



**EuroHPC**  
Joint Undertaking

# GUIDE FOR DIRECT INTEGRATION WITH A LOW-LEVEL CMS DATABASE

## CONTENT

1. Batch Data Transfer Using Sqoop.....	2
2. Stream Data Ingestion with Apache Flume .....	3
3. Custom MapReduce or Spark Job.....	4
4. Apache NiFi.....	5
5. PostgreSQL Foreign Data Wrappers (FDW).....	6
6. Database Connectors and JDBC/ODBC .....	7
7. ETL (Extract, Transform, Load) Tools .....	8
8. Change Data Capture (CDC).....	9
9. Data Virtualization .....	10
10. HDFS client libraries.....	11
Reference .....	14

Integrating Hadoop Distributed File System (HDFS) with a PostgreSQL database can be approached through various methods, allowing bidirectional data flow and utilizing the strengths of both systems. These methods range from simple data import/export operations to more complex data federation and real-time replication techniques. Below is a comprehensive list of approaches:

## 1. Batch Data Transfer Using Sqoop

Apache Sqoop is a tool designed for efficiently transferring large datasets between Hadoop and structured data stores, such as relational databases. Data can be imported from PostgreSQL to HDFS and exported from HDFS to PostgreSQL. Sqoop uses MapReduce to import and export data, ensuring parallel processing and fault tolerance.

### Advantages:

- Optimized for transferring large volumes of data.
- Provides capabilities for parallel processing.
- Integrates well with the Hadoop ecosystem.

### Disadvantages:

- Not intended for real-time data integration.
- Setup can be complex, requiring a fair understanding of the Hadoop ecosystem.
- Relies on MapReduce, which can add overhead compared to newer data processing frameworks.

### Example:

To import data from a PostgreSQL database to HDFS using Sqoop, the command might look like this::

```
sqoop import \  
--connect jdbc:postgresql://hostname:port/dbname \  
--username dbuser \  
--password dbpass \  

```

```
--table tablename \  
  
--target-dir /user/hdfs/tablename \  
  
--m 1
```

## 2. Stream Data Ingestion with Apache Flume

Apache Flume is a service for efficiently collecting, aggregating, and moving large amounts of streaming data into HDFS. Flume's architecture is flexible, allowing for various sources and destinations. For PostgreSQL, a custom Flume sink can be written to push data to a PostgreSQL database, or a custom source that reads from a PostgreSQL instance.

### Advantages:

- Suitable for continuous data ingestion, ideal for event logs and data streaming.
- Highly adaptable with the ability to build custom components.
- Can handle large amounts of data with high throughput.

### Disadvantages:

- Out-of-the-box support for PostgreSQL as a sink may be limited, potentially requiring custom development.
- May have a steeper learning curve for setting up streaming pipelines compared to batch processes.
- Limited support for transactions or complex data transformations.

### Example:

To stream data from a source to HDFS and then to PostgreSQL, you can create a Flume configuration file as follows:

```
# Define source, channel and sink  
  
agent.sources = r1  
  
agent.channels = c1  
  
agent.sinks = k1  
  
# Configure the source
```

```
agent.sources.r1.type = netcat agent.sources.r1.bind =  
localhost agent.sources.r1.port = 44444
```

```
# Configure the channel agent.channels.c1.type = memory  
agent.channels.c1.capacity = 1000  
agent.channels.c1.transactionCapacity = 100
```

```
# Configuring the HDFS sink
```

```
agent.sinks.k1.type = hdfs  
agent.sinks.k1.hdfs.path = hdfs://hostname:port/path/to/dir agent.sinks.k1.hdfs.fileType =  
DataStream
```

```
# Connecting the source, sink and channel
```

```
agent.sources.r1.channels = c1 agent.sinks.k1.channel =  
c1
```

### 3. Custom MapReduce or Spark Job

Developers can write custom MapReduce or Spark jobs that execute SQL queries against PostgreSQL to extract data and store it in HDFS or vice versa. This method provides maximum control over the integration process but requires significant development effort.

MapReduce is a programming model for processing large data sets with a parallel, distributed algorithm on a cluster, while Spark is an open-source distributed processing system used for big data workloads.

#### **Advantages:**

- Provides detailed control over data processing and integration.
- Spark jobs run in-memory and are faster than traditional MapReduce jobs
- Scalable for processing large volumes of data.

