

Oracle and Hadoop, Let Them Talk Together

Laurent Leturgez
PREMISEO
Lille - France

Keywords:

Oracle, Hadoop, Data, JDBC, ODBC, Hive, Mapper, Reducer, YARN, HDFS, Parquet, ORC, Sequence files, delimited files, Sqoop, Spark, Oracle BigData Connectors, Oracle BigData SQL, SQL, Analytics, Flow, Gluent

Introduction

We all know that Oracle has been created in 1977 and then, has evolved to become one of the most powerful rdbms.

A couple of years ago, hadoop was on the hype and is now one of the main systems for big data platforms.

Their goal is the same ... storing data, but not the same way and not for the same kind of workload.

But Oracle and Hadoop systems can interact together and exchange data: export/import, data transfer, offload.

After a short introduction on how Oracle and Hadoop are different but complementary, a presentation of many use cases will give the attendee some real cases of Oracle and hadoop integration.

Then, after presenting use cases, several products will be presented to give the attendee some clues on how they can implement data exchanges between Oracle and Hadoop:

- Sqoop and Spark for hadoop
- HS datasources (ODBC and JDBC) and Big Data SQL solutions for Oracle
- Gluent data platform

At the end of the presentation, the attendee will have a good view on how he can use hadoop and/or Oracle products to exchange data together.

IT services changed the way they managed data

Nowadays, IT services are submitting to a huge transformation. In this transformation, the way data is managed is deeply modified.

10-15 years ago, when a new project was initiated, the database management system chosen was the one mainly used in the team or by the company. It was a product approach.

More recently, the approach has shifted to a solution approach. Each project will use a kind of data (relational, non-relational, document oriented etc.) and with this kind of data, the solution that fit with it.

As a result, IT services use the right tool for the right problem.

Hadoop is one of this new products that offers many solutions to store and manage data. Usually, Hadoop clusters manage very large datasets and process these datasets with distributed algorithms.

With this transformation, legacy systems continue being used and the data it manage can enrich Hadoop systems and vice versa.

In this presentation, many products that can be used for this purpose and specially between Oracle and Hadoop will be discribed.

Hadoop, What it is, how it works, and what it can do.

Hadoop is an open source framework based upon the following main components:

- HDFS (Hadoop Distributed File System) used to store huge amounts of data (petabytes scale)
- Map Reduce algorithm: a distributed programming model made to process this amount of data.
- A cluster Manager YARN (Yet Another Resource Negotiator) used to allocate the cluster resources to various applications effectively.

Around this core, Hadoop has lots of products to do specific tasks: data movement, orchestration, security etc.

This is called the Hadoop Ecosystem.

Hadoop is open source but some enterprise editions exists: Cloudera, HortonWorks Data Platform, MAPR converged platform. All the content of this presentation is based on a Cloudera CDH 5.9 platform and an Oracle 12.1.0.2 single instance platform (No RAC).

Hadoop is designed to:

- Manage structured and unstructured data
- Run analytics workload
- Run on commodity servers
- Scale to tens of petabytes and thousands of CPUs.

Hadoop clusters have a specific architecture: the shared nothing architecture. This is a key point to understand its design and why Hadoop cluster can scale at the levels mentioned above.

On the base of what is a Hadoop cluster, we can ask ourselves how we can connect our Oracle databases to our Hadoop clusters, how exchange data between those two systems, which engine use to query data, how can we reuse Oracle Data in a Hadoop Workload etc. and which solution to use for this purpose?

Sqoop, to move data between Oracle and Hadoop

Sqoop is a command line tool provided by Hadoop framework to move data between Oracle and HDFS/Hive. It can be used for many purpose: data offloading, import and export data into and from Hadoop, data archive etc.

