

A Comparative Study of Apache Spark and Ray for Scalable Big Data Analytics and Machine Learning in Python

Kolydakis Manos

*Electrical & Computer Engineering
National Technical University of Athens
Athens, Greece
01320156*

Kouriannidis Iasonas

*Electrical & Computer Engineering
National Technical University of Athens
Athens, Greece
03120439*

Abstract—As data-intensive applications become more prevalent, the need for efficient and scalable big data processing frameworks has never been greater. This work presents a comparative analysis of two prominent distributed systems - Apache Spark and Ray - focusing on their performance in executing Extract, Transform, Load (ETL) processes and handling a range of machine learning workloads, from basic algorithms to more sophisticated AI models. To conduct our evaluation, we leverage datasets of varying sizes, sourced and synthetically generated, which are stored in a Hadoop Distributed File System (HDFS) deployed across a two-node cluster. We assess both frameworks based on critical criteria such as execution efficiency, scalability under load, and resilience to varying resource constraints. The results reveal distinct advantages for each platform: Apache Spark demonstrates strong proficiency in large-scale ETL operations, while Ray stands out in distributed machine learning tasks, offering tight integration with modern ML tools through Ray Datasets. Together, they represent complementary strengths in the evolving landscape of big data processing.

Index Terms—Big Data, Apache Spark, Ray, Distributed Computing, Machine Learning, Scalability, Python Frameworks, ETL, Benchmarking

I. INTRODUCTION

In the era of big data, organizations and researchers alike face the ongoing challenge of processing, analyzing, and drawing insights from vast and continuously growing datasets. To address this, distributed computing frameworks have emerged as essential tools for enabling scalable and efficient data processing across clusters of machines. Among the most widely adopted frameworks is Apache Spark, known for its robust support for batch processing, in-memory computation, and rich ecosystem for Extract, Transform, Load (ETL) operations. More recently, Ray has gained attention for its flexible and scalable architecture tailored toward distributed machine learning (ML) and reinforcement learning workloads.

A. HDFS and Hadoop

Hadoop is an open-source framework designed to store and process large-scale data sets across distributed computing environments. At the core of Hadoop is the Hadoop Distributed File System (HDFS), which is a scalable, fault-tolerant file

system that runs on commodity hardware and manages the storage of massive amounts of data. HDFS enables a single Hadoop cluster to scale from a handful of nodes to thousands, making it ideal for handling big data applications. The system is designed to provide high throughput access to data by distributing files into large blocks and replicating them across multiple nodes to ensure fault tolerance and data reliability. HDFS follows a master-slave architecture, where the NameNode manages the file system namespace and metadata, while multiple DataNodes handle the storage and retrieval of data blocks. This design supports parallel data processing by locating computation tasks close to where data resides, reducing network congestion and improving performance. Hadoop and HDFS together facilitate efficient batch processing, complex data analytics, and support a variety of applications, from data lakes to AI and machine learning workflows. Their portability across hardware platforms and compatibility with multiple programming languages make Hadoop and HDFS widely adopted solutions for organizations managing vast, diverse data sets [1].

B. Ray

Ray is an open-source unified framework designed to scale AI and Python applications, particularly in machine learning, without requiring expertise in distributed systems. It provides a compute layer that simplifies parallel processing and manages complex workflows such as data preprocessing, distributed training, hyperparameter tuning, reinforcement learning, and model serving. Ray offers Pythonic primitives for scaling Python applications and integrates smoothly with popular cloud platforms and cluster managers like Kubernetes, AWS, GCP, and Azure. Its architecture includes Ray Core for general-purpose distributed computing, Ray AI Libraries tailored for specific ML tasks, and Ray Clusters that can autoscale based on workload demands. Ray enables data scientists and ML engineers to seamlessly scale workloads from a laptop to large clusters using the same Python code, reducing the gap between development and production. Additionally, it handles essential distributed system tasks like orchestration,

scheduling, fault tolerance, and auto-scaling automatically, making it a powerful tool for building and deploying scalable machine learning applications [2].

