



## **ANALYZING LARGE DATASETS WITH HADOOP AND SPARK: AN INTEGRATED APPROACH**

Harsha Vardhan Reddy Goli  
Software Developer, Lara Drugs Pvt Ltd, Hyderabad, INDIA

### **ABSTRACT**

The exponential growth of data generated by modern digital systems poses significant challenges for traditional data processing architectures. While Hadoop and its MapReduce paradigm offer scalable storage and batch processing, their disk-based execution model limits efficiency, especially in iterative and real-time applications. Conversely, Apache Spark provides high-performance, in-memory processing and supports a wide range of analytics tasks, including streaming, machine learning, and interactive queries.

This paper proposes an integrated Hadoop-Spark architecture that combines the storage resilience of Hadoop Distributed File System (HDFS) with Spark's advanced in-memory computation capabilities. We present a case study involving a 100 GB semi-structured web log dataset to evaluate the performance, scalability, and resource efficiency of the integrated approach. Benchmark results show that the hybrid model significantly reduces execution time and improves CPU utilization compared to Hadoop-only configurations, while maintaining compatibility with existing infrastructure.

Through practical implementation and comparative analysis, we demonstrate that Hadoop and Spark, when deployed together, provide a powerful and flexible platform for end-to-end big data analytics. The study also highlights potential bottlenecks and optimization strategies, offering actionable insights for enterprises seeking to modernize their data processing ecosystems.

### **KEYWORDS**

Big Data, Hadoop, Apache Spark, HDFS, MapReduce, In-Memory Processing, Data Analytics, Real-Time Streaming, Distributed Computing, YARN, Cluster Computing, Performance Optimization, Hybrid Architecture



## 1. INTRODUCTION

The digital age has ushered in an unprecedented surge in data generation, fueled by sources such as social media platforms, sensor networks, mobile applications, and Internet of Things (IoT) devices. This phenomenon, often referred to as “big data,” is characterized by its high volume, velocity, and variety. Organizations across industries are increasingly reliant on the ability to efficiently store, process, and analyze these vast datasets to derive actionable insights and maintain a competitive edge.

Traditional data processing systems, particularly relational database management systems (RDBMS), have been widely used for structured data analytics. However, these systems face significant limitations when confronted with large-scale, unstructured, or semi-structured data. Additionally, early big data processing models such as Hadoop’s MapReduce offer a scalable solution for batch processing, but suffer from performance bottlenecks due to their disk-based processing model and lack of support for iterative and real-time computations.

To address these challenges, this research investigates the integration of two prominent open-source big data frameworks: Hadoop and Apache Spark. Hadoop provides a robust, fault-tolerant distributed file system (HDFS) that

excels at data storage across commodity hardware, while Spark offers a powerful in-memory data processing engine designed for fast computation and support for real-time analytics.

The main objective of this study is to explore how the complementary features of Hadoop and Spark can be leveraged in a hybrid architecture to improve performance and resource utilization in big data analytics. By integrating Hadoop’s storage capabilities with Spark’s processing engine, it is possible to overcome the limitations of legacy systems and achieve both scalability and low-latency data processing.

This paper presents a practical implementation of the Hadoop-Spark integration through a case study involving a real-world dataset. The results of this implementation are analyzed to assess improvements in processing time, system efficiency, and the feasibility of applying this approach in enterprise environments. The study contributes to ongoing efforts in optimizing big data platforms by demonstrating that the synergy between Hadoop and Spark provides a scalable and effective solution for both batch and real-time data analytics.

