



# K-Means Algorithm

Implementation of sequential and  
parallel version using CUDA

**Marco Trambusti - 7135350**

# Introduction

## Purpose of the Assignment

- Implement the K-Means algorithm in C++ in both sequential and parallelized versions using CUDA.
- Dataset: 300147 2D points related to age and total spending of supermarket customers.
- Performance comparison between sequential and parallel version using speedup and efficiency .



## K-Means Algorithm

**Def.:** The K-Means algorithm is a widely used clustering method for partitioning a data set into K distinct clusters

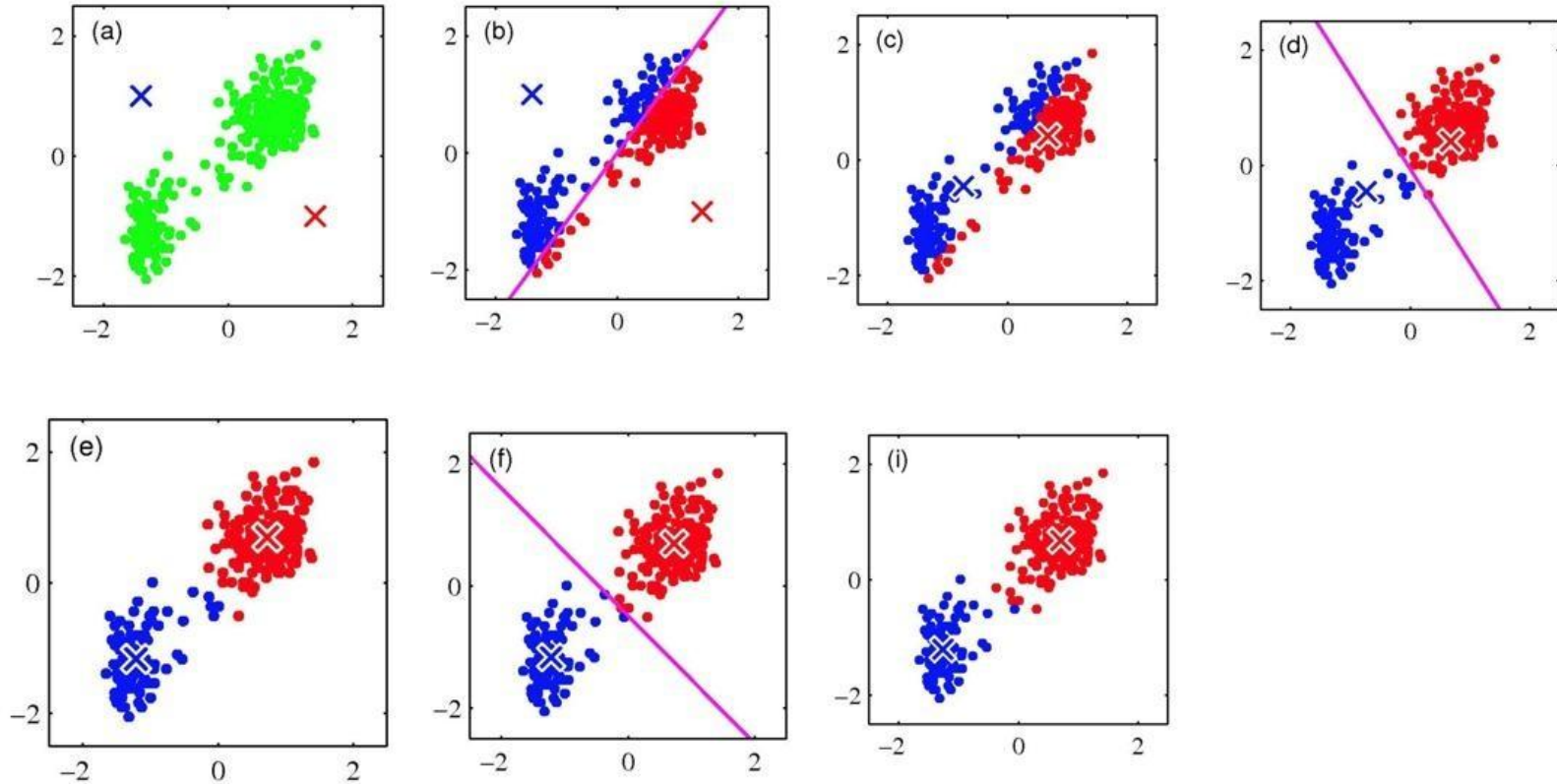
**Steps:**

1. **Initialization:** Random selection of K initial cluster centers.
2. **Cluster Assignment:** Each data point is assigned to the nearest centroid.
3. **Centroid Update:** Centroids are recalculated as the mean of the points assigned to each cluster.
4. **Iteration:** The assignment and update steps are repeated until the centroids no longer change significantly or a maximum number of iterations is reached.