C. Apache Spark

Apache Spark is an open-source distributed processing system designed for big data workloads, known for its fast analytic queries enabled by in-memory caching and optimized query execution. Supporting APIs in Java, Scala, Python, and R, Spark facilitates multiple workloads including batch processing, interactive queries, real-time analytics, machine learning, and graph processing. Originally developed in 2009 at UC Berkeley’s AMPLab, Spark aimed to overcome limitations of Hadoop MapReduce by reducing job steps and reusing data in-memory, making it significantly faster, especially for iterative tasks like machine learning. Unlike Hadoop, Spark does not have its own storage but works with systems like HDFS, Amazon S3, and others, often running on Hadoop clusters using YARN for resource management. Spark’s ecosystem includes core components like Spark SQL for fast queries, Spark Streaming for real-time data processing, MLlib for scalable machine learning, and GraphX for graph computations. Its cloud deployment is simplified through Amazon EMR, enabling quick cluster setup, auto-scaling, and cost savings, making Spark a popular choice across industries for scalable and efficient big data analytics [3].

II. EXPERIMENTAL SETUP

A. Hardware and OS

For the experiments, we used a distributed cluster of local virtual machines. The cluster consisted of three VMs, each with Ubuntu Server 22.04 LTS, 4 CPUs, 4 GB of RAM, and 50 GB of disk space.

B. Hadoop

The Hadoop configuration focused on key parameters to optimize cluster operation. We defined `dfs.replication` as 1, indicating a single copy of each data block for simplicity. The `yarn.resourcemanager.hostname` was set to `o-master` to identify the resource manager, and the web application address used the master’s IP on port 8088 for monitoring. Memory allocation parameters such as `yarn.nodemanager.resource.memory-mb` and `yarn.scheduler.maximum-allocation-mb` were configured to use 2GB at maximum. Additionally, auxiliary services like `mapreduce_shuffle` and `spark_shuffle` were enabled to facilitate shuffle operations for MapReduce and Spark within YARN.

C. Ray

The Ray cluster was configured to align resource settings with those used in the Spark environment for a fair comparison. Each node was allocated 4 CPU cores, mirroring Spark’s configuration. The Ray object store memory was set to approximately 1.2GB to match the Spark executor memory allocation.

D. Spark

The Spark configuration was set up to integrate with the Hadoop YARN cluster as the resource manager. Spark event logging was enabled, with logs stored in HDFS under the directory `/spark.eventLog` on the master node to support job history and monitoring. The Spark master was configured to run in YARN client mode, allowing dynamic resource allocation and management through YARN. Memory allocations were set to 1 GB each for the Spark driver and executors, with a single CPU core allocated per executor, balancing resource usage across the cluster. The Spark History Server was started to provide a web interface for monitoring completed Spark applications.

III. DATASETS

A. Data Discovery

The primary dataset utilized in this project is the Reddit Hyperlinks dataset [4], obtained from the Stanford SNAP repository [5]. This dataset contains records of subreddit interactions captured in a tab-separated values (TSV) format, detailing hyperlinks shared between subreddits along with various metadata and sentiment features. To facilitate processing, the original TSV data was converted into CSV format, ensuring compatibility with the data ingestion pipelines used in the project.

B. Data Generation

To support experiments requiring larger volumes of data, synthetic datasets were generated by replicating and modifying entries from the original dataset. A Python-based data generation script was developed to expand the dataset to a specified target size (e.g., in megabytes or gigabytes). The script samples from existing data attributes and ensures uniqueness in critical fields such as post identifiers, thereby creating enlarged datasets that preserve the statistical properties of the original data.

IV. EVALUATION TASKS

To evaluate the performance of Hadoop, Spark, and Ray across different workloads, we selected representative tasks that stress different aspects of distributed data processing systems. These tasks cover graph analytics, data preparation, and machine learning, providing a balanced assessment of computation, communication, and memory management. In particular, we focus on three categories of evaluation tasks: graph operations, extract-transform-load (ETL) workloads, and machine learning algorithms.

A. Graph Operations

Graph analytics workloads test the ability of distributed systems to handle irregular data structures and communication-heavy computations. We evaluated two well-known graph algorithms: *triangle counting* and *PageRank*.

1) *Triangle Counting*: The Triangle Count algorithm calculates how many triangles each node in a graph participates in. A triangle is formed when three nodes are all directly connected to each other, creating a fully connected triplet (also called a 3-clique). This technique is commonly applied in social network analysis to identify tightly connected groups and to measure how cohesive these groups are. It also provides insight into the overall structure and stability of a network and is often used in the computation of network metrics such as the Local Clustering Coefficient [6].

2) *PageRank*: PageRank is an algorithm designed to measure the relative importance of nodes within a graph. Originally developed to rank web pages, it can be applied to any network where nodes are connected by edges, such as social networks, citation networks, or transportation systems.

