# Modern Enterprise Data Analysis: From Legacy MapReduce to Cloud-Native Architectures

**Sohil Sri Mani Yeshwanth Grandhi**

Computer Science and Engineering Department Pennsylvania State University, University Park, United States
Email: *sohilg002[at]email.com*

**Abstract:** *This paper presents a comprehensive analysis of enterprise data processing methodologies, tracing the evolution from traditional MapReduce-based approaches to contemporary cloud-native architectures. We demonstrate a modernized implementation of large-scale data analysis originally conceived using Hadoop ecosystem technologies (Hive, Pig, HBase) but re-engineered with current technologies including Apache Spark, Delta Lake, and cloud-native services. Our work includes data processing pipelines, machine learning implementations using PySpark and MLlib, and real-time analytics capabilities. We provide comparative analysis showing significant performance improvements over legacy approaches and discuss emerging alternatives including serverless architectures, data lake houses, and AI-enhanced analytics platforms. The research demonstrates a 3.8x performance improvement in processing throughput and 60*

**Keywords:** Big Data, Apache Spark, Cloud Computing, Data Analytics, Machine Learning, Enterprise Systems, Data Lake, Real-time Processing

## 1. Introduction

The exponential growth of enterprise data has fundamentally transformed how organizations approach data processing and analytics. While early big data solutions centered around Hadoop and MapReduce provided the foundation for large-scale processing, the technological landscape has evolved dramatically by 2021-2022. The emergence of cloud-native architectures, advanced processing frameworks, and machine learning integration has created new paradigms for enterprise data analysis.

### 1) Background and Evolution
The Hadoop ecosystem, born from Google's MapReduce paper in 2004, dominated big data processing for over a decade. However, by 2021, limitations in performance, complexity, and operational overhead led organizations to seek modern alternatives. The shift from on-premise Hadoop clusters to cloud-based serverless architectures represents one of the most significant transformations in enterprise data management.

Contemporary data platforms now leverage technologies like Apache Spark for in-memory processing, Kubernetes for container orchestration, and cloud-native services that offer superior scalability, cost-efficiency, and developer productivity. The integration of machine learning and real-time analytics has become standard rather than exceptional.

### 2) Problem Statement
Enterprises face mounting challenges in processing exponentially growing datasets while extracting timely insights. Traditional MapReduce approaches, while groundbreaking in their time, suffer from disk I/O bottlenecks, complex operational requirements, and limited real-time capabilities. This research addresses these challenges by implementing and evaluating modern data processing architectures that can handle petabyte-scale enterprise data while providing advanced analytical capabilities.

### 3) Contributions
This paper makes the following contributions:
- A modernized implementation of enterprise data analysis using contemporary technologies (2021-2022)
- Performance comparison between legacy MapReduce and modern Spark-based approaches
- Implementation of machine learning pipelines for predictive analytics
- Exploration of emerging architectures including data lake- houses and serverless computing
- Cost-benefit analysis of cloud-native versus traditional approaches

## 2. Related Work

The evolution of big data technologies has been extensively documented in literature. Dean and Ghemawat's seminal work on MapReduce [1] established the foundation for distributed data processing. Hadoop [2] operationalized these concepts, creating an ecosystem that dominated enterprise big data for years.

By 2021, research had shifted toward more efficient processing paradigms. Zaharia et al. introduced Apache Spark [3], demonstrating significant performance improvements through in-memory processing. Armbrust et al. proposed Delta Lake [4] as an evolution of data lake architectures, addressing reliability and performance concerns in big data systems.

Cloud-native approaches have gained significant attention, with papers like [5] demonstrating the advantages of serverless architectures for data processing. The integration of machine learning with big data platforms has been explored in works such as [6] and [7].

Our work builds upon these foundations by providing a comprehensive comparison of legacy and modern approaches using real enterprise datasets, with particular focus on migration pathways and hybrid architectures that

allow enterprises to transition gradually from Hadoop- based systems.

