

Clustering algorithms on GPU for video processing

Jonathan Fin, Mat: 256178, jonathan.fin@studenti.unitn.it, *GitRepo*:
<https://github.com/Nathanoj02/GPU-Computing-2025-256178>

Abstract—This work presents a GPU-accelerated implementation of the k-means clustering algorithm, tailored for high-dimensional video data. The study evaluates different initialization strategies, distance metrics, and optimization techniques, with particular focus on OpenACC for parallelization. Experimental results show that the squared Euclidean distance achieves the best trade-off between clustering accuracy and computational efficiency, while GPU implementations consistently outperform the CPU baseline across multiple resolutions. The optimized GPU algorithm demonstrates near real-time performance at high resolutions, and a calibration strategy further reduces iteration counts for static video sequences. Additionally, comparisons of alternative approaches highlight the advantages of calibration over k-means++ for video data. The findings confirm the effectiveness of GPU acceleration for clustering tasks and suggest promising directions for further optimization, including CUDA-based low-level implementations and multi-GPU scaling for larger datasets.

Index Terms—K-means, Clustering, GPU, OpenACC, CUDA, Parallelization, Video Processing

I. INTRODUCTION

The aim of this project is to develop an efficient algorithm for performing clustering on large datasets using a GPU, leveraging OpenACC to parallelize the most computationally expensive steps.

Clustering is a fundamental operation in unsupervised learning, widely used for tasks such as data analysis, pattern recognition, and dimensionality reduction. Efficient clustering algorithms are crucial for extracting meaningful patterns from large datasets while minimizing computational time.

In this report, the focus is on the k-means algorithm, specifically applied to clustering video data. Video datasets are particularly challenging due to their high dimensionality and temporal complexity, making GPU acceleration essential for practical performance. The project investigates strategies for parallelizing k-means on the GPU, evaluates their effectiveness in terms of speedup and scalability, and explores computational optimizations as well as heuristics for choosing the initial centroids in video data.

II. PROBLEM STATEMENT

Clustering is a fundamental task in unsupervised machine learning, aiming to group data points into **clusters** such that points within the same cluster are more similar to each other than to points in other clusters. It is widely used in fields such as data analysis, image and video processing, pattern recognition, and recommendation systems.

Among clustering algorithms, **k-means** is one of the most popular due to its simplicity and efficiency. The algorithm partitions a dataset into k clusters by iteratively assigning

each data point to the nearest **centroid** and then updating the centroids as the mean of the points in each cluster. Despite its simplicity, k-means can be computationally intensive for large or high-dimensional datasets, especially when applied to video data where each frame contributes multiple dimensions.

K-means clustering is widely used in unsupervised learning and has several standard optimizations:

- **Centroid Initialization:** k-means++ is a widely adopted method to improve convergence by carefully selecting initial centroids based on data distribution;
- **Distance Computation:** efficient calculation of distances between data points and centroids is critical. GPU implementations often parallelize this computation for large datasets;
- **Cluster Evaluation:** metrics for selecting the optimal number of clusters (k) are essential for practical applications, especially in complex datasets such as videos.

A. Parallelization

The focus of this project is on **parallelization** strategies for k-means on the GPU. Initial CPU implementations serve as a baseline for performance evaluation. Early GPU implementations revealed performance bottlenecks, which motivated profiling and further optimization. Key improvements include optimized memory allocation (e.g., using `malloc` efficiently), exploring centroid selection heuristics, and calibrating the algorithm for video datasets.

Given the foundational nature of this work, the primary objective is to experiment with different parallelization strategies and computational optimizations, benchmark their performance, and explore heuristics for cluster selection in video data, rather than to propose novel clustering algorithms.