**Disadvantages:**

- Requires significant development effort and expertise in MapReduce or Spark..
- Integration with PostgreSQL is not native and requires JDBC for connectivity.
- Development and maintenance can be more complex compared to using specialized integration tools.

**Example:**

Running Spark in Scala to transfer data from PostgreSQL to HDFS might look like this:

```
val spark=SparkSession.builder.appName("PostgreSQLtoHDFS").getOrCreate(
)

val jdbcDF = spark.read

    .format("jdbc")

    .option("url", "jdbc:postgresql://hostname:port/dbname")

    .option("dbtable", "tablename")

    .option("user", "dbuser")

    .option("password", "dbpass")

    .load()

jdbcDF.write.parquet("/user/hdfs/tablename")
```

This Spark job creates a DataFrame by reading from a PostgreSQL table and then writes the DataFrame to HDFS in parquet format.

## 4. Apache NiFi

Apache NiFi is a data logistics platform that allows the automation of data movement between systems. It can be used to create data flows that read from or write to both HDFS and PostgreSQL. NiFi processors for HDFS and SQL databases can be configured to work with PostgreSQL.

### **Advantages:**

- NiFi's user interface provides easy configuration and monitoring of data flows.
- It offers data provenance, which is valuable for tracking and debugging data flow issues.
- It is designed to handle data of all shapes, sizes, and volumes.

### **Disadvantages:**

- As a graphical tool, it may not be as efficient for complex data transformation logic as code-based tools.
- The learning curve for setting up advanced configurations and optimizations can be steep.
- There may be performance issues with some processors, especially under heavy load

### **Example:**

A NiFi data flow to move data from PostgreSQL to HDFS can be set up as follows:

- GetJDBC Processor: Configured to connect to the PostgreSQL database and execute a SQL query to fetch the data.
- PutHDFS Processor: Configured to write the fetched data to HDFS.

The processors are connected with a data flow link, and each processor is configured with the necessary credentials and connection properties. NiFi handles the rest, including retrying failed operations and buffering data if necessary.

## **5. PostgreSQL Foreign Data Wrappers (FDW)**

The Foreign Data Wrapper feature in PostgreSQL allows it to act as an SQL-based federated database system. HDFS can be accessed by using an FDW, which connects to external data sources. This integration allows PostgreSQL to directly query data stored in HDFS.

### Advantages:

- FDW allows data to be accessed across different systems using PostgreSQL's SQL capabilities.
- It integrates seamlessly into existing PostgreSQL installations.
- This can be a low-overhead solution for combining SQL queries on remote data sources with local processing.

### Disadvantages:

- Performance is often limited by the capabilities of the FDW and the network latency between PostgreSQL and the data source.
- It can be complex to set up and execute.
- Write operations may be limited or unsupported depending on the specific implementation of the FDW.

### Example:

To access HDFS data from PostgreSQL using FDW:

- Install and configure the Hadoop\_FDW extension in PostgreSQL.
- Create a foreign server that connects to HDFS.
- Map external tables in HDFS to foreign tables in PostgreSQL.

```
CREATE EXTENSION hadoop_fdw;
```

```
CREATE SERVER hdfs_server FOREIGN DATA WRAPPER hadoop_fdw OPTIONS  
(host 'hdfs-hostname', port '8020');
```

```
CREATE FOREIGN TABLE hdfs_table (...) SERVER hdfs_server OPTIONS  
(filename '/path/to/hdfs/file');
```

```
SELECT * FROM hdfs_table;
```

## 6. Database Connectors and JDBC/ODBC

Various Hadoop components, such as Apache Hive or Apache Drill, can access relational databases using JDBC (Java Database Connectivity) or ODBC (Open Database Connectivity). These components can query PostgreSQL databases directly and store the results in HDFS.

**Advantages:**

- Provide a standardized way to connect to PostgreSQL from various applications and services.
- JDBC/ODBC is supported by many tools, including Hadoop ecosystem tools like Hive and Spark.
- They are generally easy to set up and use with existing database drivers.

**Disadvantages:**

- Performance can suffer due to the overhead of the JDBC/ODBC bridge, especially with large volumes of data.
- They may not support all features of a direct database connection.
- Driver compatibility issues can arise, and maintaining up-to-date drivers is necessary.

