

PC-2023/24 Parallel K-Means

Author

E-mail address

francesco.marchinil@edu.unifi.it

Abstract

This report delves into the utilization of OpenMP as an API for parallelizing C++ code, with a specific emphasis on the k-means algorithm employing Euclidean distance. The analysis compares a sequential version with a parallel version of the algorithm, focusing on the efficiency of parallelization achieved through OpenMP pragma directives. The report also provides insights into the parallelization process, highlighting the key areas of focus and the utilization of OpenMP directives.

Future Distribution Permission

The author(s) of this report give permission for this document to be distributed to Unifi-affiliated students taking future courses.

1. Introduction

This report describes the use of *OpenMP* [1] as a API for parallelizing C++ code. We will analyse and compare two implementations of the k-means algorithm with Euclidean distance: a sequential version and a parallel version. We will focus on the parallelization efficiency of our *OpenMP* pragma directive code, going to estimate the speedup achieved in terms of execution time. The algorithm analyses 1,000,000 data points ranging from 0 to 1,000 in an attempt to identify 10 clusters. The code repository is available on *GitHub* by clicking on the following link: parallel-k-means.

2. K-Means algorithm

This section explains how the code works, focusing on the most relevant parts for the future parallelization. The implementation of the k-means algorithm can be divided into two main

parts:

- setup and data generating process;
- clusters finding and points association.

The initial part of the code is not parallelizable as it is solely responsible for setting up the environment for subsequent operations. On the other hand, the second part is of our interest; this is where we are going to act with `pragma` directives to make the code multi-threaded. Our solution uses two for loops: the first one iterates through all points in the dataset, while the second one iterates through all clusters. In the first loop, we calculate the partial distance between each point and its respective cluster. This allows us to avoid an unnecessary loop to calculate the average distance of all points within each cluster.

2.1. Parallelization

The code's second portion is divided into two main operations:

- assignment of each point to its nearest centroid;
- calculate the average distance between the points of each centroid and assign the new centroid.

2.1.1 First section

The first operation consist on calculate the Euclidean distance between points and all clusters, assigning all points to the nearest cluster. This operation is implemented through the use of two for loops, the first dealing with scrolling through all

the points in the randomly generated dataset, the second scrolling through all the clusters by going to associate the point in question with the nearest cluster.

```

1 for (int i = 0; i < points.size(); i++)
2 {
3     double min_distance = INFINITY;
4     int centroid_index = -1;
5     for (int j = 0; j < centroids.size();
6         j++)
7     {
8         double distance =
9             euclidean_distance(points[i],
10                                centroids[j]);
11         if (distance < min_distance)
12         {
13             min_distance = distance;
14             centroid_index = j;
15         }
16     }
17     if (centroid_index !=
18         points[i].clusterId)
19     {
20         points[i].clusterId =
21             centroid_index;
22         changed = true;
23     }
24     update(cumulate_centroids.x)
25     update(cumulate_centroids.y)
26     update(cumulate_centroids.z)
27     update(cumulate_centroids.num_points)
28 }

```

Focusing on the external for loop, we can easily see that the operation of calculating the nearest cluster for each point is independent of each other, so we can simply parallelise the loop using the *OpenMP* directive, paying attention to the variables shared by the various threads: `#pragma omp parallel for shared (points, centroids, cumulate_centroids, changed).`

In addition, we must make the operations on newly declared shared variables atomic, to avoid race condition. This can be achieved by adding the `#pragma omp atomic write` and `#pragma omp atomic` directives.

```

1 #pragma omp parallel for shared(points,
2   centroids, cumulate_centroids, changed)
3 for (int i = 0; i < points.size(); i++)
4 {
5     double min_distance = INFINITY;
6     int centroid_index = -1;

```

```

6     for (int j = 0; j < centroids.size();
7         j++)
8     {
9         double distance =
10             euclidean_distance(points[i],
11                                centroids[j]);
12         if (distance < min_distance)
13         {
14             min_distance = distance;
15             centroid_index = j;
16         }
17     }
18     if (centroid_index !=
19         points[i].clusterId)
20     {
21         #pragma omp atomic write
22         points[i].clusterId =
23             centroid_index;
24         #pragma omp atomic write
25         changed = true;
26     }
27     #pragma omp atomic
28     update(cumulate_centroids.x)
29     #pragma omp atomic
30     update(cumulate_centroids.y)
31     #pragma omp atomic
32     update(cumulate_centroids.z)
33     #pragma omp atomic
34     update(cumulate_centroids.num_points)
35 }

```

2.1.2 Second section

The second operation involves computing the average distance between the internal points of each cluster. In this case, a for loop is used to achieve this goal. It scrolls through all clusters and calculates the mean distance between all points in the cluster, assigning the new value to them.

```

1 for (int i = 0; i < centroids.size(); i++)
2 {
3     average_distance(centroids[i].x)
4     average_distance(centroids[i].y)
5     average_distance(centroids[i].z)
6 }

