

## Advanced ML (Machine Learning) Techniques for Optimizing ETL Workflows with Apache Spark and Snowflake

Arjun Mantri

Independent Researcher, Seattle, WA, USA

### ABSTRACT

The optimization of ETL (Extract, Transform, Load) pipelines using Apache Spark and Snowflake. Apache Spark is a powerful open-source distributed data processing platform, while Snowflake is a cloud-native data warehousing solution. It discusses the challenges and solutions in tuning Spark configurations using machine learning techniques and optimizing Snowflake's architecture for cost efficiency and performance. Experimental results demonstrate significant performance gains and cost savings through these optimizations.

### \*Corresponding author

Arjun Mantri, Independent Researcher, Seattle, WA, USA.

Received: July 05, 2023; Accepted: July 20, 2023; Published: July 20, 2023

**Keywords:** ETL, Apache Spark, Snowflake, Machine Learning, Data Warehousing, Performance Optimization

### Introduction

The rapid growth of big data has necessitated efficient data processing frameworks. Apache Spark has emerged as a leading platform due to its in-memory processing capabilities and support for complex analytics. Snowflake, on the other hand, offers a scalable and flexible cloud data warehousing solution.

Traditional ETL processes involve transforming data before loading it into the data warehouse to limit data size and ensure optimal querying performance. However, Snowflake's low storage costs and powerful SQL (structured query language) capabilities, combined with significant performance improvements provided by query pushdown, enable a transition to a more modern and effective ELT (Extract, Load, Transform) model. In this model, all data is loaded into Snowflake first, and any necessary transformations are performed directly within Snowflake. This approach allows for seamless integration with existing Spark programming environments using Scala, Python, or Spark SQL [1,2]. This paper aims to address the optimization challenges in ETL pipelines using these technologies.

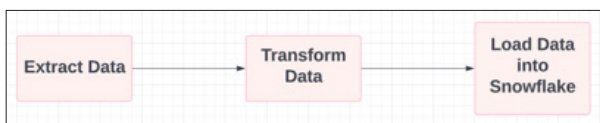


Figure 1: ETL Process

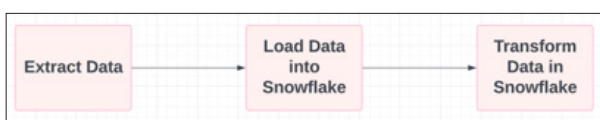


Figure 2: ELT Process

### Background

#### Apache Spark

Apache Spark is designed for large-scale data processing, leveraging distributed memory abstraction to process data efficiently. It supports various data sources and programming languages, making it versatile for different data processing needs [3].

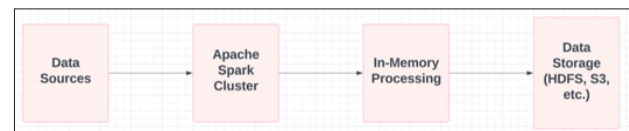


Figure 3: Apache Spark Architecture

#### Snowflake

Snowflake's architecture separates storage from computing, allowing for independent scaling of resources. It supports structured and semi-structured data, providing robust data sharing and governance capabilities [2].