**Example:**

To move data from HDFS to PostgreSQL using Hive with JDBC:

- Set up a Hive table that maps the data in HDFS.
- Use the Hive JDBC driver to connect to Hive from a Java application.
- Execute the SQL INSERT INTO command to move data from the Hive table to PostgreSQL.

```
ConnectionhiveConnection=DriverManager.getConnection(  
"jdbc:hive2://hive-      host:10000/default",      "user",  
"password");
```

```
Statement stmt = hiveConnection.createStatement();
```

```
stmt.executeUpdate("INSERT INTO postgresql_table SELECT * FROM hdfs_table");
```

## 7. ETL (Extract, Transform, Load) Tools

Traditional ETL tools like Talend, Informatica, and Pentaho can be used to integrate data between HDFS and PostgreSQL. These tools provide a graphical interface for designing data pipelines and can handle complex transformations and batch processing.



### **Advantages:**

- ETL tools often come with a wide range of built-in connectors for various data sources and destinations, including HDFS and PostgreSQL.
- They provide a visual interface for designing data workflows, which can be easier for users who are not familiar with programming.
- Advanced ETL tools offer robust data transformation capabilities, handling complex data manipulation tasks.

### **Disadvantages:**

- Commercial ETL tools can be expensive due to licensing costs.
- They may be overkill for simple data transfer tasks, adding unnecessary complexity and overhead.
- ETL processes can be resource-intensive, especially when dealing with large data sets.

### **Example:**

A typical ETL process using a tool like Talend or Informatica would include:

- Creating a job or workflow in the ETL tool.
- Configuring a source component to connect and extract data from HDFS.
- Optionally adding transformation components to process the data.
- Configuring a destination component to load the data into PostgreSQL.

The ETL tool's interface will guide the user through setting up each component and mapping data fields between the source and the destination.

## **8. Change Data Capture (CDC)**

CDC tools can capture changes made at the database level in real time and apply these changes to another system. Tools like Debezium can capture row-level changes in PostgreSQL and stream these changes to HDFS using Kafka as an intermediary data store.

**Advantages:**

- CDC allows real-time data integration, capturing changes as they happen.
- It minimizes the load on the source database since it captures only incremental changes.
- CDC is suitable for event-driven architectures and microservices patterns.

**Disadvantages:**

- Setting up CDC can be complex and requires a deep understanding of the internals of the source database.
- It may require additional infrastructure, such as Apache Kafka, for buffering and distributing change events.
- Handling edge cases and ensuring data consistency can be challenging.

**Example:**

Setting up CDC using Debezium and Kafka would involve:

- Configuring Debezium to monitor PostgreSQL for changes.
- Streaming change events to a Kafka topic.

Consuming the Kafka topic and writing changes to HDFS, possibly using a Kafka-HDFS connector or a custom consumer job.

```
KafkaConsumer<String,String>consumer=new KafkaConsumer<>(props);  
consumer.subscribe(Arrays.asList("postgres_changes"));  
  
for(ConsumerRecord<String,String>record: consumer.poll(Duration.ofMillis(100)))  
{  
  
    writeToHdfs(record.key(), record.value());  
  
}
```

## 9. Data Virtualization

CDC tools can capture changes made at the database level in real-time and apply these changes to another system. Tools like Debezium can capture row-level changes in PostgreSQL and stream these changes to HDFS using Kafka as an intermediary data store.

**Advantages:**

- CDC allows real-time data integration, capturing changes as they happen.
- It minimizes the load on the source database since it captures only incremental changes.

- CDC is suitable for event-driven architectures and microservices patterns.

**Disadvantages:**

- Setting up CDC can be complex and requires a deep understanding of the internals of the source database.
- It may require additional infrastructure, such as Apache Kafka, for buffering and distributing change events.
- Handling edge cases and ensuring data consistency can be challenging.

**Example:**

Setting up CDC using Debezium and Kafka would involve:

- Configuring Debezium to monitor PostgreSQL for changes.
- Streaming change events to a Kafka topic.

Consuming the Kafka topic and writing changes to HDFS, possibly using a Kafka-HDFS connector or a custom consumer job.

## 10. HDFS client libraries

HDFS client libraries like libhdfs for C or PyArrow for Python can be used to develop custom applications that interact with both HDFS and PostgreSQL. These applications can programmatically move data between these two systems.

**Advantages:**

- Provides programmatic control over data operations, offering flexibility to implement custom logic.
- Enables integration of HDFS data operations into existing applications.
- Useful for creating tailored solutions that meet specific requirements.