```

All operations within the loop are independent of each other. Therefore, as before, we can parallelise the cycle. The 'for' loop can be parallelized by adding the `#pragma omp parallel for shared(centroids, cumulate_centroids)` directive. Here is the final code:

```

1 #pragma omp parallel for shared(centroids,
   cumulate_centroids)
2 for (int i = 0; i < centroids.size(); i++)
3 {
4     average_distance(centroids[i].x)
5     average_distance(centroids[i].y)
6     average_distance(centroids[i].z)
7 }

```

3. Speedup

The performance achieved with the sequential version will now be compared to that of the parallel version in terms of execution time. In order to reach this task, a Python script will be utilized to execute both sequential and parallel versions of k-means varying the number of threads used. The script calculates the speedup achieved at each iteration and generates a graph displaying the results in terms of speedup for threads. In order to make the experiments replicable, we set the seed for random data generation equal to 42. All tests were performed on an 11th generation Intel i7 [2] CPU with 4 cores and 8 threads.

The speed up is computed as follow:

$$S_P = t_s/t_p \quad (1)$$

Where t_s is the execution time of sequential program and t_p the execution time of the parallel program. In this case, we run the parallel code using between one and sixteen threads to evaluate performance at each step. As expected given the CPU used, the best approach is achieved by utilising eight threads: $S_P = 4.2836$. All result are available in Fig. 1 and on Table 3.

Considering the 8 threads of the 11th Intel i7, we will measure our speedup in terms of efficiency using the following formula:

$$E_P = S_P/P \quad (2)$$

Where P is number of CPU units used. The achieved efficiency is: $E_P = 0.535$; highlighting an efficiency of 54% on our machine. The second graph (Fig.2) displays the results of the same experiment, indicating only the execution times of the parallel algorithm as the number of threads

Thread #	Speedup
1	1
2	1.6849
3	2.5299
4	3.3668
5	2.8529
6	3.162
7	3.3559
8	4.2836
9	3.7095

Table 1. Speedup increases as the number of threads used increases.

used changes. The third graph (Fig. 3) compares the sequential and parallel programs based on their execution time in 10 random trials, each with a different dataset. The x-axis shows the iterations required for the algorithm to converge, and the y-axis shows the execution time in seconds.

4. Conclusion

We can conclude that by using OpenMP, we can achieve a good degree of parallelism with just a few `pragma` instructions. We recommend experimenting with further tests on different machines with diffrent numbers of CPUs and threads, highlighting any differences in speedup achieved.

References

- [1] Openmp, 2012. <https://www.openmp.org/>.
- [2] Intel. i7 11th generation, 2021. <https://www.intel.com/>.

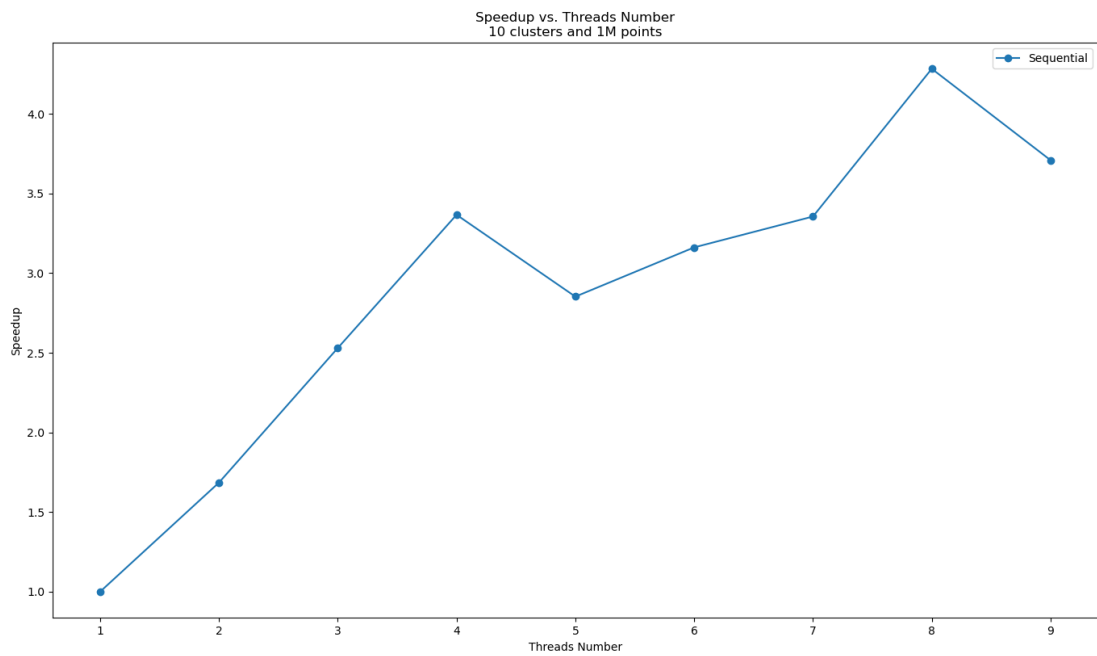


Figure 1. Speedup VS Number of Threads on 10 cluster and 1M points.

