

MLS COURSEWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper presents a comparative analysis of CPU and GPU performance for information retrieval tasks, focusing on query-to-vector search using exact k-nearest neighbours and approximate nearest neighbours, and model serving systems, testing the effectiveness of queueing and batching for efficient real-world employment. We evaluate both architectures across a range of retrieval workloads. Our findings highlight the GPU's advantages in parallelism and high-throughput computation, especially for dense retrieval and large-scale model serving, while CPUs demonstrate competitive performance for low-latency, memory-efficient tasks.

1 INFORMATION RETRIEVAL ON GPU (TASK 1)

1.1 MOTIVATION

Information retrieval is the process of extracting relevant information from a large collection of unstructured data, such as text documents, web pages or multimedia files, in an efficient and accurate way (Singhal et al., 2001). It involves the use of computational algorithms, data structures, and statistical models to efficiently match user queries with relevant documents.

The typical workflow in an information retrieval system comprises several steps (Kowalski, 2007). All the collected data is processed and stored in an index. Once a query is submitted, it is processed (expanded or refined) to improve search results, and the query is matched against indexed documents (e.g., neural embeddings (Zuccon et al., 2015) can be used for this) and ranked based on the assigned relevance scores. The most relevant documents are retrieved and presented to the user in a structured format. User interactions are used to evaluate the results and collect feedback to improve future search queries.

However, information retrieval struggles with efficiently organising and searching large, high-dimensional datasets. To address this, clustering is used to group similar documents, which helps reduce the search space and retrieval time, and makes scaling to large datasets more manageable. However, as the size of datasets increases, the time and memory consumption of clustering algorithms grow significantly (Whang et al., 2012)¹. Common clustering algorithms like K-Means (Hartigan & Wong, 1979) may struggle to find a local minimum depending on the initialisation of the centroids, often yielding suboptimal results (Celebi et al., 2013). Moreover, it is suggested that the notion of distance becomes less meaningful in high-dimensional spaces, making it harder to identify meaningful clusters (Aggarwal et al., 2001).

Graphic Processing Unit (GPU) acceleration addresses such challenges by allowing parallel processing and enhancing the speed of computations such as distance calculations and centroid updates. GPUs can simultaneously process multiple data points, decreasing the time complexity associated with clustering algorithms (Asano et al., 2009). They also have optimised memory access patterns that enable high-bandwidth, low-latency data movement when threads access memory in a coalesced and structured manner, allowing for the effective management of extensive datasets and high-dimensional data. Additionally, GPUs excel at iterative optimisation tasks involved in clustering algorithms, such as the recalculation of centroids in K-Means (Li et al., 2013). Consequently, GPU-based clustering has various practical applications, such as recommendation engines of video streaming platforms like YouTube or Netflix and image and video processing (Roshdi & Roohparvar, 2015).

¹for example, K-Means has a time complexity of $O(n^2)$, where n is the number of data points.

1.2 OVERVIEW

Our GPU-accelerated system is designed to efficiently handle large-scale information retrieval and clustering tasks by leveraging parallelism and optimised memory access patterns. It supports three core algorithms: K-Nearest Neighbours (KNN) for similarity search (Peterson, 2009), K-Means for unsupervised clustering, and Approximate Nearest Neighbours (ANN) methods for fast retrieval in high-dimensional spaces (Andoni et al., 2018). Our GPU-accelerated clustering system is divided into three layers, and visualised in Figure 1.

Data Generation Layer. This layer generates synthetic data for both the database and queries, providing customisations for size, dimensionality, and the number of clusters. Data can be generated randomly or using specified distributions. The data is converted to 32 precision to balance efficiency and accuracy.

GPU Processing Layer. This module consists of distance functions using CUDA kernels wrapped in CuPy, Triton, and Pytorch. It also consists of KNN, K-Means, and ANN algorithms for clustering and finding nearest neighbours. The generated data is first loaded onto the GPU, after which the GPU processing layer produces output, which is offloaded onto the CPU again.

Evaluation Layer. This layer computes the total and query time of all implemented algorithms. It also calculates the recall rates of the ANN implementations compared to the exact KNN results.

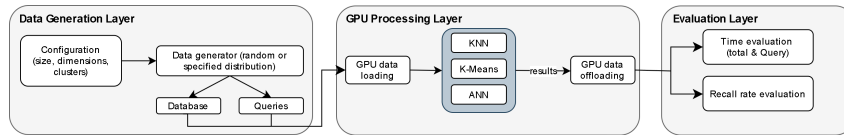


Figure 1: Workflow of the GPU-based clustering system.

We implement our system in different GPU-accelerated versions using PyTorch, CuPy, and Triton, and also include a Central Processing Unit (CPU) benchmark using Numpy. The core ideology behind our system is parallelisation, with each algorithm designed to work with 1D and 2D databases and queries.

1.3 DISTANCE FUNCTIONS

Problem Definition. Distance functions play a vital role in information retrieval, particularly in ranking and comparing document embeddings (Kusner et al., 2015), or query-document similarity scores (Yang et al., 2009). They help measure how similar or different two points (e.g., query vs. document vector) are in a given feature space and are thus crucial in clustering algorithms to group similar documents together.