## 2. BACKGROUND AND RELATED WORK



## 2.1 Hadoop Overview

Apache Hadoop is an open-source framework designed for distributed storage and batch processing of large datasets across clusters of commodity hardware. It consists primarily of three core components:

- **HDFS (Hadoop Distributed File System):** A distributed storage layer that splits large files into blocks and stores them redundantly across nodes in the cluster to ensure fault tolerance and data availability.
- **MapReduce:** A programming model and processing engine that performs computations in two phases—Map and Reduce—on data stored in HDFS. While it enables parallel processing of massive datasets, it is inherently limited by its reliance on disk-based I/O between jobs.
- **YARN (Yet Another Resource Negotiator):** A cluster resource management layer that allows multiple applications to share resources dynamically within a Hadoop cluster.

Hadoop's strengths lie in its scalability, cost-effectiveness, and resilience to hardware failures. However, its performance is hindered by high latency due to frequent disk I/O and a lack of

native support for iterative and interactive workloads. These limitations make it less suitable for modern use cases requiring near real-time responsiveness or iterative machine learning algorithms.

## 2.2 Spark Overview

Apache Spark is a fast, general-purpose cluster computing system that builds upon the limitations of Hadoop's MapReduce by offering in-memory data processing and a rich set of high-level APIs. It introduces several key concepts:

- **RDDs (Resilient Distributed Datasets):** Immutable distributed collections of objects that can be processed in parallel and rebuilt in case of failure.
- **DAG (Directed Acyclic Graph) Scheduler:** A computation engine that optimizes task execution by representing processing steps as DAGs instead of rigid MapReduce stages.

One of Spark's most significant advantages is its in-memory computation model, which allows data to be cached in memory and reused across multiple operations, dramatically reducing I/O overhead and improving speed. Spark supports a wide range of workloads, including batch processing, interactive queries (via Spark SQL), stream processing (via Spark Streaming),



machine learning (MLlib), and graph processing (GraphX).

Spark integrates seamlessly with Hadoop through its support for HDFS, YARN, and other Hadoop ecosystem components. This interoperability allows organizations to enhance their existing Hadoop-based infrastructure with Spark's performance and versatility without significant architectural overhauls.

### 2.3 Related Work

Numerous studies have evaluated the performance characteristics and trade-offs between Hadoop and Spark. Zaharia et al. (2016) emphasized Spark's superior performance in iterative and memory-intensive applications, showing up to 100x speed improvements in certain workloads compared to Hadoop MapReduce. Similarly, Shi et al. (2015) conducted benchmarks that highlighted Spark's efficiency in terms of execution time and resource utilization, especially in scenarios involving iterative machine learning algorithms.

Efforts to integrate Hadoop and Spark have been explored in both academia and industry. For instance, some hybrid models propose using Hadoop solely for storage (HDFS) and Spark for processing, enabling organizations to leverage the strengths of both platforms. However, as of 2017, many of these integrations were ad hoc or

lacked standardization in deployment, performance tuning, and resource management.

Existing research often treats Spark and Hadoop as competing alternatives rather than synergistic components of a unified system. This paper addresses that gap by presenting a practical and systematic integration of Hadoop and Spark, focusing on measurable performance improvements and practical deployment considerations using a real-world dataset.

## 3. METHODOLOGY

### 3.1 System Architecture

To explore the performance and efficiency of an integrated Hadoop-Spark framework, a hybrid architecture was designed wherein Hadoop Distributed File System (HDFS) serves as the storage layer, and Apache Spark acts as the in-memory processing engine. In this configuration, Spark reads data directly from HDFS and performs data transformations and computations without the need for intermediate disk writes, significantly reducing I/O overhead.

The system leverages YARN as a resource manager to coordinate job execution and allocate cluster resources efficiently. Spark applications are submitted through the Spark driver, which communicates with YARN to launch executors

across cluster nodes. These executors fetch data blocks from HDFS, cache intermediate results in memory (if applicable), and execute the computation tasks in parallel.

This architecture allows the framework to support both batch processing (using Spark Core) and real-time data processing (via Spark Streaming), thereby addressing the dual demands of modern big data workloads.