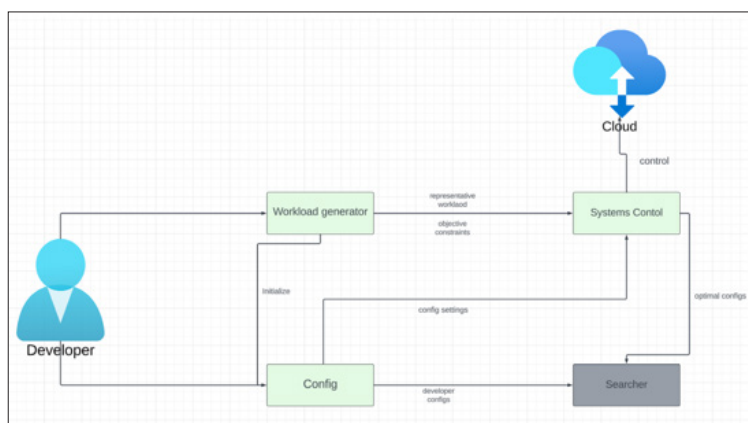
Figure 4: Snowflake Architecture

**Characteristics of the Included Studies**  
**Table 1: Characteristics of the Included Studies**

Study Reference	Title	Focus	Key Finding
Snowflake Blog	Snowflake Spark Pushing Query Processing	Query processing in Snowflake and Spark	Discusses advancements in query processing for Snowflake and Spark
George, A. Shaji	Deciphering the Path to Cost Efficiency and Sustainability in the Snowflake Environment	Cost efficiency and sustainability in Snowflake	Discusses strategies for achieving cost efficiency and sustainability in Snowflake environments
Wang, Guolu, Jungang Xu, and Ben He	A novel method for tuning configuration parameters of spark based on machine learning	Tuning Spark configuration parameters using machine learning	Introduces a machine learning-based method for optimizing Spark configuration parameters to improve performance
Sornlertlamvanich, Virach, Petchporn Chawakitchareon, and Aran Hansuebsai, eds	Information Modelling and Knowledge Bases XXIX	Information modeling and knowledge bases	Covers various topics in information modeling and knowledge bases
Snowflake Blog	Benchmarking Snowflake Versus Spark with LTI Mosaic Decisions	Benchmarking Snowflake vs. Spark	Snowflake outperforms Spark in terms of performance, stability, and ease-of-use
Wang, Kewen, and Mohammad Maifi Hasan Khan	Performance prediction for Apache Spark platform	Performance prediction for Spark	Discusses methods for predicting the performance of the Apache Spark platform
Armbrust, Michael, et al.	Spark SQL: Relational data processing in spark	Relational data processing in Spark	Describes Spark SQL and its capabilities for relational data processing
Cheng, Guoli, Shi Ying, and Bingming Wang	Tuning configuration of Apache Spark on public clouds by combining multi-objective optimization and performance prediction model	Tuning Spark configuration on public clouds	Combines multi-objective optimization and performance prediction models to tune Spark configurations on public clouds
Dünner, Celestine, et al.	Understanding and optimizing the performance of distributed machine learning applications on Apache Spark	Optimizing distributed machine learning on Spark	Focuses on understanding and optimizing the performance of distributed machine learning applications on Spark
Astsatryan, Hrachya, et al.	Performance optimization system for Hadoop and spark frameworks	Performance optimization for Hadoop and Spark	Proposes a system for optimizing the performance of Hadoop and Spark frameworks

**Challenges in ETL Optimization**  
**Spark Configuration Tuning**

Tuning Spark configurations is complex due to the large parameter space and interactions among parameters. Traditional manual tuning methods are inefficient and often fail to achieve optimal performance [3].



**Figure 5: Spark Configuration Tuning Model**

**Cost Management in Snowflake**

Snowflake's ease of scaling can lead to uncontrolled resource consumption and high costs. Proper governance and optimization techniques are essential to manage these costs effectively [4].

## Proposed Solutions

### Machine Learning-Based Spark Tuning

A novel method using machine learning for Spark configuration tuning is proposed. This method involves binary and multi-class classification models to predict the execution time of Spark applications under different configurations [3].

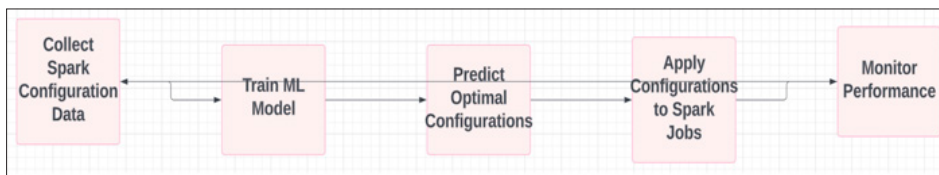


Figure 6: Machine Learning Model for Spark Tuning

### Snowflake Optimization Techniques

Several techniques are recommended for optimizing Snowflake, including right-sizing virtual warehouses, query optimization, and efficient data storage through compression and clustering [2].

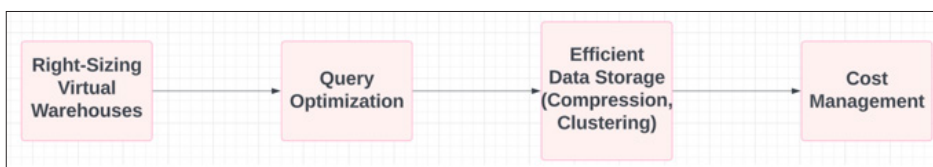


Figure 7: Snowflake Optimization Technique

## Experimental Results

### Spark Performance Gains

Experiments show that the proposed machine learning-based tuning method achieves an average performance gain of 36% compared to default configurations [3].

### Cost Savings in Snowflake

Implementing the recommended optimization techniques in Snowflake results in significant cost savings and improved performance. For instance, query optimization and result caching can reduce compute costs substantially [2].