**Disadvantages:**

- Developing custom applications requires significant coding effort and ongoing maintenance.
- Performance optimization can be complex and requires deep understanding of both HDFS and PostgreSQL.
- Potential for bugs and issues due to custom code that may not be as thoroughly tested as standard tools.

**Example:**

Python script, using `psycopg2` for PostgreSQL and `PyArrow` for HDFS, might look like this :

```
import psycopg2

import pyarrow as pa

import pyarrow.hdfs as hdfs

conn=psycopg2.connect(database="dbname",user="user", password="password",
host="hostname", port="port")

cursor = conn.cursor()

hdfs_connect = hdfs.connect(host='hdfs-hostname', port=8020)

cursor.execute("SELECT * FROM some_table") rows =
cursor.fetchall()

with hdfs_connect.open('/path/to/hdfs/file', 'wb') as file: file.write(rows)
```

Each of these approaches has its advantages and disadvantages in terms of complexity, performance, scalability, and real-time processing capabilities. The choice of integration method will largely depend on the specific use case, the volume and speed of data, and the desired level of connectivity between systems.

The pros and cons of each approach to integrating HDFS with a PostgreSQL database are summarized in the following table:

#	Approach	Advantages	Disadvantages
1	Batch data transfer using Sqoop	Efficient for large batches; parallel processing; integrates with the Hadoop ecosystem.	Not suitable for real-time processing; complex setup; overhead of MapReduce.
2	Data streaming with Apache Flume	Good for streaming data; Adaptive; can handle high-throughput data streams.	Requires a custom sink for PostgreSQL; limited transactional support; steep learning curve.
3	MapReduce / Spark jobs	Very adaptable; can perform complex processing; Scalable.	Requires significant development effort; not natively integrated with PostgreSQL.
4	Apache NiFi	Visual control of data flow; supports broad data sources and destinations; It's good for data ancestry.	Can be complex to configure; potential performance bottlenecks.
5	PostgreSQL Foreign Data Packaging (FDW)	Enables SQL queries on HDFS data; integrates with the PostgreSQL query optimizer.	Limited by the performance of FDW; additional setup and maintenance overhead.
6	Database and JDBC/ODBC connectors	Standard data access interface; widely supported by data tools; easy integration with BI tools.	May introduce execution overhead; limited by the performance of the JDBC / ODBC driver.
7	ETL Tools	Graphical interface for pipeline design; supports complex transformations; batch processing.	It can be resource intensive; commercial tool licensing costs; potential learning curve.
8	Change of data capture (CDC)	Real-time data synchronization; minimal impact on the source system; Suitable for event-driven architectures.	Complex setup; requires additional infrastructure such as Kafka; potential data consistency challenges.
9	Data virtualization	Real-time integration; avoids data duplication; unified queries between data sources.	Can be complex to implement; potential impact on performance; trading platform licensing costs.
10	HDFS client libraries	Programmatic data access control; supports complex logic; Suitable for custom applications.	Development overheads; requires custom code support; potential performance issues.

## Reference

1. WebHDFS FileSystem APIs; 2022;  
<https://learn.microsoft.com/en-us/rest/api/datalakestore/webhdfs-filesystem-apis>
2. Cloudera; Apache Nifi; <https://www.cloudera.com/products/open-source/apache-hadoop/apache-nifi.html>
3. Abhinav Agarwal; Use NiFi to extract and parse data from HTTP endpoints; 2023;  
<https://www.projectpro.io/recipes/use-nifi-extract-and-parse-data-from-http-endpoints-and-store-data-persistent-storage>
4. Flume 1.11.0 User Guide;  
<https://flume.apache.org/releases/content/1.11.0/FlumeUserGuide.html>
5. David Taylor; Apache Flume Tutorial: What is, Architecture & Hadoop Example;  
<https://www.guru99.com/create-your-first-flume-program.html>
6. Prateek Majumder; Apache Sqoop: Features, Architecture and Operations; 2023;  
<https://www.analyticsvidhya.com/blog/2022/09/apache-sqoop-features-architecture-and-operations/>
7. Foreign Data Wrappers; [https://wiki.postgresql.org/wiki/Foreign\\_data\\_wrappers](https://wiki.postgresql.org/wiki/Foreign_data_wrappers)
8. PostgreSQL Documentation-Create foreign data wrapper;  
<https://www.postgresql.org/docs/current/sql-createforeigndatawrapper.html>