



UNIVERSITÀ
DEGLI STUDI
FIRENZE

DIEF

Dipartimento di
Ingegneria Industriale

K-means

Parallel Computing Course

**Antonio Castellucci
Michela Crulli**

1. Introduction

- What is K-means
- What we have done

2. Implementation

- ImagesHandler
- C++ implementation
- OpenMP implementation
- CUDA implementations

3. Dataset

4. Results

- Test platform
- Tests

5. Conclusions

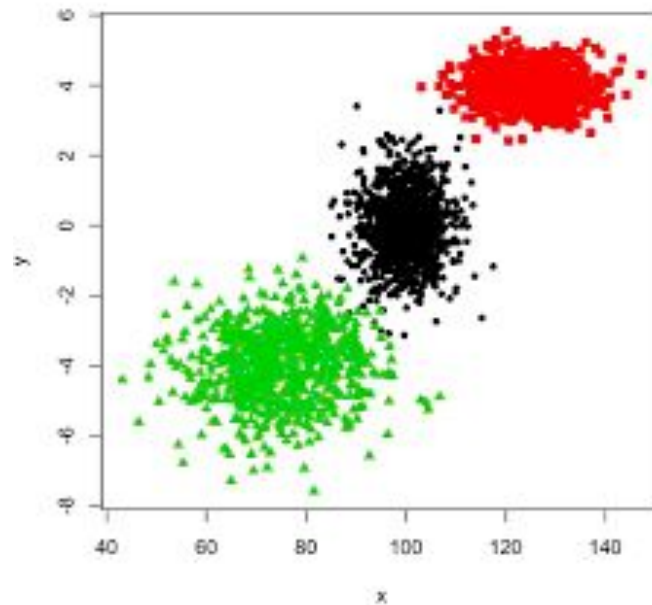
In this project, we realize a program which do **k-means** clustering with three different approaches:

- a **sequential** implementation
 - C++
- three **parallel** implementation
 - OpenMP
 - CUDA naive
 - CUDA with reduction



- What is K-means?

K-means clustering aims to **partition** input data into **k cluster** in which each data belongs to the cluster with the **nearest mean**.

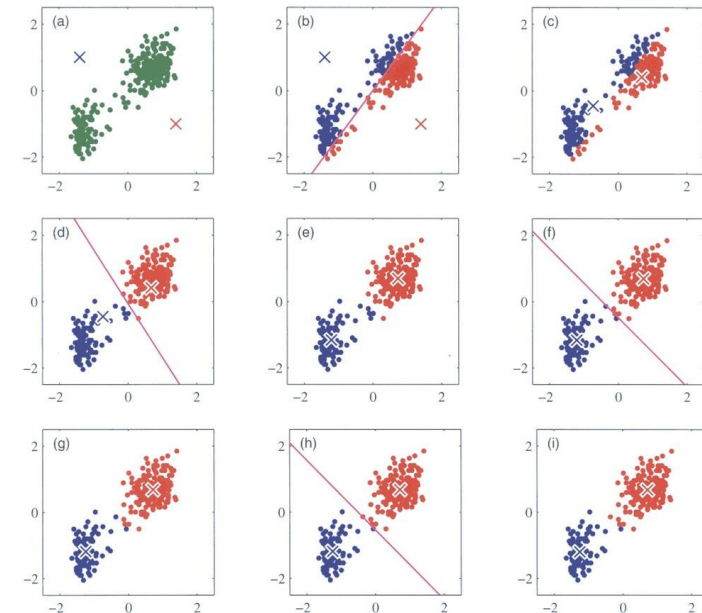


- Steps:

1. Decide number of clusters
2. Distance from each point to each centroids
3. Re-calculate of each centroid with the minimization of:

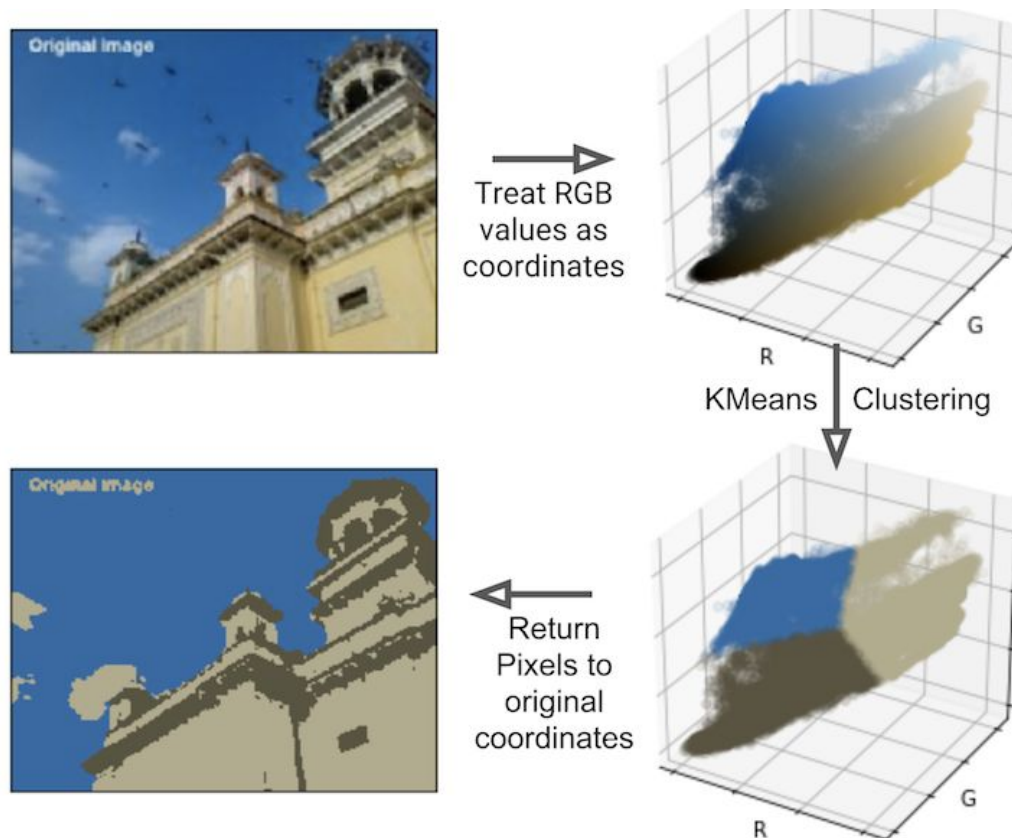
$$J = \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