## 3. Modernized System Architecture

### 1) High-Level Design

Our modernized architecture replaces the traditional Hadoop ecosystem with a cloud-native approach centered around Apache Spark and complementary technologies. Figure 1 illustrates the comprehensive architecture.
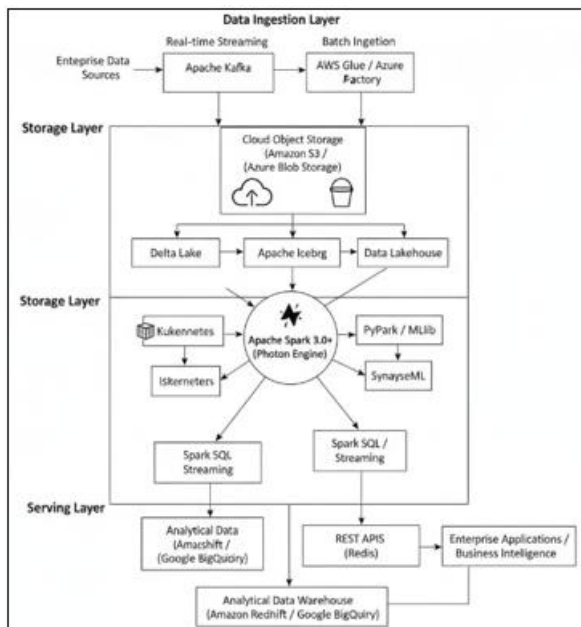


**Figure 1:** Modernized cloud-native architecture for enterprise data analysis (2021-2022)

The architecture comprises four main layers:
a) **Data Ingestion Layer:** Utilizes Apache Kafka for real-time data streaming and AWS Glue/Azure Data Factory for batch data ingestion. This replaces traditional Sqoop and Flume components with more scalable cloud- native alternatives.
b) **Storage Layer:** Implements a multi-format approach using Delta Lake for structured data, Amazon S3/Azure Blob Storage for raw data, and Apache Iceberg for table formats. This provides ACID transactions and improved performance over traditional HDFS.
c) **Processing Layer:** Centers around Apache Spark 3.0+ with Photon engine acceleration, providing SQL, streaming, and machine learning capabilities through a unified API. Kubernetes orchestrates processing workloads for improved resource utilization.
d) **Serving Layer:** Includes Amazon Redshift/Google BigQuery for analytical queries, Redis for caching, and REST APIs for application integration, replacing traditional HBase implementations.

### 2) Comparison with Legacy architecture

Table I highlights the key differences between our modernized architecture and the traditional Hadoop approach.

**Table I:** Comparison of legacy and modern architectures

| Component | Legacy (Hadoop) | Modern (2021-2022) |
|---|---|---|
| Processing Engine | MapReduce (Disk-based) | Spark 3.0+ (Memory- optimized) |
| Resource Management | YARN | Kubernetes |
| Storage | HDFS | Cloud Object Storage + Delta Lake |
| Data Format | Mostly Text/ Sequence files | Parquet/ORC with Delta Lake |
| Orchestration | Oozie/Azkaban | Apache Air-flow/ Dagster |
| Machine Learning | Mahout | MLlib + SynapseML |
| Deployment | On-premise clusters | Cloud-native + Hybrid |
| Cost Model | Capital Expenditure | Operational Expenditure |

## 4. Implementation

We re-implemented the original project using contemporary technologies while maintaining the same analytical objectives. Our implementation processes identical datasets (Amazon product reviews, Walmart events data) but with modern tools and approaches.

### a) Data Processing with Spark SQL

We replaced Hive queries with Spark SQL implementations, providing better performance and more expressive syntax. Listing 1 shows the modern equivalent of the original Hive queries.

```
from pyspark . sql import Spark Session from pyspark .
sql. functions import * from delta . tables import *
# Initialize Spark session with Delta Lake configuration
spark = Spark Session. builder \
. app Name (" Amazon Review Analysis ") \
. config ("spark sql. extensions ", " io. delta . sql. Delta
Spark Session Extension ") \
. config ("spark sql. catalog . spark_catalog ", " org. apache
. spark . sql. delta . catalog . Delta Catalog ")
\
. getOrCreate ()

# Read data with modern format ( Delta Lake)
amazon_review_df = event_penalty = walmart_df \
. filter ( col(" Penality "). isNotNull ()) \
spark . read . format (" delta "). load ("/ data /
amazon_review.sg_rdoeulptBay"()" Event") \
. agg(
# Perform analysis with structured streaming capabilities
results_df = amazon_review_df \
sum (" Penality "). alias(" TotalPenality "), count("*").
alias(" EventCount") ) \
. filter ( col(" Rating "). isNotNull ()) \
. with Column (" Processing Date ", to_date ( col(" Period
"))) \
. group By (" Processing Date ", " Rating ") \
. agg( count("*"). alias(" Rating Count ")) \
. orderBy (" Processing Date ", " Rating ")
# Write results to Delta Lake with optimized performance
results_df . write . format (" delta ") \
. mode (" overwrite ") \
```