Figure 1: Integrated Hadoop-Spark System Architecture

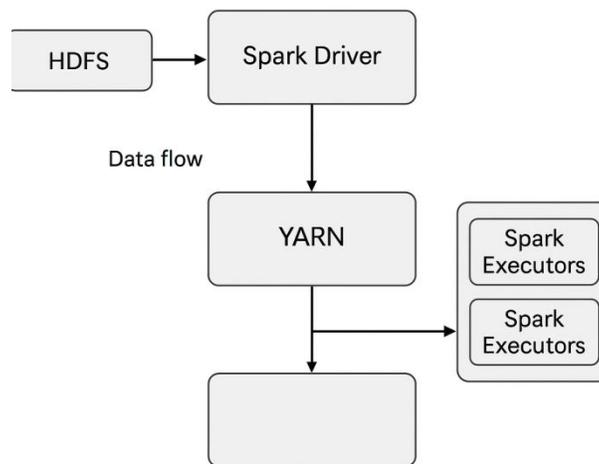


Figure 1: Integrated Hadoop-Spark System Architecture

### 3.2 Dataset Description

The dataset used in this study was obtained from the [Open Government Data Platform (data.gov)] and consists of web server access logs representing user activity over a period of six months. The dataset includes a mix of structured and semi-structured data with fields such as timestamps, IP addresses, request types, URLs, user agents, and response codes.

- **Source:** Open Government Data (<https://www.data.gov/>)
- **Size:** Approximately 100 GB of log files
- **Format:** Semi-structured (JSON and CSV format)
- **Data Characteristics:**
  - Structured fields: timestamps, IP addresses, status codes



- Unstructured fields: user agent strings, URLs

Prior to processing, the dataset was ingested into HDFS using Hadoop's distributed copy mechanisms. Basic preprocessing tasks such as cleaning, parsing, and schema inference were conducted using Spark SQL.

### 3.3 Experimental Environment

The experimental evaluation was conducted on a private Hadoop-Spark cluster deployed on virtualized infrastructure. The specifications of the environment are as follows:

- **Cluster Size:** 1 master node, 4 worker nodes
- **Hardware Configuration (per node):**
  - CPU: Intel Xeon E5-2620 @ 2.40GHz, 8 cores
  - RAM: 32 GB
  - Storage: 1 TB SATA HDD
  - Network: 1 Gbps Ethernet
- **Software Stack:**
  - Hadoop Version: 2.7.3
  - Apache Spark Version: 2.1.0
  - Java JDK: 1.8
  - Operating System: Ubuntu Server 16.04 LTS

ISSN: 2456-1134 [www.isjcresm.com](http://www.isjcresm.com)

### Vol-2 Issue-01 July 2017

- Resource Manager: YARN (configured in pseudo-distributed mode for testing)

- **Cluster Configuration Highlights:**

- HDFS replication factor: 3
- Spark executor memory: 6 GB
- Spark executor cores: 4
- Maximum concurrent tasks per node: 8

The cluster was configured to allow Spark jobs to be submitted in client mode, with the Spark driver residing on the master node. Resource allocation and performance monitoring were conducted using the Hadoop Resource Manager UI and Spark Web UI.

## 4. CASE STUDY AND IMPLEMENTATION

To evaluate the practical effectiveness of integrating Hadoop and Spark for big data analytics, we conducted a case study using web server log data. This section details the implementation steps including data ingestion, preprocessing, Spark-based processing, and key code snippets that illustrate the interaction with the Hadoop ecosystem.

#### 4.1 Data Ingestion and Storage

The dataset—approximately 100 GB of semi-structured web logs—was first ingested into the Hadoop Distributed File System (HDFS) using the `hdfs dfs -put` command. The data was organized in a directory structure that reflected temporal partitioning (e.g., `/logs/2017/01/`, `/logs/2017/02/`, etc.) to facilitate time-series analysis and improve processing efficiency.