**Stopping criterion used:** maximum number of iterations.

## Algorithm Parameters

- **Number of clusters (K):**
  - Specified a priori
  - Defined in an indicative manner (alternatively could have been used methods like Elbow, Silhouette or Gap Static)
- **Initialization of centroids:**
  - Randomly among the data points
- **Choice of distance metric:**
  - squared Euclidean distance
- **Maximum number of iterations**



Example of K-Means

## Sequential Implementation

- The project is in C++
- Chosen CUDA for the parallel implementation.
- 2D points are read from a csv file and normalized
- Clusters are initialized randomly;

The points in the dataset are represented using a struct **Point**, composed by:

- x, y coordinates (double)
- cluster (int);
- minDist (double)
- distance(Point p) (method)

Initially, the point does not belong to any cluster, and was arbitrarily set to -1.  
As a result, minDist was set to the maximum possible value.

# Sequential Implementation

## Kmeans method: Cluster assignment part

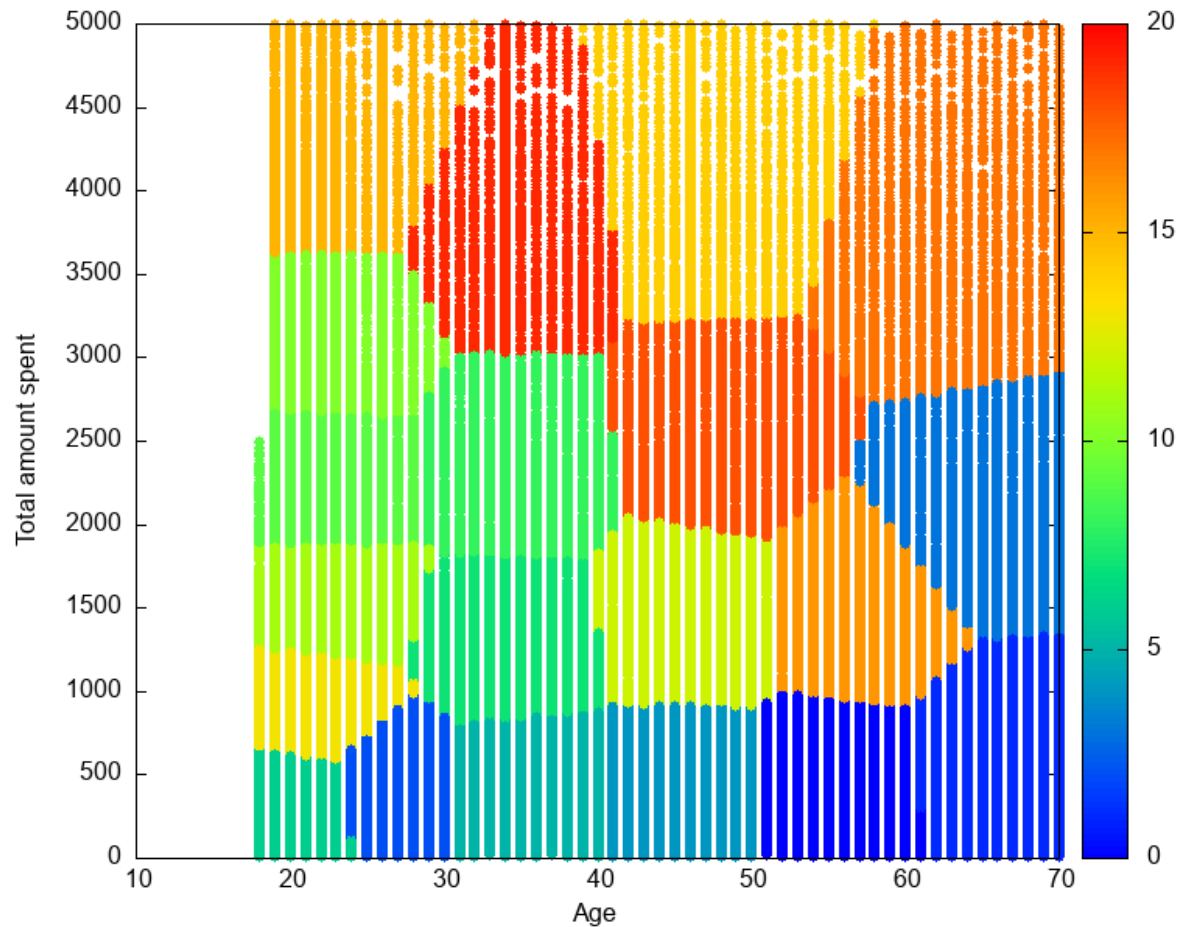
```
1 std::vector<double> KMeansHelper::kmeans_sequential(std::vector<Point>&
2 points, std::vector<Point>& centroids, const int k, const int epochs) {
3     double clusterAssignTime = 0.0;
4     double centroidsUpdateTime = 0.0;
5
6     for (int epoch = 0; epoch < epochs; ++epoch) {
7         std::vector<float> sumX(k, 0.0f), sumY(k, 0.0f);
8         std::vector<int> nPoints(k, 0);
9
10        // Assign points to clusters
11        double startClusterAssign = omp_get_wtime();
12        for (auto& point : points) {
13            for (int i = 0; i < k; ++i) {
14                const float dist=point.distance(centroids[i]);
15                if (dist < point.minDist) {
16                    point.minDist = dist;
17                    point.cluster = i;
18                }
19            }
20            const int clusterId = point.cluster;
21            nPoints[clusterId]++;
22            sumX[clusterId] += point.x;
23            sumY[clusterId] += point.y;
24            point.minDist = std::numeric_limits<float>::max();
25        }
26        ...
```