The algorithm assigns a numerical score to each node based on the structure of the graph: a node receives a higher score if it is linked to by other important nodes. Conceptually, PageRank models a “random walker” moving through the graph: at each step, the walker either follows an outgoing edge from the current node or jumps to a random node. Nodes that are more frequently visited by this process are considered more central or influential [7].

B. ETL

Extract-Transform-Load (ETL) tasks are essential for preparing raw data into a form suitable for analysis. In our benchmark, ETL workloads consisted of extracting text data from storage, performing quality checks, and transforming attributes through feature engineering. Operations included handling missing or invalid values, categorizing posts by length, readability, and sentiment, and computing composite measures such as engagement, complexity, and quality scores. The final stage involved cleansing invalid entries and exporting the processed dataset for downstream use. This benchmark highlights the importance of ETL in ensuring data reliability, scalability, and interpretability, while also stressing system I/O and computational efficiency.

C. Machine Learning

For the machine learning category, we selected the *k-means clustering* algorithm, one of the most widely used unsupervised learning methods. K-means partitions a dataset into clusters by iteratively minimizing the distance between points and their assigned centroids [8]. This algorithm is both computation and memory-intensive, as it requires repeated passes over the dataset and communication of updated centroids across nodes. Evaluating k-means allows us to assess how well each framework supports iterative, data-parallel machine learning workloads, which are central to modern analytics pipelines.

V. BENCHMARKS AND RESULTS

In this section, we present performance benchmarks for various workloads using Ray and Spark. We measure execution time and memory usage across different dataset sizes and number of workers.

A. Graph Operations

TABLE I: Triangle Counting Benchmark Results

Workers	Time (s)			Memory (MB)		
	Dataset (GB)	Ray	Spark	Dataset (GB)	Ray	Spark
3	1	111.02	64.71	1	159.66	456.58
3	5	2796.29	761.85	5	158.25	1509.84
2	5	2960.04	823.36	5	158.62	1463.75

Counting Triangles, Execution Time

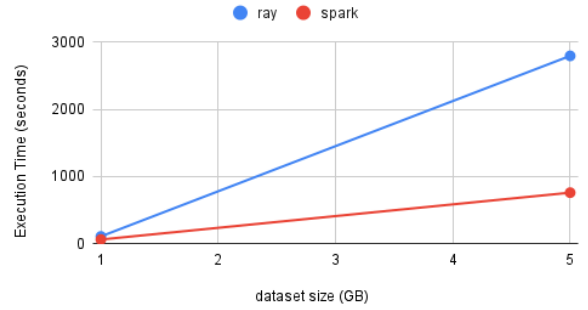


Fig. 1: Spark is significantly faster than Ray for triangle counting. As we can see in the graph, not only Ray is slower, but it also has a greater slope which means their difference in execution time increases with larger datasets.

Counting Triangles, Peak Memory

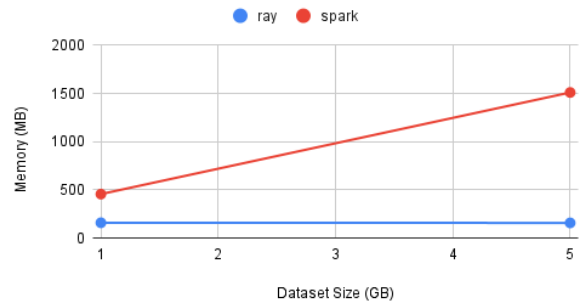


Fig. 2: Spark uses more memory (peak) than Ray for triangle counting. Spark often achieves faster runtimes by caching intermediate data, such as RDDs or DataFrames, in memory, reducing redundant recomputation but raising peak memory usage [9], [10]. Ray, on the other hand, favors streaming and chunked data processing without extensive in-memory caching. Its Ray Data feature is designed for scalable, memory-efficient streaming and handles overflow by spilling to disk when necessary [11].

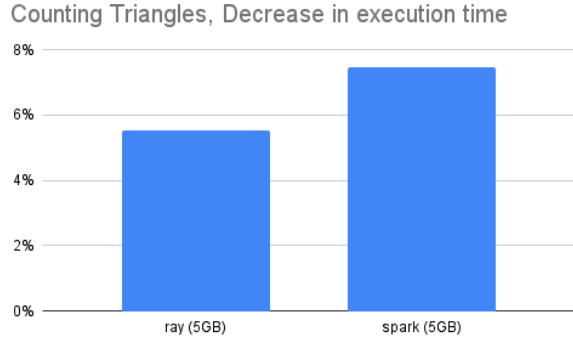


Fig. 3: Spark and Ray are equally good at scaling triangle counting with more workers.