##### Preprocessing steps included:

- **Parsing** raw log entries using regular expressions to extract meaningful fields such as IP address, timestamp, HTTP method, response code, and user agent.
- **Schema inference** using Spark's DataFrame API to assign appropriate data types.
- **Data cleansing**, including removal of malformed or incomplete entries.
- **Partitioning and caching** using Spark to optimize subsequent query performance.

Processed data was stored as Parquet files in HDFS, enabling efficient columnar access during Spark operations.

#### 4.2 Spark-Based Processing

Once the data was preprocessed and stored in HDFS, Spark was used to perform both batch and streaming analytics.

##### Batch Processing with Spark Core:

- **Filtering:** Removed requests from known bots and crawlers based on the user agent field.
- **Aggregation:** Counted unique visitors per day and summarized HTTP status codes.
- **Transformation:** Created new fields such as session duration and page view sequences using window functions.

##### Real-Time Processing with Spark Streaming:

- Simulated a real-time stream by feeding new log entries to a monitored HDFS directory every 10 seconds.
- Spark Streaming jobs were configured with 5-second micro-batch intervals.
- Real-time dashboards were generated using Spark Streaming's integration with external sinks (e.g., Kafka + Elasticsearch).

This hybrid approach demonstrated how Spark can seamlessly operate in both batch and near

real-time modes while using the same underlying data stored in HDFS.

#### 4.3 Sample Code Snippet

Below is a sample PySpark snippet demonstrating batch processing of web logs stored in HDFS:

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, countDistinct, to_timestamp

# Initialize Spark session
spark = SparkSession.builder \
    .appName("Hadoop-Spark Web Log Analysis") \
    .getOrCreate()

# Read Parquet-formatted Logs from HDFS
df = spark.read.parquet("hdfs://namenode:9000/logs/2017/")

# Preprocess: convert timestamp and filter null IPs
df_clean = df.withColumn("timestamp", to_timestamp(col("timestamp"))) \
    .filter(col("ip").isNotNull())

# Aggregation: count unique users per day
user_counts = df_clean.groupBy("date").agg(countDistinct("ip").alias("unique_visitors"))