## Sequential Implementation

Kmeans method: Centroid Update Code

```
1   for (int i = 0; i < k; ++i) {  
2       if (nPoints[i] > 0) {  
3           double newX = sumX[i] / nPoints[i];  
4           double newY = sumY[i] / nPoints[i];  
5           centroids[i].x = newX;  
6           centroids[i].y = newY;  
7       }  
8   }
```





Example of results plotted using GnuPlotted

## Parallel Implementation

### From Array of Structs to Struct of Array

in order to make parallelization with CUDA more efficient by using the GPU, it is optimal to use a struct of arrays to represent points

```
1 struct Points {
2     std::vector<float> flat_coord;
3     std::vector<int> clusters;
4     Points(const int numPoints,
5           const std::vector<Point> &points) :
6         flat_coord(numPoints * 2), clusters(numPoints, -1)
7         for (int i = 0; i < points.size(); i++) {
8             flat_coord[i * 2] = points[i].x;
9             flat_coord[i * 2 + 1] = points[i].y;
10        }
11    }
12 };
```

## Cluster Assignment

- **Parallelism Exploited:**
  - The cluster assignment function implements the K-means algorithm on GPUs using CUDA to maximize parallelism.
- **Thread-Based Processing:**
  - Each thread is identified and processes a group of data points independently.
  - Defined constant POINTS\_PER\_THREAD=10 optimizes parallelization and maintains speedup.



# Cluster Assignment

```

1 __global__ void assignClusters(const float *d_datapoints,
2   int *d_clust_assn, float *d_centroids, int *d_nPoints,
3   float *d_sumX, float *d_sumY, int numPoints, const int k) {
4
5   int idx = blockIdx.x * blockDim.x + threadIdx.x;
6   extern __shared__ float s_centroids[];
7   int *s_nPoints=reinterpret_cast<int*>(&s_centroids[2*k]);
8   float *s_sumX = reinterpret_cast<float *>(s_nPoints + k);
9   float *s_sumY = s_sumX + k;
10  if (threadIdx.x < k) {
11      s_nPoints[threadIdx.x] = 0;
12      s_sumX[threadIdx.x] = 0.0f;
13      s_sumY[threadIdx.x] = 0.0f;
14      s_centroids[2*threadIdx.x]
15          = d_centroids[2*threadIdx.x];
16      s_centroids[2*threadIdx.x+1]
17          =d_centroids[2*threadIdx.x+1];
18  }
19  __syncthreads();
20  if (idx < static_cast<int>(ceil(static_cast<float>(
21      numPoints)/POINTS_PER_THREAD))) {

```

- centroids are loaded into the shared memory at the beginning and variables nPoints, sumX, sumY are defined in the shared memory to speed up the memory access operations for these variables.
- After the initialization, the threads are synchronized

# Cluster Assignment

```

22     for (int i = 0; i < POINTS_PER_THREAD; ++i) {
23         int pointIdx = idx * POINTS_PER_THREAD + i;
24         if (pointIdx < numPoints) {
25             float minDist = FLT_MAX;
26             int bestCluster = -1;
27             float px = d_datapoints[2 * pointIdx];
28             float py = d_datapoints[2 * pointIdx + 1];
29
30             #pragma unroll 5
31             for (int c = 0; c < k; ++c) {
32                 float cx = s_centroids[2 * c];
33                 float cy = s_centroids[2 * c + 1];
34                 float dist = (px-cx)*(px-cx)+(py-cy)*(py-cy);
35
36                 if (dist < minDist) {
37                     minDist = dist;
38                     bestCluster = c;
39                 }
40             }
41             d_clust_assn[pointIdx] = bestCluster;
42             atomicAdd(&s_nPoints[bestCluster], 1);
43             atomicAdd(&s_sumX[bestCluster], px);
44             atomicAdd(&s_sumY[bestCluster], py);
45         }
46     }
47 }
48 __syncthreads();
49 if (threadIdx.x < k) {
50     atomicAdd(&d_nPoints[threadIdx.x], s_nPoints[threadIdx.x]);
51     atomicAdd(&d_sumX[threadIdx.x], s_sumX[threadIdx.x]);
52     atomicAdd(&d_sumY[threadIdx.x], s_sumY[threadIdx.x]);
53 }
54 }

```

- The inner loop, responsible for iterating over the centroids (up to 20), is unrolled to reduce GPU overhead.
- The outer loop does not require unrolling as the value is negligible
- Post centroid assignment, threads synchronize again.
- The relevant variables on the device, including nPoints, sumX, and sumY, are updated to reflect the latest computation results

