

VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



SUMMARY of MASTER'S THESIS PROPOSAL

**Building an On-Premises Data Lakehouse
for Big Data Storage and Analysis**

Advisor: Dr. Phan Trong Nhan

Author: Dinh Thanh Phong - 2270243

Ho Chi Minh City, June 2024



Contents

1	Advantages of a Data Lakehouse	2
2	Several existing cloud Data Lakehouse Architecture	3
2.1	Reference Lakehouse Architecture by Microsoft Azure	3
2.2	Reference Lake House Architecture on AWS	4
3	The Need for On-Premises Data Lakehouse	6
4	Proposed On-Premises Data Lakehouse Architecture	7
4.1	Proposed Architecture	7
4.2	Technologies Used	7
5	Demo PhongDinhCS-data_lakehouse has been built and published on GitHub	9
5.1	Example Scenario Applying Data Lakehouse Architecture	9
5.2	GitHub Repository PhongDinhCS and Docker-compose	9
5.3	Components and Trial Running of PhongDinhCS-data_Lakehouse	14
6	Next Steps for Completing the On-Premises Data Lakehouse Architecture	20

1 Advantages of a Data Lakehouse

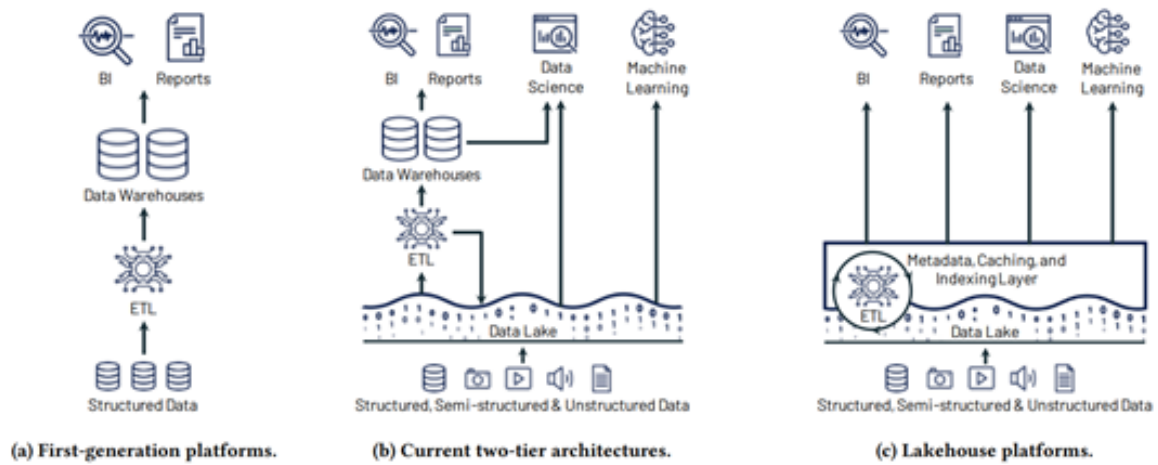


Figure 1: Comparing Data Lakehouse Architecture with Existing Data Architectures

Benefits of the Data Lakehouse: *Simplicity, Flexibility, and Low Cost*

The data lakehouse architecture offers several benefits:

	Data Warehouse	Data Lake	Lakehouse
Data	Rational data from transactional systems, operational databases & business applications	All data including structured, semi-structured, and unstructured.	Query Every kind of data Image, audio, video, and others.
Data Access	SQL Only	Open API, SQL, Python	Open API, SQL, Python
Data Format	Proprietary Format	Open format	Open Format
Governance	Fine-grained Security	Weak Governance and security	Fine-grained Security
Reliability	High with ACID Transactions	Low Quality- Data Swamps	High with ACID Transactions.
Performance	Fast query results using local language	Faster query results, decoupling of computing and storage	Faster and deeper insights without data movement.
Scalability	Scaling becomes expensive	Low cost of scaling Regardless of the data-type	Low cost of scaling regardless of the data-type

Figure 2: Comparison of Features: Data Warehouse, Data Lake, and Data Lakehouse

Unified System: By integrating data warehouse features directly into the low-cost storage of a data

lake, it eliminates the need for separate data warehouse layers.

Simplified Access: Data teams can access and utilize data more efficiently without navigating multiple systems.

Updated Data: Ensures that the most complete and up-to-date data is available for data science, machine learning, and business analytics.

2 Several existing cloud Data Lakehouse Architecture

2.1 Reference Lakehouse Architecture by Microsoft Azure

Organization of the Reference Architectures

The reference architecture is structured around the workflows of Source, Ingest, Transform, Query and Process, Serve, Analysis, and Storage:

Source

The architecture distinguishes between semi-structured and unstructured data (sensors and IoT, media, files/logs) and structured data (RDBMS, business applications). SQL sources (RDBMS) can also be integrated into the lakehouse and Unity Catalog without ETL through lakehouse federation. Additionally, data can be loaded from other cloud providers.