### Case Study: Snowflake Performance and Agility

#### Performance

Snowflake offers a data processing capacity that is typically 200% of the Apache Spark analytics engine. In terms of performance and total cost of ownership (TCO), Snowflake runs faster and outperforms Spark by a significant margin across the ETL cycle. If its other features align with business needs, Snowflake becomes a natural and preferred choice to be integrated and used with Mosaic Decisions [5].

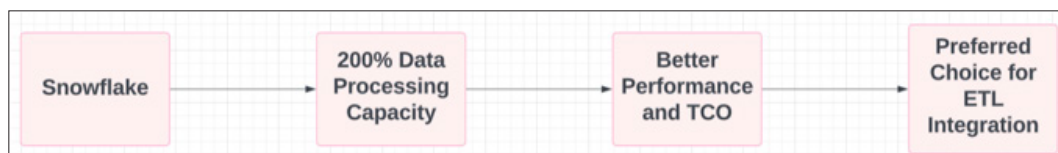


Figure 8: Spark Performance Comparison

#### Agility

Given that Snowflake is a true SaaS (Software as a service) solution, it is simple to get started with; it requires no hardware or software to install, configure, and manage; and it even takes care of the maintenance operations for its components. On the other hand, Spark is a technology built for analytics experts and could prove a challenge for less tech-savvy users. Additionally, data pipelines executed on a Spark cluster took about five minutes to start running, thus delaying the overall processing, whereas on Snowflake all data executions began instantaneously [5].

#### Concurrency

When there are too many concurrent users, it becomes necessary for the system to scale up to cater to the users' needs. Here both Spark interactive clusters as well as Snowflake virtual warehouses offer an auto-scaling capability. However, Snowflake performed 3x better even when using only 25% of the resources while the Spark cluster struggled to manage 100+ concurrent users. The combination of LTI (Larsen & Toubro Infotech) Mosaic Decisions and Snowflake is a win-win solution for enterprises, as it harnesses and complements the capabilities of each product [5].

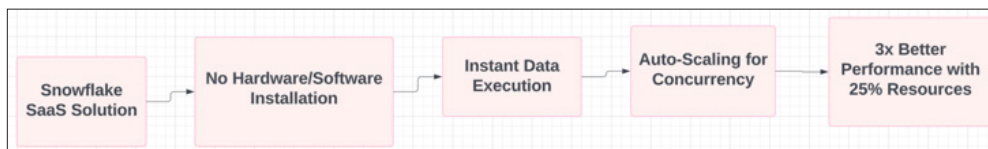


Figure 9: Agility and Concurrency

### Query Pushdown

Some of the more resource-intensive processing can be performed in Snowflake instead of Spark, resulting in significant performance improvements. Snowflake's ability to push query processing down from Spark into Snowflake helps transition from a traditional ETL process to a more flexible and powerful ELT model. Snowflake uses a virtual warehouse to process the query and copies the query result into AWS (Amazon Web Services) S3. The connector retrieves the data from S3 and populates it into Data Frames in Spark. This allows for optimal performance by balancing compute capacity and IO bandwidth against S3 [1].

### Case Study: Apache Spark Performance Optimization

#### Experimental Setup

Experiments were performed on a cluster of 13 nodes, where one node served as the master node, and the other 12 nodes served as worker nodes. The master node had 4 Central Processing Unit (CPU) cores and 6 GB of memory. Among the 12 worker nodes, 6 nodes had 8 CPU cores and 8 GB of memory, and the other 6 nodes had 12 CPU cores and 16 GB of memory. Apache Spark was run using its standalone cluster mode on top of Hadoop Distributed File System (HDFS) with a default 64 MB block setting [6].

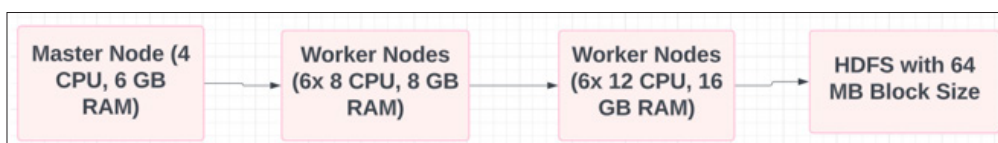


Figure 10: Experimental Setup

#### Data Collection and Analysis

For data collection, Spark event logs generated by the Apache Spark platform were used to record execution profiles and performance metrics. These logs were saved in JSON format. By analyzing these logs, metrics such as StageTime and TaskTime were calculated. Job startup time, job cleanup time, stage startup time, and stage cleanup time were also recorded. I/O cost details for each task were calculated using shuffle write and read metrics to construct the I/O profile [6].

