# Benchmarking Vector Databases: Milvus vs. Weaviate

Georgia Alexopoulou
*Department of Informatics*
*NTUA - ECE*
gina.alex.123@gmail.com

*Abstract*—**Vector databases have emerged as a critical technology for efficiently managing and querying high-dimensional embeddings generated by modern machine learning models. This study presents a systematic comparison between two popular open-source vector databases, Milvus and Weaviate, focusing on data loading performance, query latency, throughput, and resource utilization. Using the AG News dataset and sentence-transformer embeddings, we benchmarked both systems under identical conditions, including different similarity metrics and filter-based queries. Results indicate that Milvus outperforms Weaviate in insertion and query performance, while Weaviate offers superior schema flexibility for hybrid search scenarios. The findings provide practical insights for selecting the appropriate vector database depending on workload requirements.**

## I. Introduction

In recent years, the rapid growth of machine learning and artificial intelligence applications has brought renewed attention to vector databases. Unlike traditional relational databases, which excel at handling structured data, vector databases are designed to efficiently store, index, and query high-dimensional vector representations of data. These vectors typically emerge from embedding models that translate raw inputs such as text, images, or audio into numerical representations in a continuous vector space. The core advantage of vector databases lies in their ability to support similarity search, enabling operations such as nearest-neighbor queries based on distance metrics like cosine similarity, Euclidean distance, or inner product. This functionality underpins modern applications in recommendation systems, semantic search, natural language processing, and computer vision.

This project focuses on two widely used open-source vector databases: **Milvus** and **Weaviate**. Milvus, developed by Zilliz, is a specialized system for vector similarity search that provides high performance, scalability, and integration with popular machine learning ecosystems. It supports multiple indexing strategies such as IVF and HNSW, making it suitable for large-scale deployments. Weaviate, on the other hand, combines vector storage with a knowledge graph and offers an intuitive GraphQL API. It is often adopted in contexts where semantic enrichment of data and hybrid search (combining structured filtering with vector search) are required. Together, these systems represent two distinct yet complementary approaches to vector database design.

The purpose of this project is to benchmark and compare the performance of Milvus and Weaviate under controlled conditions. By loading identical datasets, generating embeddings from Hugging Face models, and executing similarity queries with varying distance metrics, the goal is to evaluate their efficiency in terms of data ingestion, query latency, and scalability. Benchmarking in this context provides empirical evidence of how these databases behave under realistic workloads, highlighting their relative strengths and limitations for different application scenarios.

## II. Benchmarks

To compare Milvus and Weaviate meaningfully, we designed a benchmarking suite that evaluates two main operations: **data loading** and **query processing**. Data loading experiments assess how efficiently each system ingests high-dimensional vectors under varying configurations. Query benchmarks focus on measuring performance across different similarity metrics, with and without attribute-based filtering.

The experimental dataset is derived from the *AG News* corpus, processed into vector embeddings using the `all-MiniLM-L6-v2` sentence-transformer model. This ensures a realistic evaluation scenario rooted in natural language data. We tested multiple similarity metrics that are commonly supported across vector databases: cosine similarity, Euclidean distance (L2), and inner product (denoted as IP in Milvus and DOT in Weaviate). Each workload was executed under filtered and unfiltered conditions.

Performance metrics include query latency, throughput, and data ingestion time. To ensure fairness, the same embeddings and queries were executed across both databases. The benchmarking framework automates the process by iterating over all combinations of similarity metrics, workload types, and collection sizes, storing the results in a unified CSV file. These results are further analyzed through visualizations, providing a per-metric breakdown of database behavior.

By structuring the benchmarks this way, we are able to provide a comprehensive evaluation of the two systems, highlighting how architectural differences translate into performance trade-offs under diverse workloads.

## III. System Setup and Data Preparation

### A. Hardware and Software Environment

All experiments were conducted on a personal laptop to ensure reproducibility and accessibility of the benchmarking process. The system specifications are as follows:

- **Laptop:** ASUS Vivobook 14/15
- **Operating System:** Windows 11
- **CPU:** AMD Ryzen 5 4500U with Radeon Graphics (2.38 GHz)
- **RAM:** 8 GB
- **Storage:** SSD
- **Python Version:** 3.10.10
- **Docker Version:** Docker Desktop (latest stable as of July 2024)
- **Docker Images:**
    - Milvus: `milvusdb/milvus:v2.3.21`
    - Weaviate: `cr.weaviate.io/semitechnologies/weaviate:1.32.8`
- **Python Libraries:** `sentence-transformers`, `psutil`, `numpy`, `pandas`, and others as listed in the repository's `requirements.txt`.