These kinds of scenarios can be represented by the figures below:

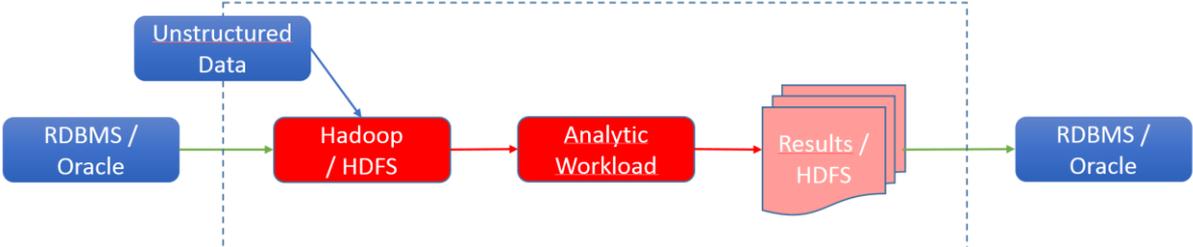


Illustration. 1: Sqoop to enrich analytic workloads with multiple data sources

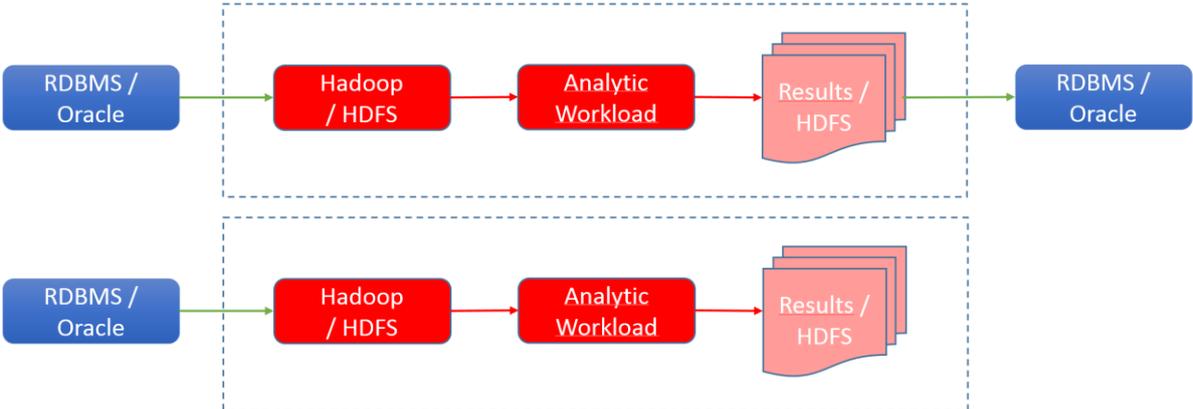


Illustration. 2: Sqoop to offload analytic workloads on Hadoop

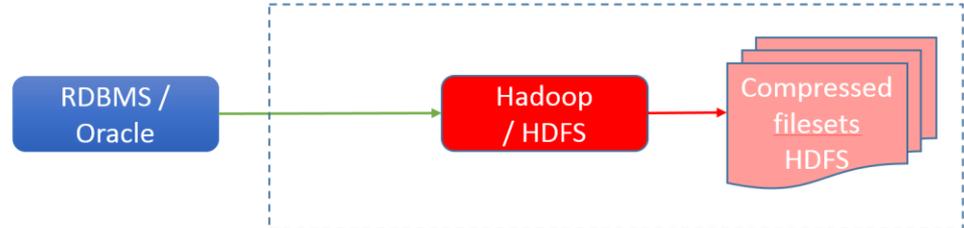


Illustration. 3: Sqoop to archive data into Hadoop

On the Oracle side, it can import and export tables or partitions data completely or partially (based on a SQL statement).

On the Hadoop side, data can be located in a delimited file, a Hive table, compressed or not.

Sqoop runs Map Reduce jobs in conjunction with a JDBC data source to import or export data.

Sqoop can import AVRO files, Parquet files, Sequence files, delimited files (compressed or not), Hive tables, or Hbase.