#### Algorithm Performance

##### Logistic Regression

Logistic Regression is an iterative algorithm with 10 stages, where there is no shuffle Input/Output (I/O) cost. For Logistic Regression, the input data size was 50 GB. As logistic regression is a computing-intensive job, this minimizes the effect of I/O and leads to better prediction accuracy. For memory cost prediction, all four simulation setups achieved 100% accuracy and correctly predicted the value of 11.5 GB needed to create the initial RDD (Resilient Distributed Dataset) [6].

##### K-Means

The same data set used for testing the Logistic Regression algorithm was used for the K-Means clustering algorithm. K-Means is an iterative algorithm with 22 stages. In the third stage of the job, the I/O cost involves shuffle write, and the I/O cost involves shuffle read in the fourth stage. Later stages follow the same pattern. However, as the time cost for I/O operations is small and has minimal effect on total execution time, the model still achieved high prediction accuracy for execution time. For memory cost prediction, the model had 100% accuracy, correctly predicting the value of 11 GB [6].

##### PageRank

PageRank is an iterative algorithm with 13 stages where there are I/O read and write costs for each stage. In the evaluation, 25 GB of input data was used. For time prediction, the accuracy was above 80% for simulating one task per core but dropped to around 70% for simulating two tasks per core. This result may be due to the fluctuation in execution time across tasks within a stage. For I/O write cost prediction, the accuracy was above 90% for large-scale simulation [7,8].

#### Evaluation of Spark SQL

The performance of Spark SQL was evaluated on two dimensions: SQL query processing performance and Spark program performance. Spark SQL's extensible architecture not only enables a richer set of functionalities but also brings substantial performance improvements over previous Spark-based SQL engines. For Spark application developers, the DataFrame API (application programming interface) can bring substantial speedups over the native Spark API while making Spark programs more concise and easier to understand. Applications that combine relational and procedural queries run faster on the integrated Spark SQL engine than by running SQL and procedural code as separate parallel jobs [7].

#### Spark SQL Code Example

To illustrate the use of Spark SQL, consider the following example where we create a DataFrame from a CSV (comma separated values) file, perform some transformations, and execute SQL queries: This example demonstrates how to load data into a DataFrame, create a temporary view, and execute SQL queries using Spark SQL. The `createOrReplaceTempView` method registers the DataFrame as a temporary view, allowing SQL queries to be run against it [8].

```
import org.apache.spark.sql.SparkSession
// Create a SparkSession
val spark = SparkSession.builder()
  .appName("Spark SQL Example")
  .config("spark.some.config.option", "some-value")
  .getOrCreate()
// Load data from a CSV file into a DataFrame
val df = spark.read.format("csv")
  .option("header", "true")
  .option("inferSchema", "true")
  .load("path/to/csvfile.csv")
// Create a temporary view
df.createOrReplaceTempView("data")
// Execute SQL queries
val result = spark.sql("SELECT * FROM data WHERE age > 30")
result.show()
```

Figure 11: Spark SQL Code Example

### Compression in Spark

In Big Data infrastructures, reducing the size of data chunks is crucial to minimize I/O operation delays and save space on local disks. However, finding the optimal trade-off is challenging, as a high compression factor may underload I/O but overload the CPU, while a weak compression factor may underload the CPU but overload I/O. The ideal configuration ensures that both I/O

and CPU are fully utilized without waiting for each other. Memory and I/O consume a significant portion of Big Data processing resources. Many scientific studies have focused on memory optimizations at various levels, including hardware, kernel memory, and middleware layers. This paper aims to optimize resources at the application level using several compression algorithms within the Apache Hadoop and Spark frameworks. The goal is to reduce the size of files to be processed, loaded into memory, and written back to the disk. The degree of data size or I/O reduction depends on the compression ratio, which is the size of compressed data divided by the size of uncompressed data. The compression ratio relies on the data and the compression algorithm. A lower ratio means less memory and I/O usage. Data compression in Hadoop and Spark frameworks increases storage space and improves job performance. If the input or intermediate output of the map phase is compressed, the framework chooses a decompression algorithm before processing according to the file extension. Below table shows different file formats [7,9,10].