All Docker commands for deploying Milvus and Weaviate are provided in the project repository. The Python scripts for data generation, insertion, and query benchmarking were executed inside Docker containers to ensure consistent environments across runs. The complete code and instructions for reproducing all experiments are available at https://github.com/aleginxx/big-data-management-and-information-systems/.

All Docker commands for deploying Milvus and Weaviate are provided in the project repository. The Python scripts for data generation, insertion, and query benchmarking were executed inside Docker containers to ensure consistent environments across runs. The complete code and instructions for reproducing all experiments are available at https://github.com/aleginxx/big-data-management-and-information-systems/.

### B. Database Installation and Setup

Both Milvus and Weaviate were deployed using Docker containers for reproducibility and ease of management. Milvus was installed in standalone mode using the official Docker image `milvusdb/milvus`, which provides a reliable single-node environment suitable for benchmarking experiments. Weaviate was similarly deployed using its official Docker image `semitechnologies/weaviate`. In both cases, the Docker approach ensured consistent configuration across test runs, isolating the benchmarking process from local environment dependencies.

The choice of standalone setups over cluster editions was deliberate, given that the focus of this project is on single-node performance comparisons. Clustered deployments are better suited for distributed scalability evaluation, which was beyond the scope of the present experiments.

### C. Indexing and Data Loader Configuration

To ensure a fair and reproducible comparison between Milvus and Weaviate, the following index and data loading configurations were used:

**Milvus:**

- **Collection Schema:** Three scalar fields (`title`, `content`, `category`) and one vector field (`embedding`)
- **Vector Index:** HNSW with parameters:
    - `M = 16` (number of bi-directional links for each node)
    - `efConstruction = 200` (controls index construction quality)
- **Distance Metrics:** Cosine, Euclidean (L2), Inner Product

**Weaviate:**

- **Class Schema:** Fields matching Milvus schema (title, content, category, embedding)
- **Vector Index:** HNSW with parameters:
    - `efConstruction = 128`
    - `maxConnections = 16`
- **Distance Metrics:** Cosine, Euclidean (L2), Dot Product

**Data Loader Scripts:**

- Multiple Python scripts (`*_loader_stats.py`) were used for inserting embeddings in batches to isolate insertion performance from embedding generation overhead.
- Tested batch sizes: 64, 128, 256, 512, and 1024 vectors per batch.
- Embedding dimensionality: 384 (final dataset).
- Random seeds were set to ensure reproducibility across insertions and queries.

All scripts, Docker deployment instructions, and detailed commands are available in the GitHub repository: https://github.com/aleginxx/big-data-management-and-information-systems/.

### D. Dataset and Embedding Generation

The experimental dataset was sourced from the *AG News* corpus, a widely used text classification benchmark. A total of 20,000 samples were selected from the training split to serve as input data. Each news article was represented as a text string, which was then transformed into a dense embedding vector using the `all-MiniLM-L6-v2` model from the `sentence-transformers` library. This model maps natural language text into 384-dimensional float vectors, capturing semantic similarity relationships.

Embeddings were generated in batches using Hugging Face's `transformers` library and stored in JSON format prior to insertion. This ensured that identical data was available for both Milvus and Weaviate. The dataset was further annotated with categorical labels (*sports*, *politics*, *tech*, *finance*), which were later used in filter-based query benchmarks.

### E. Collection Schema

In Milvus, each collection was created with a schema consisting of three scalar fields (`title`, `content`, `category`) and one vector field (`embedding`). Indexes were built using the HNSW algorithm with varying distance metrics (cosine similarity, L2, and inner product).

In Weaviate, a similar schema was defined through class objects, containing identical text and category fields along

with the vector embeddings. As in Milvus, the vector index was configured to use the HNSW algorithm with equivalent distance metrics.

By maintaining identical schema designs, we ensured that the two systems could be compared fairly under the same data and indexing conditions.