III. STATE OF THE ART

The state-of-the-art solutions for clustering large-scale datasets leverage GPU acceleration to reduce computational time, particularly for high-dimensional data. Libraries such as cuML [1] provide GPU-accelerated implementations of clustering algorithms, including k-means, and are optimized for dense datasets. These libraries typically support multiple precision operations and include routines for distance computation, centroid initialization, and iterative refinement. Notably, cuML's Python API mirrors that of scikit-learn, enabling seamless integration into existing Python-based machine learning workflows.

In addition to cuML, there are other notable implementations of GPU-accelerated k-means clustering. For instance, the

repository [2] presents a highly-optimized CUDA implementation of the k-means algorithm. This approach, documented in a conference paper [3], focuses on optimizing the kernel execution and memory access patterns to enhance performance on NVIDIA GPUs.

These existing solutions provide a foundation for GPU-accelerated k-means clustering. However, they often require adaptation to specific use cases, such as clustering video data. This project takes inspiration from these foundations to implement a GPU-accelerated k-means algorithm tailored for video data, incorporating optimizations for memory allocation, centroid initialization, and distance computation, and evaluating its performance against CPU-based implementations.

IV. METHODOLOGY AND CONTRIBUTIONS

The work in this project can be divided into several stages, starting from the CPU baseline implementation and progressively moving towards optimized GPU-based solutions tailored for video clustering. Each stage was designed to address specific computational challenges, improve efficiency, and investigate algorithmic trade-offs.

A. CPU Baseline

A CPU implementation of the k-means algorithm was first developed as a reference baseline. Several **distance metrics** were explored to measure similarity between data points and centroids, allowing for a fair comparison of their computational cost and impact on clustering quality. This stage served as the ground truth for validating GPU-based implementations.

Algorithm 1 report a simple k-means clustering algorithm, where the *nearest* centroid is calculated using the various distance functions.

Algorithm 1 K-means Clustering

Input: Dataset $X = \{x_1, x_2, \dots, x_n\}$, number of clusters k
Output: Cluster assignments and centroids $C = \{c_1, c_2, \dots, c_k\}$
Initialize centroids C (e.g., random or with K-means++)
repeat
 foreach $x_i \in X$ **do**
 Assign x_i to nearest centroid c_j
 foreach cluster j **do**
 Update c_j as mean of all points assigned to it
until *convergence*;

B. Cluster Evaluation Metrics

To identify the optimal number of clusters k , two widely used evaluation methods were employed:

- **Elbow Method:** evaluates the total within-cluster sum of squares (WCSS) across different values of k , selecting the point where the marginal gain in performance decreases;
- **Silhouette Method:** measures cluster cohesion and separation by comparing intra-cluster distance to inter-cluster distance, with higher silhouette scores indicating better clustering quality.

These methods ensured that the chosen value of k reflected meaningful structures in the dataset, rather than being arbitrarily selected.

C. GPU Implementations

The first GPU implementation of k-means revealed that performance was slightly better than the CPU baseline but never achieved more than a $2\times$ speedup.

Memory Bound Behaviour: the computational structure of the k-means algorithm indicates a **memory-bound** nature. Each iteration involves loading all data points and centroid coordinates from global memory to compute pairwise distances, followed by storing updated assignments. The arithmetic intensity is therefore very low, meaning that execution time is primarily limited by memory bandwidth.

Using *NVIDIA Nsight Systems* (*nsys*), profiling confirmed bottlenecks primarily in memory transfers. To address these, multiple GPU kernels and algorithmic variations were explored:

- **K-means with Random Initialization:** centroids chosen randomly from the dataset;
- **K-means++ Initialization:** centroids initialized with the k-means++ heuristic to improve convergence;
- **Centroids from Image Pixels:** Initial prototypes selected directly from random pixels of the input video frames;
- **Tiling:** introduced memory tiling strategies to improve memory coalescing and reduce global memory latency.
- **Convergence Skipping:** convergence checks skipped for a fixed number of iterations to reduce synchronization overhead.

