# Comparative Analysis of Apache Sqoop and Apache Spark for Efficient Data Transfer between Relational Databases and Hadoop Distributed File System (HDFS)

**Sainath Muvva**

**Abstract**

**With the growing adoption of big data technologies like Hadoop, many companies are overhauling their data infrastructure. A crucial aspect of this transition is the ability to transfer both transactional and analytical data from traditional relational database management systems (RDBMS) into the new ecosystem. This migration enables advanced data processing and facilitates deeper analytical insights. This paper focuses on exploring the various tools available for importing data from relational databases into the Hadoop Distributed File System (HDFS). It delves into the underlying mechanisms of these tools and highlights the key distinctions between them.**

**Keywords: HDFS, Sqoop, Spark, SQL Loaders**

## Introduction

Data Engineers routinely face the challenge of importing data from diverse sources into big data ecosystems. This process, known as data ingestion, often involves using both off-the-shelf third-party tools and custom-built solutions. Relational Database Management Systems (RDBMS) are particularly common data sources, requiring daily ingestion for processing through established data pipelines.

The method of data ingestion can vary depending on the specific RDBMS. For instance, extracting data from Teradata might involve using 'bteq' [1] to export data to a file, placing that file on a shared drive accessible to the Hadoop cluster, then running a Hadoop job to load the data into a Hive text table, and finally converting it to ORC or Parquet file formats. This multi-step process can be complex and time-consuming.

This paper narrows its focus to two specific tools: Apache Sqoop and Apache Spark. Both utilize JDBC (Java Database Connectivity) to establish connections with relational databases and facilitate data transfer to the Hadoop Distributed File System (HDFS). These tools aim to streamline the data ingestion process, offering more efficient alternatives to manual, multi-step procedures.

## Apache Sqoop

Sqoop is a specialized tool in the Hadoop ecosystem for bidirectional data transfer between relational databases (RDBMS) and Hadoop Distributed File System (HDFS). It offers data transformation capabilities and utilizes MapReduce for efficient, parallel, and fault-tolerant data operations.

By streamlining the data transfer process, Sqoop enables developers to focus on data transformation and post-import tasks. It's compatible with various databases like MySQL, Oracle, Teradata, and SQL Server, and supports multiple file formats including Text, ORC, Parquet, and Avro.

When initiating an import, Sqoop first extracts metadata from the source RDBMS. For tables with integer primary keys, it determines the minimum and maximum values to establish import boundaries. The data is then split into chunks based on the specified number of mappers (default is four). If the primary key isn't an integer, users can specify an alternative column for data splitting. Without an integer primary key or alternative column, Sqoop uses a single mapper.

Each mapper operates as a separate database session, so it's crucial to balance the number of mappers to avoid overwhelming the RDBMS. These mappers work concurrently, writing data to HDFS as they process their assigned chunks.

It's worth noting that Sqoop functions solely as a mapper and doesn't support complex operations like joins. The number of mappers should be carefully optimized to balance performance and database load.
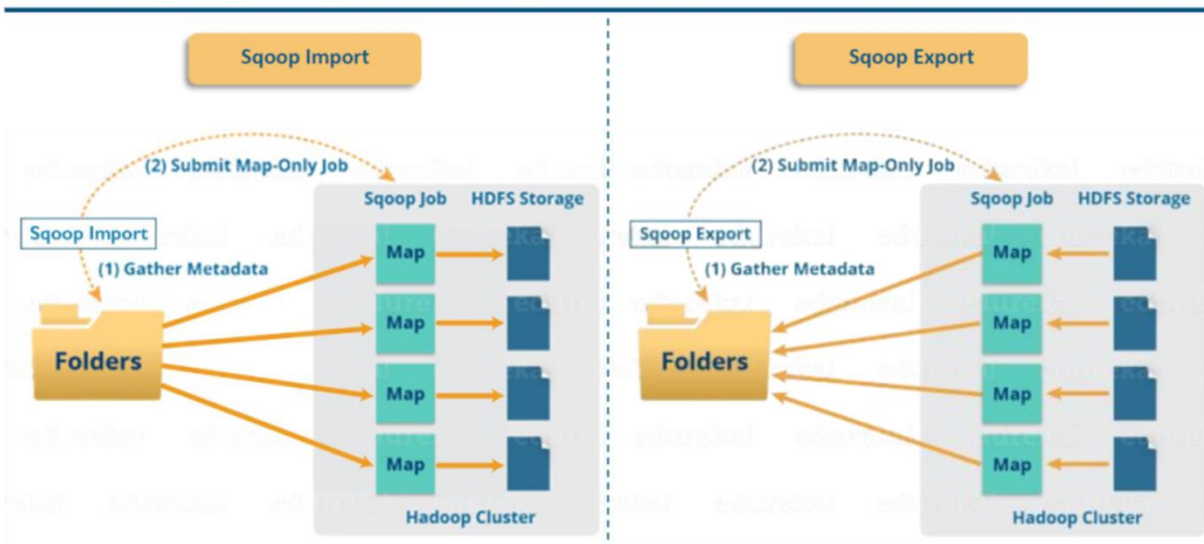