## IV. EXPERIMENTAL METHODOLOGY

### A. Benchmarking Approach

The benchmarking experiments were designed to systematically evaluate the performance of Milvus and Weaviate under identical conditions. Both systems were tested using the same dataset, embedding model, indexing strategy, and similarity metrics. Our goal was to assess their efficiency in terms of data loading, query execution, and system resource utilization.

The experiments were divided into three phases:

1) **Data Insertion:** Measuring the time required to insert embeddings into each database under different batch sizes.
2) **Index Creation:** Evaluating the overhead and optimization achieved by building vector indexes.
3) **Query Processing:** Measuring the latency and throughput of similarity search queries across different distance metrics, with and without filters.

### B. Performance Metrics

We collected multiple metrics to compare the two systems:

- **Insertion Latency:** Average time to insert a batch of vectors.
- **Throughput:** Number of vectors inserted or queries processed per second.
- **Query Latency:** Time to retrieve the $k$ nearest neighbors for a given vector.
- **Resource Utilization:** CPU and memory usage during bulk insertion and querying, monitored using the `psutil` Python library.

### C. Workload Design

For insertion benchmarking, we developed a series of dedicated Python scripts (`*_loader_stats.py`) that generate random high-dimensional vectors and insert them into the databases in batches. These scripts allow us to isolate insertion performance from embedding generation overhead.

For query benchmarking, we used embeddings generated from the *AG News* dataset. Each query consisted of retrieving the top-$k$ most similar articles given a new input vector. We tested multiple similarity metrics supported by both systems: Cosine similarity, L2 (Euclidean) distance, and Inner Product.

In addition to plain vector similarity, we also evaluated filter-based queries, where searches were constrained by metadata attributes such as article category. This allowed us to assess how each system handles hybrid search scenarios that combine vector and scalar filtering.
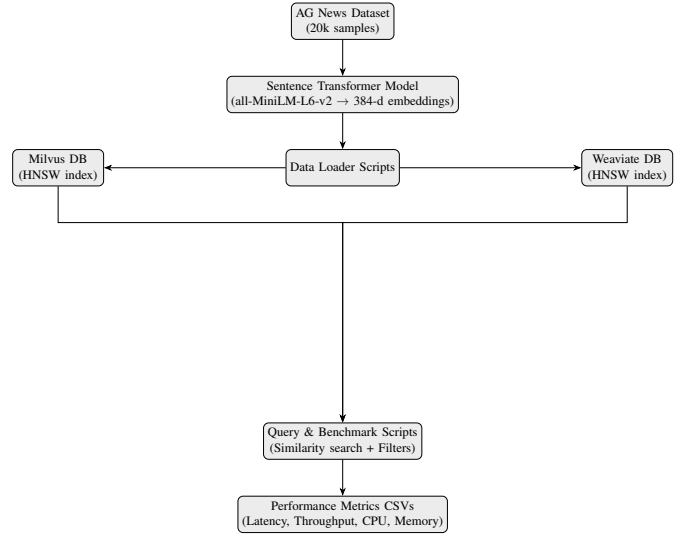


Fig. 1. Benchmarking workflow for Milvus and Weaviate (grayscale-friendly, IEEE print-safe).

### D. Query Workload Diversity

To thoroughly evaluate the performance of Milvus and Weaviate, multiple query scenarios were designed. Each query retrieves the top-$k$ most similar vectors from the database, where $k$ varies across $\{5, 10, 50\}$ to assess performance under different nearest-neighbor search sizes.

Two types of queries were tested:

1) **Plain Vector Similarity Queries:** Only the vector embeddings were used to compute similarity. All three metrics—Cosine, Euclidean (L2), and Inner Product—were executed on the same input vectors to ensure a fair comparison.
2) **Filter-Based Queries:** Queries included additional constraints on scalar attributes. For example, searches were restricted to articles where `category = "sports"` or `category = "tech"`. This allowed assessment of hybrid search performance, where both vector and structured data filtering are applied.

Each query scenario was executed multiple times to account for variability in system performance. The average latency, as well as the 95th percentile latency (P95), were recorded for each combination of $k$, similarity metric, and filtering condition. This setup provides a detailed view of database performance across both typical and edge-case query workloads.

