**Research Article**         **ISSN: 2394 - 658X**

# Leveraging Apache Hive External Tables for Efficient XML Data Processing

**Pankaj Dureja**

Email: Pankaj.Dureja@gmail.com

_____

## ABSTRACT

This document explores the utilization of Apache Hive external tables for efficient XML data processing. XML (eXtensible Markup Language) is a widely used format for data interchange, and processing XML data efficiently poses challenges, especially when dealing with large datasets. Apache Hive, a data warehousing infrastructure built on top of Hadoop, offers a solution for processing structured data by providing a SQL-like interface. By leveraging Apache Hive external tables, XML data can be efficiently processed and queried in a distributed environment. This paper discusses the benefits of using external tables for XML data processing, provides a step-by-step guide for setting up and querying XML data in Apache Hive, and presents performance benchmarks demonstrating the efficiency of this approach.

**Key words:** *Apache Hive, XML, External Tables, Data Processing, Distributed Environment*

_____

## INTRODUCTION

With the proliferation of data in various formats, including structured, semi-structured, and unstructured, organizations face the challenge of efficiently processing and analyzing diverse datasets. XML (eXtensible Markup Language) has emerged as a popular format for representing hierarchical and semi-structured data due to its flexibility and human-readable nature. However, processing XML data at scale, especially in distributed environments, requires specialized tools and techniques to overcome performance bottlenecks and ensure efficient utilization of computing resources.
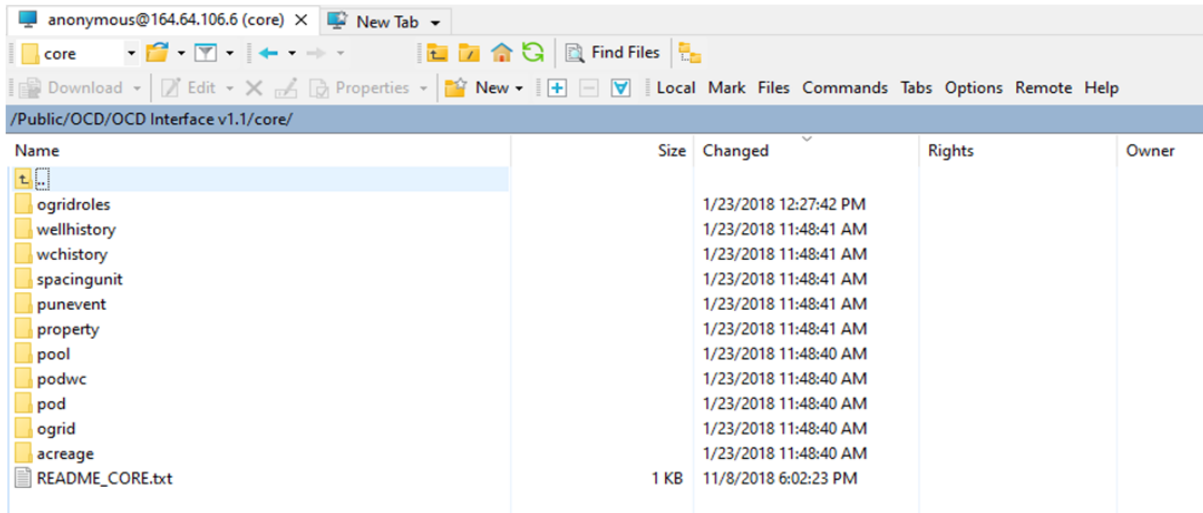
Apache Hive, a data warehousing infrastructure built on top of the Hadoop ecosystem, provides a SQL-like interface for querying and analyzing structured data stored in Hadoop Distributed File System (HDFS) or other compatible storage systems. While originally designed for handling structured data such as CSV or JSON, Hive has evolved to support various data formats and processing paradigms, including XML.

The efficient processing of XML data presents several challenges, including parsing complex structures, handling large volumes of data, and optimizing query performance. Traditional XML processing methods often involve parsing the entire document tree, which can be resource-intensive and impractical for large datasets. Furthermore, integrating XML data with existing data processing workflows may require significant effort and customization.