Design Choice. In this work, we consider four distance functions to calculate the similarity between two points in a feature space. For conciseness, we extensively discuss the implementations using CuPy.

First, we utilise the Euclidean distance, or L2 distance, which is defined as follows: $d_{L2}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$, where x and y denote the document vector and query vector, respectively, and n represents the dimensionality of both vectors. It measures the straight-line distance between two points, derived from the Pythagorean theorem. As the differences are squared, it is sensitive to large feature difference, and works well when magnitude matters or when features are normalized or comparable in range.

The second distance function is the Manhattan distance, or L1 distance, which measures the sum of absolute differences across dimensions. It is defined as follows: $d_{L1}(x, y) = \sum_{i=1}^n |x_i - y_i|$. Compared to L2 distance, it is more robust to outliers as it does not exaggerate large differences. It is, therefore, more useful when changes in individual features are important and in sparse data or high-dimensional settings where many features are zero.

Third, we make use of the negative dot product. The dot product measures similarity, projecting one vector onto another, and as such, it would return a high value for vectors that are close to each other.

108 However, our sorting algorithms look for the lowest distances, and taking the negative transforms
 109 the dot product into a form of distance (lower is better). It is defined as follows: $d_{\text{dot}}(x, y) =$
 110 $-\sum_{i=1}^n x_i y_i$, which measures the directional similarity and magnitude of two vectors. It is often
 111 utilised for applications with embedding spaces.

112 Last, the cosine distance is used, which is formulated as follows: $d_{\text{cos}}(x, y) = 1 - \frac{x \cdot y}{\|x\| \|y\|}$. Cosine
 113 distance measures the angle between two vectors, regardless of their magnitude. It focuses on ori-
 114 entation and similarity of direction rather than distance. Cosine itself is also a measure of similarity
 115 by default, and subtracting cosine from 1 results in a distance function. It can be used for comparing
 116 text embeddings, or anything where scale is not informative, with applications in text retrieval and
 117 natural language processing.

118 **Implementation Details.** For implementation, we wrap CUDA kernels in Python using CuPy’s
 119 `RawKernel`. For certain applications, it is desirable to compare a database to a matrix of queries
 120 instead of a single vector, which may result in slow execution times as the system has to loop over
 121 each query. To address this, our distance functions are able to parallelise vector to vector, matrix
 122 to vector and matrix to matrix operations. For memory access, each CUDA thread is responsible
 123 for computing the distance between one pair of vectors, where threads are arranged in a 2D grid
 124 depending on the number of vectors in the database and query matrix, respectively. An overview of
 125 the algorithm is provided in Algorithm 1 in the Appendix. For a general CPU-based comparison,
 126 we implement the same distance functions using Numpy operations.

128 1.4 TOP-K RETRIEVAL

130 **Problem Definition.** Exact KNN search becomes computationally expensive for large-scale
 131 datasets, with complexity $O(N)$ per query. For a dataset with N points in D dimensions, finding
 132 exact nearest neighbours requires computing and sorting N distances for each query point. This be-
 133 comes impractical for large N or high D . ANN approximates KNN, presenting a trade-off between
 134 perfect accuracy and dramatically improved speed, making ANN algorithms essential for real-world
 135 applications.

136 **Design Choice.** Two ANN algorithms are implemented in this work based on two different under-
 137 lying principles: K-Means and Locality-Sensitive Hashing (LSH). For K-Means, the data is divided
 138 into K clusters ($K > K_1$) using K-Means clustering, then the algorithm finds the K_1 closest cluster
 139 centroids to the query. Next, for each elected cluster, it finds the K_2 nearest points, after which
 140 the K closest points from this candidate pool are returned. This reduces the search complexity to
 141 approximately $O(K + K_1 K_2)$ per query, making it particularly efficient for naturally clustered data.

142 The second implementation utilises LSH, which uses random projections to hash similar points
 143 into the same buckets with high probability. The key principle is that nearby points in the original
 144 space are likely to share similar hash values. This reduces complexity from $O(N)$ to approximately
 145 $O(L + nC)$, where L is the number of hash tables and nC is the average number of candidates.

146 **Implementation Details.** Both of our designs for ANN distinguish between total and query time:
 147 K-Means ANN requires initial clustering of the dataset, while LSH needs hash table generation
 148 and hash code computation. However, given a database, this only needs to be done once, and
 149 as such is not taken into account when comparing performance with KNN. For K-Means ANN,
 150 query time involves finding K_1 nearest centroids and examining K_2 points from their clusters. One
 151 challenge we encountered was balancing speed and recall with the hyperparameters K_1 and K_2 ,
 152 both determine the number of points the algorithm is searching through. We set $K_1 = \lceil 0.8K \rceil$ and
 153 $K_2 = \lceil 2K \rceil$, which was determined empirically. Similarly, for LSH, the hyperparameters include
 154 L , the number of hash tables determining the chance that similar points hash to the same bucket as
 155 the query, and the number of hash functions per table determining how specific each hash bucket is.
 156 We set $L = 6$ and the number of hashes to 2. An overview of ANN K-Means and ANN LSH can be
 157 found in Figures 2 and 3 in the Appendix, respectively.