The tiling optimization was tested to improve memory coalescing, but experimental results confirmed that tiling is the wrong tool for this job. Unlike convolutional or stencil-based algorithms, k-means does not benefit from spatial locality since pixels are treated as independent with no neighborhood dependencies. Therefore, tiling introduced overhead without measurable speedup. The best-performing version was obtained using k-means++ initialization combined with optimized memory transfers, and so was adopted as the primary GPU solution.

Algorithm 2 shows the k-means++ initialization used in the project, while algorithm 3 shows a more detailed GPU implementation of the k-means algorithm.

Algorithm 2 K-means++ Initialization (Distance-Max Variant)

Input: Dataset $X = \{x_1, x_2, \dots, x_n\}$, number of clusters k , dimensionality d
Output: Initial centroids $C = \{c_1, c_2, \dots, c_k\}$
Choose first centroid c_1 randomly
for $j = 2$ **to** k **do**
 For each $x \in X$, compute its distance to the nearest centroid in C
 Let c_j be the point $x \in X$ with the maximum such distance
 Add c_j to C

Algorithm 3 GPU-Accelerated K-means Clustering

Input: Dataset X of n pixels with d dimensions, number of clusters k , stability threshold ϵ , maximum iterations T

Output: Cluster assignments

Initialize centroids C using K-means++

for $t = 1$ **to** T **do**

 Save current centroids C^{old}

 // Assignment step (parallel over pixels)

for each pixel $x_i \in X$ **do in parallel**

 Assign x_i to the nearest centroid in C

 // Update step (parallel accumulation)

for each cluster j **and dimension** d **do in parallel**

 Compute the mean of all pixels assigned to j

 Update centroid c_j with this mean

 // Convergence check

 Compute maximum centroid shift $\Delta = \max_j \|c_j - c_j^{old}\|^2$

if $\Delta \leq \epsilon$ **then**

break

for each pixel $x_i \in X$ **do in parallel**

 Replace x_i with its assigned centroid

In the most optimized version of the algorithm, centroid initialization is done in parallel to achieve minimum computation time. All the instructions inside the for loop are done *in GPU memory*, to avoid costly memory transfers just to perform convergence checks.

The parallel *Update step* was originally handled by having all threads **atomically** update global sum and count arrays for each cluster, necessary to ensure safety but introducing some bottlenecks in the execution. To optimize throughput, a **2-step reduction** approach was introduced. Threads compute partial sums **locally**, which are then combined in parallel using OpenACC's *reduction* clause, eliminating atomic contention and improving performance.

D. GPU Calibration for Video Clustering

To further adapt the algorithm to video data, a **calibration** strategy was introduced. A subset of frames (e.g., 10 frames) was clustered first to identify representative prototypes. These calibrated centroids were then reused as initial prototypes throughout the entire video. This approach reduced the variance introduced by random initialization, stabilized clustering results across frames, and improved overall performance in video analysis tasks.

V. SYSTEM DESCRIPTION AND EXPERIMENTAL SET-UP

The experiments were executed on the university's high-performance computing (HPC) cluster using the configuration detailed in Table I and Table II. Compilation was performed with the NVIDIA HPC compiler (*nvc*) using the flags `-O2`, `-acc` and `-gpu=cc80`. These flags enabled aggressive optimizations and GPU-specific targeting for compute capability 8.0. The use of the `-acc` flag enabled OpenACC directives, which allow GPU parallelization to be expressed at a higher

level. By doing so, the focus of this project remains on investigating algorithmic challenges and optimization strategies, rather than dealing with low-level implementation details of GPU programming.

To ensure reliability and minimize the impact of scheduling noise, each experiment was repeated 30 times with at least 3 warm-up cycles, with the results averaged across runs. This methodology reduces variability caused by fluctuations in operating system priorities and ensures consistent benchmarking conditions.

A. System Description

All the systems specifications are described in Table I and Table II.