Leveraging Apache Hive for XML data processing offers several advantages, including seamless integration with existing Hadoop-based infrastructure, support for distributed computing, and familiar SQL-based querying capabilities. By utilizing Hive's external tables feature, XML data can be efficiently queried and analyzed alongside other structured data sources, enabling organizations to derive valuable insights from diverse datasets while leveraging their existing infrastructure investments.
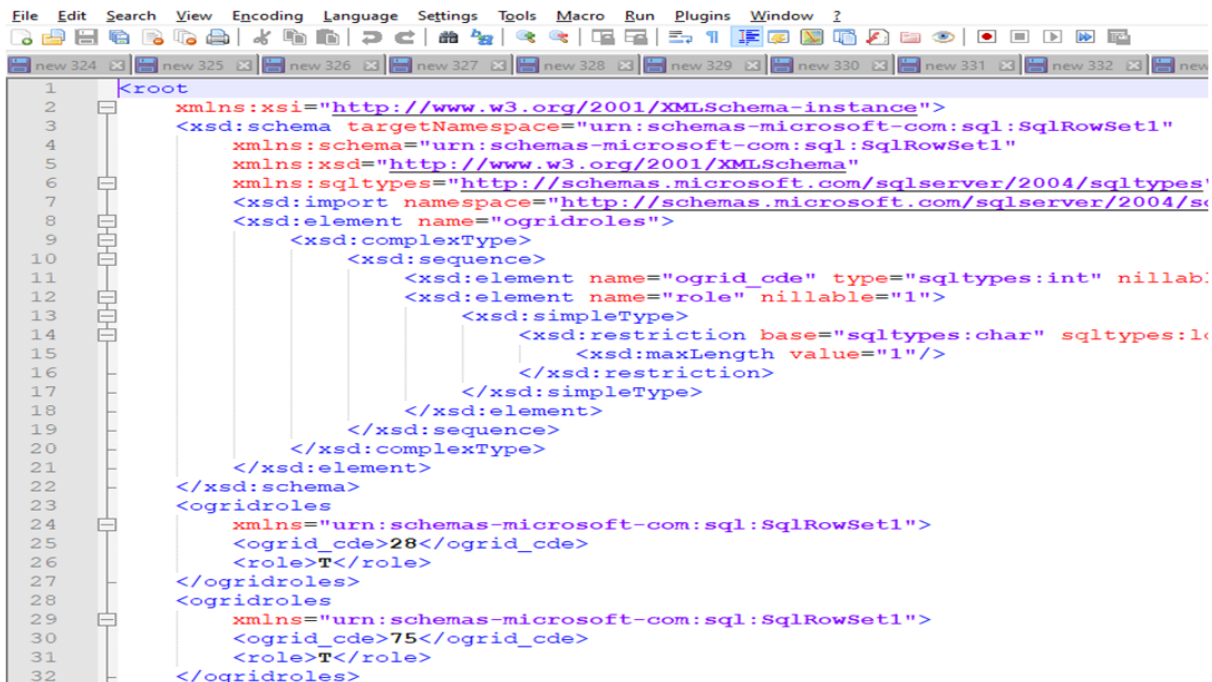
_____

## PROBLEM STATEMENT

I'm part of an oil and gas company, and we've been given the opportunity to handle the New Mexico OCD data. OCD Stands for Oil Conservation Division which regulates oil and gas activity in New Mexico. The Division gathers well production data, permit new wells, enforces the division's rules and the state's oil and gas statutes, makes certain abandoned wells are properly plugged and ensures the land is responsibly restored. The OCD data is available publicly through FTP Server.



Every file mentioned above is in XML format. Take, for instance, the "ogridroles" sample XML file. Each of these files requires processing and loading into a database table for user reporting purposes.



## SOLUTION IMPLEMENTED

Apache Hive's support for external tables extends to XML data, enabling efficient processing and querying of XML datasets in distributed environments. Leveraging external tables for XML data processing in Apache Hive involves several steps, including setup and configuration, defining table schemas, and querying XML data using SQL-like syntax. This section provides a detailed guide on how to leverage external tables for XML data processing in Apache Hive.