4. Iterate from step two until in two iteration clusters don't change



- What we have done

K-means applied to **images** to have a better and immediate view



- Input output structure

- scriptable programs with unique interface between three versions

```
./fastCudaKmeans ./testImages/car.jpg RGB 3 500
```



car.jpg

Execution

```
selected file: ./testImages/car.jpg  
mode: RGB  
number of clusters: 3 clusters  
iterations: 500 iterations
```

image processing...

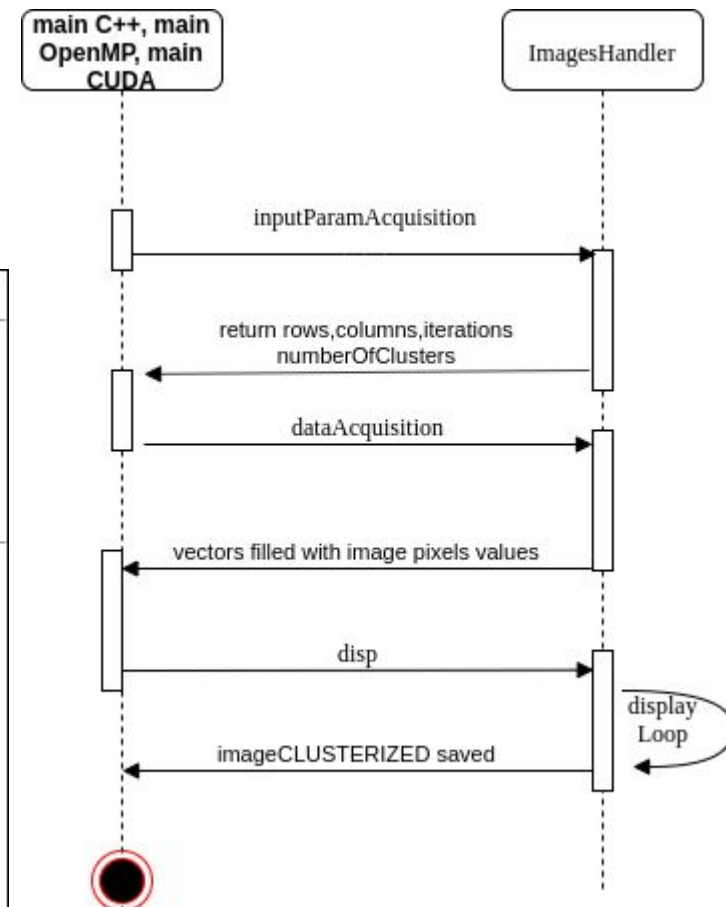
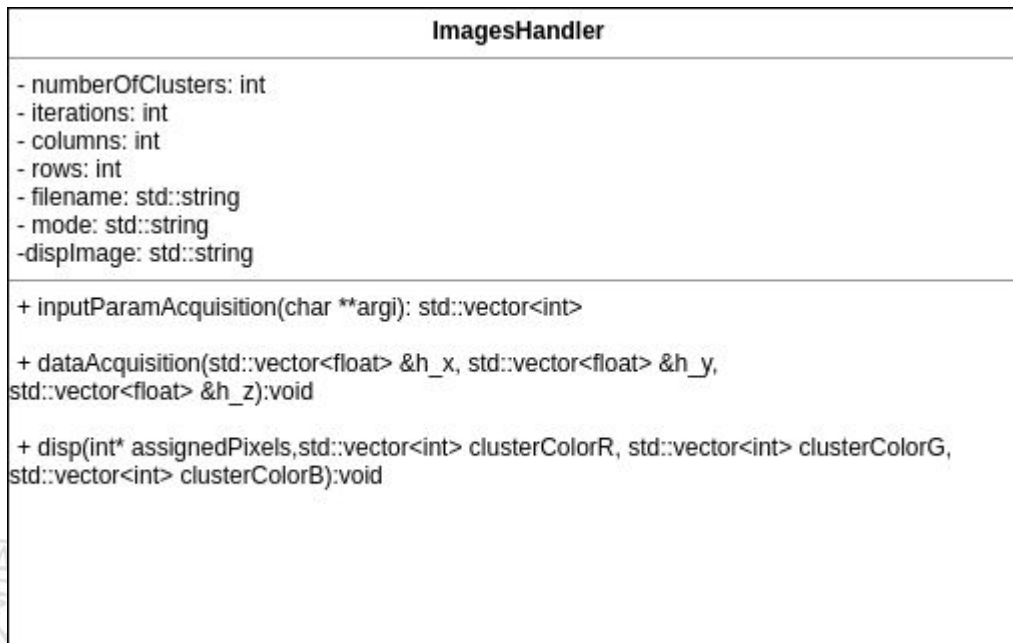
PROCESSING TIME: 0.811982 s



carCLUSTERIZED.jpg

● ImagesHandler

- What is? It's a **common class** that define interfaces for pixels values acquisition and processed data handling
- Why? For the separation of concerns and code reusability