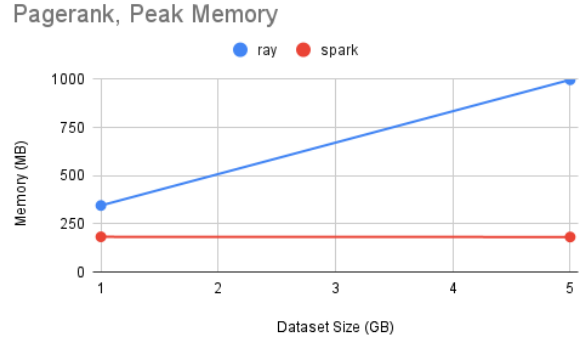


Fig. 6: For PageRank, Spark also used less memory. Spark's memory usage remains stable, while Ray's memory increases with dataset size.

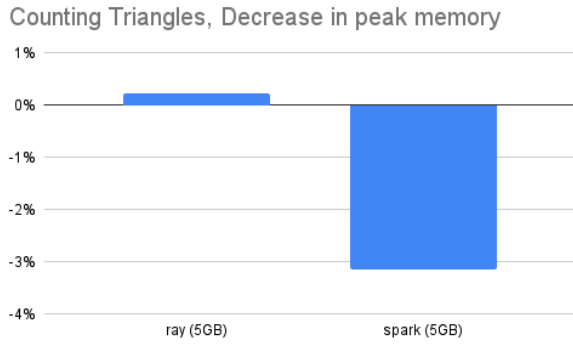


Fig. 4: Scaling didn't significantly impact peak memory used by Ray or Spark for triangle counting.

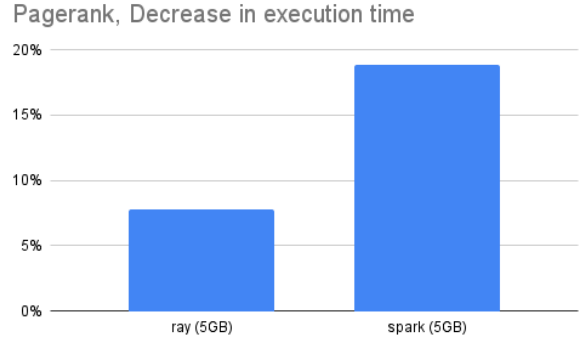


Fig. 7: For PageRank, Spark is faster at scaling, showing significant speedup with more workers.

TABLE II: PageRank Benchmark Results

Workers	Time (s)			Memory (MB)		
	Dataset (GB)	Ray	Spark	Dataset (GB)	Ray	Spark
3	1	93.65	69.32	1	345.43	183.16
3	5	467.54	144.78	5	996.92	182.16
2	5	506.98	178.43	5	995.28	182.48

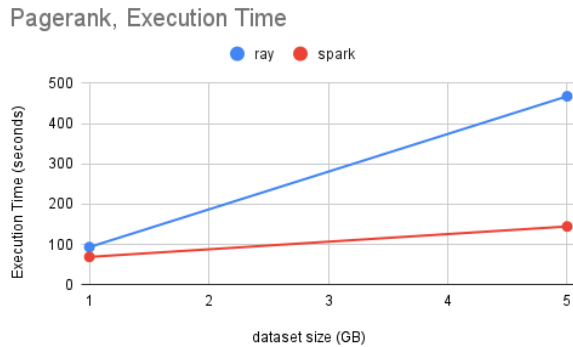


Fig. 5: For PageRank, Spark is faster than Ray.

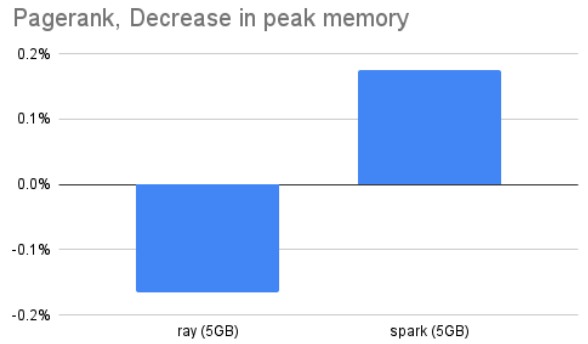


Fig. 8: Scaling didn't significantly impact peak memory used by Ray or Spark for PageRank.