System Component	Specification
CPU	Intel Xeon Silver 4309Y
Cores / Frequency	8 cores @ 2.80 GHz
GPU Accelerator	NVIDIA A30
GPU Memory	24 GB HBM2e
Compute Capability	8.0

TABLE I
HARDWARE SPECIFICATIONS

Software Component	Version
Operating System	Rocky Linux 9.3
CUDA Toolkit	12.5
C++	11.4.1
NVIDIA HPC Compiler (NVC)	24.7-0

TABLE II
SOFTWARE ENVIRONMENT SPECIFICATIONS

B. Dataset description

The dataset used for experimentation consists of frames extracted from a single high-definition video source at 1080p resolution. To analyze the scalability and robustness of the k-means clustering algorithm across different input sizes, the same video was also downsampled to 720p, 480p, and 240p resolutions. This multi-resolution approach enables the evaluation of performance trade-offs as a function of input dimensionality, allowing for a thorough comparison of execution time and clustering quality across varying dataset sizes. By using the same video content across all resolutions, the experiments maintain consistency in data distribution while highlighting the computational effects of scaling.

VI. EXPERIMENTAL RESULTS

A. Evaluation metrics

Elbow method and **average silhouette method** were run on the first frame of the test video to determine the best k to

Video	Width (pixels)	Height (pixels)	Total pixels per frame
walking_1080.mp4	1080	1920	$\sim 2.07M$
walking_720.mp4	720	1280	$\sim 922k$
walking_480.mp4	480	854	$\sim 410k$
walking_240.mp4	240	426	$\sim 102k$

TABLE III
DATASET

adopt in following calculations. The most suitable k value is usually low, so the tests were run on k ranging from 1 to 10 included.

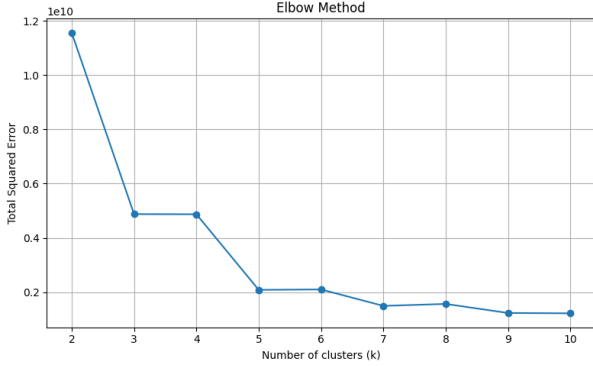


Fig. 1. Elbow method

The Elbow method (Figure 1) initially suggested $k = 3$ as a potential candidate, but its visual nature makes it inherently subjective. The Silhouette method (Figure 2) quantitatively confirmed $k = 3$ as optimal, since it achieved the maximum average silhouette score (~ 0.634) among all tested values. This indicates that clusters are well-separated while maintaining internal cohesion — a desirable property for color-based video segmentation, where three dominant regions typically correspond to background, foreground, and transitional areas. Therefore, $k = 3$ was adopted in all subsequent experiments.

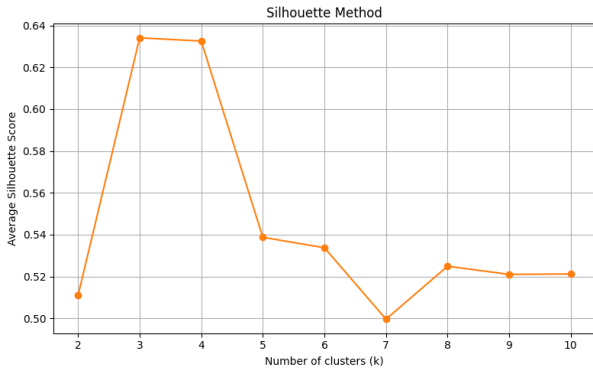


Fig. 2. Average silhouette method

B. Distance Metrics Evaluation

To determine the most efficient and effective distance metric for k-means clustering on video frames, several distance

measures were tested on CPU. The total execution times were recorded for each metric and are summarized in Figure 3.