- C++ implementation

Algorithm 1 c++

```
1: Input: data, number_of_cluster, number_of_iterations
2: Output: means, data.assignments
3: for i=0 to number_of_cluster
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   for point = 0 to data.size()
8:      $best\_distance = d(data[point], means[0])$ 
9:      $best\_cluster = 0$ 
10:    for cluster = 1 to number_of_cluster
11:       $dist = d(data[point], means[cluster])$ 
12:      if  $dist < best\_distance$ 
13:         $best\_distance = dist$ 
14:         $best\_cluster = cluster$ 
15:      end if
16:    end for
17:     $data.assignments[point] = best\_cluster$ 
18:  end for
```

```
19:   for point = 0 to data.size()
20:      $cluster = data.assignments[point]$ 
21:      $new[cluster].x += data[point].x$ 
22:      $new[cluster].y += data[point].y$ 
23:      $new[cluster].z += data[point].z$ 
24:      $counts[cluster] += 1$ 
25:   end for
26:   for cluster = 0 to number_of_cluster
27:      $count = counts[cluster]$ 
28:      $means[cluster].x = new[cluster].x / count$ 
29:      $means[cluster].y = new[cluster].y / count$ 
30:      $means[cluster].z = new[cluster].z / count$ 
31:   end for
32: end for
```



- OpenMP implementation

Algorithm 2 openMP

```
1: Input: data, number_of_cluster, number_of_iterations
2: Output: means, data.assignments
3: for i=0 to number_of_cluster
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   #pragma omp parallel for
8:   for point = 0 to data.size()
9:     best_distance = d(data[point], means[0])
10:    best_cluster = 0
11:    for cluster = 0 to number_of_cluster
12:      dist = d(data[point], means[cluster])
13:      if dist < best_distance
14:        best_distance = dist
15:        best_cluster = cluster
16:      end if
17:    end for
18:    data.assignments[point] = best_cluster
19:  end for
```

```
20:   for point = 0 to data.size()
21:     cluster = data.assignments[point]
22:     new[cluster].x += data[point].x
23:     new[cluster].y += data[point].y
24:     new[cluster].z += data[point].z
25:     counts[cluster] += 1
26:   end for
27:   #pragma omp parallel for
28:   for cluster = 0 to number_of_cluster
29:     count = counts[cluster]
30:     means[cluster].x = new[cluster].x/count
31:     means[cluster].y = new[cluster].y/count
32:     means[cluster].z = new[cluster].z/count
33:   end for
34: end for
```

- CUDA Naive implementation

Algorithm 3 *cudaKmeans*

```

1: Input: data, k=number_of_cluster,number_of_iterations
2: Output: means, data.assignments
3: for i=0 to k
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   assign_cluster(data, means, k, new_sums)
8:   cudaDeviceSynchronize()
9:   compute_new_means(means, new_sums)
10:  cudaDeviceSynchronize()
11: end for

```

Algorithm 5 *assign_cluster*

```

1: Input: data, means, new_sums, k
2: Output: data
3: index = blockIdx.x * blockDim.x + threadIdx.x
4: if index  $\geq$  data_size return
5: end if
6: best_distance = d(data[index], means[0])
7: best_cluster = 0
8: for cluster = 1 to k
9:   distance = d(data[index], means[cluster])
10:  if distance < best_distance
11:    best_distance = distance
12:    best_cluster = cluster
13:  end if
14: end for
15: data.assignments[index]=best_cluster
16: atomicAdd(new_sums.x[best_cluster], data.x[point])
17: atomicAdd(new_sums.y[best_cluster], data.y[point])
18: atomicAdd(new_sums.z[best_cluster], data.z[point])
19: atomicAdd(counts[best_cluster], 1)

```

Algorithm 4 *compute_new_means*

```

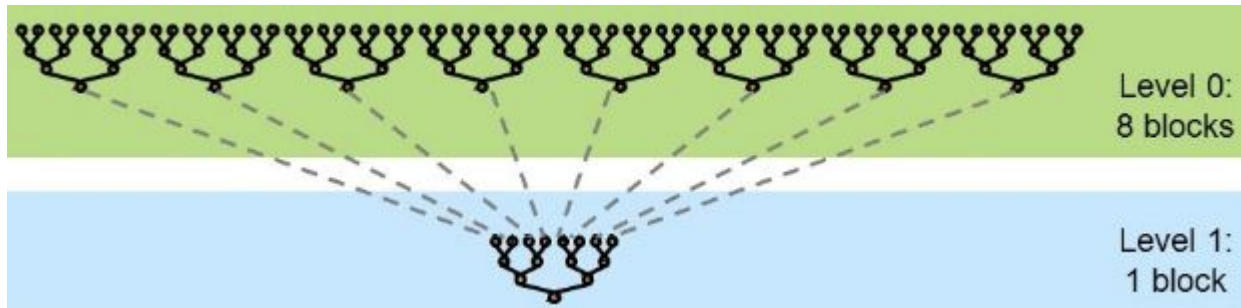
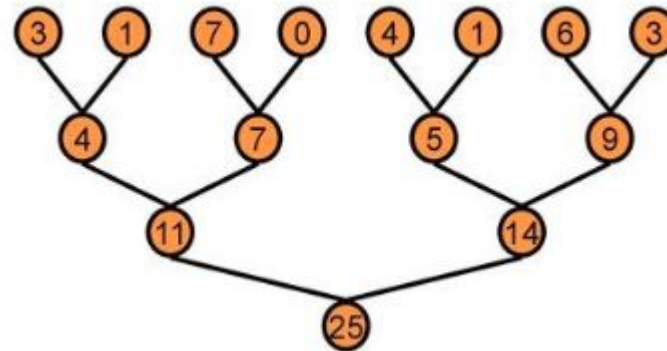
1: Input: means, new_sums
2: Output: means
3: cluster = threadIdx.x
4: means.x[cluster] = new_sum.x[cluster] / counts[cluster]
5: means.y[cluster] = new_sum.y[cluster] / counts[cluster]
6: means.z[cluster] = new_sum.z[cluster] / counts[cluster]

```

- Reduction

Mark Harris tricks:

- Unrolling loops
- Multiple adds/threads



- fastCuda 1.0

Algorithm 6 fastCudaKmeans 1.0

```

1: Input: data, k=number_of_cluster, number_of_iterations
2: Output: processed data
3: for i=0 to k
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   assign_cluster(data, means, k)
8:   cudaDeviceSynchronize()
9:   for cluster = 0 to k
10:    populate(data, clusterData, cluster)
11:    cudaDeviceSynchronize()