**Fig 1. Sqoop Import and ExportExecution [3]**

**Sqoop Import Command Example:**

```
sqoop import \
  --connect jdbc:mysql://[Database_Server]/[Database_Name] \
  --username [User] \
```

```
--password [Password] \
--table [Table_to_Import] \
--target-dir[HDFS_Destination_Path] \
--where "[Optional_SQL_Filter_Condition]"
```

**Apache Spark**

Apache Spark is a versatile, distributed computing platform designed for large-scale data processing and analysis. It can operate independently or in conjunction with external resource management systems like YARN, Kubernetes, or Mesos.

At its core, Spark utilizes Resilient Distributed Datasets (RDDs), which are fault-tolerant collections of elements distributed across a cluster. These RDDs allow Spark to perform parallel operations on data using multiple executors, enabling efficient querying and transformation of large datasets.

Building upon RDDs, Spark introduced DataFrames, which add a structured schema to the distributed data collections. This enhancement provides a more intuitive and optimized way to work with structured data. DataFrames are flexible and can interact with various data sources, including file systems, relational and NoSQL databases, and real-time data streams.

This architecture allows Spark to handle a wide range of data processing tasks, from batch processing to real-time analytics, making it a powerful tool in the big data ecosystem.

This Spark code snippet demonstrates how to establish a connection with a MySQL database and retrieve data:

```
valdf = spark.read.format("jdbc")
.option("url", "jdbc:mysql://db1.zaloni.com/customer")
.option("driver", "com.mysql.jdbc.driver")
.option("dbtable", "customerProfile")
.option("user", "*****")
.option("password", "******")
.load()
```

Like Sqoop, Spark offers mechanisms for parallel data extraction from databases, enhancing performance through distributed processing. Spark's approach to partitioning data for efficient retrieval is comparable to Sqoop's split functionality.

In Spark, the 'partitionColumn' parameter serves a similar purpose to Sqoop's '--split-by' option. This column is typically used to divide the data into manageable chunks for parallel processing.

To define the range of data to be extracted, Spark uses 'lowerBound' and 'upperBound' parameters. These set the minimum and maximum values of the partition column, analogous to how Sqoop determines the boundaries of its splits.

The 'numPartitions' parameter in Spark is crucial for parallelization. It determines how many segments the data range is divided into, each handled by a separate task. This parameter also indirectly controls the maximum number of simultaneous JDBC connections to the database.

It's important to note that while 'numPartitions' sets an upper limit on concurrent database connections, the actual number may be lower. This depends on the available Spark executors for the job, which can be influenced by cluster resources and configuration.