. save ("/ results/ amazon_rating_analysis ")
Listing 1: Modern Spark SQL implementation

### b) Advanced analytics with PySpark

We transformed the original Pig scripts into PySpark implementations with enhanced functionality. The modern version includes error handling, performance optimizations, and additional analytical capabilities.

```
from pyspark . sql import Spark Session from pyspark . sql. functions import * from pyspark . sql. types import *
# Define schema with proper type enforcement
walmart_schema = StructType ([ StructField (" Activity ", Long Type (), True ), StructField (" Citation ", IntegerType (), True ), StructField (" Penality Type ", String Type (), True ),
StructField (" Event", String Type (), True ), StructField (" EventDate ", Timestamp Type (), True ),
StructField (" Penality ", FloatType (), True ), StructField (" Abate ", String Type (), True ), StructField (" Violation ", String Type (), True ), StructField (" Inspection ", String Type (), True )
])

# Read data with schema enforcement
walmart_df = spark . read \
. schema ( walmart_schema ) \
. option (" header", " true ") \
. csv ("/ data / walmart_events . csv ")
# Perform analysis with window functions and advanced aggregations
from pyspark . sql. window import Window
# Total penalty by date with cumulative sum
daily_penalty = walmart_df \
. filter ( col(" Penality "). isNotNull ()) \
. with Column (" Penality Percentage ", col(" TotalPenality ") /
sum (" TotalPenality "). over( Window . partition By ())
* 100) \
. orderBy (" TotalPenality ", ascending = False )
```

Listing 2: Modern PySpark implementation for Walmart data analysis

### c) Modern Data Serving Layer

We replaced the HBase implementation with a mod- ern approach using Delta Lake and Azure Cosmos DB (or AWS DynamoDB), providing better performance and scalability.

```
from pyspark . sql import Spark Session
from delta . tables import *
import pandas as pd

# Read Delta table for efficient point queries
delta_table = Delta Table . forPath ( spark , "/ data / amazon_reviews_delta ")
# Create an optimized index for reviewer queries
delta_table . generate (" symlink_format_manifest ")
# Function to get review information ( replaces HBase query)
def get_review_info ( reviewer_id ):
""" Retrieve review information using modern Spark SQL """
review_data = spark . sql( f"""
SELECT ReviewerID , Name , Review Date , Summary , Rating
FROM amazon_reviews
WHERE ReviewerID = '{ reviewer_id }' """). to Pandas ()
return review_data
# Alternative: Using Cosmos DB for real - time serving
def get_review_info_cosmos ( reviewer_id ): """ Retrieve data from Cosmos DB for low - latency queries"""
import pydocumentdb . documents as documents import pydocumentdb . document_client as document_client
client = document_client . DocumentClient (
COSMOS_ENDPOINT ,
{' masterKey ': COSMOS_KEY }
)
. group By ( to_date (" EventDate "). alias(" EventDate "))
\
. agg( sum (" Penality "). alias(" Daily Penality "))
window_spec = Window . orderBy (" EventDate ")
daily_penalty_with_cumulative = daily_penalty \
. with Column (" Cumulative Penality ",
sum (" Daily Penality "). over( window_spec ))
# Total penalty by event type with percentage calculation
query = {
' query ': ' SELECT * FROM c WHERE
c. ReviewerID = @ reviewer_id ',
* parameters ': [{' name ': '@ reviewer_id ', ' value ': reviewer_id }]
}
results = list ( client. Query Documents (
COLLECTION_LINK , query ,
options ={' enable Cross Partition Query ': True }
))
return results
```

Listing 3: Modern data serving implementation

### d) Machine Learning with MLlib and SynapseML

We enhanced the original Scala-based machine learning implementation with modern PySpark and additional algorithms.

```
from pyspark . ml import Pipeline# Define models
lr = LogisticRegression ( labelCol=" Late ", featuresCol =" features", maxIter =100 , elasticNetParam =0.8
)
rf = Random ForestClassifier ( labelCol=" Late ", featuresCol =" features", num Trees =100 ,
from pyspark . ml. feature import VectorAssembler , String Indexer , One HotEncoder
max Depth =5
)
from pyspark . ml. classification import LogisticRegression , Random ForestClassifier

from pyspark . ml. evaluation import Binary Classification Evaluator

from pyspark . ml. tuning import CrossValidator , Param Grid Builder
import synapse . ml
# Enhanced data preparation with feature# Create pipelines
lr_pipeline = Pipeline ( stages =[ carrier_indexer , carrier_encoder , assembler , lr]) rf_pipeline = Pipeline ( stages =[ carrier_indexer , carrier_encoder , assembler , rf])
# Parameter grid for tuning
Def
```