```

means: #Cluster (k)

x	250	100		
y	250	100		
z	250	100		
ass	0	0		

↑ ↑

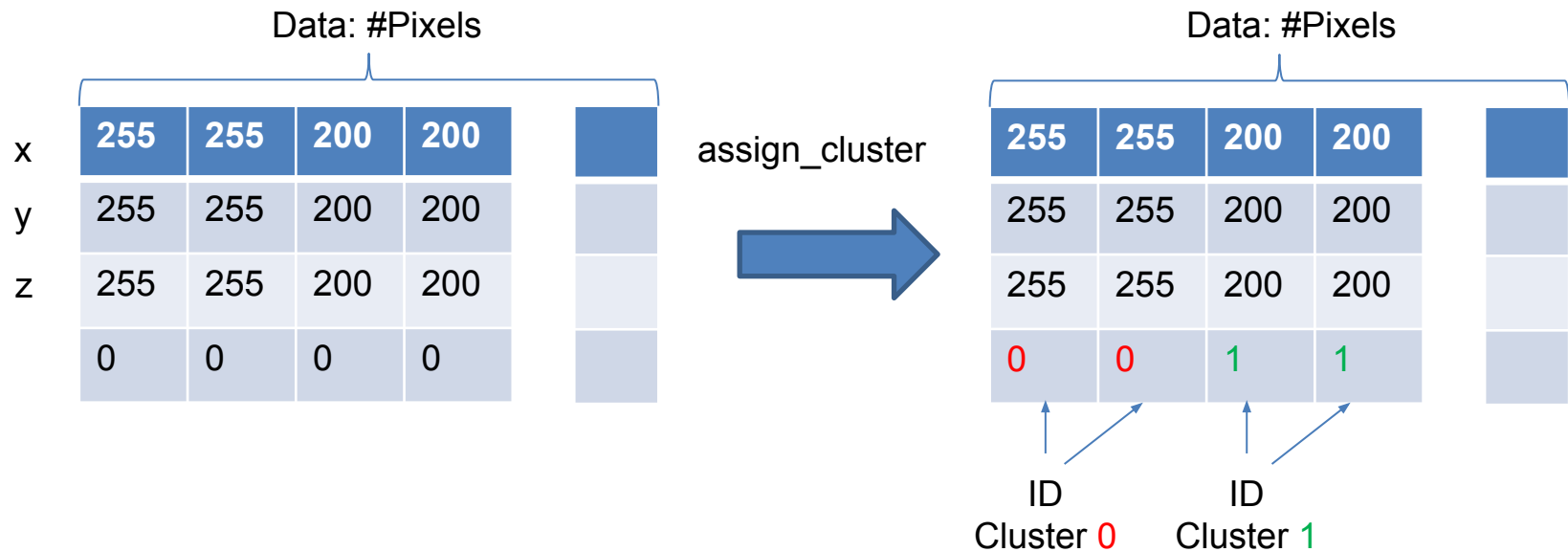
ID ID

Cluster 0 Cluster 1

data: #Pixels

x	255	255	200	200
y	255	255	200	200
z	255	255	200	200
	0	0	0	0

- `assign_cluster(data, means, k)`



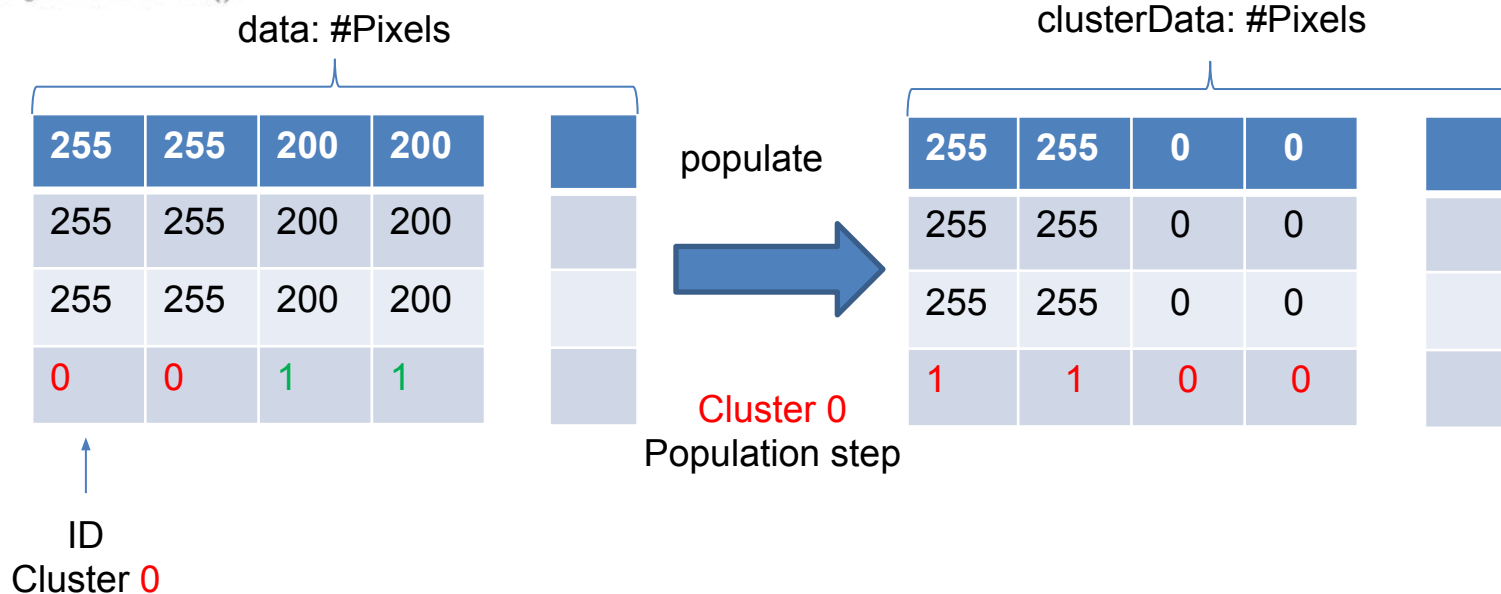
• fastCuda 1.0

Algorithm 6 fastCudaKmeans 1.0

```

1: Input: data, k=number_of_cluster,number_of_iterations
2: Output: processed data
3: for i=0 to k
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   assign_cluster(data, means, k)
8:   cudaDeviceSynchronize()
9:   for cluster = 0 to k
10:    populate(data, clusterData, cluster)
11:    cudaDeviceSynchronize()