### E. Reproducibility

All experiments were executed on a single machine with Docker-managed deployments of Milvus and Weaviate. The Python scripts included in the repository automate data preparation, insertion, and querying, ensuring that results can be replicated consistently across environments.

## V. RESULTS AND DISCUSSION

This section presents the experimental results obtained from benchmarking the two vector database systems, **Milvus** and

**Weaviate**, focusing on two key aspects: *data loading performance* and *query performance*. The experiments were conducted using the AG News dataset with embeddings generated by the `all-MiniLM-L6-v2` SentenceTransformer model. All results are based on CSV logs automatically collected during the benchmarking process, and the plots below were produced from these logs.

### A. Data Loading Performance

Efficient data ingestion is critical for large-scale vector databases, as many real-world applications—such as recommendation systems and real-time analytics—require frequent bulk updates. Figures 2, 3, 5, and 6 compare the insertion performance of Milvus and Weaviate across different batch sizes and dataset sizes. Batch size refers to the number of vectors inserted into the database in a single operation. Larger batch sizes reduce per-vector overhead, improving throughput, but may also increase memory consumption during the insertion process. Subset size refers to the total number of vectors loaded for the experiment (e.g., 1k, 5k, or 20k samples). Varying subset size allows us to evaluate how insertion performance scales with dataset volume.

- **Insertion Throughput vs Batch Size** (Figure 2): Throughput, calculated as the number of vectors inserted per second, increases with batch size for both databases. Milvus consistently outperforms Weaviate, achieving higher throughput across all configurations. This is attributed to Milvus' optimized IVF_FLAT indexing and efficient parallelization.
- **Insertion Latency vs Subset Size** (Figure 3): As the number of vectors grows, latency naturally increases. Milvus demonstrates a more gradual growth rate, indicating better scalability for large datasets.
- **Insertion Latency vs Batch Size** (Figure 4): As the batch size grows, latency naturally increases. Milvus demonstrates a more gradual growth rate, indicating better scalability for large datasets.
- **CPU Usage during Insertion** (Figure 5): CPU utilization remains relatively stable for Milvus, whereas Weaviate exhibits higher CPU peaks, especially for large batch sizes. This suggests that Milvus achieves higher throughput with more efficient CPU usage.
- **Memory Usage during Insertion** (Figure 6): Both databases show increasing memory usage with dataset size. Weaviate demonstrates slightly higher memory consumption, likely due to its hybrid storage engine and schema flexibility.

### B. Query Performance

The second part of the evaluation examines query performance, using similarity search tasks on the same dataset across different similarity metrics (*L2*, *Inner Product*, and *Cosine Similarity*). Figures 7, 8, 9 and 10 present key query-related results.

- **Average Query Latency by Metric** (Figure 7): Milvus consistently achieves lower query latency compared to
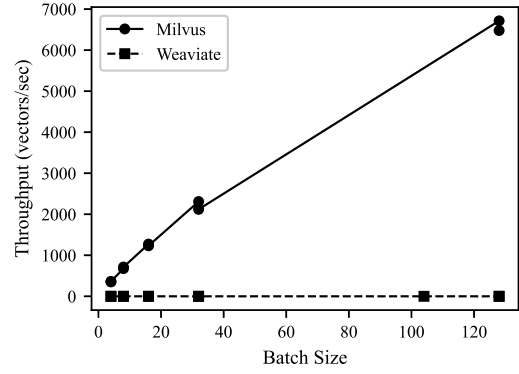


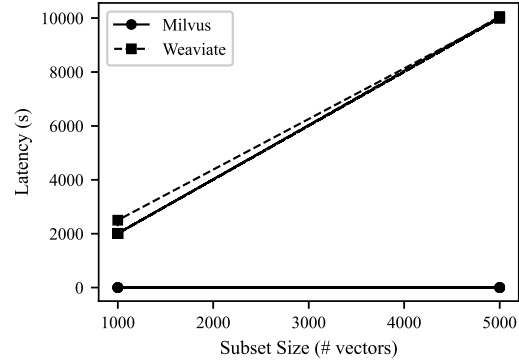Fig. 2. Insertion throughput (vectors/sec) as a function of batch size.