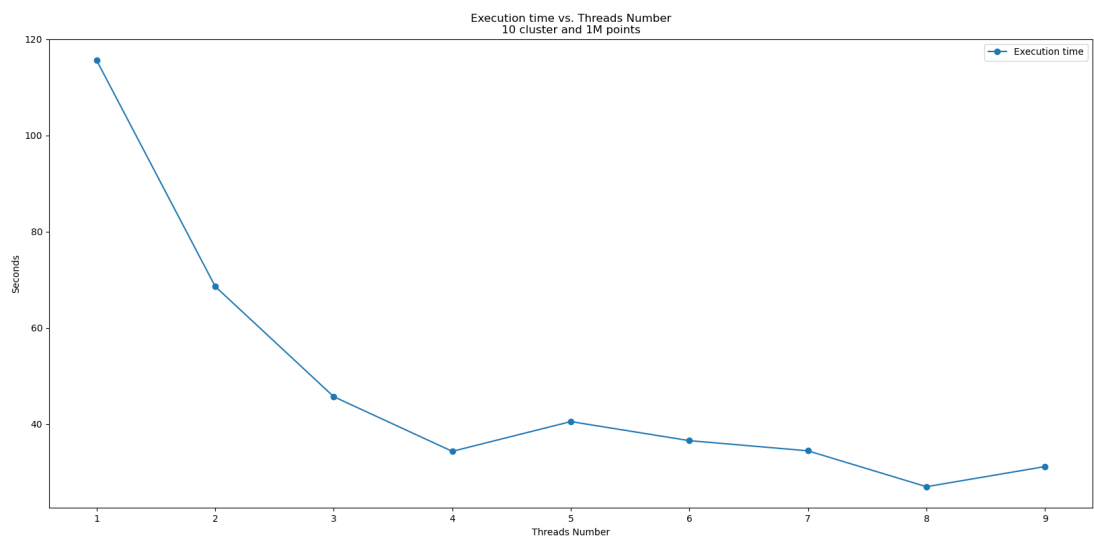


Figure 2. Execution time VS Number of Threads on 10 cluster and 1M points.

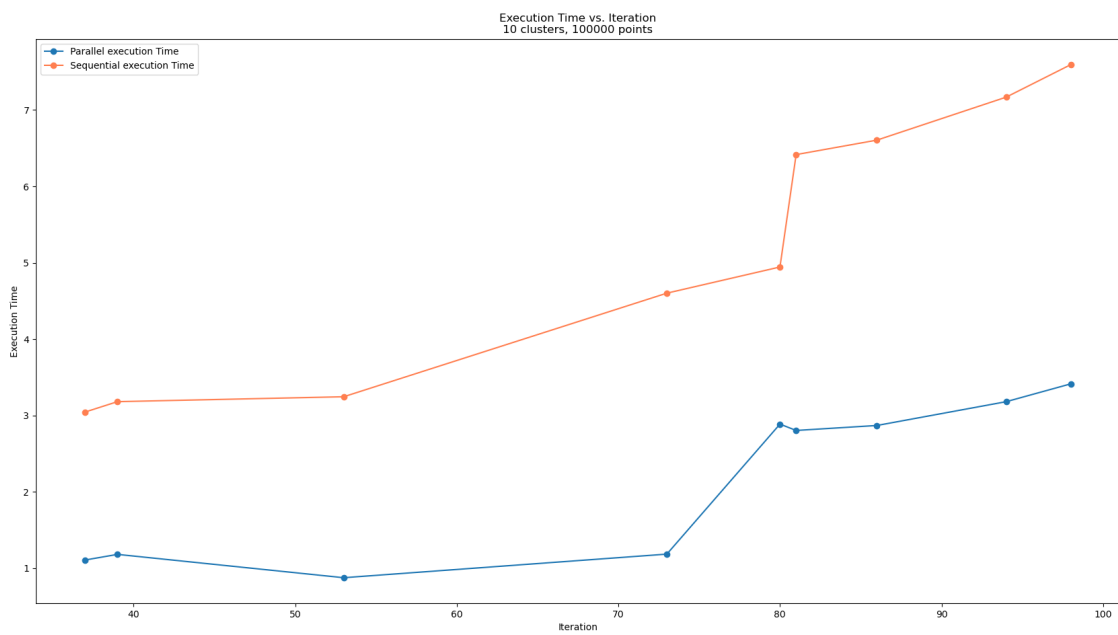


Figure 3. Parallel vs sequential execution times on 10 clusters and 100000 points using different dataset.