Ingest

Data can be ingested into the lakehouse via batch or streaming:

- Files transferred to cloud storage can be directly loaded using Databricks Auto Loader.
- For batch ingestion of data from enterprise applications into Delta Lake, the Databricks lakehouse relies on partner ingestion tools with specific adapters for these systems.
- Streaming events can be ingested directly from streaming systems like Kafka using Databricks Structured Streaming. Streaming sources can include sensors, IoT, or change data capture processes.

Storage

Data is typically stored in cloud storage systems where ETL pipelines use the medallion architecture to organize and store data as Delta files/tables.

Transform and Query and Process

The Databricks lakehouse utilizes its engines, Apache Spark and Photon, for all transformations and queries. Due to its simplicity, the declarative framework DLT (Delta Live Tables) is a good option for building reliable, maintainable, and testable data processing pipelines. Supported by Apache Spark and Photon, the Databricks Data Intelligence Platform supports both types of workloads: SQL queries through SQL warehouses and SQL, Python, and Scala workloads through workspace clusters. For data science (ML Modeling and Gen AI), the Databricks AI and Machine Learning platform provides specialized ML runtimes for AutoML and coding ML jobs. All data science workflows and MLOps are best supported by MLflow.

Serve

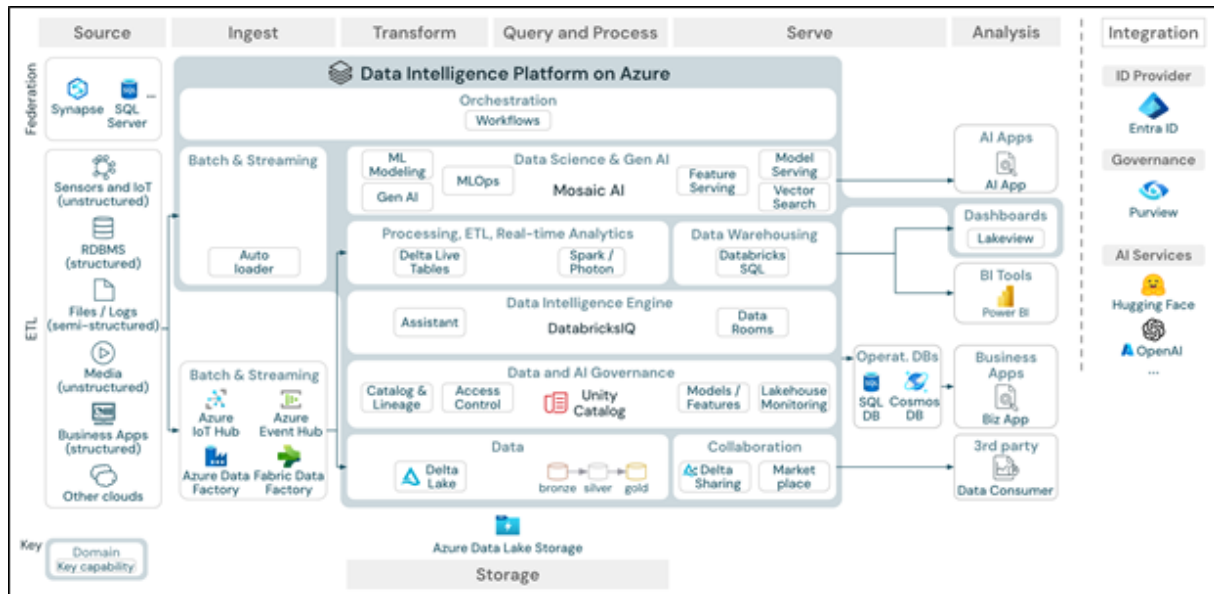


Figure 3: Data Lakehouse Architecture proposed by Microsoft using Azure Data Lake Storage

For DWH and BI use cases, the Databricks lakehouse offers Databricks SQL, a data warehouse supported by SQL warehouses and serverless SQL warehouses. For machine learning, model serving provides enterprise-level, large-scale, real-time model serving capabilities hosted in the Databricks control plane. Operational databases: External systems, such as operational databases, can be used to store and serve the final data products for user applications. Collaboration: Business partners gain secure access to the data they need through Delta Sharing. Based on Delta Sharing, the Databricks Marketplace is an open forum for exchanging data products.

Analysis

The final business applications reside within this workflow. Examples include custom clients like AI applications connected to Databricks Model Serving for real-time inference or applications accessing data pushed from the lakehouse to an operational database. For BI use cases, analysts typically use BI tools to access the data warehouse. SQL developers can also use the Databricks SQL Editor (not shown in the diagram) to query and create dashboards. The Data Intelligence Platform also provides dashboards for building data visualizations and sharing insights.

2.2 Reference Lake House Architecture on AWS

The diagram below illustrates the reference Lake House architecture on AWS.