```



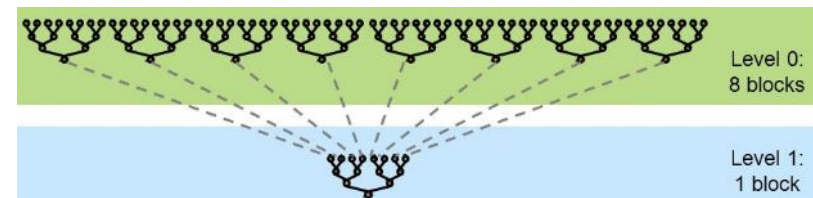
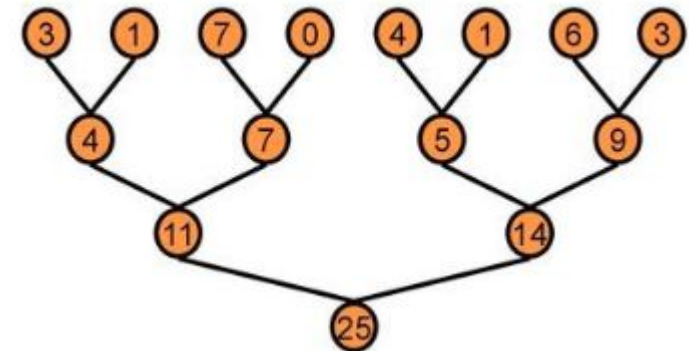
- fastCuda 1.0

Algorithm 6 fastCudaKmeans 1.0

```

1: Input: data, k=number_of_cluster,number_of_iterations
2: Output: processed data
3: for i=0 to k
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   assign_cluster(data, means, k)
8:   cudaDeviceSynchronize()
9:   for cluster = 0 to k
10:    populate(data, clusterData, cluster)
11:    cudaDeviceSynchronize()
12:    n = number_of_pixels
13:    do
14:      blocksPerGrid = n/BLOCKSIZE
15:      reduction(clusterData, tmp, n)
16:      cudaDeviceSynchronize()
17:      clusterData = tmp
18:      n = blocksPerGrid
19:    while n > BLOCKSIZE

