



**INDIAN INSTITUTE OF INFORMATION TECHNOLOGY
DESIGN AND MANUFACTURING KANCHEEPURAM**

**PROJECT REPORT
ON
K-MEAN CLUSTERING ALGORITHM
USING CUDA**

**SUBMITTED BY
AMAR KUMAR
(CED17I029)
TO
DR. NOOR MAHAMMAD**

Problem Statement

A Pizza company wants to open its delivery centres across a very big city.

The challenges they will face is that they need to analyze areas from where pizza is being ordered frequently.

They need to figure out the locations for the pizza stores within all these areas in order to keep the distance between the store and delivery points minimum.

Resolving these challenges includes a lot of analysis and mathematics. We would learn here about how clustering can provide a meaningful and easy method of sorting out such real life challenges.

My Solution to this problem

I have used an unsupervised machine learning algorithm to solve this problem. I have used the K-Means clustering algorithm for solving the above problem.

My dataset contains x and y coordinates of locations of different people/their houses across the city and some randomly generated pizza centers(According to my algorithm, the number of pizza centers can vary from 1 to 100).

We have to find the optimum location for these pizza centers so that each house can get the delivery fastest.

Step1 : For this, I have calculated the distance of every house from every pizza center. The house which is nearest to a particular pizza center, will be served by that particular pizza center in the first iteration.

Step2 : Now the location of the pizza center will change according to the location of houses which come under it. X coordinate of the pizza center will be calculated by taking the mean of x coordinates of the location of houses. Similarly Y coordinate of the pizza center will be calculated by taking the mean of y coordinates of the location of the houses.

Step3 : For further iteration repeat from Step1

Note 1 : More the number of iterations, more optimal will be the location of the pizza center. But the location of the pizza center will not vary after a certain number of iterations as it will start converging to an optimal location of the pizza center such that a house can get its delivery faster.

Note 2 : This implementation of k-mean clustering algorithm in CUDA is done on Google Colab.

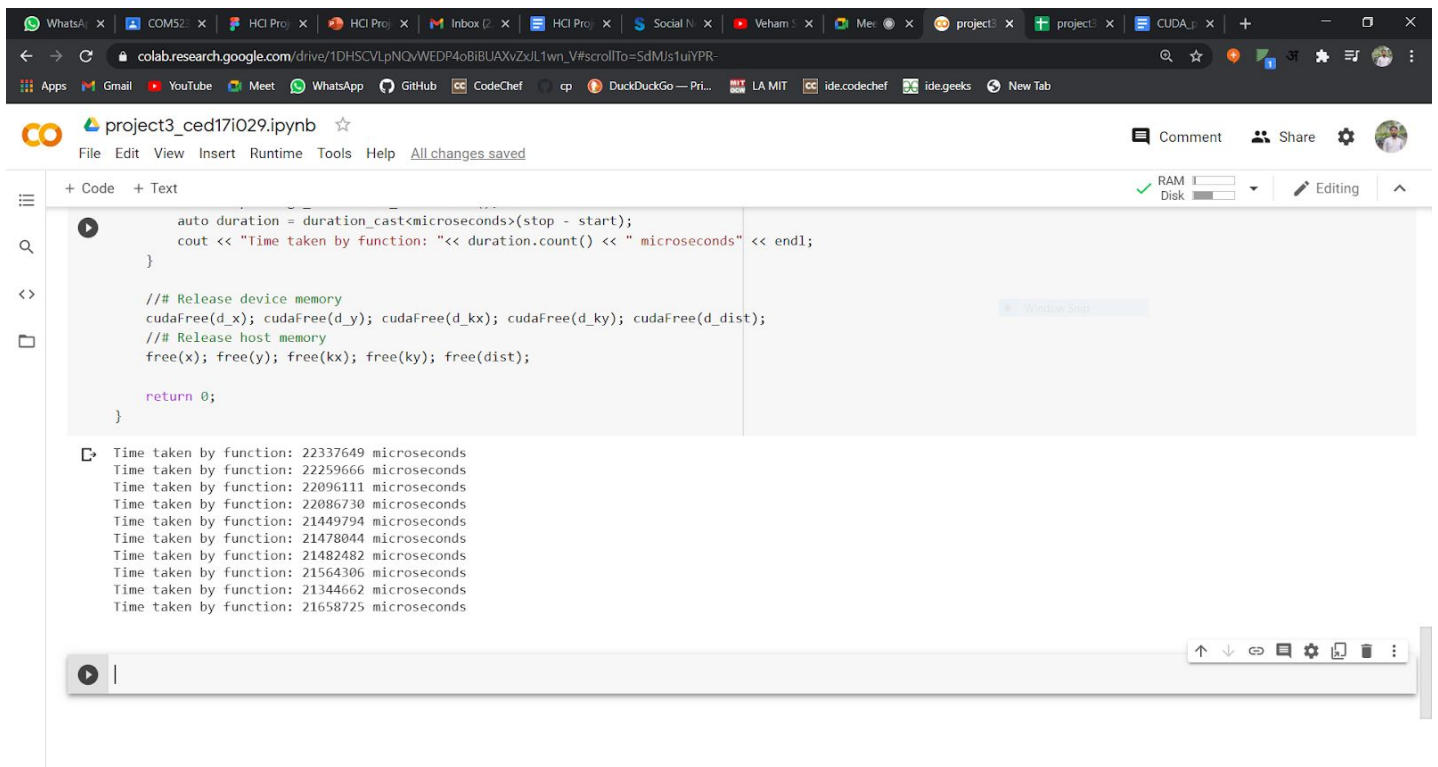
How did i do parallel programming

I have used CUDA technique to parallelize my program.