## Centroids Update

```
1 __global__ void updateCentroids(float *d_centroids,  
2   int *d_nPoints, float *d_sumX, float *d_sumY, int k) {  
3  
4     int idx = blockIdx.x*blockDim.x+threadIdx.x;  
5  
6     if (idx < static_cast<int>(ceil(static_cast<float>(k) /  
7       CLUSTERS_PER_THREAD))) {  
8       for (int i = 0; i < CLUSTERS_PER_THREAD; ++i) {  
9         int cIdx = idx*CLUSTERS_PER_THREAD+i;  
10        if (cIdx < k) {  
11          if (d_nPoints[cIdx] > 0) {  
12            d_centroids[2 * cIdx] =  
13              d_sumX[cIdx]/d_nPoints[cIdx];  
14            d_centroids[2 * cIdx + 1] =  
15              d_sumY[cIdx]/d_nPoints[cIdx];  
16          }  
17        }  
18      }  
19    }  
20 }
```

Each thread identified by an index calculated from the current block and thread and independently processes a group of centroids (defined in the constant CLUSTERS\_PER\_THREAD=5)

## Results and Performance Evaluation

### Setup:

- Tests conducted on an Intel Core i7-6700K processor (4 cores and 8 threads), and a NVIDIA GeForce GTX 970 GPU, based on Maxwell architecture, features 1664 CUDA cores, a base clock speed of 1050 MHz, 4 GB of GDDR5 memory with a 256-bit memory interface.
- Data sets of various sizes with a commensurate number of clusters and epochs.
- The results are averaged across five runs.



## Results and Performance Evaluation

### Duration:

Average execution time of each configuration.

- operations on larger data sets take significantly longer.

Small Data Set (500 points):

- 0.0011 seconds (1 thread)
- 0.0637 seconds (32 threads): Duration does not decrease with additional threads indicates management overhead.

Large Data Set (301,487 points):

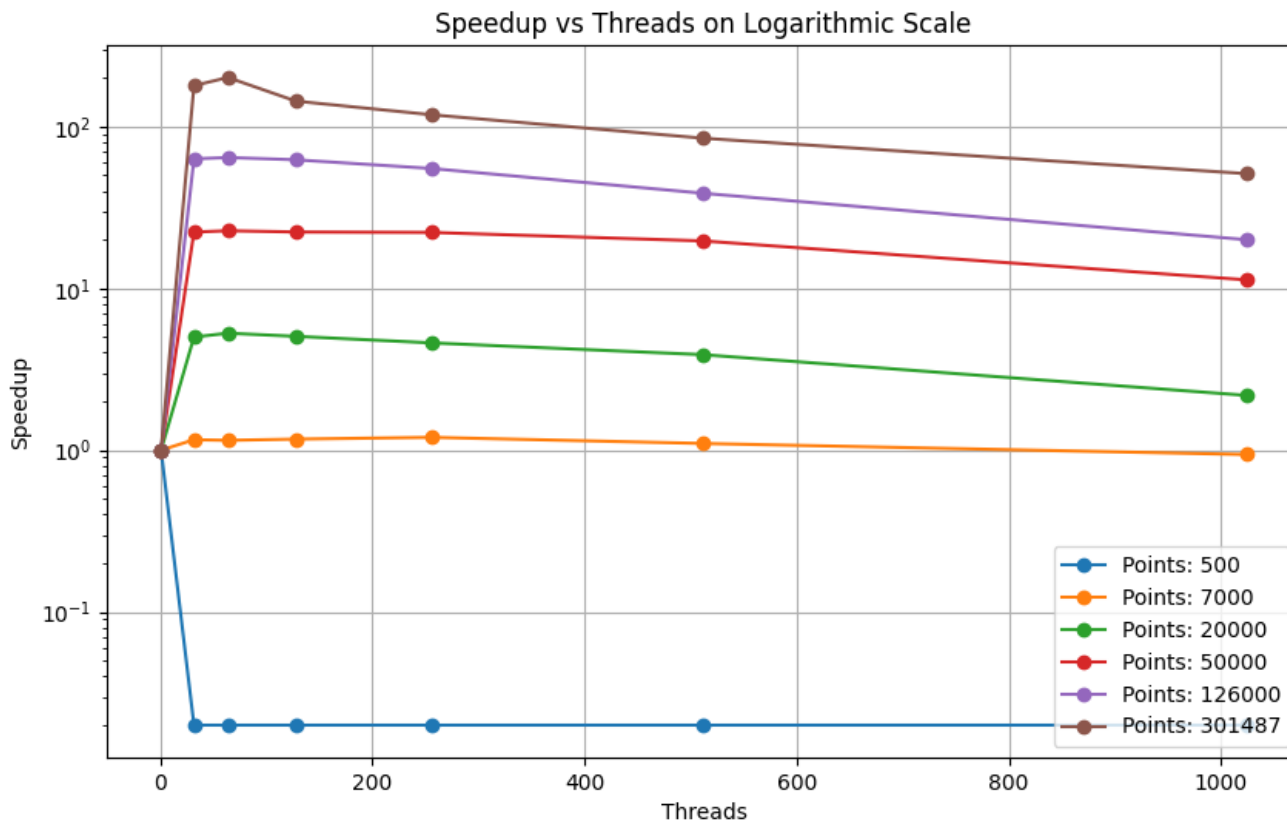
- 22.8738 seconds (1 thread)
- 0.1279 seconds (32 threads) - shows effectiveness of parallelization.





# Results and Performance Evaluation

**Speedup:** Ratio of sequential to parallel execution time.



## Speedup:

### Small Data Sets (500 & 7000 points):

- Speedup  $\sim 0.02$ ; minimal improvement with additional threads due to overhead.