B. ETL

TABLE III: ETL: Total (for Extract, Transform and Load) Benchmark Results

Workers	Time (s)			Memory (MB)		
	Dataset (GB)	Ray	Spark	Dataset (GB)	Ray	Spark
3	1	50.20	37.77	1	230.49	160.33
3	5	261.73	176.43	5	228.31	162.47
2	5	263.91	201.82	5	232.18	160.16

ETL, Execution Time

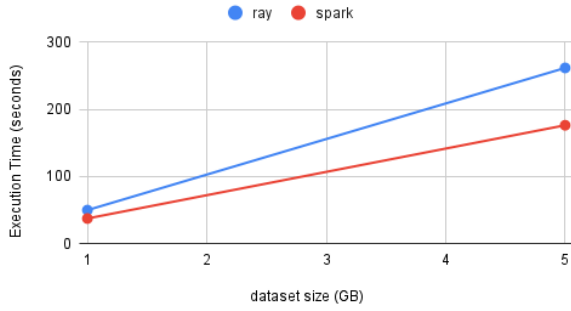


Fig. 9: ETL: Execution time comparison. Spark is faster than Ray.

ETL, Peak Memory

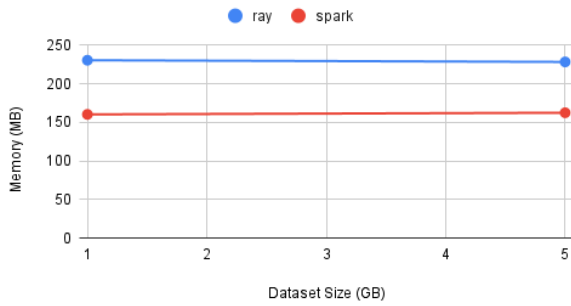


Fig. 10: ETL: Memory usage comparison. Ray uses more memory than Spark. Both remain stable for larger datasets.

ETL, Decrease in execution time

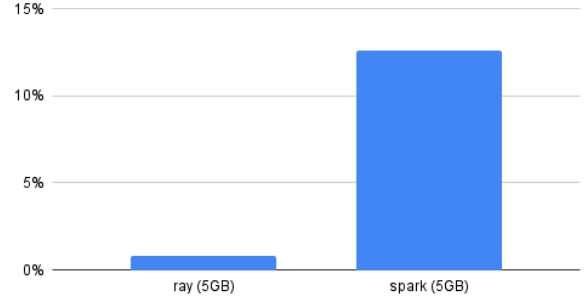


Fig. 11: ETL: Spark is faster at scaling, showing significant speedup with more workers. Ray's speedup is less pronounced.

ETL, Decrease in peak memory

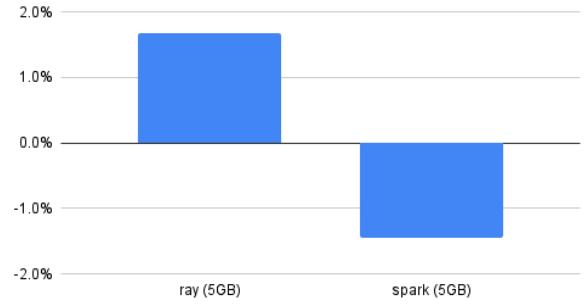


Fig. 12: ETL: Scaling didn't significantly impact peak memory used by Ray or Spark.

TABLE IV: ETL Extract: Execution time for Ray and Spark

Workers	Dataset (GB)	Ray	Spark
3	1	15.06	22.60
3	5	78.52	99.34
2	5	79.17	127.62

ETL Extract, Execution Time

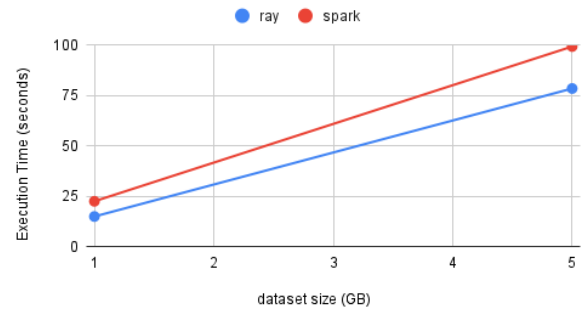


Fig. 13: ETL Extract phase: Similar results, but Ray is faster.

