_t *result, int n, int k) {
  int tid = blockDim.x * blockIdx.x + threadIdx.x;
  if (tid < k) {
    int count = 0;
    double sum_x = 0.0;
    double sum_y = 0.0;
    for (int j = 0; j < n; ++j) {
      if (assignment[j] == tid) {
        sum_x += points[j].x;
        sum_y += points[j].y;
        count++;
    } }
    if (count != 0.0) {
      centroids[tid].x = sum_x / count;
      centroids[tid].y = sum_y / count;
    }
    result[(iter + 1) * k + tid].x = centroids[tid].x;
    result[(iter + 1) * k + tid].y = centroids[tid].y;
} }
```

$ make run-small
Executing k-means clustering with 20 iterations, 1000 points, and 5 centroids...
Time Elapsed (kernel): 0.001543 s
Time Elapsed (total): 0.001608 s

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

# Task 2.4
# Optimize and Parallelize K-means on the GPU

■ **Optimize the algorithm and parallelize it to utilize the GPU as much as possible. Keep the data layout optimizations from task 1.4 in mind.**

→ Similar findings here: optimal coalescing

→ Distribute updating of centroids over multiple thread blocks to prevent serialization

→ Updating of centroids could be done in shared memory

- **Check you results for correctness and evaluate the actual performance with the performance estimated with your model in task 3.1. How close is you implementation to the theoretical peak performance for that hotspot based on the model? Evaluate reasons for potential differences between the modeled and the achieved performance.**

    → Exact same results

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Performance Comparison (large data set)**

  → $t_{data} = t_{total} - t_{kernel} = 0.607147$ s - $0.600960$ s ≈ 6.2 ms

  → Modelled: $t_{data} = 2\alpha + \frac{2(n+k+k(no\_iter+1)) * 8\ bytes + n * 4\ bytes}{12\ GB/s} + \frac{2k(no\_iter+1) * 8\ bytes}{12\ GB/s} \approx$

  $$\frac{2(1000000+5000+10000(50+1)) * 8\ bytes + 1000000 * 4\ bytes}{12\frac{GB}{s}} \approx 2.35\ ms$$

  → $t_{kernel} = 0.600960$

  → Modelled: $t_{kernel} = \max(t_{compute}, t_{memory})$

  → $t_{compute} = \frac{arithmetic\ operations}{P_{max}} = \frac{50 * 1000000 * 5000 * 7\ Flop}{6656\ GFlop/s} = 262.9\ ms$

  → $t_{memory} = \frac{data\ transfers(LOAD,STORE)}{b_s} = \frac{50 * 1000000 * 5000 * 4 * 8\ bytes}{865\ GB/s} = 9.25\ s$

- **Lots of caching**

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University

- **Compare your obtained performance results of the k-means algorithm on the GPU to the ones on the CPU (task 1.4). Is this problem suitable for the GPU (meaning: is the algorithm accelerated by the usage of the GPU)? Justify your answer.**

  → $t_{CPU}$ = 124.785246 s

  → $t_{GPU}$ = 0.607147 s

  → Speedup of ~231

- **Is this comparison fair? Justify your answer.**

  → No, use all 48 cores on CPU node

  → $t_{CPU}$ = 4.007033 s ->   436.7 Gflop/s  ($P_{max,CPU}$ = 2150 Gflop/s)

  → $t_{GPU}$ = 0.607147 s -> 2882.3 Gflop/s  ($P_{max,GPU}$ = 6656 Gflop/s)

  → Speedup of ~9

**Parallel and Data-centric Programming**
Chair for High Performance Computing, RWTH Aachen University