### Medium Data Set (20,000 points):

- Speedup: 5.00 (32 threads), 5.28 (64 threads).
- Indicates a sub-linear trend.

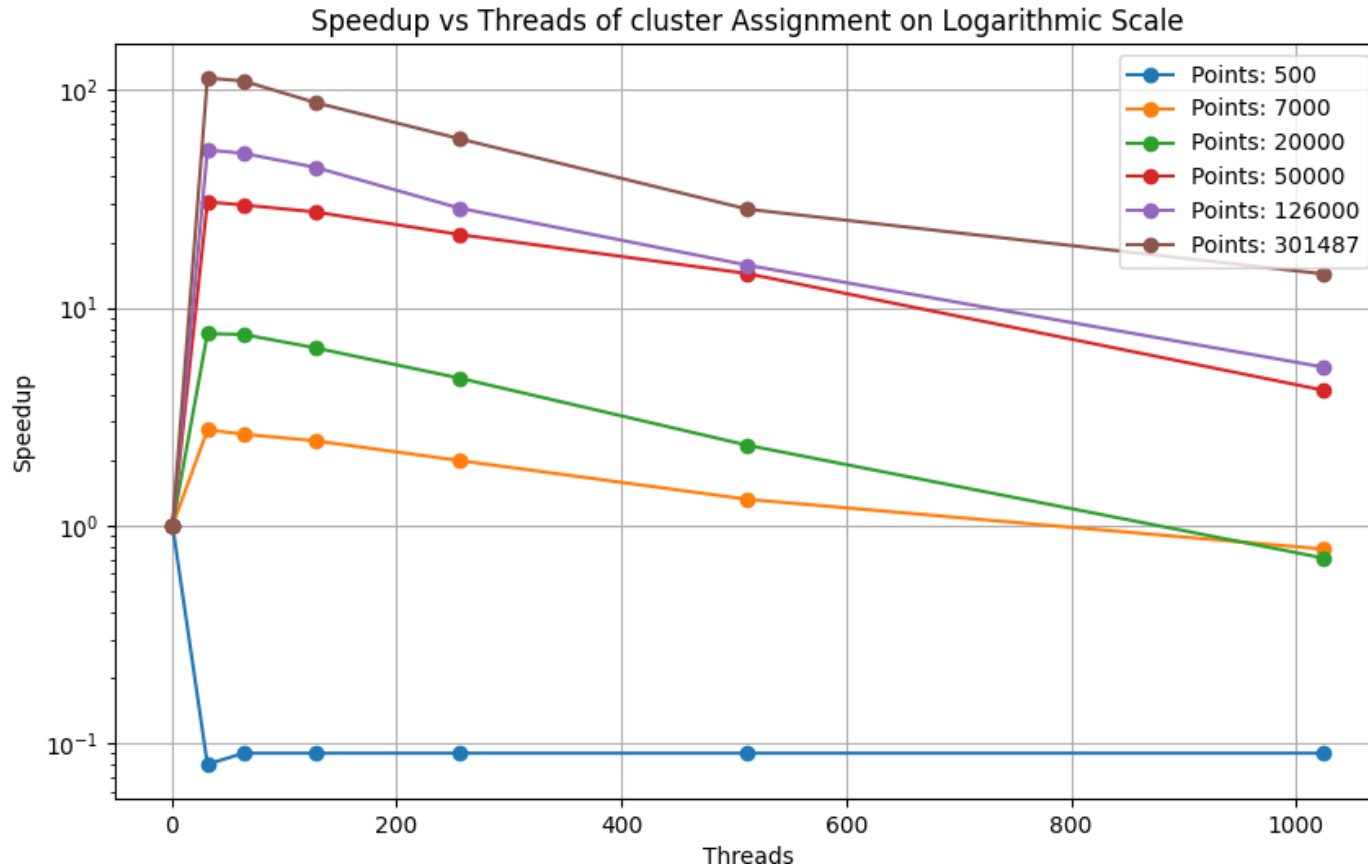
### Large Data Sets (50,000, 126,000, 301,487 points):