Data Ingestion Layer

The ingestion layer in the AWS Lake House architecture uses various AWS services specifically designed to ingest data from different sources into the storage layer. The key services include:

- **Operational Database Sources:** AWS Data Migration Service (DMS) imports and replicates data from operational RDBMS and NoSQL databases to Amazon S3 or Amazon Redshift.

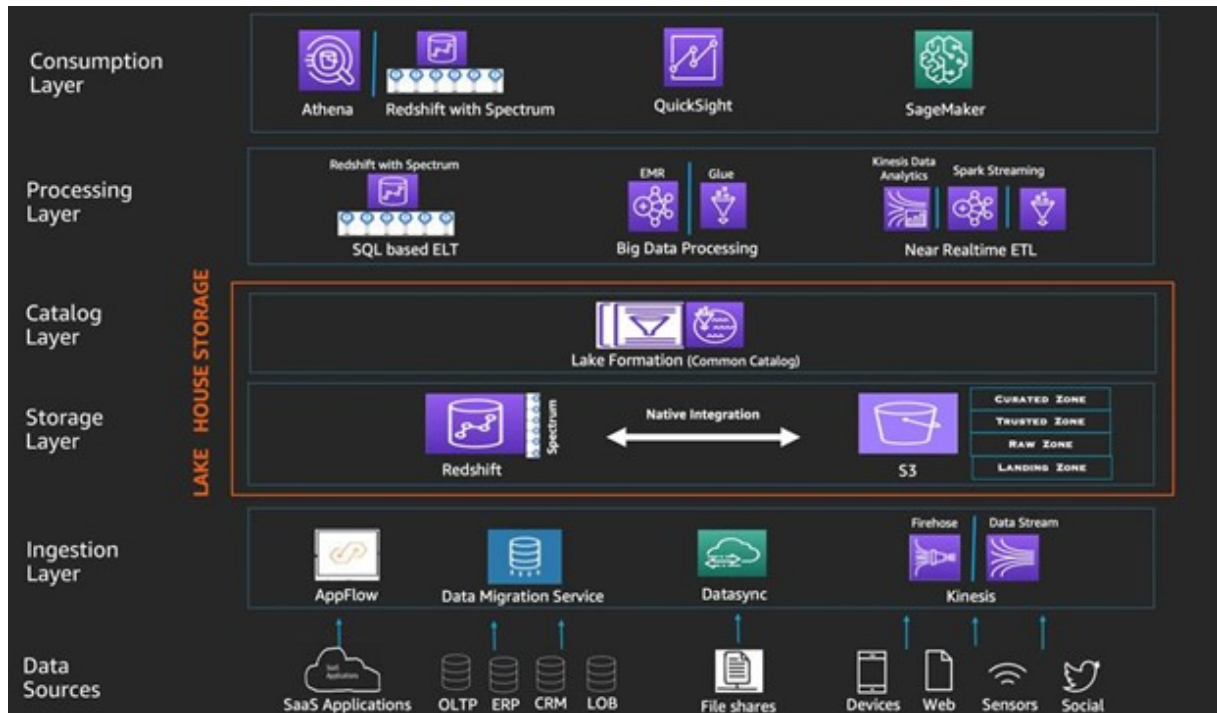


Figure 4: Proposed Data Lakehouse Architecture by AWS

- **SaaS Applications:** Amazon AppFlow ingests data from SaaS applications like Salesforce and Google Analytics into Amazon S3 or Redshift.
- **File Shares:** AWS DataSync transfers data from NAS devices to Amazon S3, handling large volumes and ensuring data integrity.
- **Streaming Data Sources:** Amazon Kinesis Data Firehose processes and transfers streaming data to Amazon S3 or Redshift.

Lake House Storage Layer

AWS Lake House integrates Amazon Redshift and Amazon S3 for storage, supporting both structured data in Redshift and diverse data types in S3. The key components include:

- **Amazon Redshift Spectrum:** Allows Redshift to query data stored in S3 using a unified SQL interface, supporting large-scale analytics without data movement.
- **Storage Zones in S3:** Data is organized into landing, raw, trusted, and curated zones to efficiently manage the lifecycle and stages of data processing.

Data Processing Layer

The processing layer includes various AWS services designed for different data processing needs:

- **SQL-based ELT:** Amazon Redshift and Redshift Spectrum support robust ELT pipelines for transforming structured data using SQL.

- **Big Data Processing:** AWS Glue and Amazon EMR provide large-scale ETL capabilities and big data processing using Apache Spark.
- **Near-real-time ETL:** Services like Kinesis Data Analytics and Spark Streaming (on AWS Glue or EMR) enable near-real-time data processing.

Data Consumption Layer

The data consumption layer democratizes data access across the organization through:

- **Interactive SQL:** Amazon Redshift (with Spectrum) and Athena allow interactive SQL queries on data in S3 and Redshift.
- **Machine Learning:** Amazon SageMaker provides tools for data exploration, feature engineering, model training, and deployment integrated with the Lake House storage.
- **Business Intelligence:** Amazon QuickSight offers serverless BI capabilities to create interactive dashboards and integrates with ML models from SageMaker.