158 For conciseness, we restrict the implementation details to CuPy. Our initial implementation us-
 159 ing standard CuPy utilities for distance calculations (`cp.linalg.norm`, `cp.dot`) and top-K
 160 selection (`cp.argmaxpartition`) proved inefficient for large-scale nearest neighbour search. We
 161 optimised these calculations by developing custom CUDA kernels for all distance metrics (cosine,
 L2, dot, and L1) to maximise GPU utilisation. Furthermore, the LSH implementation leverages

vectorised operations through `cp.tensor_dot` for efficient hash computation, while both methods use a custom quickselect-based kernel for top-K selection, significantly improving performance over `cp.argmaxpartition`. It performs an efficient block-wise top-K selection using a custom CUDA kernel: it splits the input into blocks, finds the K smallest distances within each block in parallel, and then does a final global top-K selection from all block-level candidates. This improves speed by massively reducing the number of values considered in the final top-K step.

Although no Out-Of-Memory (OOM) issues were encountered on the GPU, being able to handle datasets that do exceed GPU memory is essential for real-world problems. Therefore, we also implemented batched versions of each framework, since the large size of the database is a common cause OOM issues. We estimate the available GPU memory, after which the database is loaded onto the GPU in batches based on the available memory, after which batched distance calculations are performed. Similarly, for the ANN algorithms, everything that includes the database is done in batches. For testing, we simulate that the GPU has 4 GBs of VRAM.

Evaluation Each function was run 13 times, with the first 3 runs used as warm-up and discarded. Final results are averages over the remaining 10 runs. We compare three implementations: NumPy KNN, CuPy KNN, and CuPy LSH ANN. Metrics include total execution time, execution time excluding data transfer (to isolate algorithm performance), and for ANN, separate query time and recall (for correctness). Note that the batched implementation does not report data loading time due to its integrated loading and computation process. Tests were conducted on four datasets to evaluate performance across varying dimensionalities and sizes. The largest tested file was 400000×1024 (N, D), limited by disk space and available compute time. Below we present results for the cosine distance function and only cover ANN LSH, while detailed results for the other distance functions and the K-Means ANN can be found in the appendix.

Table 1: Results using the cosine distance.

(N,D)	CuPy										
	KNN CPU		KNN				ANN LSH			CuPy batched	
	total (s)	w/o data (s)	total (s)	w/o data (s)	total (s)	w/o data (s)	query (s)	recall	total (s)	total (s)	query (s)
(4000,2)	0.0004 ± 0.0000	0.0004 ± 0.0000	0.0023 ± 0.0006	0.0006 ± 0.0020	0.0044 ± 0.0021	0.0039 ± 0.0021	0.0033 ± 0.0020	1.00	0.0026 ± 0.0005	0.0037 ± 0.0010	0.0021 ± 0.0006
(4 mil.2)	0.5653 ± 0.0254	0.5609 ± 0.0254	0.0116 ± 0.0005	0.0050 ± 0.0004	0.0162 ± 0.0008	0.0118 ± 0.0004	0.0113 ± 0.0004	1.00	0.0293 ± 0.0024	0.3882 ± 0.0266	0.0877 ± 0.0084
(4000,65536)	0.5454 ± 0.0142	0.5403 ± 0.0137	0.0119 ± 0.0009	0.0049 ± 0.0004	0.0157 ± 0.0005	0.0116 ± 0.0005	0.0111 ± 0.0005	1.00	0.3499 ± 0.0180	5.1337 ± 0.1269	1.4016 ± 0.0344
(400000,1024)	0.8271 ± 0.0811	0.4978 ± 0.0636	4.5820 ± 0.4674	0.0126 ± 0.0008	4.3943 ± 0.3285	0.0293 ± 0.0012	0.0288 ± 0.0011	0.80	2.7219 ± 0.5790	10.8444 ± 0.8801	2.8530 ± 0.6470

Notes: We report the total time the algorithm takes, the time without converting the data to float32 and onto the GPU if applicable (w/p is the abbreviation for without), and the query time and recall for the ANN algorithms. For the batched versions, no times without data loading are available.

Results and Analysis LSH performed better on high-dimensional, uniformly distributed data, while k-means ANN showed superior performance on naturally clustered datasets [see Section 1.7]. Memory usage was higher for LSH due to multiple hash tables, but it offered more consistent query times across different data distributions. The trade-off between recall rate and query time could be adjusted through the involved hyperparameters (K_1 , K_2 for k-means; number of hash tables for LSH), allowing flexibility based on application requirements.