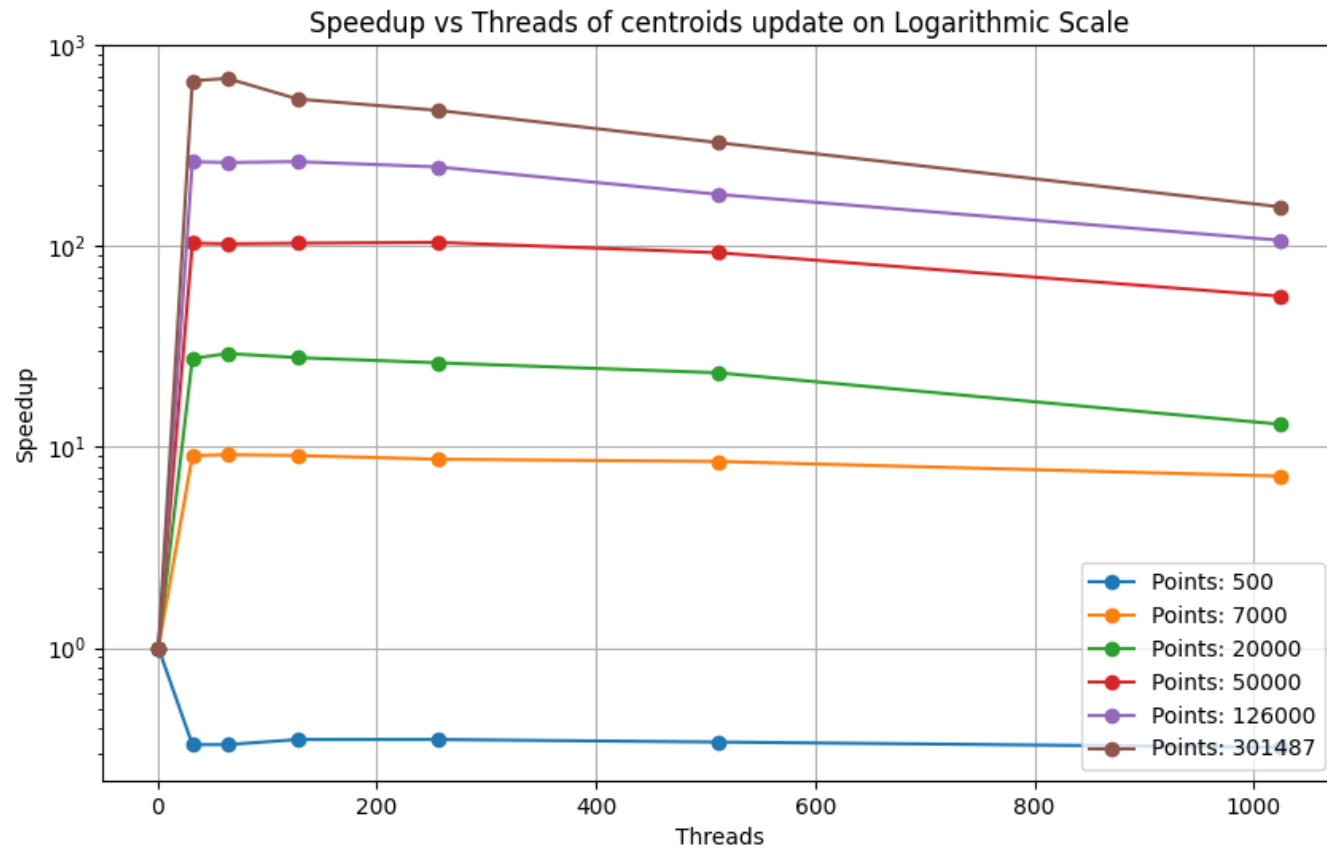
- 50,000 points: Speedup = 22.27 (32), 22.68 (64).
- 126,000 points: Speedup = 63.05 (32), 64.30 (64).
- 301,487 points: Speedup = 178.82 (32), peaks at 200.73 (64), then drops with 128 threads.

Beyond a certain number of threads, speedup almost stabilizes, suggesting saturation