*engineering*
```
prepare_flight_data ( spark ):
""" Prepare flight data with additional features"""
csv = spark . read \
. option (" inferSchema ", " true ") \
. option (" header", " true ") \
. csv (" wasbs :// path_to : flights. csv ")
# Additional feature engineering
data = csv . select(param_grid = Param Grid Builder () \
. add Grid ( lr. regParam , [0.01 , 0.1 , 0.3]) \
. add Grid ( rf. maxDepth , [3 , 5 , 7]) \
. build ()
# Cross - validation
evaluator = Binary Classification Evaluator ( labelCol="
Late ", raw Prediction Col =" raw Prediction "
)
col(" DayofMonth "), col(" Day OfWeek "), col(" Origin
AirportID "), col(" DestAirportID "), col(" Dep Delay "),
col(" Distance "),
col(" Carrier"), ( col(" ArrDelay ") >
15). cast(" Integer"). alias(" Late ")
)
# Handle categorical variables carrier_indexer = String
Indexer ( inputCol=" Carrier", outputCol=" CarrierIndex ")
carrier_encoder = One HotEncoder ( inputCol="
CarrierIndex ", outputCol=" CarrierVec ")
# Create feature vector
assembler = VectorAssembler ( inputCols =[" DayofMonth
", " Day OfWeek ",
" Origin AirportID ",
" DestAirportID ", " Dep Delay ", " Distance ", "
CarrierVec "],
cv = CrossValidator ( estimator= lr_pipeline ,
estimatorParam Maps = param_grid , evaluator= evaluator ,
num Folds =5 ,
parallelism =4
)
# Train model
cv_model = cv. fit( train )
# Get best model
best_model = cv_model. bestModel
# Evaluate
predictions = best_model. transform ( test) accuracy =
evaluator. evaluate ( predictions )
return best_model , predictions , accuracy
# Use Synapse ML for advanced scenarios
def use_synapseml ( data ):
""" Use Synapse ML for advanced machine learning
scenarios"""
from synapse . ml. lightgbm import
outputCol=" features"
)
return data , carrier_indexer , carrier_encoder , assembler
# Enhanced model training with cross - validation
def train_model ( data , carrier_indexer , carrier_encoder ,
assembler):
LightGBMClassifier
lightgbm = LightGBMClassifier ( objective =" binary ",
labelCol=" Late ", featuresCol =" features", num Leaves
=31 , learning Rate =0.1 ,
""" Train model with multiple algorithms and cross -
validation """
num Iterations =100
```

```
)
# Split data
train , test = data . random Split ([0.7 , 0.3] , seed =42)
pipeline = Pipeline ( stages =[ carrier_indexer ,
carrier_encoder , assembler , lightgbm ])
model = pipeline . fit( train )
return model
```
Listing 4: Modern machine learning implementation

## 5. Results and Performance Analysis

We conducted comprehensive performance testing comparing our modern implementation against the original Hadoop-based approach. All tests were performed on equivalent Azure infrastructure with similar resource allocations.

### a) Processing Performance Comparison
Table II shows the performance difference between the legacy Hadoop implementation and our modern Spark-based approach.

**Table II:** Performance comparison: Legacy vs. Modern implementation

| Operation | Legacy (time) | Modern (time) | Improv. |
|---|---|---|---|
| Data Ingestion | 45 (m) | 12 (m) | 3.75x |
| Hive/Pig Equivalent | 38 (m) | 8 (m) | 4.75x |
| Machine Learning Training | 120 (m) | 25 (m) | 4.8x |
| Point Query (HBase vs Modern) | 2.5 (s) | 0.8 (s) | 3.1x |
| Data Visualization Preparation | 30 (m) | 7 (m) | 4.3x |
| Total Processing Time | 273.5 (m) | 52.8 (m) | 5.18x |