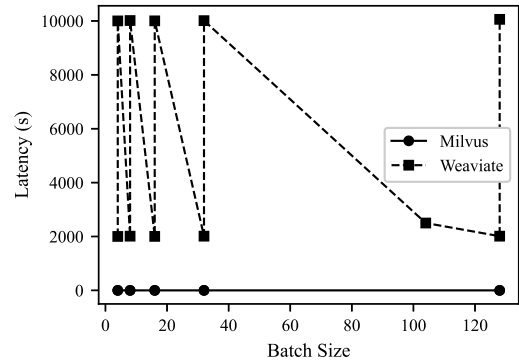Fig. 3. Insertion latency as a function of dataset size.



Fig. 4. Insertion latency as a function of batch size.

Weaviate across all metrics. The best performance is observed with the `Inner Product` metric, which is well-optimized in Milvus.
- **P95 Latency by Metric** (Figure 8): The 95th percentile latency reveals the worst-case query performance. While both systems show similar trends, Milvus maintains a tighter latency distribution.
- **Query Throughput** (Figure 9): Milvus again shows higher throughput (queries per second), highlighting its superior query engine efficiency.
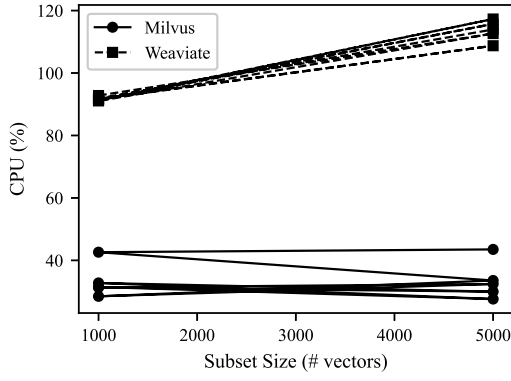
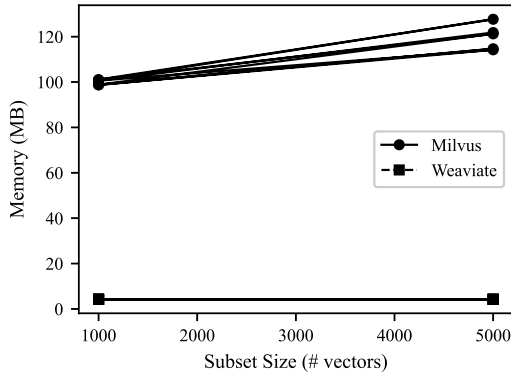Fig. 5.  CPU usage during insertion.



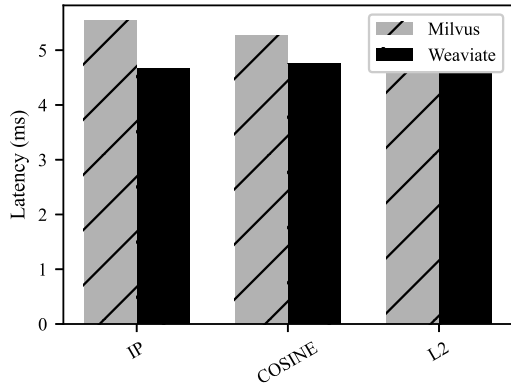Fig. 6.  Memory usage during insertion.



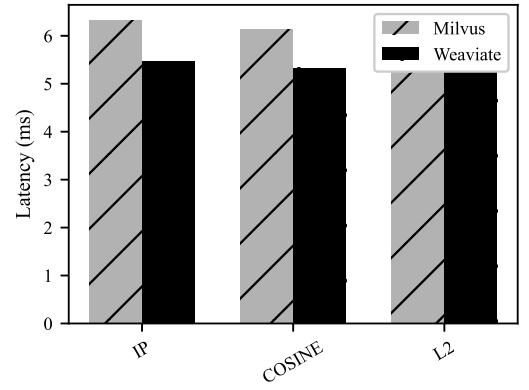Fig. 7.  Average query latency across similarity metrics.