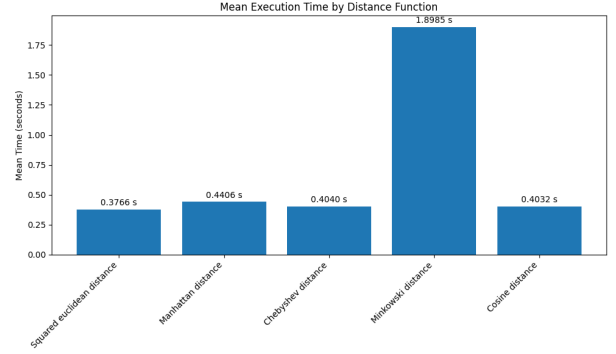


Fig. 3. Execution times of different distance metrics on CPU

From the chart, *squared Euclidean distance* emerges as the best choice, balancing both computational efficiency and clustering quality. The squared form is preferred over the standard Euclidean distance because it avoids costly square root operations while preserving the relative distances necessary for clustering. This distance metric is also the one producing the most visually accurate clusters in the sample video frames.

C. GPU Implementation Performance

The CPU baseline was compared with multiple GPU implementations to assess speedup and performance improvements. Figure 4 reports execution times for each GPU variant, including initial and optimized kernels. A dashed line indicates the CPU baseline.

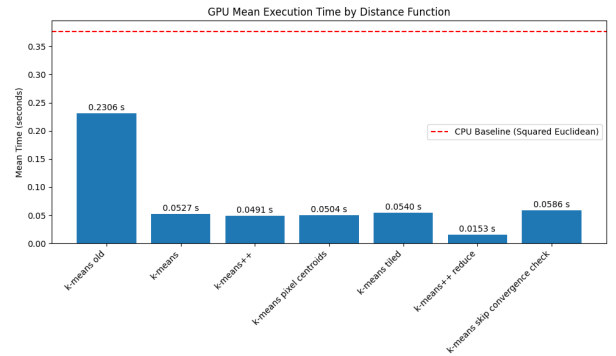


Fig. 4. Execution times of GPU implementations

The charts highlight significant **speedup** in GPU implementations compared to CPU execution. The initial GPU kernel was slower than expected due to memory access patterns, while optimized kernels with improved memory transfers and centroid heuristics achieved a speedup of approximately $24.6\times$ over CPU. The corresponding processing rate reaches approximately 65.4 FPS for standard 1080p video frames, compared to the CPU's ~ 2.7 FPS.

D. Multi-Resolution Video Performance

To evaluate the scalability of the GPU implementations, the algorithm was tested across multiple video resolutions: 1080p, 720p, 480p, and 240p.

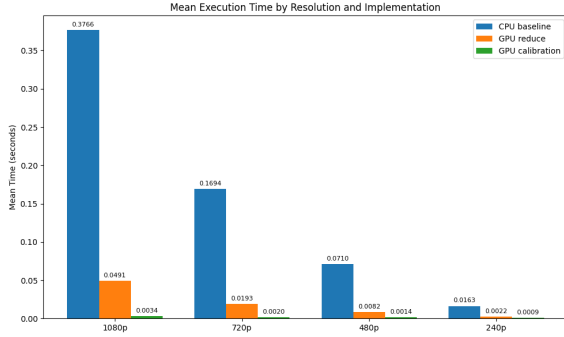


Fig. 5. Execution time comparison across video resolutions for CPU, GPU, and GPU with calibration

Figure 5 provides a comprehensive view of the algorithm’s scalability and the impact of the final “Calibration” optimization.