By leveraging these parameters, Spark can efficiently distribute the workload of data extraction across multiple executors, potentially significantly reducing the time required for large-scale data transfers from relational databases [4].
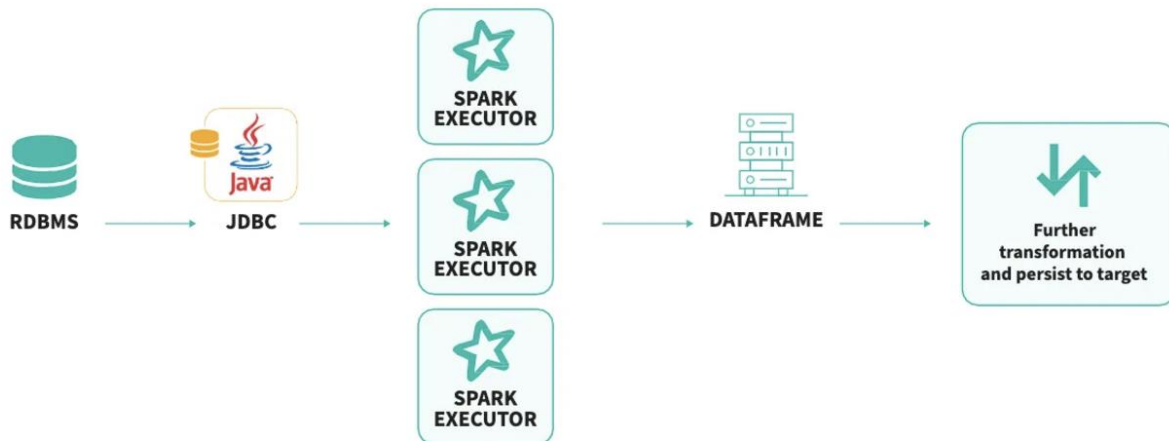


**Fig 2. Spark Execution flow**

**Pros and Cons of using Spark over Sqoop**

Pros:
1. Data Persistence Flexibility
- With Spark, persisting data is optional, unlike Sqoop which requires data to be persisted into HDFS
- Users can perform transformations and only persist the transformed data if needed
- Offers multiple target systems for data persistence (AVRO, Parquet, JSON, relational databases, NoSQL databases, streams)

2. Interactive Data Processing
- Allows interactive work with data using Jupyter Notebooks or Spark Shell
- Users can examine data before deciding to persist it
- Sqoop only submits MapReduce jobs without interactive capabilities

3. Deployment Flexibility

- Can run in standalone mode or with resource managers (YARN/Mesos/Kubernetes)
- Ability to start lightweight transient Spark clusters
- No competition for resources in long-running Hadoop clusters

4. Federated Data Queries
- Can join data from different sources in memory
- No need to persist data from various sources into a common location first
- Can work with multiple data sources simultaneously

Cons:
1. Data Type Mapping
- May require extending or creating new implementations of JDBCDialect
- Need to handle conversion of SQL data types to Catalyst data types

2. Performance Tuning Complexity
- Requires careful configuration of numPartitions
- Proper selection of PartitionColumn is crucial for parallelism
- Need to manage coalesce or repartition functions to optimize file writing and JDBC connections

**Challenges**

The primary challenge when using Spark or Sqoop for data ingestion from databases lies in managing the high number of concurrent connections. As organizations deal with larger data volumes, the increased parallelism can significantly impact database performance.

This issue has led many companies to reconsider their data ingestion strategies. A recommended approach is to set up dedicated database servers specifically for Sqoop or Spark jobs. This separation ensures that the performance degradation caused by these high-intensity operations doesn't affect the primary servers, allowing business operations to continue uninterrupted.

For cloud-based operations, it's advisable to deploy specialized SQL loaders with enhanced memory and CPU capabilities. These loaders are designed explicitly for Sqoop or Spark ingestion jobs. To keep data current, companies typically establish transaction log (T-log) replication between the primary databases and these cloud-based SQL loaders [5].