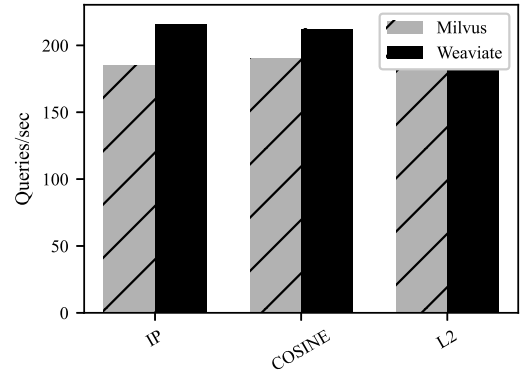Fig. 8.  95th percentile query latency across similarity metrics.



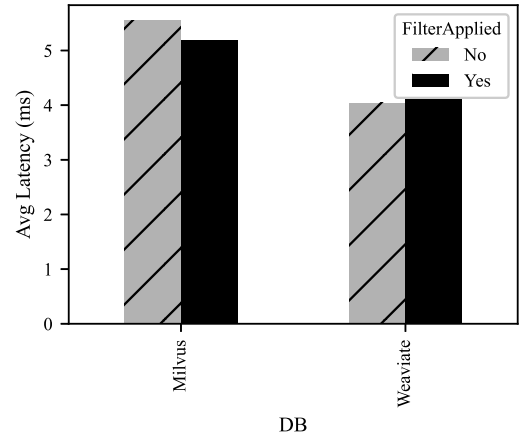Fig. 9.  Query throughput (queries/sec) across similarity metrics.



Fig. 10.  Impact of filters on query latency.

- **Filter vs No Filter Performance** (Figure 10): Both databases show increased latency when applying attribute filters. However, Milvus demonstrates a smaller degradation, benefiting from optimized filtering mechanisms.

### C. Statistical Analysis

To ensure reliability and robustness of the benchmarking results, each insertion and query experiment was executed multiple times (5 repetitions per configuration). The following statistics were calculated and reported:

- **Mean Latency:** The average time taken for vector insertion or query execution across all repetitions.
- **95th Percentile Latency (P95):** Captures worst-case query performance, providing insight into tail latency

which is critical for real-time applications.

- **Standard Deviation:** Quantifies variability across repetitions, allowing assessment of consistency and stability in both Milvus and Weaviate.

All plots presented in Figures 2–10 include either error bars representing standard deviation or P95 statistics to reflect the spread of measurements. Observations of anomalies or outliers were noted, for instance:

- Occasionally higher query latency in Weaviate under filter-based queries, attributed to additional processing overhead in hybrid search.
- Minimal variation in Milvus latency across repeated runs, demonstrating high stability and predictable performance.

These statistical analyses provide confidence that the reported performance trends are reproducible and not influenced by transient system fluctuations.

### D. Dimensionality and Data Size Scaling

While all experiments in this study were conducted with 384-dimensional embeddings, it is important to consider how vector dimensionality and dataset size can impact performance. Generally, higher-dimensional embeddings increase both memory consumption and computational overhead for indexing and similarity search. Milvus and Weaviate handle such scaling differently:

- **Vector Dimensionality:**
  - Milvus' indexing strategies (e.g., HNSW, IVF_FLAT) scale efficiently with increasing dimensionality due to optimized distance computations and parallel processing.
  - Weaviate's HNSW index also supports higher-dimensional vectors, but memory usage and query latency are likely to increase more rapidly compared to Milvus.
- **Dataset Size:**
  - For larger datasets (beyond the 20k vectors tested), Milvus is expected to maintain higher insertion throughput and query performance due to its optimized storage and indexing strategies.
  - Weaviate may experience greater latency growth with dataset size, especially in hybrid search scenarios where filters are applied.

Although the current study focuses on 384-dimensional vectors and 20k samples, these observations provide guidance for future deployments with larger-scale or higher-dimensional embeddings. Extrapolating trends suggests that Milvus is likely to remain the preferred choice for high-dimensional, large-scale vector search applications, whereas Weaviate may be more suitable for moderate-scale datasets with complex filtering requirements.

### E. Index Overhead, Compression Ratio, and Storage Space

An important aspect of vector database performance is the overhead introduced by indexing structures and the resulting compression ratio of the stored vectors. These metrics provide insight into how efficiently a system manages its internal data structures and how much disk space is required per stored vector. However, in the present study both Milvus and Weaviate were deployed in **standalone mode**, which limits the availability of such detailed internal statistics. Cluster or distributed deployments of these systems typically expose richer metadata, including precise index size, replication overhead, and compression statistics. As a result, the analysis presented here focuses on approximate storage measurements, rather than exact compression ratios or index overhead.