Central Data Catalog and Metadata Management

AWS Lake Formation serves as the central catalog for metadata management across the entire data lake and data warehouse, supporting schema-on-read, granular permissions, and integration with various components of the processing and consumption layers.

3 The Need for On-Premises Data Lakehouse

We acknowledge the tremendous advantages and utilities provided by Public Cloud platforms. However, building an on-premises Data Lakehouse architecture not only optimizes data management efficiency but also ensures compliance with data security and privacy regulations. Specifically, this helps businesses meet strict requirements for the protection of citizen data in various countries, such as GDPR in Europe, HIPAA in the healthcare sector in the United States, as well as cybersecurity regulations in Vietnam and China.

At the same time, the concept of a Data Lakehouse is still relatively new in recent years, with no specific standards and limited research applied to specific fields. For instance, in the financial sector, development processes are still in the research phase. Particularly in Vietnam, where businesses are still struggling to transition from data warehouses to data lakes, the Data Lakehouse concept remains fairly novel.

Therefore, researching and proposing an on-premises Data Lakehouse architecture for businesses that require high security to ensure both efficiency and safety, without relying on external parties, is a highly urgent and practical need that many companies are currently concerned with and carefully considering.

4 Proposed On-Premises Data Lakehouse Architecture

To build a fully functional data lakehouse pipeline as mentioned above, the author proposes a solution with the following key components and specific technologies:

4.1 Proposed Architecture

- **Ingestion Layer:** The layer for ingesting data from various data sources.
- **Storage Layer:** The layer for storing data with support for structured, semi-structured, and unstructured data.
- **Catalog Layer:** This layer manages metadata and schema enforcement for structured and semi-structured data, ensuring data consistency and discoverability across the lakehouse.
- **Processing Layer:** The layer for processing data for ETL tasks, data analysis, and machine learning.
- **Serving Layer:** The layer for serving data to BI tools, reporting, and analytics.
- **Monitoring and Management Layer:** The layer for monitoring and managing to ensure the system operates stably.

4.2 Technologies Used

- **Ingestion Layer:**
 - Apache Kafka: This technology is used to ingest real-time data from various sources such as log files, databases, and API services. Kafka provides high fault tolerance and reliability in data transmission.
 - The data ingestion process involves collecting and importing data into the data lakehouse. This process can be automated or done manually using custom scripts or command-line tools, such as CRM applications, relational databases, and NoSQL databases.
- **Storage Layer:**
 - Hadoop HDFS (Hadoop Distributed File System): HDFS provides flexible and scalable storage, supporting structured, semi-structured, and unstructured data.
 - Delta Lake: Built on top of Apache Spark, Delta Lake provides storage with ACID support, time travel (the ability to revert to previous data states), and efficient metadata management.
- **Catalog Layer:**
 - Hive Metastore: The Hive Metastore is on top of Hadoop HDFS and Delta Lake Table. It ensures data consistency, facilitates data discovery, and integrates seamlessly with other components in the data lakehouse.