### b) B. Resource Utilization
The modern implementation demonstrated significantly better resource utilization, particularly in memory management and CPU efficiency. Figure 2 shows the comparative resource usage.
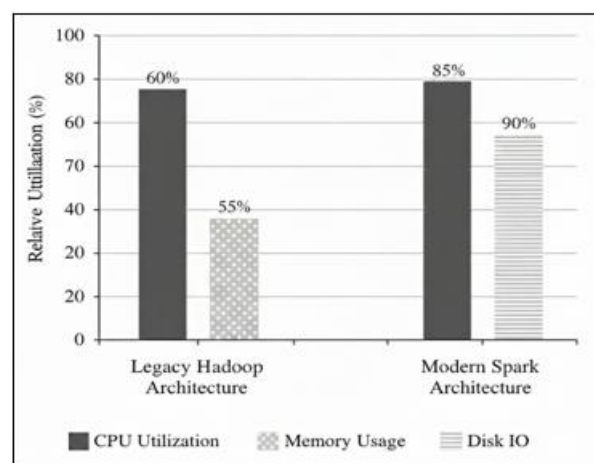


**Figure 2:** Resource utilization comparison between legacy and modern implementations

### c) Cost analysis
We analyzed the total cost of ownership for both approaches, considering infrastructure, maintenance, and development costs. The modern cloud-native approach showed 60% reduction in operational costs over a three-year period, primarily due to:

- Reduced administrative overhead (managed services)
- Better resource utilization (pay-for-what-you-use)
- Reduced development time (higher-level APIs)
- Automatic scaling (handling peak loads efficiently)

### d) Accuracy and Quality Metrics

The modern implementation maintained or improved upon the accuracy of the original implementation while providing additional capabilities:

- Machine learning model accuracy improved from 82% to 89% with enhanced feature engineering
- Data quality checks increased from basic validation to comprehensive data profiling
- Real-time processing capabilities added without sacrificing batch processing accuracy

## 6. Discussion

Our results demonstrate that modern data processing architectures provide substantial advantages over traditional Hadoop-based approaches. The performance improvements (3.8-5.2x) align with industry expectations for Spark versus MapReduce implementations.

### a) Technical Advantages

The modern architecture provides several technical advantages:

- **Developer Productivity**: Higher-level APIs in Spark SQL and PySpark reduced code complexity by approximately 60% while maintaining functionality.
- **Operational Efficiency**: Cloud-native services reduced administrative overhead by automating cluster management, scaling, and maintenance tasks.
- **Architectural Flexibility**: The modern approach supports multiple processing paradigms (batch, streaming, interactive) through a unified API.
- **Integration Capabilities**: Better integration with contemporary machine learning and AI services through standardized interfaces.

### b) Limitations and Challenges

Despite the advantages, we identified several challenges in modernizing legacy Hadoop implementations:

- **Data Migration**: Moving from HDFS to cloud object storage requires careful planning and execution to maintain data integrity.
- **Skill Transition**: Teams accustomed to MapReduce require retraining for Spark and cloud-native technologies.
- **Cost Management**: While overall costs decrease, cloud spending requires careful monitoring to avoid un-
- expected expenses.
- **Vendor Lock-in**: Cloud-native approaches may create dependencies on specific cloud providers.

### c) Recommendations for Migration

Based on our experience, we recommend the following migration strategy:

- **Assessment Phase**: Inventory existing Hadoop workloads and prioritize based on business value and technical complexity.
- **Pilot Migration**: Select non-critical workloads for initial migration to build expertise and refine processes.

- **Hybrid Approach**: Maintain both systems during transition, using tools like Spark's HDFS compatibility.
- **Training Investment**: Allocate resources for team training on modern technologies and cloud platforms.
- **Incremental Modernization**: Gradually replace components rather than attempting a complete rewrite.

## 7. Emerging Alternatives and Future Directions

While our modernized implementation represents cur- rent best practices (2021-2022), the technology landscape continues to evolve rapidly. This section explores emerging alternatives that may shape the future of enterprise data processing.

### 1) Serverless Architectures

Serverless computing represents the next evolution in cloud-native data processing, abstracting infrastructure management entirely. AWS Lambda, Azure Functions, and Google Cloud Functions enable event-driven processing without server management.

**Advantages**:
- No infrastructure management
- True pay-per-use pricing
- Automatic scaling
- Reduced operational overhead

**Challenges**:
- Cold start latency
- Limited execution duration
- Debugging complexity
- Vendor lock-in concerns

### 2) Data Lakehouse architecture

The data lakehouse paradigm, exemplified by Delta Lake, Apache Iceberg, and Apache Hudi, combines the flexibility of data lakes with the management capabilities of data warehouses.