Spark, to process Hadoop and Oracle Data Together

Spark is an open source distributed computing framework, it's fault tolerant by design. It's usually integrated into Hadoop Cluster because Spark works with a cluster manager. This cluster manager can be YARN or another one (MESOS for example).

There are many components in the Spark framework, and specifically one for data abstraction named Spark SQL.

Spark is built around a data structure named RDD (Resilient Distributed Dataset). This data structure and its evolved data structures (DataFrame, and DataSets) can be filled with data coming from an Oracle Database.

This can be done by using one of the API language available: Scala, Java, Python, R. Once done, the power of Spark (for example a Machine Learning algorithm) can be used on data coming from Oracle, or other compatible data sources (HDFS, Hive, S3 etc.).

ODBC and Heterogeneous Services to process Hadoop Data with Oracle SQL Engine

One of the main powerful components of Oracle is its SQL engine.

It is possible to use it on Oracle tables, and Hadoop table structures managed by Hive or Impala.

To do that, we have to set up an ODBC heterogeneous data source that will use ODBC drivers provided by Hadoop provider (Cloudera for example).

Once done, a simple database link to these Hive or Impala tables can be created in an Oracle Schema.

This can be illustrated by the following figure:

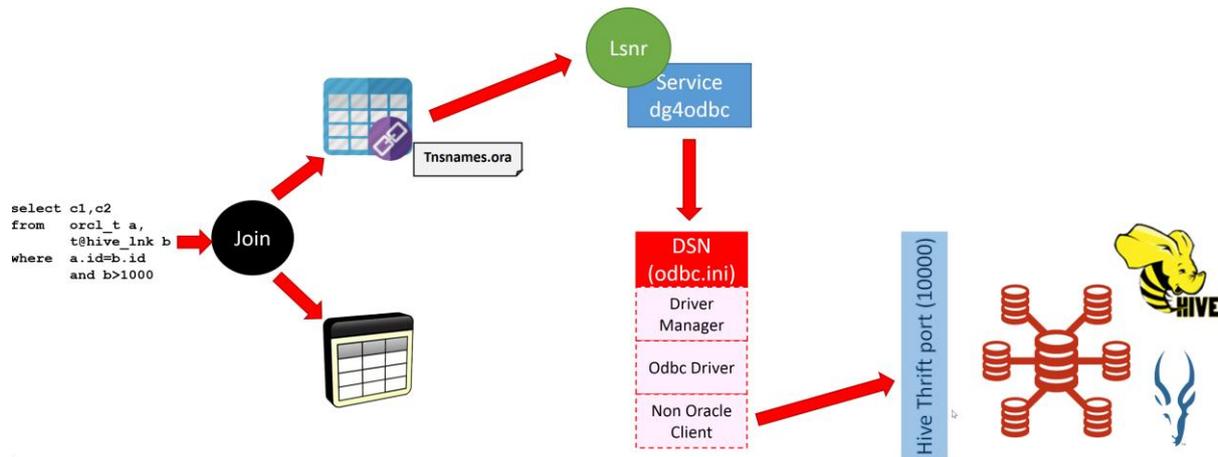


Illustration. 4: Using ODBC Heterogeneous datasource with Hadoop

Oracle Big Data Connectors, various components to move data between Oracle and Hadoop

Oracle Big Data Connectors is a set of tools provided by Oracle to address a Hadoop cluster into an Oracle Database.

In this part, I will only talk about three of them:

- Oracle Datasource for Apache Hadoop
- Oracle Loader for Hadoop
- Oracle SQL Connector for HDFS

Oracle Datasource for Apache Hadoop enables to create external tables in Hive with data stored in an Oracle table.

This is secured and optimized.

Secured because an integration with a Wallet or a Kerberos authentication is possible.

Optimized because filter predicates are pushed down to Oracle, column projection is pushed down too, and partition pruning is enabled.