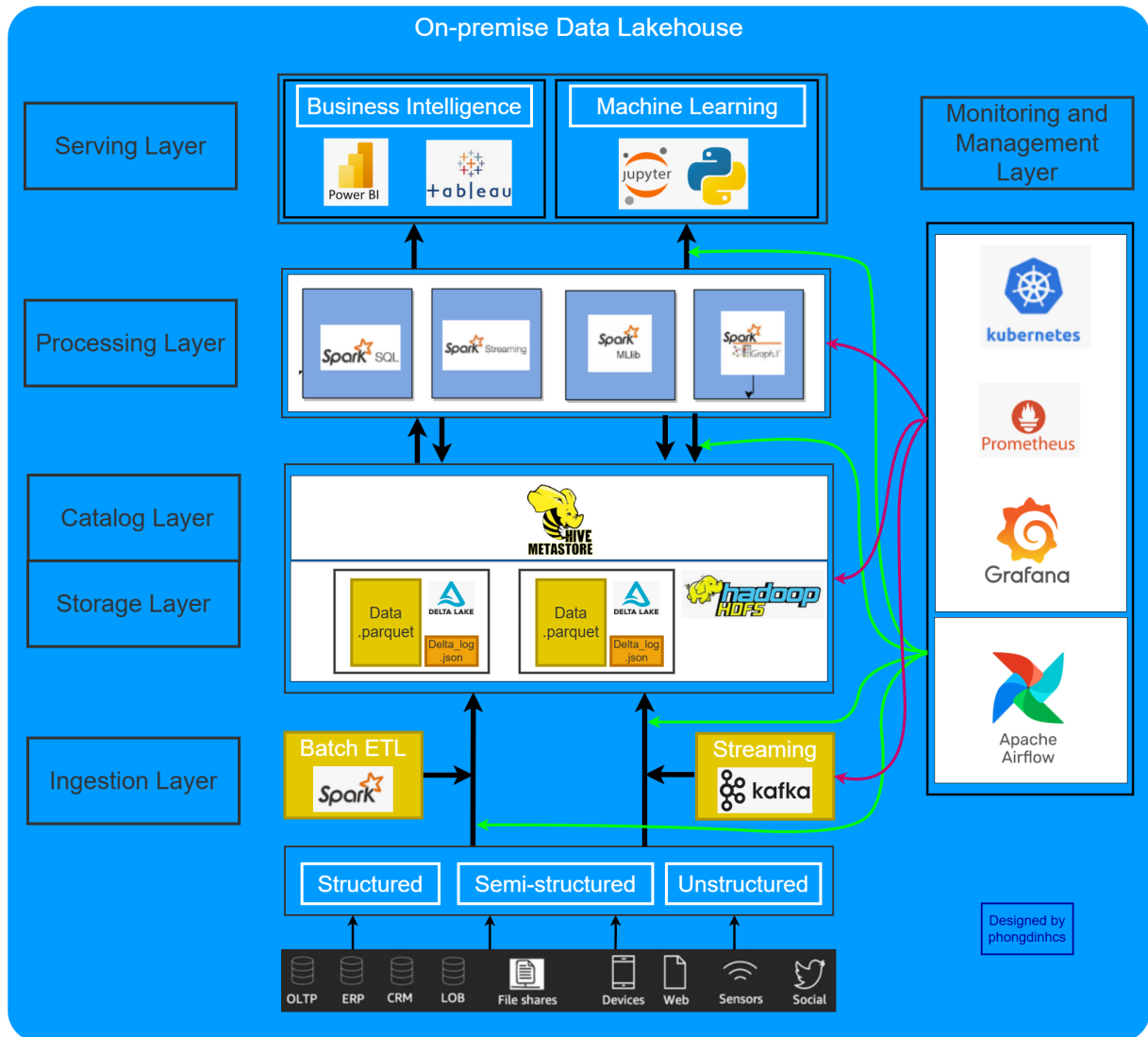


Figure 5: Proposed On-Premise Data Lakehouse Architecture

- **Processing Layer:**

- Apache Spark: Apache Spark is used to perform ETL (Extract, Transform, Load) tasks, data analysis, and machine learning applications. Spark offers powerful and efficient parallel data processing.

- **Serving Layer:**

- Power BI: Power BI allows access and visualization of data from the data lakehouse, providing easy-to-understand reports and dashboards for end users.
- Jupyter Notebook: This is an integrated environment for data scientists to explore, analyze data, and develop machine learning models.

- **Monitoring and Management Layer:**

- Kubernetes with Prometheus and Grafana: Prometheus and Grafana are deployed on Kubernetes to manage and deploy containers. Prometheus is used to monitor system performance metrics, while Grafana displays these metrics in charts.
- Apache Airflow: Used to manage and orchestrate ETL pipelines, ensuring tasks are executed in the correct order and at the specified times.

5 Demo PhongDinhCS-data_lakehouse has been built and published on GitHub

5.1 Example Scenario Applying Data Lakehouse Architecture

In this report, I present an example of real-time data updates in the stock finance domain. The model includes four information tables:

- Stock Symbol Table: This table is pre-built and directly stored as a Delta Table within HDFS.
- Current Stock Price Table: Updated with simulated data from a PostgreSQL source on an online server. PySpark updates this table every 10 seconds and writes the data to a Delta Table in HDFS.
- Transaction Table: PySpark records updates to the Current Stock Price Table, capturing historical stock prices.
- Market Information Update Table Based on Electronic News Reports: This table uses a Python script to crawl data from online news articles. A Kafka producer sends this data to a Kafka topic, and a Kafka consumer retrieves the data and writes updates to an HTML store containing two pieces of information: timestamp and the raw HTML. The purpose of crawling this data in its raw form is to assess early risk if a stock is frequently mentioned in the news or to store these HTML records for future machine learning or deep learning purposes.

5.2 GitHub Repository PhongDinhCS and Docker-compose

The demo architecture has been built and is available on GitHub:

https://github.com/PhongDinhCS/PhongDinhCS-data_lakehouse

In this section, I will describe the functionality of the complete-version-1.0 of the data lakehouse architecture, available in the following GitHub repository:

https://github.com/PhongDinhCS/PhongDinhCS-data_lakehouse/tree/complete-version-1.0

In the 'complete-version-1.0' repository, I have set up several containers, including Kafka, Hadoop, Hive Metastore, PostgreSQL, and Delta Spark, which have been successfully tested with Flask API, Jupyter Notebook, and Power BI. All containers were built from official images on Docker Hub.