1.5 ALTERNATIVE IMPLEMENTATIONS (PYTORCH AND TRITON)

For Pytorch, we mainly make use of built-in operations. For the Euclidean and Manhattan distance, `torch.cdist` is used. This function computes batched pairwise distances between two sets of vectors and leverages matrix multiplication for the Euclidean distance for efficiency. Similar to the CuPy implementation, Pytorch is implemented in such a way that the distance functions also handle matrix-to-matrix operations. `torch.topk(distances, k=K, largest=False, sorted=False)` is used, which retrieves the indices of the K smallest values along the last dimension of the distances tensor, which typically represents distances from a query point to all other points (e.g., in KNN search). By setting `largest=False`, it selects the smallest values (nearest neighbours), and `sorted=False` returns them in an arbitrary order to improve performance. The ANN algorithms for Pytorch are similar to the CuPy implementations but use the Pytorch distance function.

Similar to CuPy, Triton also utilises kernels to implement the distance functions. In our Triton kernels, each program instance (i.e., thread) is responsible for computing the distance between one row/vector in a database and one row/vector in a query matrix. In case when the query is a vector, the kernel launches N threads, one per row, using `tl.program_id(0)` to identify the row index. Each thread processes its row in chunks of columns, using `tl.make_block_ptr` to efficiently load data with boundary checks and padding, calculates the distance between the row and query for this chunk, and uses `tl.advance` to move the pointer to the next chunk. Furthermore, Triton leverages `triton.autotune` for efficiency. Triton’s autotune feature automatically searches for the best kernel configuration, such as block sizes and tiling parameters, based on the input shapes and hardware characteristics. Triton benchmarks each one on the fly and caches the fastest for future use. This helps maximise performance without manual tuning, adapting the kernel to different GPUs or input sizes efficiently.

For the ANN implementations, they mimic the CuPy implementation with a block-wise top-K selection function using a kernel. Last, batched versions of Pytorch and Triton are also implemented in similar fashion to CuPy.

1.6 TESTING FRAMEWORK

To facilitate benchmarking and iterative development, a reusable `Testing` class was implemented as part of a wider testing framework. It supports customisable configuration via `argparse` and allows dynamic specification of runtime parameters, such as trials to run (e.g., $N=8000$ $D=1024$), the distance function(s) to use, the number of warm-up runs, among others.

Internally, `Testing` leverages a `Benchmark` class responsible for collecting performance metrics during execution. `Benchmark` supports runtime checkpoints, such as query-time capture, which contribute to flexible timing breakdowns (e.g., isolating total runtime from core algorithm query time). `Testing` is iterable, and emits standardised test cases across runs and frameworks (Algorithm 4). It achieves this by embodying an iterator abstraction, which provides fine-grained control over the iterated elements (Algorithm 5).

For each trial, performance data is logged and a CSV summary is subsequently generated. When a specific framework file (e.g., PyTorch) is executed, a corresponding report is produced. When the central testing program is run, it sequentially executes all framework files, aggregating results into a unified report for cross-comparison.

1.7 INFLUENCE OF DIFFERENT DATA DISTRIBUTIONS

The effectiveness of our k-means based ANN implementation varies significantly with data distribution patterns. Specifically, we focused on two distinct variations: well-clustered data (Figure 3) and more spread-out data (Figure 4). For well-clustered distributions ($cluster\ std. = 0.3-0.5$), a smaller K_1 (4–8) yields high recall due to clearly defined cluster boundaries. In contrast, spread-out data ($cluster\ std. = 0.8-1.0$) requires larger K_1 values (8–16) to achieve comparable recall, owing to increased cluster overlap. This direct correlation between cluster separation and the optimal K_1 facilitates more efficient hyperparameter tuning based on data characteristics.

2 MODEL SERVING SYSTEMS (TASK 2)

2.1 MOTIVATION

Machine learning model serving is the process of deploying trained models to make predictions on new, incoming data, often in real time. However, challenges occur when handling online inference requests. Many online applications (e.g., search engines) require sub-second response times. The number of individual requests may fluctuate over time, making it difficult to maintain consistent latency. Furthermore, processing single requests at a time leads to the underutilisation of hardware accelerators, which are designed to perform optimally with parallel computations. In addition, real-time deployment can be expensive, as continuously spinning up infrastructure for individual requests increases operational costs. Last, high request volumes can overwhelm the system if not properly managed, leading to increased response times or dropped requests.

To address such shortcomings, efficient request queuing and batch processing are necessary. Queuing temporarily holds incoming requests in a buffer. As such, it allows the server to organise and group them for more efficient processing. Batching combines multiple requests into a single batch and processes them together. It takes advantage of vectorised computation on CPUs/GPUs, significantly improving throughput, and reduces overheads across multiple requests. Together, queuing ensures smooth intake of requests, while batching maximises computation per cycle, striking a balance between latency and throughput. For example, OpenAI APIs internally use intelligent batching to serve large-scale inference, ensuring responsiveness despite handling thousands of concurrent users.

2.2 OVERVIEW

The system implements a Retrieval-Augmented Generation (RAG) pipeline deployed via a FastAPI web server, supporting both standard request handling and batched processing for improved throughput. Requests are submitted through one endpoint, `/rag-batched` queues requests and processes them in mini-batches for efficiency. The core logic revolves around embedding the query using a transformer model, retrieving top-K similar documents based on embedding similarity, and generating answers using a text-generation pipeline. Figure 2 provides a high-level overview of the request queue and batch processing pipeline.