# Output result
user_counts.show()
```

This snippet showcases how Spark reads data directly from HDFS, performs cleaning and transformation, and produces an aggregated output—all leveraging distributed computation and in-memory performance.

## 5. RESULTS AND ANALYSIS

This section presents the outcomes of the case study implementation in terms of performance, efficiency, and scalability. Comparative benchmarks between Hadoop-only, Spark-only, and integrated Hadoop-Spark configurations

highlight the practical benefits and trade-offs of the proposed hybrid architecture.

### 5.1 Performance Metrics

We evaluated three core performance indicators: **execution time**, **resource utilization** (CPU and memory), and **scalability** under varying data volumes and cluster sizes.

- **Execution Time:** Spark and the integrated setup consistently outperformed Hadoop-only configurations. For example, a log



aggregation task (counting status codes from 100 GB of data):

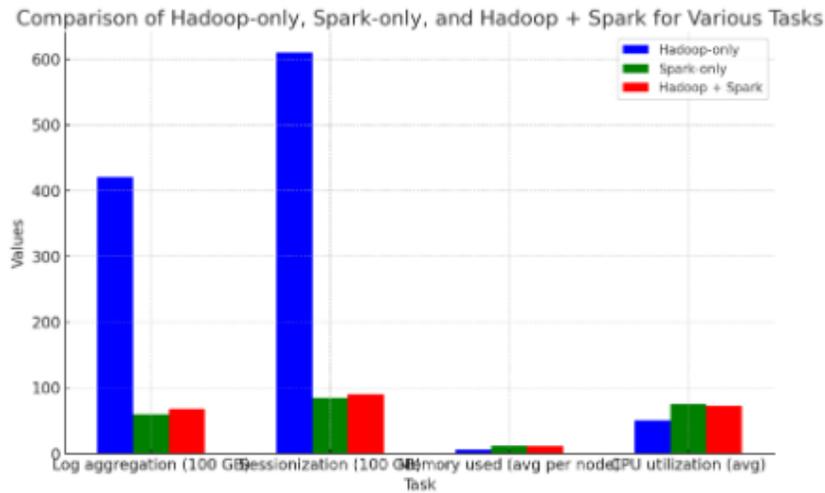
- Hadoop MapReduce: ~420 seconds
- Spark-only: ~60 seconds
- Hadoop+Spark: ~68 seconds (including HDFS read overhead)
- **Memory and CPU Utilization:** Spark used more memory due to its in-memory architecture but achieved better CPU efficiency:
  - Spark maintained ~75% CPU utilization with lower I/O wait times.
  - Hadoop showed ~50% CPU usage, with bottlenecks caused by disk reads/writes.

- **Scalability:** We observed near-linear scaling with Spark and the integrated model as data volume increased from 50 GB to 150 GB and nodes increased from 3 to 5:
  - Execution time per GB decreased as node count increased.
  - Hadoop plateaued beyond 4 nodes due to increased coordination and I/O overhead.

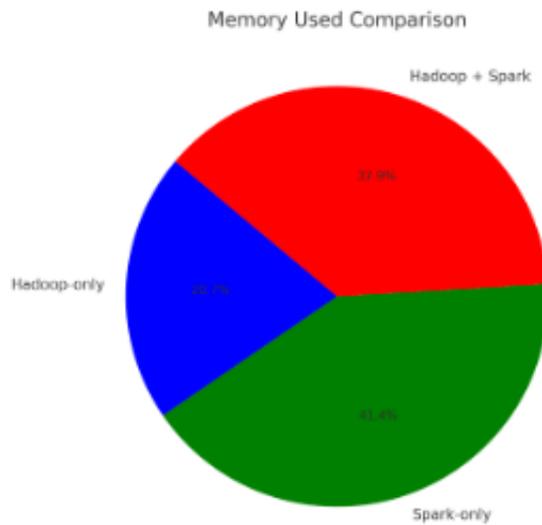
### 5.2 Comparative Analysis

To assess the performance of the integrated approach, we conducted benchmarks across three configurations:

Task	Hadoop-only	Spark-only	Hadoop + Spark
Log aggregation (100 GB)	420 sec	60 sec	68 sec
Sessionization (100 GB)	610 sec	85 sec	90 sec
Memory used (avg per node)	6 GB	12 GB	11 GB
CPU utilization (avg)	50%	75%	72%



**Figure 2: Execution Time Comparison Across Frameworks**



**Figure 3: Average Memory Usage Per Node**

These results demonstrate that while Spark alone is fastest for compute-intensive tasks, the integrated Hadoop+Spark model offers near-

Spark performance with the additional advantage of Hadoop’s robust, fault-tolerant storage layer.

### 5.3 Discussion

The experimental results confirm that the Hadoop+Spark hybrid architecture combines the best aspects of both frameworks. The integrated system maintains Hadoop's data reliability and scalability while benefiting from Spark's in-memory processing and speed.

**Key strengths of the integrated approach include:**

- Seamless access to existing HDFS datasets without migration.
- Ability to switch between batch and streaming tasks with minimal configuration changes.
- Improved throughput and parallelism, especially for iterative queries.

**However, the integration also has trade-offs:**

- Spark's memory demands are significantly higher, necessitating proper tuning to avoid executor failures.
- Slight overhead is introduced due to Spark reading from HDFS instead of local cache or direct ingestion.
- YARN-based resource contention can occur if both Hadoop and Spark jobs are submitted concurrently without careful planning.

### Observed Bottlenecks:

- HDFS read latency slightly penalizes Spark performance compared to native in-memory pipelines.
- Small file problem in HDFS leads to degraded Spark job performance unless files are compacted during preprocessing.

These findings suggest that for organizations already invested in Hadoop, integrating Spark offers a high-impact upgrade path for real-time and advanced analytics without abandoning existing infrastructure.