Image versions:

- apache/kafka:3.8.0
- apache/hadoop:3
- postgres:11
- apache/hive:3.1.3
- deltaio/delta-docker:0.8.1_2.3.0

The docker-compose.yml file includes:

```
1 version: '3'
2
3 services:
4   kafka:
5     container_name: kafka
6     restart: unless-stopped
7     image: apache/kafka:3.8.0
8     ports:
9       - 9092:9092
10    volumes:
11      - ./kafka/config:/mnt/shared/config
12      - kafka-bin:/opt/kafka
13    networks:
14      lakehouse:
15        ipv4_address: 172.18.0.99
16
17    namenode:
18      container_name: namenode
19      restart: unless-stopped
20      image: apache/hadoop:3
21      hostname: namenode
22      command: ["hdfs", "namenode"]
23      ports:
24        - 9870:9870
25        - 8020:8020
26        - 9000:9000
27      env_file:
28        - ./hadoop/config
29      volumes:
30        - ./hadoop/dfs/name:/hadoop/dfs/name
31        - hadoop-bin:/opt/hadoop
32      networks:
33        lakehouse:
34          ipv4_address: 172.18.0.98
35
36    datanode1:
```

```
37     container_name: datanode1
38     restart: unless-stopped
39     image: apache/hadoop:3
40     command: ["hdfs", "datanode"]
41     env_file:
42     - ./hadoop/config
43     volumes:
44     - ./hadoop/dfs/data1:/hadoop/dfs/data
45     networks:
46     lakehouse:
47     ipv4_address: 172.18.0.97
48
49 datanode2:
50     container_name: datanode2
51     restart: unless-stopped
52     image: apache/hadoop:3
53     command: ["hdfs", "datanode"]
54     env_file:
55     - ./hadoop/config
56     volumes:
57     - ./hadoop/dfs/data2:/hadoop/dfs/data
58     networks:
59     lakehouse:
60     ipv4_address: 172.18.0.96
61
62 resourcemanager:
63     container_name: resourcemanager
64     restart: unless-stopped
65     image: apache/hadoop:3
66     hostname: resourcemanager
67     command: ["yarn", "resourcemanager"]
68     ports:
69     - 8088:8088
70     env_file:
71     - ./hadoop/config
72     networks:
73     lakehouse:
74     ipv4_address: 172.18.0.94
75     depends_on:
76     - namenode
77     - datanode1
78     - datanode2
79
80 nodemanager:
81     container_name: nodemanager
82     restart: unless-stopped
83     image: apache/hadoop:3
84     command: ["yarn", "nodemanager"]
85     env_file:
```

```
86     - ./hadoop/config
87 networks:
88     lakehouse:
89         ipv4_address: 172.18.0.93
90     depends_on:
91         - resourcemanager
92
93 postgres:
94     image: postgres:11
95     restart: unless-stopped
96     container_name: postgres
97     hostname: postgres
98     environment:
99         POSTGRES_DB: 'metastore_db'
100        POSTGRES_USER: 'hadoop'
101        POSTGRES_PASSWORD: 'hadoop'
102     ports:
103         - '5432:5432'
104     volumes:
105         - postgresql:/var/lib/postgresql
106     networks:
107         lakehouse:
108             ipv4_address: 172.18.0.91
109
110 metastore:
111     image: apache/hive:3.1.3
112     depends_on:
113         - postgres
114     restart: unless-stopped
115     container_name: metastore
116     hostname: metastore
117     environment:
118         HIVE_CONF_DIR: /opt/hive/conf
119         DB_DRIVER: postgres
120         IS_RESUME: 'true'
121         SERVICE_NAME: 'metastore'
122     ports:
123         - '9083:9083'
124     volumes:
125         - ./hive/conf/hive-site.xml:/opt/hive/conf/hive-site.xml
126         - ./hadoop/core-site.xml:/opt/hive/conf/core-site.xml
127         - ./hadoop/hdfs-site.xml:/opt/hive/conf/hdfs-site.xml
128         - ./hadoop/yarn-site.xml:/opt/hive/conf/yarn-site.xml
129         - type: bind
130           source: ./hive/postgresql-42.2.27.jre6.jar
131           target: /opt/hive/lib/postgres.jar
132     networks:
133         lakehouse:
134             ipv4_address: 172.18.0.90
```

```
135
136 hiveserver2:
137     image: apache/hive:3.1.3
138     depends_on:
139         - metastore
140     restart: unless-stopped
141     container_name: hiveserver2
142     environment:
143         HIVE_SERVER2_THRIFT_PORT: 10000
144         IS_RESUME: 'true'
145         SERVICE_NAME: 'hiveserver2'
146         HIVE_CONF_DIR: /opt/hive/conf
147     ports:
148         - '10000:10000'
149         - '10002:10002'
150     volumes:
151         - ./hive/conf/hive-site.xml:/opt/hive/conf/hive-site.xml
152         - ./hadoop/core-site.xml:/opt/hive/conf/core-site.xml
153         - ./hadoop/hdfs-site.xml:/opt/hive/conf/hdfs-site.xml
154         - ./hadoop/yarn-site.xml:/opt/hive/conf/yarn-site.xml
155     networks:
156         lakehouse:
157             ipv4_address: 172.18.0.89
158
159 delta-spark:
160     container_name: delta-spark
161     restart: unless-stopped
162     image: deltaio/delta-docker:0.8.1_2.3.0
163     entrypoint:
164         - bash
165         - -c
166         - |
167             if ! hdfs dfs -test -d hdfs://namenode:8020/lakehouse; then
168                 hdfs dfs -mkdir hdfs://namenode:8020/lakehouse
169                 hdfs dfs -chmod 777 hdfs://namenode:8020/lakehouse
170             fi
171             tail -f /dev/null
172     volumes:
173         - hadoop-bin:/opt/hadoop
174         - kafka-bin:/opt/kafka
175         - ./delta-spark/volume:/opt/spark/work-dir/volume
176         - ./hive/conf/hive-site.xml:/opt/spark/conf/hive-site.xml
177         - ./hadoop/core-site.xml:/opt/spark/conf/core-site.xml
178         - ./hadoop/hdfs-site.xml:/opt/spark/conf/hdfs-site.xml
179         - ./hadoop/yarn-site.xml:/opt/spark/conf/yarn-site.xml
180     environment:
181         - DELTA_PACKAGE_VERSION=delta-core_2.12:2.3.0
182         - SPARK_CONF_DIR=/opt/spark/conf
183         - HADOOP_HOME=/opt/hadoop
```

```
184     - HADOOP_CONF_DIR=/opt/hadoop/etc/hadoop
185     - PATH=/opt/spark/bin:/opt/hadoop/bin:/opt/kafka/bin:/home/NBuser/.local/bin:/home/
NBuser/.local/lib/python3.9/site-packages:$PATH
186     - HADOOP_USER_NAME=hadoop
187     - SPARK_MODE=master
188     - SPARK_SQL_CATALOG_IMPLEMENTATION=hive
189     networks:
190       lakehouse:
191         ipv4_address: 172.18.0.92
192     ports:
193       - "8888:8888" # Jupyter Notebook
194       - "7077:7077" # Spark Master
195       - "4040:4040" # Spark Web UI
196       - "5000:5000" # Flask API
197     depends_on:
198       - namenode
199
200 networks:
201   lakehouse:
202     name: lakehouse
203     driver: bridge
204     ipam:
205       config:
206         - subnet: 172.18.0.0/16
207           gateway: 172.18.0.1
208
209 volumes:
210   hadoop-bin:
211   kafka-bin:
212   postgresql:
```