ETL Extract, Decrease in execution time

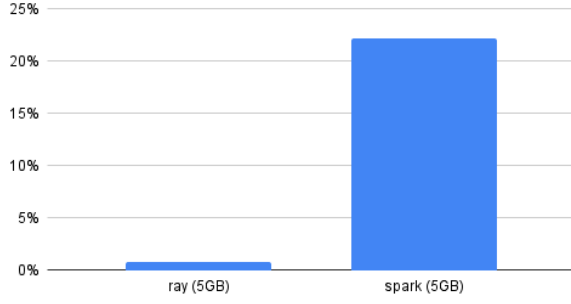


Fig. 14: ETL Extract phase: Spark is better at utilizing more workers.

TABLE V: ETL Transform: Execution time for Ray and Spark

Workers	Dataset (GB)	Ray	Spark
3	1	35.14	4.00
3	5	183.21	21.00
2	5	184.74	19.06

ETL Transform, Execution Time

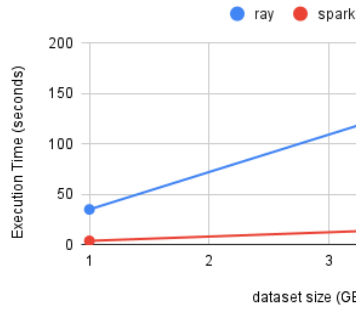


Fig. 15: ETL Transform phase: Spark is significantly faster than Ray for Transform operations.

ETL Transform, Decrease in execution time

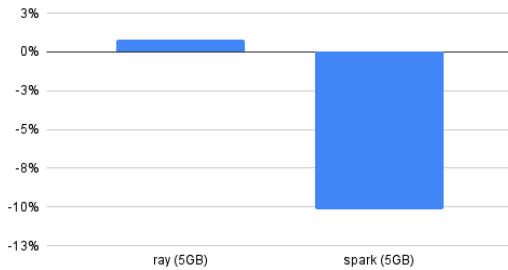


Fig. 16: ETL Transform: Spark increased memory usage when scaling.

TABLE VI: ETL Load: Execution time for Ray and Spark

Workers	Dataset (GB)	Ray	Spark
3	1	0.00	11.17
3	5	0.00	56.09
2	5	0.00	55.14

ETL Load, Execution Time

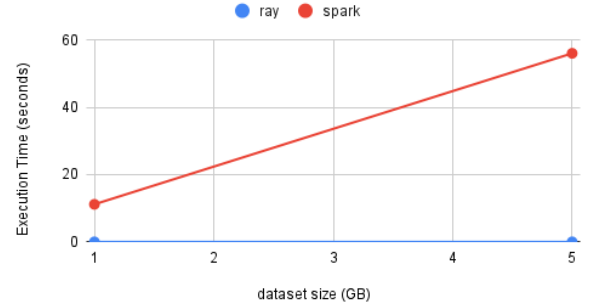


Fig. 17: ETL Load phase: Ray achieves near-zero load time due to in-memory optimizations, whereas Spark's load time increases with dataset size.

ETL Load, Decrease in execution time

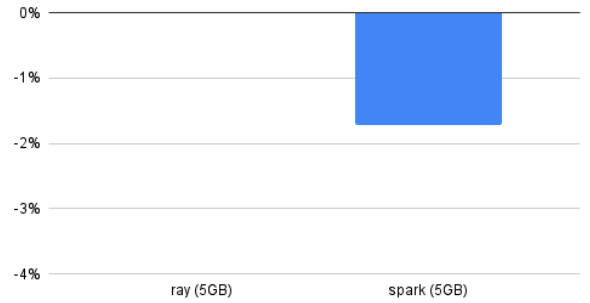


Fig. 18: ETL Load phase: Scaling didn't significantly impact peak memory used by Ray or Spark.

C. Machine Learning

TABLE VII: KMeans Benchmark Results

Workers	Time (s)			Memory (MB)		
	Dataset (GB)	Ray	Spark	Dataset (GB)	Ray	Spark
3	1	14.52	106.48	1	367.60	168.13
3	5	194.62	247.24	5	705.59	160.21
2	5	244.45	313.99	5	708.64	169.49

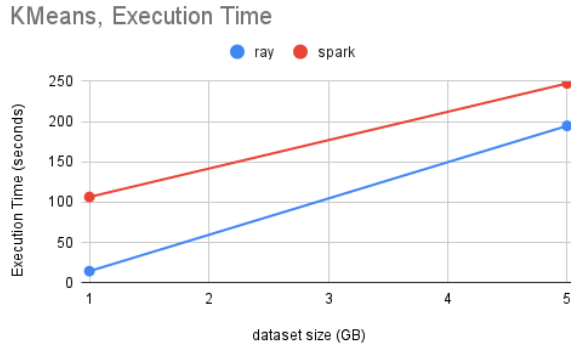


Fig. 19: KMeans: Ray is significantly faster for small datasets and remains faster than Sparks for larger datasets too.