By utilizing these dedicated resources for data ingestion tasks, organizations can maintain the performance of their primary systems while efficiently handling large-scale data transfers. This approach strikes a balance between the need for high-volume data ingestion and the importance of maintaining stable, responsive database services for day-to-day operations.
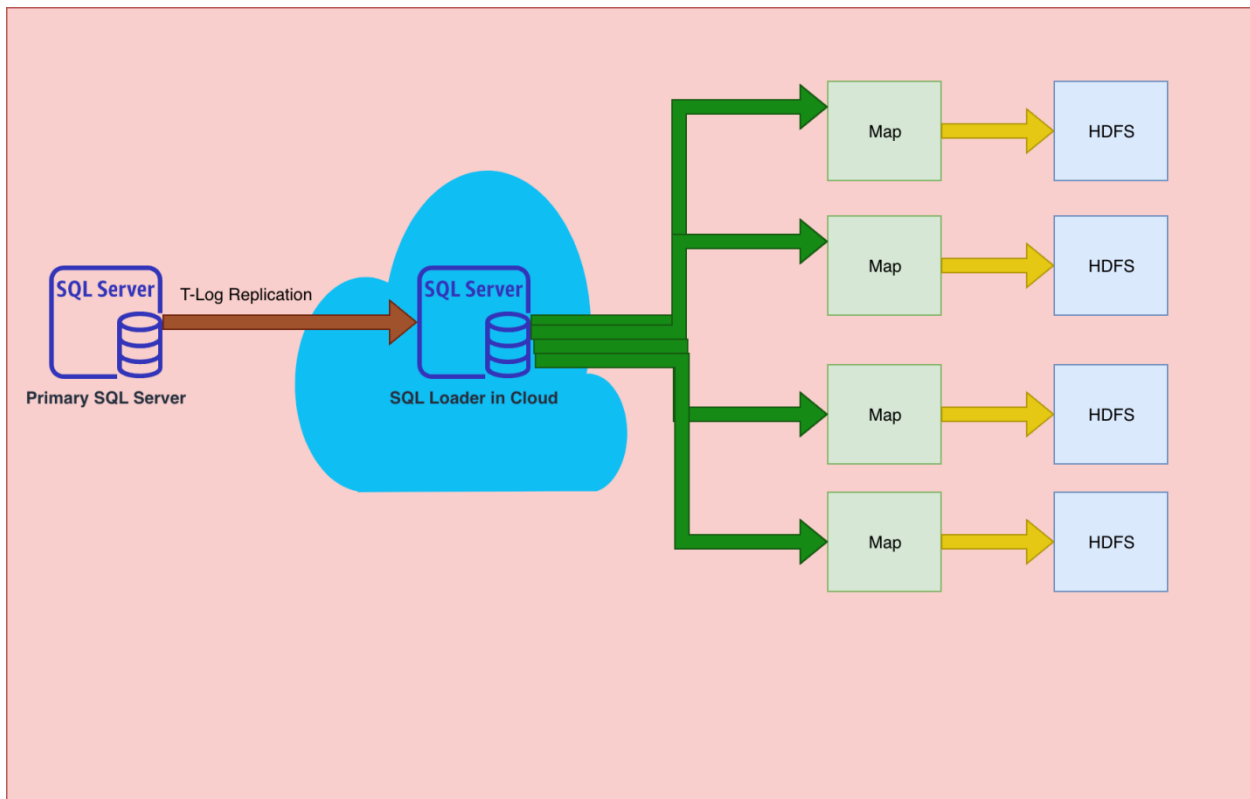
Fig 3. Sqoop/Spark using SQL Loader in Cloud

**Conclusion**

For importing data from relational databases into Hadoop Distributed File System (HDFS), Spark and Sqoop are popular choices due to their seamless integration with the Hadoop ecosystem and user-friendly nature. These tools also offer capabilities to write data to various cloud storage platforms, including Amazon S3 and Google Cloud Storage. It's worth noting that while Spark continues to evolve with ongoing development and regular updates, Sqoop's development has slowed. Despite this, Sqoop remains in use in many environments.

**Reference**:
1. https://readvitamin.wordpress.com/2007/07/27/how-to-export-data-in-teradata-sql-using-bteq/
2. https://sqoop.apache.org/docs/1.4.7/SqoopUserGuide.html#_importing_data_into_hive
3. https://avinash333.com/spark-3/
4. Nikhil Goel, "Apache Spark vs. Sqoop: Engineering a better data pipeline", https://medium.com/zaloni-engineering/apache-spark-vs-sqoop-engineering-a-better-data-pipeline-ef2bcb32b745
5. Pamela Mooney, "How to Set Up Transactional Replication", https://www.red-gate.com/simple-talk/databases/sql-server/database-administration-sql-server/how-to-set-up-transactional-replication/