Figure 12: Compression Technique

Table 2: Types of File Formats

File Format	Type	Compression Codecs/or CODES	Schema Support	Splittable	Use Case
Text	Unstructured	bzip2, deflate, gzip, Snappy, zstd	No	Yes	General use, human-readable
Avro	Structured	Snappy, gzip, deflate	Yes	Yes	Row-based storage, schema evolution
Parquet	Structured	Snappy, gzip, zstd, LZ4	Yes	Yes	Columnar storage, efficient for analytics
ORC (Optimized Row Columnar)	Structured	gzip, Snappy, LZ4	Yes	Yes	Columnar storage, optimized for BI
RC File (Resource Script)	Structured	Snappy, gzip, deflate, bzip2	Yes	Yes	Columnar storage, high compression
SequenceFile	Structured	Snappy, gzip, deflate, bzip2	Yes	No	Key-value pairs, block-level compression
JSON (JavaScript Object Notation)	Unstructured	gzip, Snappy	Yes	Yes	Web communication, flexible schema
Protocol Buffers	Structured	None	Yes	No	Streaming, RPC (Remote procedure call) (e.g., gRPC)

### Discussion

The integration of Spark and Snowflake provides a powerful solution for ETL pipelines. While Spark excels in data processing, Snowflake offers scalable and efficient data storage. The referenced optimizations enhances the performance and cost-efficiency of these platforms.

### Conclusion

Optimizing ETL pipelines using Apache Spark and Snowflake involves addressing configuration and cost management challenges. Machine learning-based tuning for Spark and strategic optimization techniques for Snowflake can lead to significant performance improvements and cost savings. Future work will focus on further refining these methods and exploring additional optimization opportunities.

### References

1. Torsten Grabs, Edward Ma (2017) Snowflake and Spark: Pushing Spark Query Processing to Snowflake. Snowflake <https://www.snowflake.com/blog/snowflake-spark-pushing-query-processing/>.
2. George A Shaji (2023) Deciphering the Path to Cost Efficiency and Sustainability in the Snowflake Environment. Partners Universal International Innovation Journal 1: 231-250.
3. Wang Guolu, Jungang Xu, Ben He (2016) A novel method for tuning configuration parameters of spark based on machine learning. 2016 IEEE 18th International Conference on High Performance Computing and Communications; IEEE 14th International Conference on Smart City; IEEE 2nd International Conference on Data Science and Systems (HPCC/SmartCity/DSS). IEEE <https://ieeexplore.ieee.org/>

- document/7828429.
4. Jaakkola H, Thalheim B, Kiyoki Y, Yoshida N, Sornlertlamvanich V (2018) Information Modelling and Knowledge Bases XXIX. IOS Press 301.
  5. Ritesh Vikram (2021) Benchmarking Snowflake Versus Spark with LTI Mosaic Decisions. Snowflake <https://www.snowflake.com/blog/benchmarking-snowflake-versus-spark-with-lti-mosaic-decisions/>
  6. Wang Kewen, Mohammad Maifi Hasan Khan (2015) Performance prediction for apache spark platform. 2015 IEEE 17th International Conference on High Performance Computing and Communications, 2015 IEEE 7th International Symposium on Cyberspace Safety and Security, and 2015 IEEE 12th International Conference on Embedded Software and Systems. IEEE <https://wangkewen.github.io/HPCC2015.pdf>.
  7. Michael Armbrust, Reynold S. Xin, Cheng Lian, Yin Huai, Davies Liu, et al. (2015) Spark sql: Relational data processing in spark. Proceedings of the 2015 ACM SIGMOD international conference on management of data [https://people.csail.mit.edu/matei/papers/2015/sigmod\\_spark\\_sql.pdf](https://people.csail.mit.edu/matei/papers/2015/sigmod_spark_sql.pdf).
  8. Cheng Guoli, Shi Ying, Bingming Wang (2021) Tuning configuration of apache spark on public clouds by combining multi-objective optimization and performance prediction model. Journal of Systems and Software 180: 111028.
  9. Celestine Dünner, Thomas Parnell, Kubilay Atasu, Manolis Sifalakis, Haralampos Pozidis (2017) Understanding and optimizing the performance of distributed machine learning applications on apache spark. 2017 IEEE international conference on big data (big data). IEEE <https://ieeexplore.ieee.org/document/8257942>.
  10. Hrachya Astsatryan, Aram Kocharyan, Daniel Hagimont, Arthur Lalayan (2020) Performance optimization system for hadoop and spark frameworks. Cybernetics and Information Technologies 20: 5-17.

**Copyright:** ©2023 Arjun Mantri. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.