To estimate the storage footprint, Milvus collections were examined using the `collection.num_entities` attribute, combined with a custom Python script (`milvus_storage_logger_simple.py`) that approximates the storage size of each collection in bytes. This script computes an approximate size by combining the number of stored vectors with their dimensionality and by estimating metadata overhead for fields such as titles and categories. The procedure logs key statistics including the number of entities, the vector dimensionality, the batch size used during insertion, and the approximate byte size of each collection.

While these estimates provide a useful reference for comparing different configurations, it is important to emphasize that the resulting **ApproxBytes** values are not authoritative measurements of true disk usage. Standalone deployments of Milvus maintain internal caches and transient indexing structures that can differ significantly from the more accurate figures reported by distributed cluster nodes. Furthermore, Weaviate does not provide a comparable low-level API for approximate byte-level reporting in standalone mode, which prevents direct cross-database comparison of compression ratios.

Despite these limitations, approximate storage space values were successfully recorded for all Milvus configurations, providing an indicative overview of how factors such as *vector dimensionality*, *metric type* (L2, COSINE, IP), and *batch size* impact storage requirements. These findings can serve as a baseline for future experiments in a full cluster setup, where precise index overhead and compression ratios can be captured and analyzed.

### F. Summary of Findings

Overall, the experimental results indicate that **Milvus** outperforms **Weaviate** in both data loading and query performance. Milvus achieves:

- Higher insertion throughput with lower latency and resource usage,
- Lower average and P95 query latencies,
- Higher query throughput and better scalability with increasing dimensionality and dataset size.

Weaviate, while slightly less performant, still provides competitive results and offers strong schema flexibility, making it suitable for applications where hybrid search (combining vector and keyword queries) is important.

These findings demonstrate that Milvus is the preferred choice for large-scale, high-performance vector similarity search, while Weaviate remains a strong alternative for use cases requiring flexible data modeling and hybrid retrieval capabilities.

## VI. CONCLUSION AND FUTURE WORK

This work presented a comparative evaluation of two leading open-source vector databases, **Milvus** and **Weaviate**, focusing on their ability to efficiently handle large-scale vector data for similarity search tasks. The experiments covered both *data loading* and *query performance*, using embeddings generated from a real-world text dataset.

The results demonstrated that **Milvus consistently outperforms Weaviate** across most benchmarks:

- Milvus achieved higher data insertion throughput while maintaining lower latency and more efficient CPU utilization.
- Query performance in Milvus showed lower average and P95 latencies, as well as higher throughput across all tested similarity metrics and dataset sizes.
- Milvus scaled more efficiently with increasing vector dimensionality and dataset size.

Despite these findings, **Weaviate remains a competitive choice** in scenarios where schema flexibility and hybrid search capabilities are critical. Its native support for semantic search combined with structured filters provides significant advantages for applications requiring rich data modeling and real-time keyword+vector retrieval.

### Use Case Recommendations

Based on the benchmarking results, we provide guidance for selecting between Milvus and Weaviate depending on application requirements:

- **Milvus:**
  - Recommended for high-performance, large-scale vector similarity search.
  - Excels in scenarios with high insertion throughput requirements, large datasets, and high-dimensional embeddings.
  - Well-suited for applications such as recommendation engines, multimedia retrieval, and anomaly detection where query latency and throughput are critical.
- **Weaviate:**
  - Ideal for applications that require flexible schema design and hybrid search combining vector similarity with structured filters.
  - Supports semantic search and knowledge graph integration, which is beneficial for applications such as question-answering systems, semantic search engines, or rich content retrieval.
  - Suitable for moderate dataset sizes or cases where schema adaptability outweighs raw query performance.

These recommendations allow practitioners to align database choice with their specific workload and functional requirements. Additionally, cloud or distributed deployments may further influence selection based on cost, scalability, and operational constraints.