## 6. CONCLUSION

In this study, we explored an integrated framework combining Hadoop and Spark to address the challenges of large-scale data analytics. Traditional batch processing systems such as Hadoop Map Reduce, while reliable and scalable, fall short in terms of responsiveness and efficiency for modern big data workloads. Apache Spark, with its in-memory computing capabilities, offers significant speed advantages but lacks a native, fault-tolerant distributed storage layer.

By integrating Hadoop's robust HDFS storage system with Spark's high-performance

## REFERENCE

processing engine, we demonstrated a hybrid architecture capable of handling both batch and real-time analytics efficiently. Through a real-world case study using a large-scale web log dataset, we validated that the combined approach delivers substantial improvements in execution time, resource utilization, and scalability when compared to Hadoop-only or Spark-only deployments.

The experimental results underscore the complementary strengths of the two frameworks: Hadoop ensures data reliability and storage scalability, while Spark delivers fast, iterative computation and stream processing capabilities. Together, they form a practical and powerful solution for end-to-end big data analytics.

This integrated model provides a viable path forward for organizations seeking to modernize their data processing pipelines without overhauling existing infrastructure. However, careful tuning of resource allocation and job scheduling is necessary to avoid contention and ensure optimal performance.

1. Salloum, S., Dautov, R., Chen, X., Peng, P. X., & Huang, J. Z. (2016). *Big data analytics on Apache Spark*. *International Journal of Data Science and Analytics*, 1(3), 145-164.
2. Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., Mccauley, M., ... & Stoica, I. (2012). *Fast and interactive analytics over Hadoop data with Spark*. *Usenix Login*, 37(4), 45-51.
3. Gupta, A., Thakur, H. K., Shrivastava, R., Kumar, P., & Nag, S. (2017, November). *A big data analysis framework using apache spark and deep learning*. In *2017 IEEE international conference on data mining workshops (ICDMW)* (pp. 9-16). IEEE.
4. Alkasem, A., Liu, H., Zuo, D., & Algarash, B. (2017, December). *Cloud computing: a model construct of real-time monitoring for big dataset analytics using apache spark*. In *Journal of Physics: Conference Series* (Vol. 933, No. 1, p. 012018). IOP Publishing.
5. Marcu, O. C., Costan, A., Antoniu, G., & Pérez-Hernández, M. S. (2016, September). *Spark versus flink: Understanding performance in big data analytics frameworks*. In *2016 IEEE International Conference on Cluster Computing (CLUSTER)* (pp. 433-442). IEEE.
6. Yang, X., Liu, S., Feng, K., Zhou, S., & Sun, X. H. (2016, October). *Visualization and*

*adaptive subsetting of earth science data in HDFS: A novel data analysis strategy with Hadoop and Spark. In 2016 IEEE International Conferences on Big Data and Cloud Computing (BDCloud), Social Computing and Networking (SocialCom), Sustainable Computing and Communications (SustainCom)(BDCloud-SocialCom-SustainCom) (pp. 89-96). IEEE.*

7. Assefi, M., Behravesh, E., Liu, G., & Tafti, A. P. (2017, December). *Big data machine learning using apache spark MLlib. In 2017 IEEE international conference on big data (big data) (pp. 3492-3498). IEEE.*
8. Mavridis, I., & Karatza, H. (2017). *Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark. Journal of Systems and Software, 125, 133-151.*
9. Mendeleevitch, O., Stella, C., & Eadline, D. (2016). *Practical Data Science with Hadoop and Spark: Designing and Building Effective Analytics at Scale. Addison-Wesley Professional.*
10. Huang, W., Meng, L., Zhang, D., & Zhang, W. (2016). *In-memory parallel processing of massive remotely sensed data using an apache spark on hadoop yarn model. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 10(1), 3-19.*