```



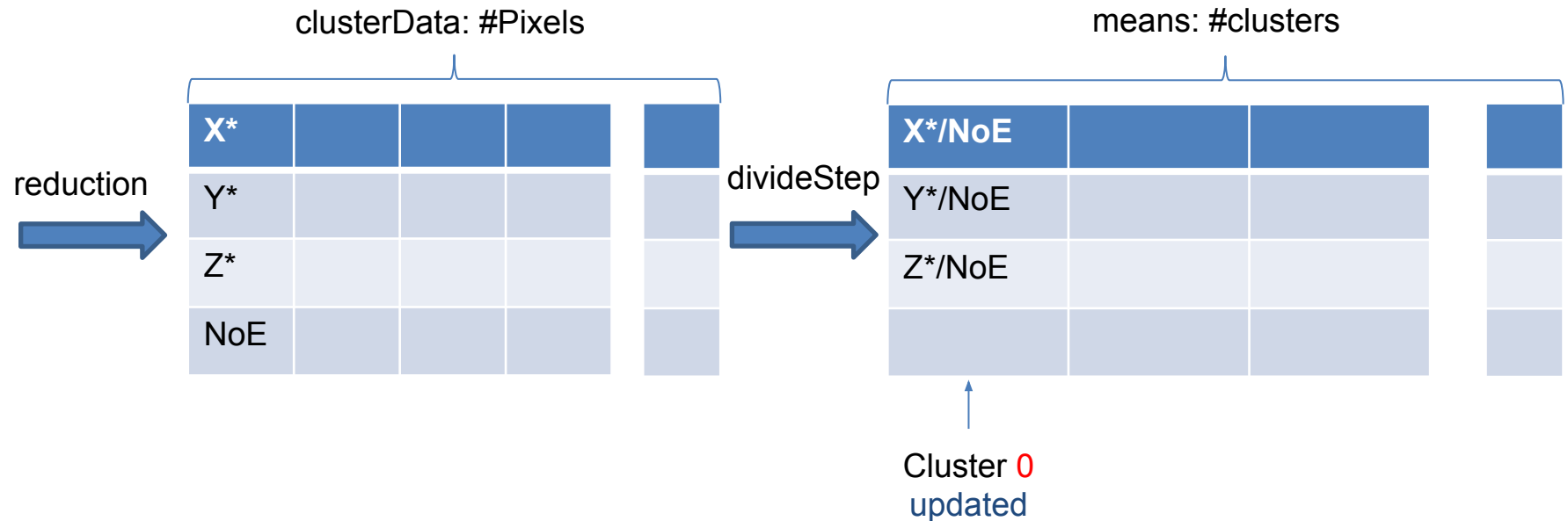
- fastCuda 1.0

Algorithm 6 fastCudaKmeans 1.0

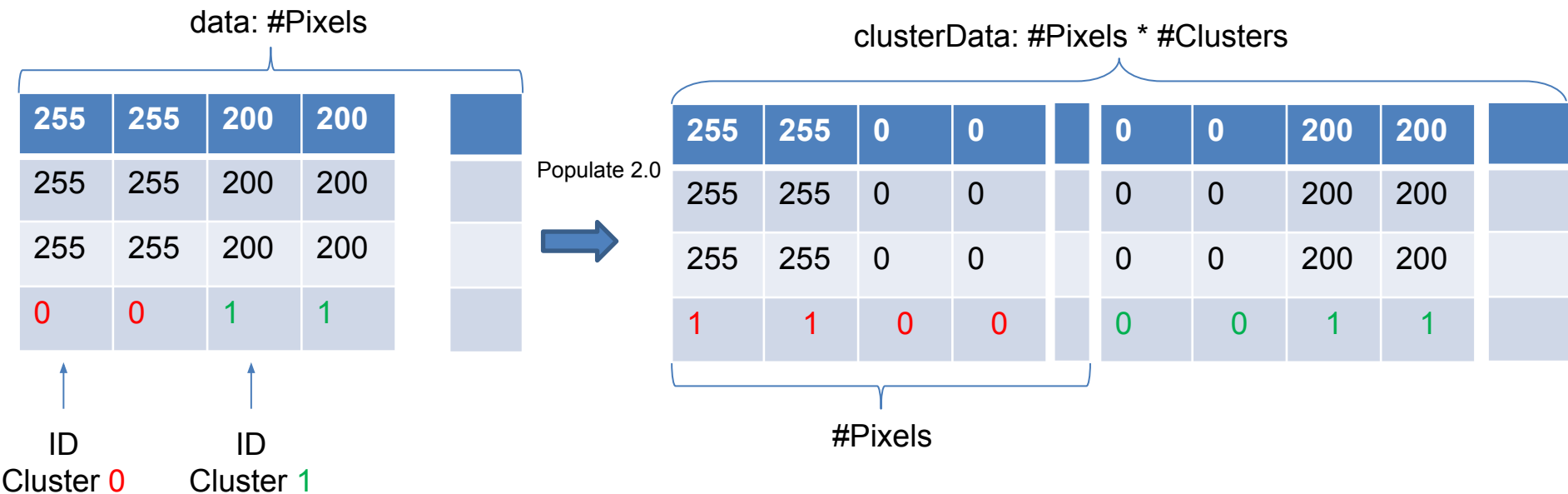
```
1: Input: data, k=number_of_cluster,number_of_iterations
2: Output: processed data
3: for i=0 to k
4:   centroids population in means[i]
5: end for
6: for iteration = 0 to number_of_iterations
7:   assign_cluster(data, means, k)
8:   cudaDeviceSynchronize()
9:   for cluster = 0 to k
10:    populate(data, clusterData, cluster)
11:    cudaDeviceSynchronize()
12:    n = number_of_pixels
13:    do
14:      blocksPerGrid = n/BLOCKSIZE
15:      reduction(clusterData, tmp, n)
16:      cudaDeviceSynchronize()
17:      clusterData = tmp
18:      n = blocksPerGrid
19:    while n > BLOCKSIZE
20:      reduction(tmp, tmp, n)
21:      cudaDeviceSynchronize()
22:      divideStep(means, tmp, cluster, k)
23:      cudaDeviceSynchronize()
24:    end for
25: end for
```



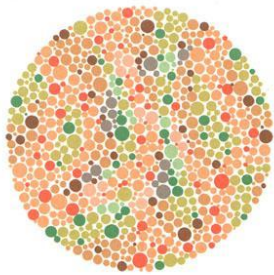
- `divideStep(means,tmp,cluster,k)`



- populate 2.0 (data, clusterData, ~~cluster~~)



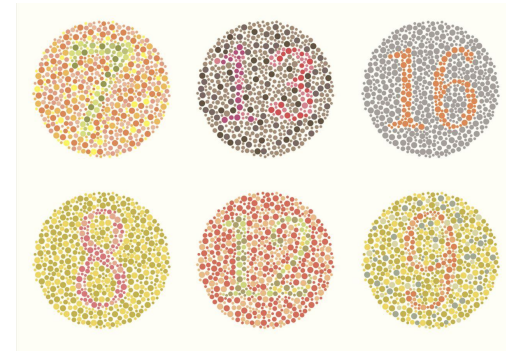
We use images of different sizes :