Most tests show optimal performance gains between 32 and 64 threads

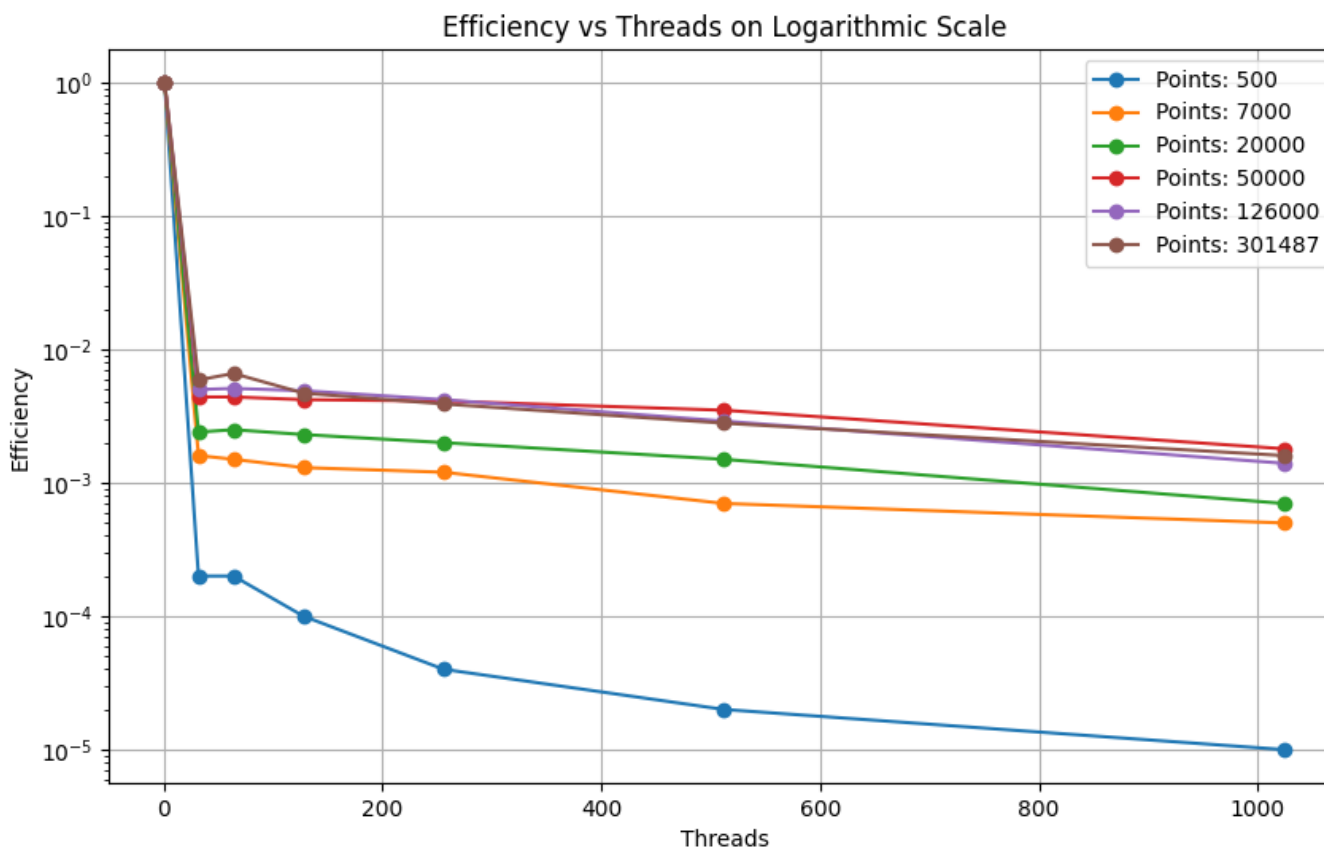






## Results and Performance Evaluation

**Efficiency:** Ratio of Speedup to Number of Threads.



## Efficiency:

### Smaller Datasets:

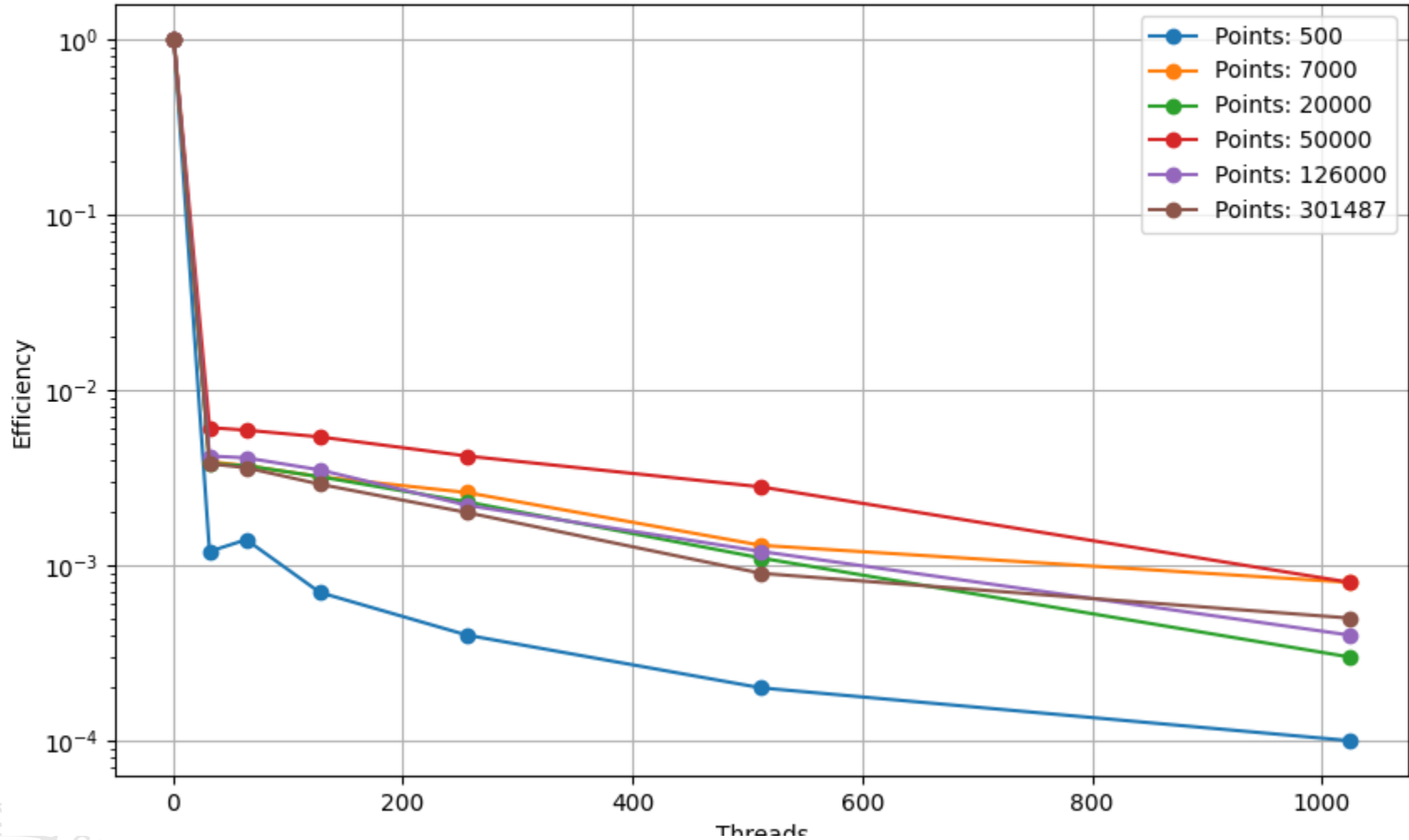
- Extremely low efficiency (e.g., 0.00001 with 1024 threads).