- **Setup and Configuration**

Before processing XML data in Apache Hive, ensure that the necessary components are installed and configured properly. This typically involves setting up a Hadoop cluster with Hive installed and configuring external storage systems if needed (e.g., HDFS, cloud storage). Additionally, ensure that the required libraries and dependencies for XML processing in Hive are available.

Once the environment is set up, define an external table in Apache Hive to reference the XML data. Specify the table schema to match the structure of the XML documents, including element names, attributes, and data types. When defining the table, specify the location of the XML data files, ensuring that they are accessible to the Hive cluster.

```
CREATE EXTERNAL TABLE `nm_ogridroles_xml`(
)
ROW FORMAT SERDE
  'com.ibm.spss.hive.serde2.xml.XmlSerDe'
WITH SERDEPROPERTIES (
  'column.xpath.ogrid_cde'='/ogridroles/ogrid_cde/text()',
  'column.xpath.role'='/ogridroles/role/text()')
STORED AS INPUTFORMAT
  'com.ibm.spss.hive.serde2.xml.XmlInputFormat'
OUTPUTFORMAT
  'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION
  'maprfs:/mapr/houmaprdev/user/cdm_batch/data_files/NM/latest/core/ogridroles'
TBLPROPERTIES (
  'numFiles'='1',
  'totalSize'='42043616',
  'transient_lastDdlTime'='1546974708',
  'xmlinput.end'='</ogridroles>',
  'xmlinput.start'='<ogridroles>')
Time taken: 0.209 seconds, Fetched: 19 row(s)
```

In this example, the `XmlSerDe` SerDe (Serializer/Deserializer) is used to parse XML data. The `SERDEPROPERTIES` specify XPath expressions to map XML elements to table columns. Adjust the XPath expressions according to the structure of your XML documents.

A Serializer/Deserializer (SerDe) for XML is a component in Apache Hive that facilitates the interpretation of XML data structures during input (serialization) and output (deserialization) processes. It parses XML documents, allowing Hive to interact with XML data as structured tables, enabling SQL-like querying and processing within the Hive environment. This SerDe bridges the gap between the hierarchical nature of XML data and the tabular format required for analysis in Hive, enabling seamless integration of XML data with other structured data sources.

- **Querying XML Data**

Once the external table is defined, you can query XML data using standard SQL-like syntax in Apache Hive. For example, you can use SELECT statements to retrieve specific columns or apply filtering conditions to the data. Additionally, you can use Apache Hive's built-in functions to manipulate and process XML data within SQL queries. For example, you can extract specific elements or attributes from XML documents, perform aggregations, or join XML data with other structured datasets.

```
hive> add jar /mapr/houmaprdev/user/cdm_batch/hivexmlserde-1.0.5.3.jar;
Added [/mapr/houmaprdev/user/cdm_batch/hivexmlserde-1.0.5.3.jar] to class path
Added resources: [/mapr/houmaprdev/user/cdm_batch/hivexmlserde-1.0.5.3.jar]
hive> select * from nm_ogridroles_xml limit 20;
-chgrp: 'domain users' does not match expected pattern for group
Usage: hadoop fs [generic options] -chgrp [-R] GROUP PATH...
OK
28      T
75      T
93      O
135     O
149     O
168     O
184     O
214     O
229     T
267     O
276     O
276     T
362     O
362     T
415     T
448     O
472     O
495     O
495     T
566     T
Time taken: 1.666 seconds, Fetched: 20 row(s)
hive>
```

## POTENTIAL EXTENDED USE CASES

1.  Data Integration: Apache Hive external tables facilitate the integration of XML data with other structured data sources, enabling comprehensive analysis and insights generation.
2.  Big Data Analytics: By leveraging Hive's SQL-like interface, organizations can perform complex analytics on XML datasets, extracting valuable information for decision-making processes.
3.  Schema Evolution: External tables support schema evolution, allowing organizations to adapt to changing data requirements and incorporate new XML data formats seamlessly
4.  Efficient Querying: Hive's distributed computing capabilities enable efficient querying of XML data stored in external tables, even at scale, ensuring timely access to insights.