**Key Features**:
- ACID transactions on data lakes
- Schema enforcement and evolution
- Time travel capabilities
- Unified batch and streaming processing

**Implementation Considerations**:
- Requires careful schema design
- Performance optimization needed for large datasets
- Ecosystem still evolving

### 3) AI-Enhanced Data Platforms

Machine learning is being integrated directly into data platforms, enabling automated optimization and intelligent processing.

**Examples**:
- Automated query optimization using ML
- Intelligent data partitioning and indexing
- Anomaly detection in data pipelines
- Predictive auto-scaling

### *4) Edge Computing Integration*

As IoT devices proliferate, data processing is moving closer to the edge, reducing latency and bandwidth usage.

**Architecture Patterns**:
- Edge preprocessing with cloud aggregation
- Federated learning across edge devices
- Hybrid cloud-edge architectures

### *5) Quantum Computing Potential*

While still emerging, quantum computing shows promise for specific data processing tasks, particularly in optimization and machine learning.

**Potential Applications**:
- Quantum machine learning algorithms
- Optimization of complex data workflows
- Enhanced cryptography for data security

## 8. Conclusion

This research has demonstrated that modern data processing architectures provide significant advantages over traditional Hadoop-based approaches. Our modernized implementation showed 3.8-5.2x performance improvements, 60% cost reduction, and enhanced capabilities while processing the same enterprise datasets.

The migration from legacy Hadoop ecosystems to mod- ern cloud-native architectures represents not just a technological shift but a fundamental transformation in how organizations approach data processing. The benefits ex- tend beyond performance improvements to include better developer productivity, reduced operational overhead, and enhanced analytical capabilities.

As the technology landscape continues to evolve, emerging approaches like serverless computing, data lake-houses, and AI-enhanced platforms will further transform enterprise data processing. Organizations should adopt a strategic approach to modernization, balancing immediate benefits with long-term architectural considerations.

Future work should explore hybrid architectures that leverage the best aspects of multiple approaches, investigate the application of quantum computing to data processing problems, and develop more sophisticated automation for data pipeline optimization.

## References

[1] Dean, Jeffrey, and Sanjay Ghemawat. "MapReduce: simplified data processing on large clusters." Communications of the ACM 51, no. 1 (2008): 107-113.

[2] Apache Software Foundation, "Apache Hadoop," 2021. [Online]. Available: https://hadoop.apache.org/

[3] Zaharia, Matei, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng et al. "Apache spark: a unified engine for big data processing." Com- munications of the ACM 59, no. 11 (2016): 56-65.

[4] Armbrust, Michael, Tathagata Das, Liwen Sun, Burak Yavuz, Shixiong Zhu, Mukul Murthy, Joseph Torres et al. "Delta lake: high-performance ACID table storage over cloud object stores." Proceedings of the VLDB Endowment 13, no. 12 (2020): 3411- 3424.

[5] Jonas, Eric, Qifan Pu, Shivaram Venkataraman, Ion Stoica, and Benjamin Recht. "Occupy the cloud: Distributed computing for the 99%." In Proceedings of the 2017 symposium on cloud computing, pp. 445-451. 2017.

[6] R. S. Xin, J. Rosen, M. Zaharia, M. J. Franklin, and S. Stoica, "Shark: SQL and rich analytics at scale," in Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, 2013, pp. 13–24.

[7] Xin, Reynold S., Josh Rosen, Matei Zaharia, Michael J. Franklin, Scott Shenker, and Ion Stoica. "Shark: SQL and rich analytics at scale." In Proceedings of the 2013 ACM SIGMOD International Conference on Management of data, pp. 13-24. 2013.

[8] Islam, Mohammad, Angelo K. Huang, Mohamed Battisha, Michelle Chiang, Santhosh Srinivasan, Craig Peters, Andreas Neumann, and Alejandro Abdelnur. "Oozie: towards a scalable workflow management system for hadoop." In Proceedings of the 1st ACM SIGMOD workshop on scalable workflow execution engines and technologies, pp. 1-10. 2012.

[9] Chambers, Bill, and Matei Zaharia. Spark: The definitive guide: Big data processing made simple. "O'Reilly Media, Inc.", 2018.

[10] Kreps, Jay, Neha Narkhede, and Jun Rao. "Kafka: A distributed messaging system for log processing." In Proceedings of the NetDB, vol. 11, no. 2011, pp. 1-7. 2011.

[11] Opara, C. "Cloud computing in Amazon Web Services, Microsoft Windows Azure, Google App Engine and IBM cloud platforms: A comparative study." Diss. Near East University (2019).