5.3 Components and Trial Running of PhongDinhCS-data_Lakehouse

The setup and execution of the PhongDinhCS-data_lakehouse containers on Docker Desktop running on Windows 10, using Windows Subsystem for Linux (WSL) for compatibility with the Linux-based containers. It is also can run on Linux Ubuntu 22.04 LTS within a VMware virtual machine. The environment enables the orchestration of multiple containers, such as Kafka, Hadoop, and Delta Spark.

The Hadoop information can be accessed through a web browser interface. It showcases the Hadoop Distributed File System (HDFS) and ResourceManager web UI, which are used to monitor the cluster's storage and job execution status.

The HiveServer2 information being checked via a web browser. It highlights the Hive Web UI, where users can manage and monitor Hive queries and configurations, demonstrating Hive's integration within the PhongDinhCS-data_lakehouse setup.

The Delta table is stored in Parquet format, accompanied by a delta_log file. Both are stored in the Hadoop Distributed File System (HDFS), showcasing Delta Lake's capabilities in managing data storage with ACID transaction properties.

Jupyter Notebook within the Delta-Spark container to test Kafka's functionality. The notebook environment allows interactive development and testing of Spark jobs, including Kafka data processing, within the PhongDinhCS-data_lakehouse framework.

Flask is used within the Delta-Spark container to implement a Web API. The API facilitates easy data sharing by exposing data stored in Delta tables through a RESTful service, which can be consumed by various client applications.

Power BI is used to retrieve data from the Web API and generate reports. It highlights the integration of Power BI with the PhongDinhCS-data_lakehouse project, showcasing its ability to visualize and analyze data sourced from Delta Lake tables via the Flask API.