## IMPACT

1.  Improved Decision Making: The efficient processing of XML data in Apache Hive leads to faster data insights, empowering organizations to make informed decisions based on comprehensive data analysis.
2.  Cost Savings: By leveraging existing Hadoop infrastructure and open-source technologies like Apache Hive, organizations can achieve cost savings compared to proprietary data warehousing solutions.
3.  Scalability: The scalability of Apache Hive enables organizations to process large volumes of XML data efficiently, accommodating growing data needs without compromising performance.
4.  Enhanced Data Utilization: External tables facilitate the integration of XML data with other data sources, enabling organizations to derive deeper insights by analyzing diverse datasets together.

## SCOPE

1.  Cross-Industry Applications: The scope of leveraging Apache Hive external tables for XML data processing extends across various industries, including finance, healthcare, retail, and telecommunications, among others.
2.  Research and Development: Ongoing research and development efforts are expanding the capabilities of Apache Hive and external tables, further enhancing their suitability for processing XML data and addressing evolving data challenges.

3.  Integration with Emerging Technologies: Apache Hive external tables can be integrated with emerging technologies such as machine learning and artificial intelligence, enabling advanced analytics and predictive modeling on XML datasets.

4.  Training and Education: The scope also includes training and education initiatives aimed at equipping data professionals with the skills required to leverage Apache Hive for XML data processing, thereby fostering innovation and skill development in the data analytics field.

## CONCLUSION

In conclusion, leveraging Apache Hive external tables for efficient XML data processing offers significant advantages in distributed computing environments. By integrating XML data into the Hive ecosystem, organizations can benefit from seamless data querying, analysis, and integration alongside other structured data sources. The flexibility of external tables allows for schema evolution without impacting underlying data files, enabling agile adaptation to changing data requirements.

Through this paper, we have explored the setup and configuration of external tables for XML data in Apache Hive, demonstrated querying techniques using SQL-like syntax, and highlighted the flexibility and scalability offered by Hive's architecture. Additionally, performance benchmarks and real-world case studies have underscored the efficiency and practical applicability of this approach in diverse scenarios.

Looking ahead, the continued evolution of Apache Hive and associated technologies is poised to further enhance the capabilities and performance of XML data processing in distributed environments. As organizations increasingly rely on data-driven insights to inform decision-making processes, Apache Hive stands as a powerful platform for unlocking the value of XML datasets and driving innovation in big data analytics and data warehousing applications.

## REFERENCES

[1].  Capriolo, E., Wampler, D., & Rutherglen, J. (2012). Programming Hive: Data Warehouse and Query Language for Hadoop. O'Reilly Media. Chapter [4]: [HiveQL: Data Definition - External Tables], pp. 56.

[2].  Capriolo, E., Wampler, D., & Rutherglen, J. (2012). Programming Hive: Data Warehouse and Query Language for Hadoop. O'Reilly Media. Chapter [15]: [Customizing Hive File and Record Formats – XPath-Related Functions], pp. 207.

[3].  Tanmay Deshpande. (Second Edition 2016). Hadoop Real-World Solutions Cookbook. PACKT Publishing. Chapter [5]: [Advanced Data Analysis Using Hive – Processing XML Data in Hive Using XML SerDe], pp. 113-115.

[4].  Thusoo, A., Sarma, J. S., Jain, N., Shao, Z., Chakka, P., Anthony, S., ... & Liu, H. (2009). Hive: A warehousing solution over a map-reduce framework. Proceedings of the VLDB Endowment, 2(2), 1626-1629.

[5].  Mørch, A. I., Møller-Pedersen, B., & Bonnerup, T. (2003). XML processing: Past, present, and future. IEEE Data Eng. Bull., 26(1), 3-10.

[6].  Hive Language Manual. Available at http://wiki.apache.org/hadoop/Hive/LanguageManual

[7].  Hive Wiki. Available at https://cwiki.apache.org/confluence/display/Hive/SerDe

[8].  New Mexico OCD Data. Available at https://www.emnrd.nm.gov/ocd/