In this way, it reduces the data volume between Oracle Instance and Hadoop cluster.

An interesting thing to notice is that it is possible to write data from Hive table into the linked Oracle table.

Oracle Loader for Hadoop enables data loading from Hadoop into an Oracle Table. This can be done by an Oracle Java Map Reduce application provided by Oracle.

The main advantage of this solution is that it can be done Online and Offline. Online mode is possible with the use of JDBC or OCI driver. Offline mode is possible with a intermediate datafile that can be a delimited file or a datapump file.

Input data in Hadoop can be AVRO files, delimited files or Key Value format if the input data is located on Oracle SQL Database.

This application is configured by a set of XML files that describes input format, output format, destination table and schema, and database connection.

Oracle SQL Connector for HDFS is, like Oracle Loader for Hadoop, a Java Map Reduce application provided by Oracle.

It creates an external table in Oracle Database. This table's data is located in local files that point to hdfs files.

This external table has the same limitation that traditional external tables: no indexing possible, no DML and only full scan.

Data in the external table is not a living data, indeed if some files are added in HDFS source, it will not be updated in real time into Oracle. The database administrator will have to refresh the external table definition.

Oracle SQL Connector for HDFS can be used to target delimited files stores in HDFS or in Hive tables, but it can be used to store datapump files on HDFS and read them through an external table (a kind of data offloading based on datapump files).

Oracle Big Data SQL to query non-relational data sources

The first goal of Oracle Big Data SQL is to provide a support for queries against non-relational data sources. These data sources can be Hive Tables, HDFS files, Oracle NoSQL, Apache HBase and other NoSQL Databases.

Product installation is made in three phases. A parcel has to be deployed on the Cloudera Cluster (it's compatible with HortonWorks data platform too), then a database bundle configuration has to be done, and finally a package is deployed on the Oracle Server (with Hadoop binaries, libraries etc.).

Oracle Big Data SQL is, like Oracle SQL Connector for HDFS, based on external tables with new access drivers that allow users to access existing hive tables or files stored in HDFS.

Users take the benefit of Oracle SQL Engine and Big Data SQL enables some features to bypass external tables limitation.

- Smart Scan for HDFS acts like a local filter locally to Hadoop nodes to ensure that only requested elements are sent to Oracle Server. This reduces network traffic and data movement between Oracle Server and Hadoop Cluster.
- Storage indexes. Like Exadata Storage indexes, Oracle Big Data SQL maintains storage indexes for HDFS to eliminate unnecessary I/O. Storage Indexes can be used for data elimination, and to improve joins.
- Predicates pushdown against Hive Partitioned table or column oriented storage format in Hadoop (ORC and Parquet).

Oracle Big Data SQL to offload read only data to Hadoop

Another feature of Oracle Big Data SQL is the capability to offload read only data to Hadoop. Technically, a FUSE filesystem is mounted on the Oracle Database Server. This filesystem is, in fact, a link to the Hadoop cluster HDFS.

Then, the datafiles' tablespaces are switched to read only and moved automatically (or manually) to this mount point.

Data located on these tablespaces can be read from any Oracle compatible tools, SQL*Plus, SQL Developer, RMAN etc.

Gluent Data Platform

Gluent Data Platform is the new platform on the hype.

With gluent, you can present data stored in Hadoop in various formats to any compatible rdbms including Oracle.

You can offload your data (one table or a contiguous range of partitions) and your workload on Hadoop.

Finally, Gluent advisor allows you to analyze your Oracle database performance statistics to produce a report.

This report helps you to identify schemas and tables that can be safely offloaded to Hadoop.

Contact address:

Laurent Leturgez

Premiseo

7 rue Jules Guesde

59320, EMMERIN - FRANCE

Phone: +33 (0)6 60 99 70 46

Fax:

Email laurent.leturgez@premiseo.com

Internet: www.premiseo.com