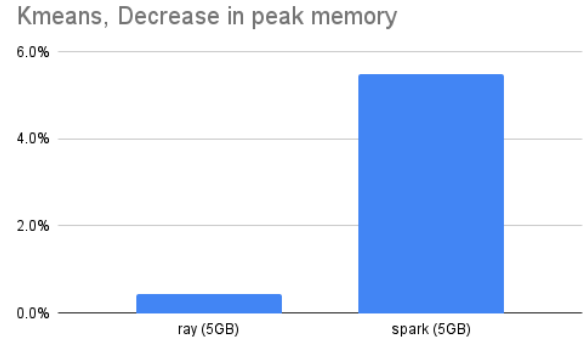


Fig. 22: KMeans: Spark used 5% less memory when scaling. Ray didn't show improvement.

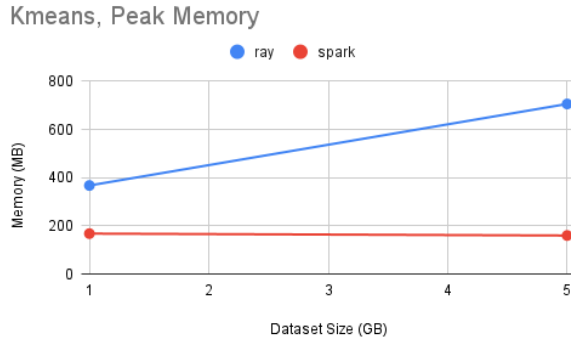


Fig. 20: KMeans: Spark used less memory than Ray.

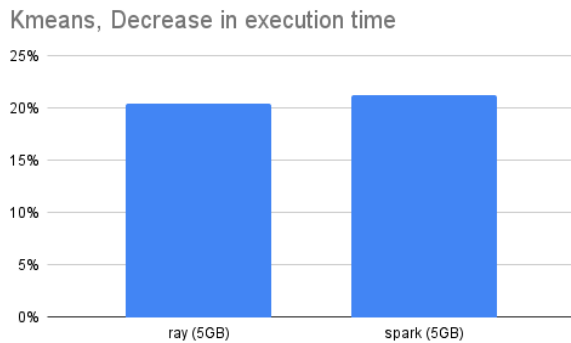


Fig. 21: KMeans: Both showed significant improvement in execution time when scaling.

VI. COMPARISON

A. Performance

In graph processing workloads, Spark consistently outperforms Ray in terms of execution time, reflecting its specialized graph processing optimizations. Memory usage, however, depends on the specific algorithm: Ray requires less memory for triangle counting, while Spark is more memory-efficient for PageRank computations due to its caching and shuffle strategies. In ETL workloads, each framework demonstrates strengths in different phases. Spark achieves significantly faster performance during the transformation phase, leveraging in-memory operations and optimized data shuffles, but it is relatively slower in data extraction and loading. In contrast, Ray provides faster execution for machine learning tasks, though this advantage comes with higher memory consumption. Overall, Spark offers superior performance for graph operations and ETL transformations, whereas Ray is more efficient for machine learning and ETL extraction and loading phases. These results suggest a general trade-off: frameworks that prioritize execution speed often incur higher memory usage, while those optimizing memory tend to sacrifice some runtime performance.

B. Scalability

Both Ray and Spark exhibit notable improvements in execution time when scaling graph operations and K-Means clustering workloads, demonstrating their ability to efficiently leverage additional computational resources. Among these tasks, Spark maintains a performance advantage for PageRank computations, likely due to its optimized graph processing libraries and in-memory caching mechanisms. In ETL workloads, the behavior of the frameworks diverges. Ray shows relatively consistent performance across extraction, transformation, and loading phases, with minimal improvement when scaling resources. Spark, in contrast, delivers mixed results: it achieves significant acceleration during the extraction phase but exhibits slower performance during transformation and loading, suggesting that the efficiency gains from parallelization can be limited by I/O and memory bottlenecks in these phases. Regarding memory usage, both frameworks

demonstrate comparable patterns when the number of workers increases, indicating that scaling primarily improves runtime performance without substantially affecting peak memory consumption. These observations highlight the importance of selecting the appropriate framework based on the specific workload characteristics and performance priorities.