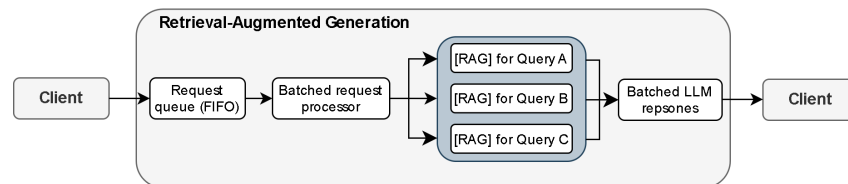


Figure 2: The implemented RAG system using queuing and batching.

2.3 SYSTEM MEASUREMENT

We evaluated the RAG system’s performance using a structured testing framework in `query_script.py`, employing three traffic patterns to simulate load conditions: (1) constant rate at 1 request/second, (2) random rate pattern with variable delays (e.g., 0.1–2.0s), and (3) Poisson-distributed inter-arrival times to approximate real-world traffic. Each test used the same query input and ran for 60 seconds, targeting both `/rag` and `/rag-batched` endpoints. We measured average and median response times, along with the 90th, 95th, and 99th percentiles to capture both central tendency and tail latency. These metrics were selected to provide a holistic view of system responsiveness while highlighting potential performance degradation under high load or outlier conditions. Metrics were collected by logging request timestamps and durations directly within the testing script, with results aggregated post-run. A challenge encountered was an inconsistent system load affecting baseline measurements; this was mitigated by isolating the test environment and repeating runs to average out anomalies.

2.4 REQUEST QUEUE AND BATCHER

Rationale for Request Queue and Batcher. In RAG systems, batching multiple requests improves GPU acceleration by minimising per-request overhead and maximising parallelisation. A queue-based system further manages load by queuing requests during traffic spikes, preventing overload. Together, batching and queuing enhance throughput and optimise memory usage.

Implementation Details. Our implementation employs a thread-safe FIFO queue for managing requests, using Python’s Future objects to handle results asynchronously. A dedicated background daemon thread² continuously monitors the queue, batching requests based on two globally set parameters: `MAX_BATCH_SIZE` (4) and `MAX_WAITING_TIME` (1.0s). This design ensures efficient

²ensures that the main program is not blocked from exiting

GPU utilisation through controlled batching while also maintaining responsive performance under varying load conditions.

Experiments and Metrics Analysis. We evaluated the RAG system’s performance under various conditions, utilising a structured testing framework implemented in `query_script.py`. The tests were conducted over 60-second intervals with 60 total requests, comparing two service endpoints: the standard `/rag` and the batched `/rag-batched` implementations. We tested under both constant rate (1 request/second) and random rate patterns (with delays between 0.1-2.0 seconds), using the query “Which animals can hover in the air?” as our test input.

We collected all the important performance metrics such as total requests processed, error counts, average response times, median response times, and various percentiles (90th, 95th, 99th). The distribution of response times from the load testing is visualised as a histogram in Figure 5.

Challenges and Solutions. A significant challenge for us was setting up the system so it could do queuing and batch processing simultaneously while maintaining memory efficiency and system stability. We addressed this challenge by implementing a batching mechanism that aggregates requests based on either time intervals or batch size. Each batch is then executed on a separate thread, with the total number of active threads capped to ensure controlled resource usage.

3 DISCUSSION OF TASK SYNERGY

GPU kernel optimisation (Task 1) and distributed ML system optimisation (Task 2) are interdependent components in modern machine learning pipelines. Kernel-level performance determines how efficiently compute resources (e.g., GPUs) are utilised, while distributed system logic governs data and task coordination across nodes. Optimising one without considering the other often leads to inefficiencies; e.g., fast kernels may idle wait for data due to poor scheduling, or alternatively, optimal data distribution may overload poorly tuned kernels.

Their synergy lies in co-design: predictable and well-optimised kernels enable effective overlapping of communication and computation; conversely, distributed systems must adapt to kernel characteristics to minimise bottlenecks, such as stragglers or bandwidth contention. Trade-offs arise between maximising GPU utilisation and ensuring system-level balance. For example, fused kernels may reduce launch overhead but introduce latency variability that disrupts distributed pipelines.