100000 pixels



2 millions pixels



1 millions pixels



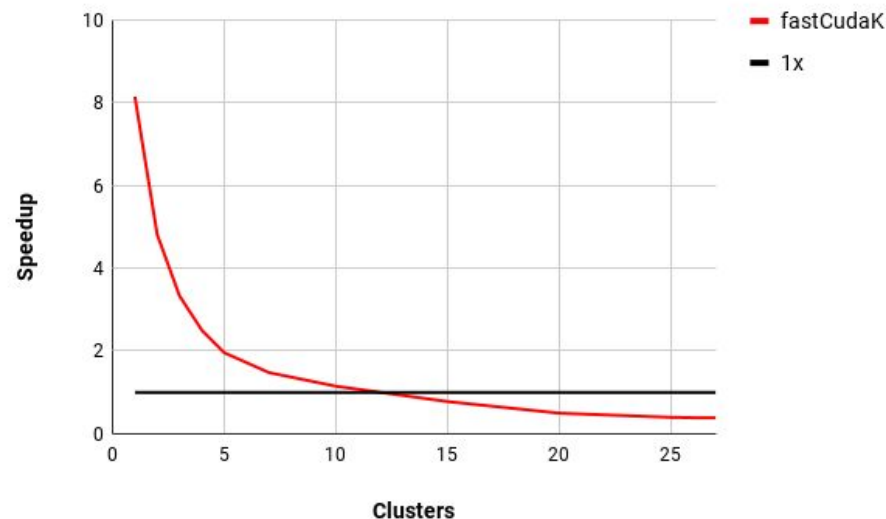
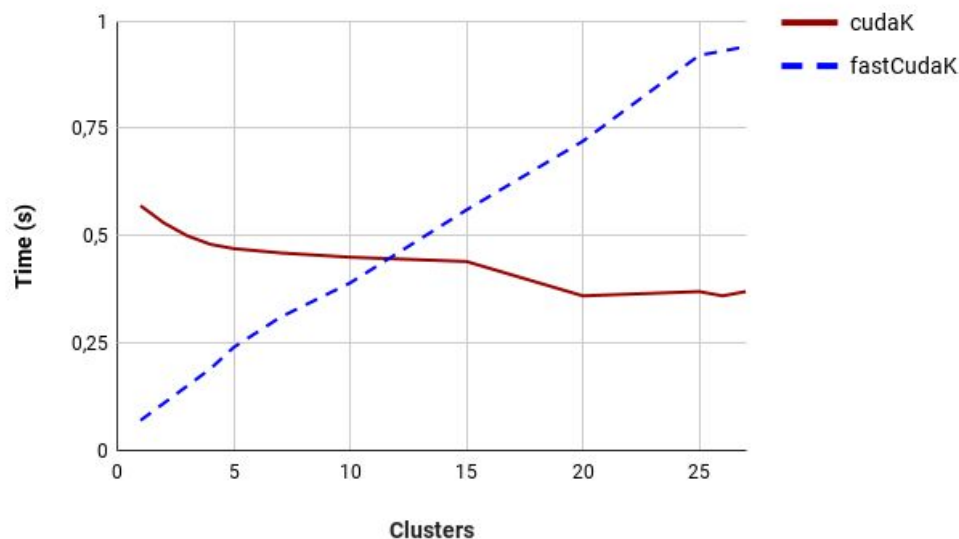
10 millions pixels

Test Platform

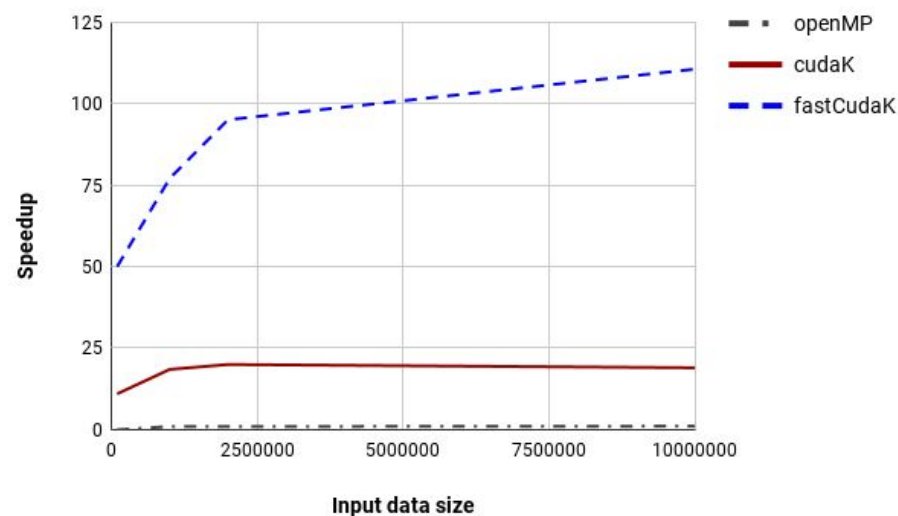
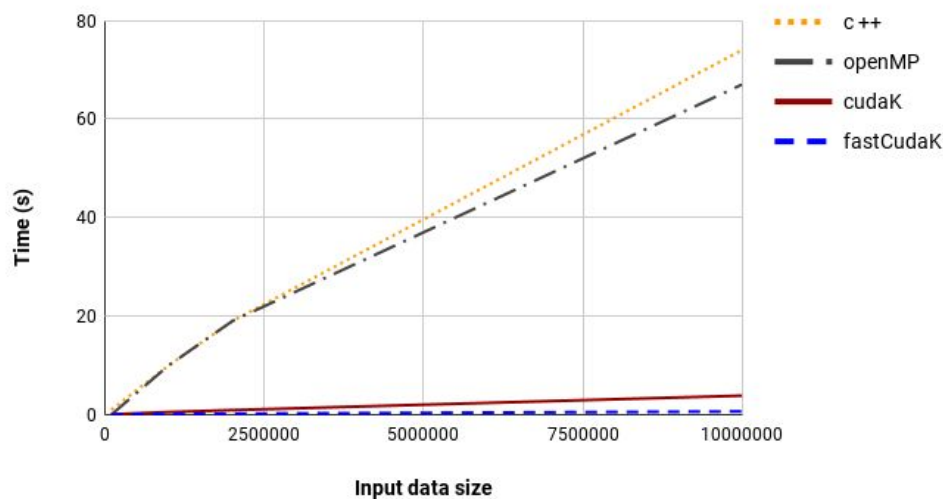
MSI GS65 with:

- i7-8750H 6 CORES 12 THREAD
- 16 GB 2400 MHz DDR4 RAM
- RTX 2060 6GB GDDR6

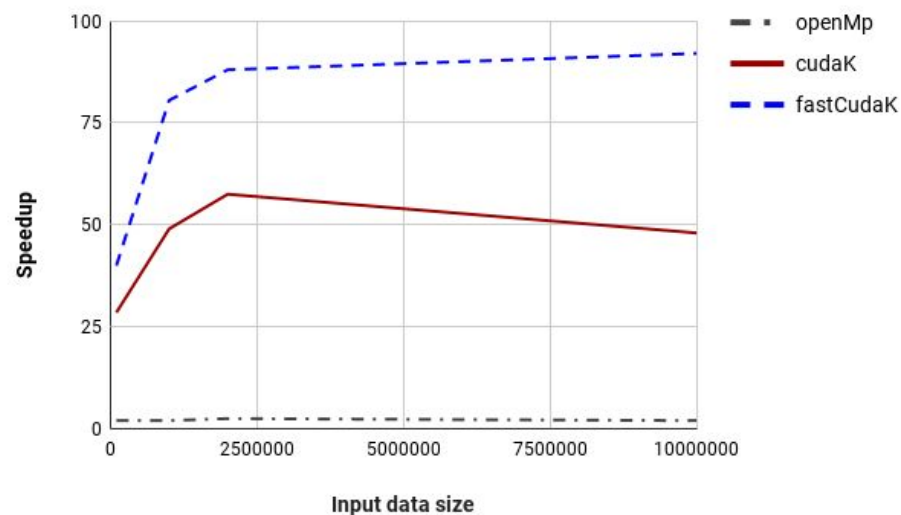
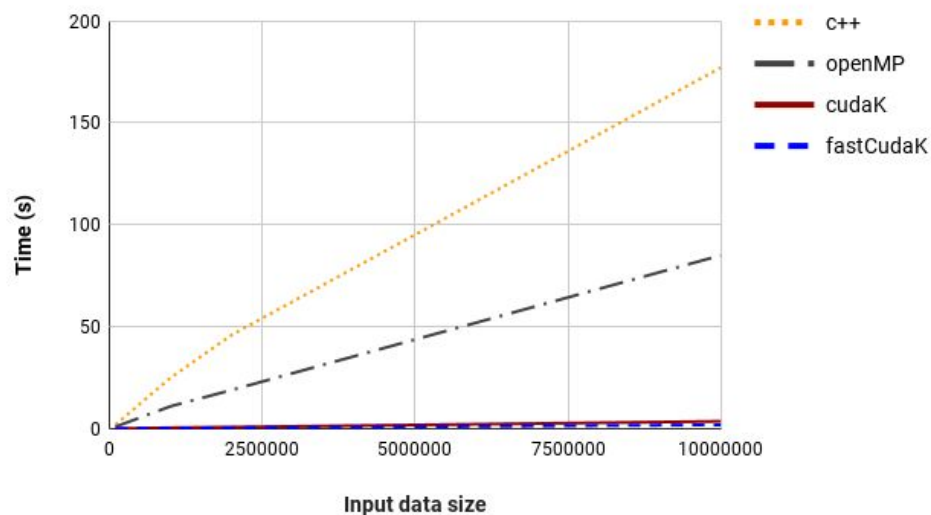
- cudaKmeans Vs fastCudaKmeans 2.0 with 2 millions pixels, 100 iterations as the number of clusters changes



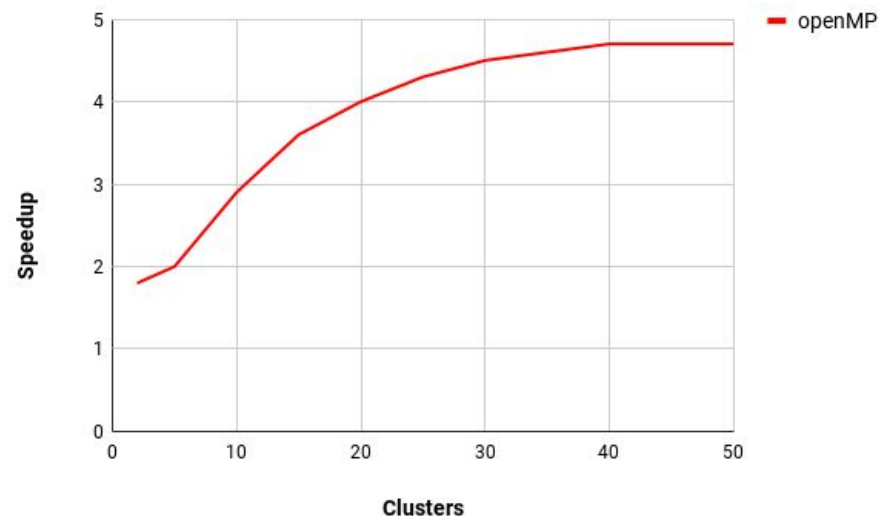
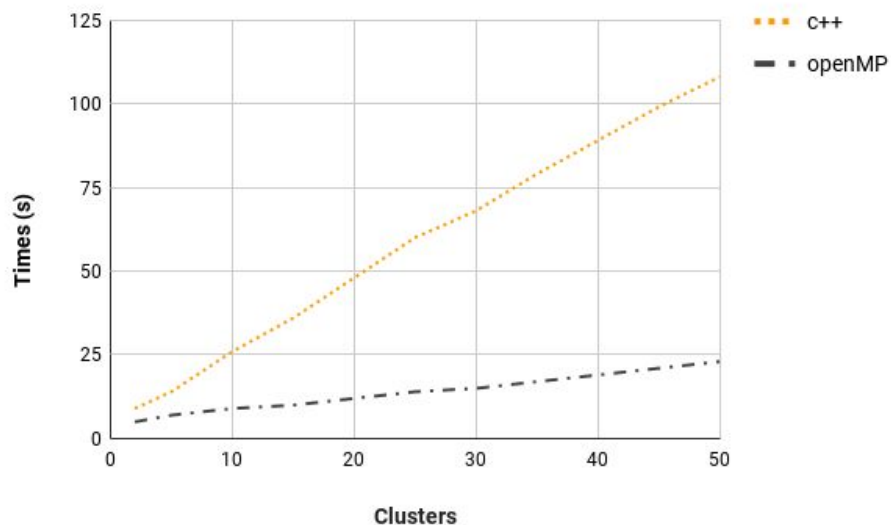
- All versions compared respect to input data size with 2 clusters and 200 iterations



- All versions compared respect to input data size with 7 clusters and 200 iterations



- C++ Vs OpenMP with 10 millions pixels compared respect to the number of clusters changes



- Understand CUDA performance characteristics
- Use peak performance metrics to guide optimization
- Understand parallel algorithm complexity theory, identify type of bottleneck e.g. memory, core computation, or instruction overhead
- Optimize the algorithm, then unroll loops
- Use template parameters to generate optimal code

Thanks for the attention!