- **Scalability:** The performance of all implementations (CPU, GPU, GPU Calibrated) scales in a near-linear fashion with the number of pixels. This is expected, as it confirms the algorithm’s $O(nk)$ complexity.
- **Calibration Impact:** The “GPU with calibration” strategy is the clear winner, achieving a massive $109.5\times$ speedup over the CPU baseline at 1080p. This optimization pushes the performance from ~ 65.4 FPS (standard GPU) to ~ 290.7 FPS (calibrated GPU), moving the solution from “usable” to “true real-time” for full HD video. This proves that the domain-specific heuristic (calibration) provides a far greater performance boost than the k-means++ initialization alone.
- **GPU Underutilization:** An important artifact is visible at lower resolutions (240p and 480p). The relative speedup of the GPU versions decreases. This indicates that the problem size (total pixels) is too small to fully saturate the GPU’s thousands of cores. This is a classic example of *Amdahl’s Law* in practice: the fixed overhead of launching kernels and managing data becomes a larger fraction of the total execution time when the parallel workload is small.

The same happens for the 2-step reduction, as when the number of pixels decreases, the improvement brought by this approach also diminishes, where 240p videos actually have better performance using the regular atomic operations, due to the overhead of launching separate reduction kernels just for a small amount of computation.

VII. OTHER WORK

Several additional optimization strategies were explored during the development of the GPU k-means implementation:

- **Shared Memory (Cache) Optimization:** attempts to exploit shared memory for caching pixel data were investigated. However, since the OpenACC implementation already manages memory efficiently, these optimizations did not yield significant performance improvements;
- **Pre-allocation of Memory for Video Frames:** for video sequences, allocating all necessary memory on the first frame and reusing it across subsequent frames was tested. This approach did not result in measurable improvements, as memory allocation overhead is negligible relative to the computational workload;
- **Video-specific Strategies:** two approaches were analyzed for video clustering:
 - *Temporal Centroid Reuse:* initializing the centroids of the current frame with the final centroids of the previous frame reduces convergence time and overhead. This method is particularly effective for longer videos where the content changes gradually;
 - *Calibration:* the calibration procedure adopted in this work, which relies on multiple random initializations and selecting the best prototypes, proved to be more accurate than k-means++ initialization. It is especially suited for static or low-motion videos, where it achieves lower clustering error, while temporal reuse is better for highly dynamic sequences.

VIII. CONCLUSIONS

This work presented the design and evaluation of a GPU-accelerated k-means clustering algorithm using OpenACC. The results show that the GPU implementations achieve significant speedups compared to the CPU baseline, while maintaining clustering accuracy. Among the distance metrics evaluated, the squared Euclidean distance provided the best balance between computational efficiency and clustering quality. The optimized GPU algorithm demonstrated consistent performance across video resolutions and achieved near real-time processing for high-resolution frames. The calibration strategy further improved convergence, reducing iteration counts and enhancing stability for static video content.

Future Work

The study demonstrates that GPU performance in k-means is limited primarily by memory bandwidth rather than raw computational power. Further optimization should therefore focus on improving memory reuse, overlapping data transfers with computation, or adopting asynchronous streams. Other future extensions of this work could focus on:

- **CUDA Implementation:** rewriting the algorithm in CUDA would allow for fine-grained control over memory management and thread scheduling, enabling further optimization beyond OpenACC’s abstraction;
- **Multi-GPU Scaling:** leveraging multiple GPUs could extend the method to much larger datasets and higher-resolution video streams, distributing the workload for even greater speedups.

REFERENCES

- [1] V. Singh, “cuml - rapids machine learning library,” 2025. [Online]. Available: <https://github.com/VishankSingh/cuml>
- [2] M. Kruliš, “cuda-kmeans: A novel, highly-optimized cuda implementation of the k-means clustering algorithm,” 2020. [Online]. Available: <https://github.com/krulis-martin/cuda-kmeans>
- [3] M. Kruliš and M. Kratochvíl, “Detailed analysis and optimization of cuda k-means algorithm,” 2020. [Online]. Available: https://jnamaral.github.io/icpp20/slides/Krulis_Detailed.pdf