Alternative designs include hierarchical optimisation (kernel-to-cluster), operator-aware scheduling, and asynchronous execution models to decouple computation and communication. Overall, integrated optimisation across both tasks is essential for scalable and efficient ML system performance.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

REFERENCES

- Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. On the surprising behavior of distance metrics in high dimensional space. In *International conference on database theory*, pp. 420–434. Springer, 2001.
- Alexandr Andoni, Piotr Indyk, and Ilya Razenshiteyn. Approximate nearest neighbor search in high dimensions. In *Proceedings of the International Congress of Mathematicians: Rio de Janeiro 2018*, pp. 3287–3318. World Scientific, 2018.
- Shuichi Asano, Tsutomu Maruyama, and Yoshiaki Yamaguchi. Performance comparison of fpga, gpu and cpu in image processing. In *2009 international conference on field programmable logic and applications*, pp. 126–131. IEEE, 2009.
- M Emre Celebi, Hassan A Kingravi, and Patricio A Vela. A comparative study of efficient initialization methods for the k-means clustering algorithm. *Expert systems with applications*, 40(1): 200–210, 2013.
- John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979.
- Gerald J Kowalski. *Information retrieval systems: theory and implementation*, volume 1. springer, 2007.
- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. From word embeddings to document distances. In *International conference on machine learning*, pp. 957–966. PMLR, 2015.
- You Li, Kaiyong Zhao, Xiaowen Chu, and Jiming Liu. Speeding up k-means algorithm by gpus. *Journal of Computer and System Sciences*, 79(2):216–229, 2013.
- Leif E Peterson. K-nearest neighbor. *Scholarpedia*, 4(2):1883, 2009.
- Akram Roshdi and Akram Roohparvar. Information retrieval techniques and applications. *International Journal of Computer Networks and Communications Security*, 3(9):373–377, 2015.
- Amit Singhal et al. Modern information retrieval: A brief overview. *IEEE Data Eng. Bull.*, 24(4): 35–43, 2001.
- Joyce Jiyoung Whang, Xin Sui, and Inderjit S Dhillon. Scalable and memory-efficient clustering of large-scale social networks. In *2012 IEEE 12th International Conference on Data Mining*, pp. 705–714. IEEE, 2012.
- Yin Yang, Nilesh Bansal, Wisam Dakka, Panagiotis Ipeirotis, Nick Koudas, and Dimitris Papadias. Query by document. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining*, pp. 34–43, 2009.
- Guido Zuccon, Bevan Koopman, Peter Bruza, and Leif Azzopardi. Integrating and evaluating neural word embeddings in information retrieval. In *Proceedings of the 20th Australasian document computing symposium*, pp. 1–8, 2015.

A VISUALISATION OF DATA DISTRIBUTIONS

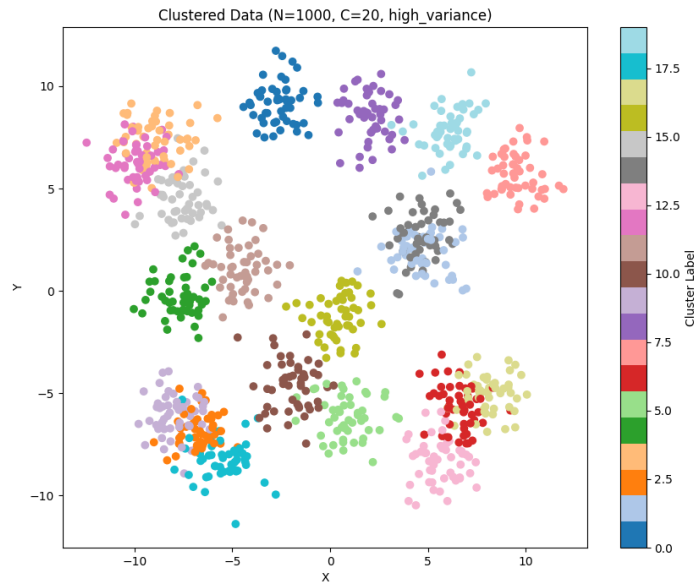


Figure 3: Clustered data with $N = 1000$, $D = 2$, and high variance

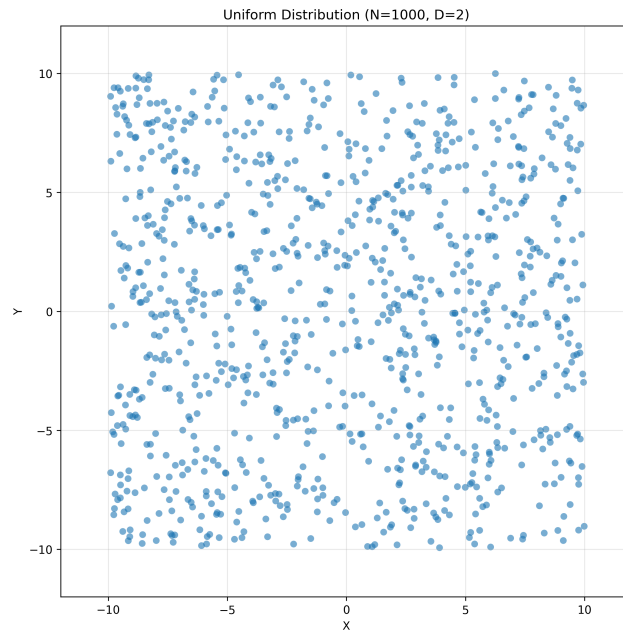


Figure 4: Uniformly distributed data with $N = 1000$, $D = 2$

B ALGORITHMS

Algorithm 1 Cosine Distance Computation

Require: Matrix $X \in \mathbb{R}^{N \times D}$, Matrix $Y \in \mathbb{R}^{K \times D}$