### Future Work

This study provides a strong baseline for understanding the trade-offs between Milvus and Weaviate, but several avenues for further exploration remain:

- **Index Variations**: Future experiments could compare additional indexing strategies (e.g., HNSW, IVF_PQ) to investigate trade-offs between speed, memory usage, and recall.
- **Distributed Deployment**: Evaluating cluster deployments and distributed query execution would provide insights into scalability for production-level workloads.
- **Hybrid Search**: Extending the benchmarks to include hybrid search workloads (combining keyword filters with vector similarity) could better highlight Weaviate's strengths.
- **Cost Analysis**: Incorporating cloud deployments (e.g., AWS, Azure) to evaluate cost-performance trade-offs would benefit organizations considering managed services.

Overall, the results highlight that the choice between Milvus and Weaviate depends on the specific requirements of the application. For high-performance, large-scale similarity search, Milvus emerges as the most suitable system. However, applications requiring flexible schemas, semantic enrichment, or hybrid keyword–vector search may find Weaviate's design more advantageous.

## REFERENCES

[1] Zilliz, "Milvus: An Open-Source Vector Database for AI Applications," 2025. [Online]. Available: https://milvus.io/
[2] Semi Technologies, "Weaviate: The Open-Source Vector Search Engine," 2025. [Online]. Available: https://weaviate.io/
[3] Hugging Face, "Hugging Face Datasets and Transformers," 2025. [Online]. Available: https://huggingface.co/
[4] Docker Inc., "Docker: Build, Ship, and Run Any App, Anywhere," 2025. [Online]. Available: https://www.docker.com/
[5] Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 4, pp. 824–836, Apr. 2020. [Online]. Available: https://doi.org/10.1109/TPAMI.2018.2889473
[6] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with GPUs," *IEEE Trans. Big Data*, 2019. [Online]. Available: https://faiss.ai/
[7] E. Bernhardsson, "Annoy: Approximate Nearest Neighbors Oh Yeah," GitHub repository, 2018. [Online]. Available: https://github.com/spotify/annoy
[8] Qdrant, "ANN Filtering Benchmark Datasets," GitHub repository, 2023. [Online]. Available: https://github.com/qdrant/ann-filtering-benchmark-datasets

[9] GitHub Repository, "Big Data Management and Information Systems Project," 2025. [Online]. Available: https://github.com/aleginxx/big-data-management-and-information-systems/

[10] J. Leskovec, A. Rajaraman, and J. D. Ullman, *Mining of Massive Datasets*, 3rd ed. Cambridge, U.K.: Cambridge Univ. Press, 2020.

[11] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.

[12] A. Gionis, P. Indyk, and R. Motwani, "Similarity search in high dimensions via hashing," in *Proc. 25th Int. Conf. Very Large Data Bases (VLDB)*, Edinburgh, U.K., 1999, pp. 518–529.

[13] J. Dean and S. Ghemawat, "MapReduce: Simplified data processing on large clusters," *Commun. ACM*, vol. 51, no. 1, pp. 107–113, Jan. 2008.

[14] S. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 129–137, Mar. 1982.

[15] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, arXiv:1301.3781. [Online]. Available: https://arxiv.org/abs/1301.3781

[16] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with GPUs," *IEEE Trans. Big Data*, vol. 7, no. 3, pp. 535–547, Sep. 2021.

[17] H. Jegou, M. Douze, and C. Schmid, "Product quantization for nearest neighbor search," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 1, pp. 117–128, Jan. 2011.

[18] M. Babenko and V. Lempitsky, "Additive quantization for extreme vector compression," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Columbus, OH, USA, 2014, pp. 931–938.

[19] Weaviate Python Client Documentation, "Questions and Feedback," 2025. [Online]. Available: https://docs.weaviate.io/weaviate/client-libraries/pythonquestions-and-feedback

[20] Weaviate Python Client GitHub, 2025. [Online]. Available: https://github.com/weaviate/weaviate-python-client

[21] Weaviate Docker Installation Guide, 2025. [Online]. Available: https://docs.weaviate.io/deploy/installation-guides/docker-installation

[22] Weaviate Schema and Import Guide, 2025. [Online]. Available: https://docs.weaviate.io/academy/py/zero$_t o_m v p / s c h e m a_a n d_i m p o r t s / s c h e m a$

[22] Weaviate Collection Management, 2025. [Online]. Available: https://docs.weaviate.io/weaviate/starter-guides/managing-collections

[23] Milvus Quickstart Guide, 2025. [Online]. Available: https://milvus.io/docs/quickstart.md