C. Ease of Use

Both Ray and Spark provide high-level APIs that simplify the development of distributed data processing pipelines, though they differ in abstraction and workflow management. Ray was much easier to set up, as it is installed as a Python library and requires minimal configuration to run, making it accessible for quick experiments and small-scale workflows. Spark, in contrast, needs substantial configuration, including cluster setup, resource management, and environment tuning, which can be time-consuming and challenging for new users. Once configured, Spark provides optimized execution and built-in task scheduling for large workloads, but the initial setup overhead is significant. Overall, Ray offers a faster and simpler way to start distributed computations, whereas Spark is more complex to initialize but provides extensive optimizations for high-performance workloads.

VII. AVAILABILITY AND REPRODUCIBILITY

All experiments were performed on publicly available datasets and using open-source versions of Ray and Spark. Scripts, configuration files, and instructions for reproducing the experiments are provided in the supplementary materials and a public repository (link omitted here). The experiments use standardized benchmarks for PageRank, ETL, triangle counting, and KMeans, ensuring that the results can be independently verified and extended by other researchers. Containerization and virtual environment specifications further improve reproducibility.

The scripts, configuration files, and instructions for reproducing the experiments are available at our GitHub repository.

VIII. CONCLUSION

This study provides a comparative evaluation of Ray and Spark across graph processing, ETL, and machine learning workloads. Spark demonstrates superior performance in graph operations and ETL transformation tasks, largely due to its optimized libraries, in-memory caching, and efficient shuffle strategies. Ray, however, exhibits advantages in machine learning workloads and the extraction and loading phases of ETL, benefiting from its flexible task-based execution model. Memory usage patterns indicate a trade-off: frameworks that prioritize execution speed often incur higher memory consumption, while memory-efficient frameworks may experience slower runtimes. Both frameworks scale effectively with additional computational resources, although Spark’s specialized optimizations allow it to maintain a performance edge in certain workloads.

REFERENCES

- [1] IBM, “What is HDFS?,” IBM. Available: <https://www.ibm.com/think/topics/hdfs>.
- [2] Ray Project, “Overview - Ray 2.48.0”. Available: <https://docs.ray.io/en/latest/ray-overview/index.html>.
- [3] Amazon Web Services, “What is Apache Spark?,”. Available: <https://aws.amazon.com/what-is/apache-spark/>.
- [4] S. Kumar, W. L. Hamilton, J. Leskovec, and D. Jurafsky, “Community interaction and conflict on the web,” in *Proc. of the 2018 World Wide Web Conference (WWW)*, pp. 933–943, 2018.
- [5] S. Kumar, W.L. Hamilton, J. Leskovec, D. Jurafsky, “Reddit Hyperlinks Dataset” Stanford SNAP. Available: <https://snap.stanford.edu/data/soc-RedditHyperlinks.html>.
- [6] Neo4j, “Triangle Count,” Neo4j Graph Data Science. Available: <https://neo4j.com/docs/graph-data-science/current/algorithms/triangle-count/>.
- [7] E. Roberts and K. Schroeder, “The Google PageRank Algorithm,” Stanford University CS54N. Available: <https://web.stanford.edu/class/cs54n/handouts/24-GooglePageRankAlgorithm.pdf>.
- [8] C. Piech, “K-Means,” Stanford University CS221. Available: <https://stanford.edu/~cpiech/cs221/handouts/kmeans.html>.
- [9] M. Zaharia, R. Xin, P. Wendell, T. Das, M. Armbrust, A. Dave, X. Meng, J. Rosen, S. Venkataraman, M. J. Franklin, A. Ghodsi, J. Gonzalez, S. Shenker, and I. Stoica, “Apache Spark: A Unified Engine for Big Data Processing,” *arXiv:1804.10563*, 2018. Available: <https://arxiv.org/abs/1804.10563>.
- [10] Stack Overflow, “Spark cache vs broadcast,” Available: <https://stackoverflow.com/questions/38056774/spark-cache-vs-broadcast>.
- [11] Ray Project, “Performance Tips for Ray Data,” Available: <https://docs.ray.io/en/latest/data/performance-tips.html>.