Ensure: Output matrix $out \in \mathbb{R}^{N \times K}$ containing pairwise cosine distances

```

492   for each thread  $(i, j)$  in parallel where  $i < N$  and  $j < K$  do
493        $dot \leftarrow 0.0$  ▷ Initialize dot product
494        $normX \leftarrow 0.0$  ▷ Initialize norm for vector  $X_i$ 
495        $normY \leftarrow 0.0$  ▷ Initialize norm for vector  $Y_j$ 
496       for  $d = 0$  to  $D - 1$  do
497            $a \leftarrow X[i \cdot D + d]$  ▷ Get component  $d$  of vector  $X_i$ 
498            $b \leftarrow Y[j \cdot D + d]$  ▷ Get component  $d$  of vector  $Y_j$ 
499            $dot \leftarrow dot + a \cdot b$ 
500            $normX \leftarrow normX + a \cdot a$ 
501            $normY \leftarrow normY + b \cdot b$ 
502       end for
503        $normX \leftarrow \sqrt{normX}$ 
504        $normY \leftarrow \sqrt{normY}$ 
505       if  $normX < 10^{-6}$  then
506            $normX \leftarrow 10^{-6}$  ▷ Handle small norm values
507       end if
508       if  $normY < 10^{-6}$  then
509            $normY \leftarrow 10^{-6}$  ▷ Handle small norm values
510       end if
511        $cos\_sim \leftarrow dot / (normX \cdot normY)$ 
512        $out[i \cdot K + j] \leftarrow 1.0 - cos\_sim$  ▷ Compute cosine distance
513   end for
514   return  $out$ 

```

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

Algorithm 2 Approximate Nearest Neighbors using K-Means

Require: Data points matrix $A \in \mathbb{R}^{N \times D}$, Query vector $X \in \mathbb{R}^D$, Number of neighbors K

Ensure: Array of indices of the top- K nearest neighbors

```

A_gpu ← Convert A to GPU array (float32)
X_gpu ← Convert X to GPU array (float32)
Reshape X_gpu to (1, D) if necessary

K1 ← ⌈0.8 · K⌉                                ▷ Number of closest centroids to consider
K2 ← ⌈2 · K⌉                                    ▷ Number of points to retrieve per cluster

labels ← Run k-means on A_gpu with K clusters    ▷ Cluster assignment for each point
centroids ← Empty matrix of shape (K, D)

for k = 0 to K − 1 do
  cluster_mask ← (labels == k)

  if cluster_mask has any True values then
    centroids[k] ← Mean of points where cluster_mask is True
  else
    centroids[k] ← 0                                ▷ Handle empty clusters
  end if
end for

nearest_cluster_indices ← Find min(K1, K) nearest centroids to X_gpu
candidate_indices_list ← Empty list
candidate_vectors_list ← Empty list

for each index i in nearest_cluster_indices do
  indices_in_cluster ← Indices of all points where labels == i

  if indices_in_cluster is not empty then
    candidate_vectors ← A_gpu[indices_in_cluster]
    k2 ← min(|candidate_vectors|, K2)              ▷ Number of points to select
    candidate_local_indices ← Find k2 nearest points in candidate_vectors to X_gpu

    Append indices_in_cluster[candidate_local_indices] to candidate_indices_list
    Append candidate_vectors[candidate_local_indices] to candidate_vectors_list
  end if
end for

if candidate_indices_list is empty then
  return Empty array
end if

candidate_indices ← Concatenate all arrays in candidate_indices_list
candidate_vectors ← Concatenate all arrays in candidate_vectors_list

final_k ← min(|candidate_vectors|, K)
top_k_local ← Find final_k nearest points in candidate_vectors to X_gpu
top_k_indices ← candidate_indices[top_k_local]      ▷ Map to original indices

return top_k_indices

```

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Algorithm 3 Approximate Nearest Neighbors using Locality-Sensitive Hashing

Require: Data points matrix $A \in \mathbb{R}^{N \times D}$, Query vector $X \in \mathbb{R}^D$, Number of neighbors K , Number of hash tables L , Number of hashes per table num_hashes

Ensure: Array of indices of the top- K nearest neighbors

```

A_gpu ← Convert A to GPU array (float32)
X_gpu ← Convert X to GPU array (float32)
Reshape X_gpu to (1, D) if necessary

bit_values ← [ $2^0, 2^1, \dots, 2^{num\_hashes-1}$ ]          ▷ For binary-to-integer conversion
hyperplanes ← Random normal matrix of shape (L, num_hashes, D)

dots_all ← A_gpu · hyperplanes                          ▷ Matrix-tensor product
codes ← (dots_all > 0)                                  ▷ Binary hash codes
codes_int ←  $\sum_{j=0}^{num\_hashes-1} codes[:, :, j] \cdot bit\_values[j]$           ▷ Integer hash values

query_dots_all ← X_gpu · hyperplanes
query_codes ←  $\sum_{j=0}^{num\_hashes-1} (query\_dots\_all > 0)[:, :, j] \cdot bit\_values[j]$ 
query_codes ← query_codes[0]

candidate_mask ←  $\exists l \in \{0, \dots, L-1\} : codes\_int[:, l] = query\_codes[l]$  ▷ Match in any table
candidates ← Indices where candidate_mask is True

if candidates is empty then
  candidates ← [0, 1, ..., N - 1]                      ▷ Fallback to full search
end if

A_candidates ← A_gpu[candidates]
dists ← distance_func(A_candidates, X_gpu)              ▷ Compute distances
dists ← Flatten dists

topk_indices ← Find indices of K smallest values in dists
final_indices ← candidates[topk_indices]                ▷ Map to original indices

return final_indices