### Larger Datasets:

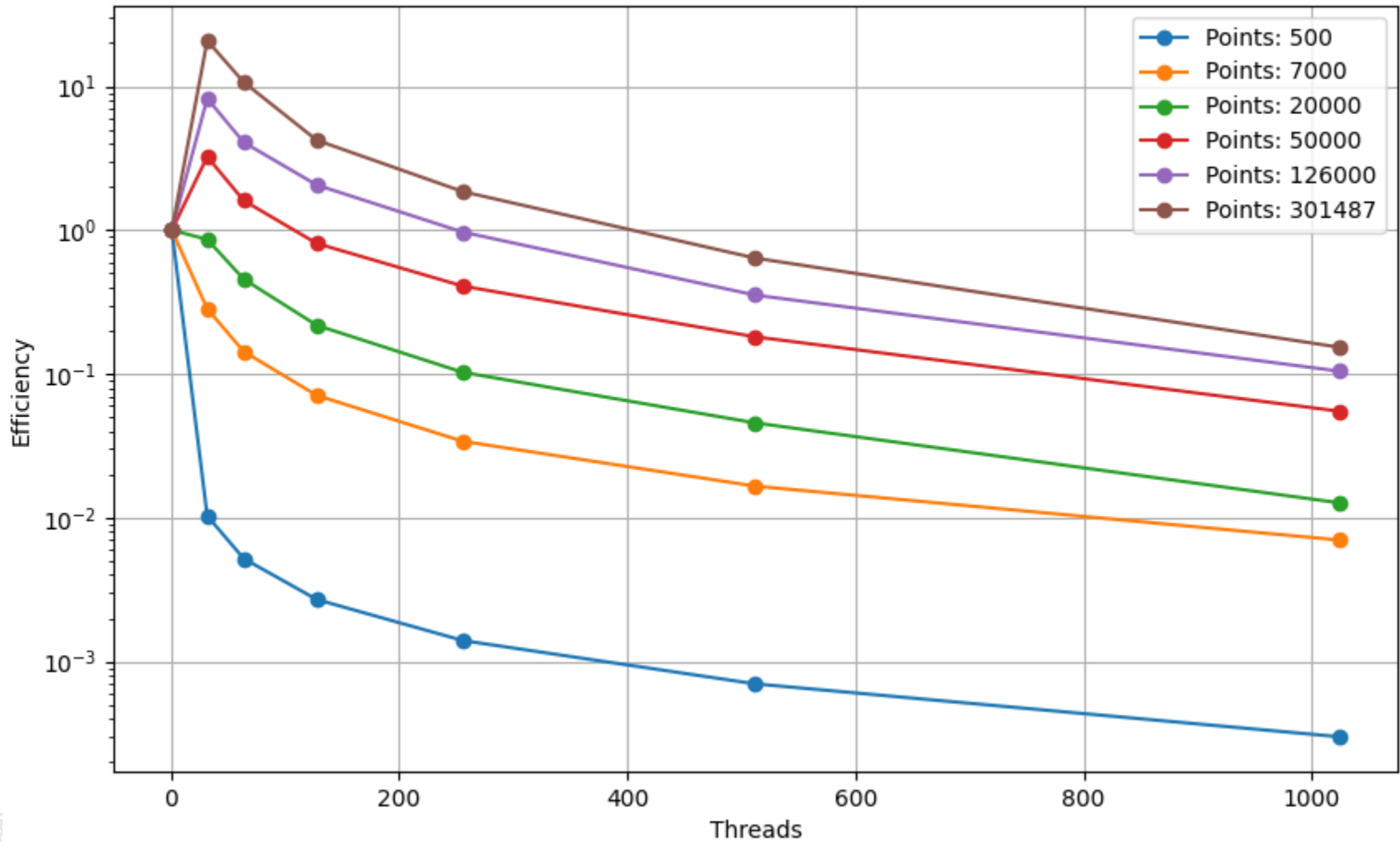
- Maximum efficiency of 0.0066 with 64 threads.
- Decreasing efficiency as thread count increases due to overheads.



Efficiency vs Threads of cluster Assignment on Logarithmic Scale



Efficiency vs Threads of centroids update on Logarithmic Scale





## Conclusions

### General Findings:

- Significant Speedup performance gains for large datasets.
- Efficiency is very low and decreases as the number of threads increases

### Performance Characteristics:

- Minimal speedup for smaller datasets due to overhead.
- Different KMEANS components respond variably to parallelization.

### Key Takeaway:

- CUDA parallelization is effective for large datasets, but efficiency diminishes with excessive threads; a balance is crucial for optimal performance.