From my solution to the problem statement, we can clearly see that we have to calculate the distance between pizza center and location of houses. So I have parallelized the program for calculating the distance by running calculateDistance function in device parallelly.

Graph and table

https://docs.google.com/spreadsheets/d/1RiWsOfH94EjRSFfPNMegSW8tQSAgpQqePwF_0k2A04/edit#gid=0



```
auto duration = duration_cast<microseconds>(stop - start);
cout << "Time taken by function: " << duration.count() << " microseconds" << endl;
}

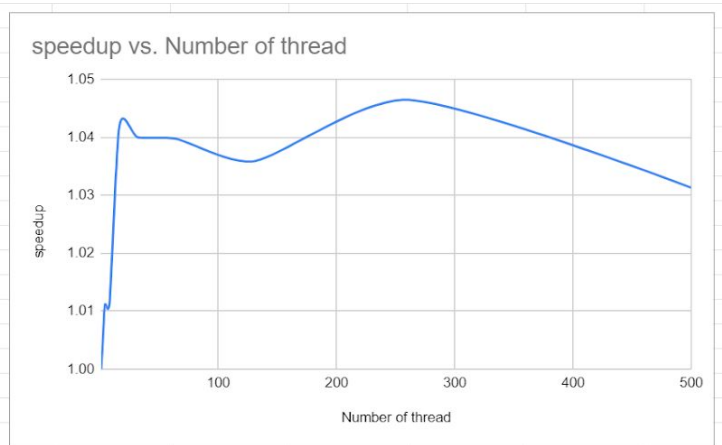
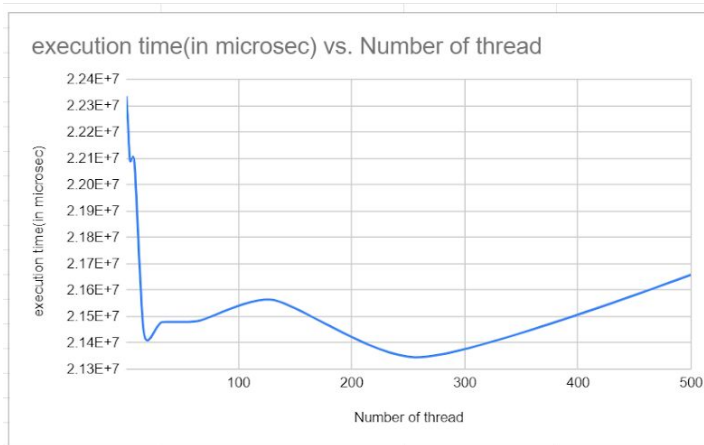
/** Release device memory
cudaFree(d_x); cudaFree(d_y); cudaFree(d_kx); cudaFree(d_ky); cudaFree(d_dist);
/** Release host memory
free(x); free(y); free(kx); free(ky); free(dist);

return 0;
}
```

Time taken by function: 22337649 microseconds
Time taken by function: 22259666 microseconds
Time taken by function: 22096111 microseconds
Time taken by function: 22086730 microseconds
Time taken by function: 21449794 microseconds
Time taken by function: 21478044 microseconds
Time taken by function: 21482482 microseconds
Time taken by function: 21564306 microseconds
Time taken by function: 21344662 microseconds
Time taken by function: 21658725 microseconds

CUDA implementation of project problem

Number of thread	execution time(in microsec)	speedup	parallelization fraction(f)
1	22337649	1	0
2	22259666	1.003503332	0.006982203006
4	22096111	1.010931245	0.0144173931
8	22086730	1.011360622	0.01283772394
16	21449794	1.041392239	0.04239682221
32	21478044	1.040022499	0.03972370564
64	21482482	1.039807644	0.03889133898
128	21564306	1.035862179	0.03489321168
256	21344662	1.046521561	0.04462784184
500	21658725	1.031346444	0.03045461808



Calculation of parallelization fraction

$T(1) = 22337649$ microseconds

Here, for $P = 256$ the execution time is minimum

$T(P) = 21344662$ microseconds

$$\text{Speedup} = \frac{T(1)}{T(P)} = \frac{22337649}{21344662} = 1.046521561$$

From Amdahl's Law,

$$\text{Speedup} = \frac{1}{(f/P) + (1-f)} \text{ Where, } f = \text{Parallelization factor } P = \text{Thread Number}$$

$$\text{So, } f = \frac{(1 - T(P)/T(1))}{(1 - (1/P))}$$

Therefore, $f = 0.04462784184$ which means that approx. 4.5% of the program is parallelized.