```

Algorithm 4 Execution flow using the `Testing` class abstraction

```

1: Testing ← Set configuration via argparse
2: args ← Parse runtime args
3: t ← Instantiate new Testing instance

4: if program_entry then
5:   for each (dist_fn, [N, D, A, X, K]) in t do
6:     knn_result ← t.run(our_knn, [N, D, A, X, K])
7:     t.run(our_kmeans, [N, D, A, X])
8:     t.run(our_ann, [N, D, A, X, K], knn_result) Recall rate comparator knn_result
9:     t.run(our_ann_lsh, [N, D, A, X, K], knn_result) Recall rate comparator knn_result
10:  end for
11: end if

```

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Algorithm 5 Iterator logic defined by `Testing` class

```
1:  $iter\_trial\_idx \leftarrow$  The current trials index
2:  $iter\_dist\_idx \leftarrow$  The current distance function index
3:  $iter\_run\_idx \leftarrow$  The current run for the given trial and distance function

4: if the set of runs for the trial and distance function have completed then
5:    $results \leftarrow benchmark.construct\_results()$ 
6:    $benchmark.reset()$  Reset internal collections
7: end if

8: if  $iter\_trial\_idx$  is equivalent to length of trial collection then
9:    $disk \leftarrow$  Write results to a CSV file
10:  Yield  $stop$  to callee
11: end if

12:  $dist\_fn \leftarrow dist\_fn\_collection.get(iter\_dist\_idx)$ 
13: Assign  $dist\_fn$  to global namespace

14:  $[N, D] \leftarrow trials\_collection.get(iter\_trial\_idx)$ 
15:  $K \leftarrow$  10 nearest neighbours

16:  $A \leftarrow$  load matrix file for  $[D, N]$ 
17:  $X \leftarrow$  load vector file for  $D$ 

18: if there are still runs to execute per  $iter\_run\_idx$  then
19:    $iter\_run\_idx \leftarrow iter\_run\_idx + 1$ 
20: else if there are still distance functions to trial per  $iter\_dist\_idx$  then
21:    $iter\_dist\_idx \leftarrow iter\_dist\_idx + 1$ 
22:    $iter\_run\_idx \leftarrow 0$ 
23: else
24:    $iter\_dist\_idx \leftarrow 0$ 
25:    $iter\_run\_idx \leftarrow 0$ 
26:    $iter\_trial\_idx \leftarrow iter\_trial\_idx + 1$ 
27: end if

28: Yield  $dist\_fn, [N, D, A, X, K]$  to callee
```

C VISUALISATION OF RESULTS FROM TASK 2

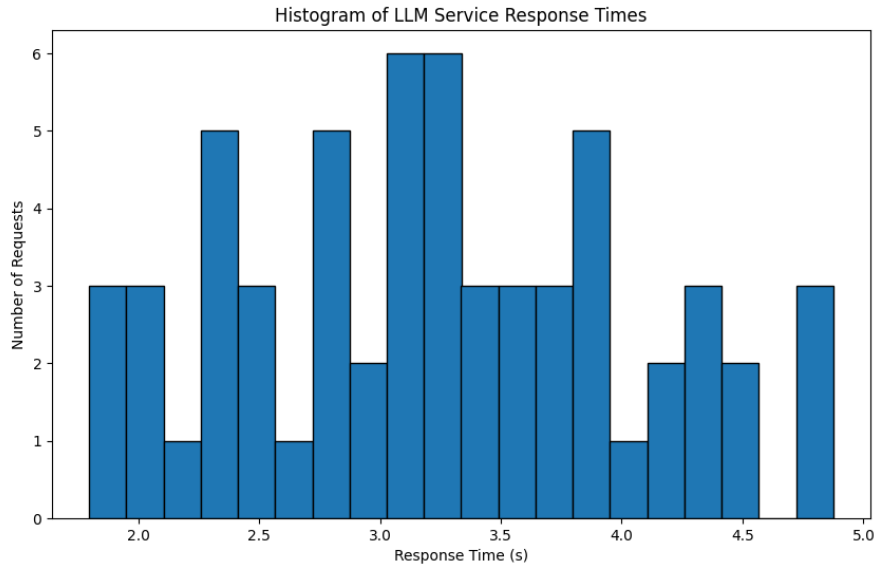


Figure 5: The distribution of response times from the load testing

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755