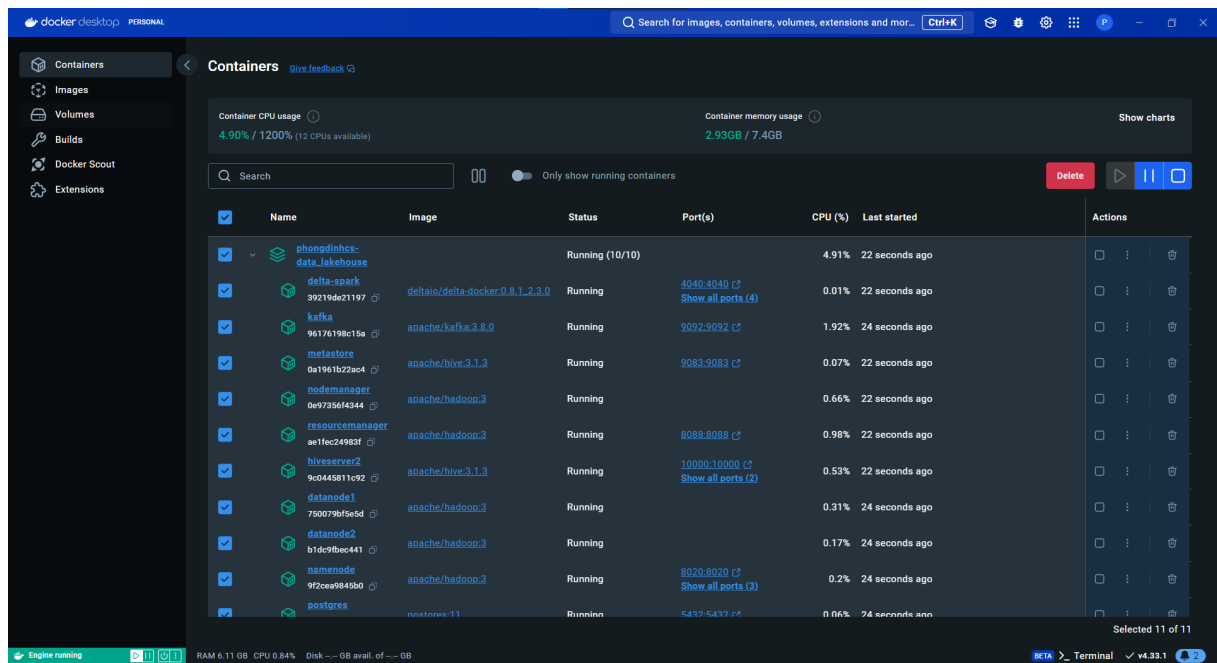


Figure 6: Running PhongDinhCS-data_lakehouse containers on Docker Desktop Windows 10 with WSL

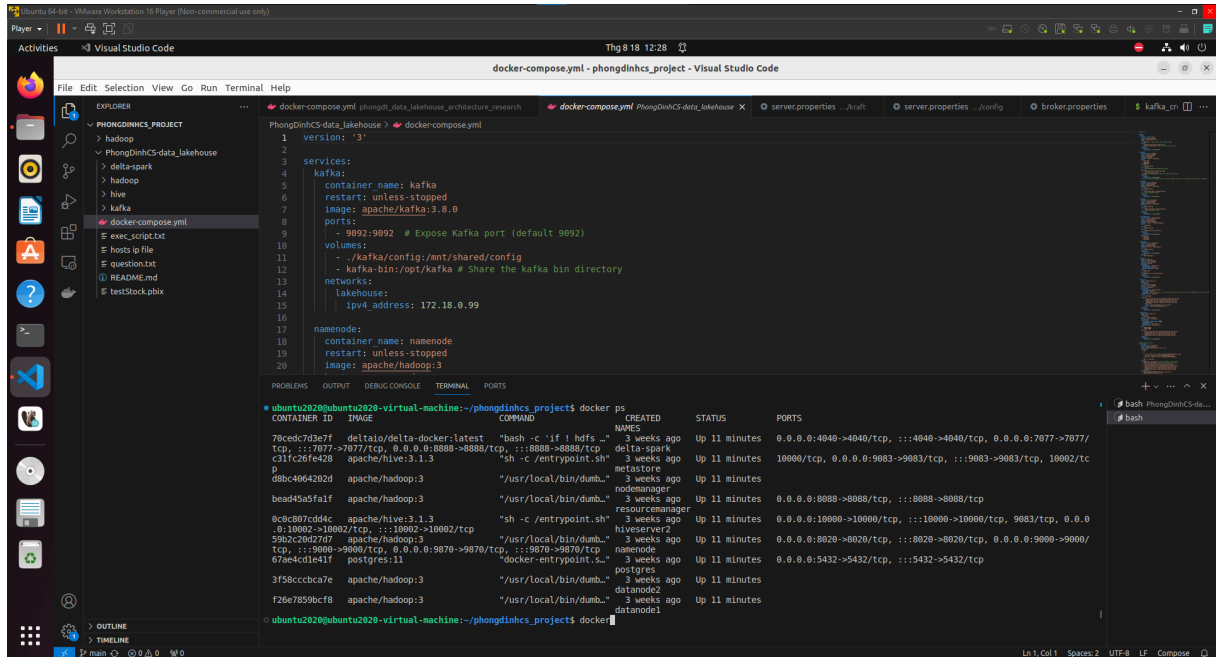


Figure 7: Running PhongDinhCS-data_lakehouse containers on Linux Ubuntu 22.04 LTS in VMware

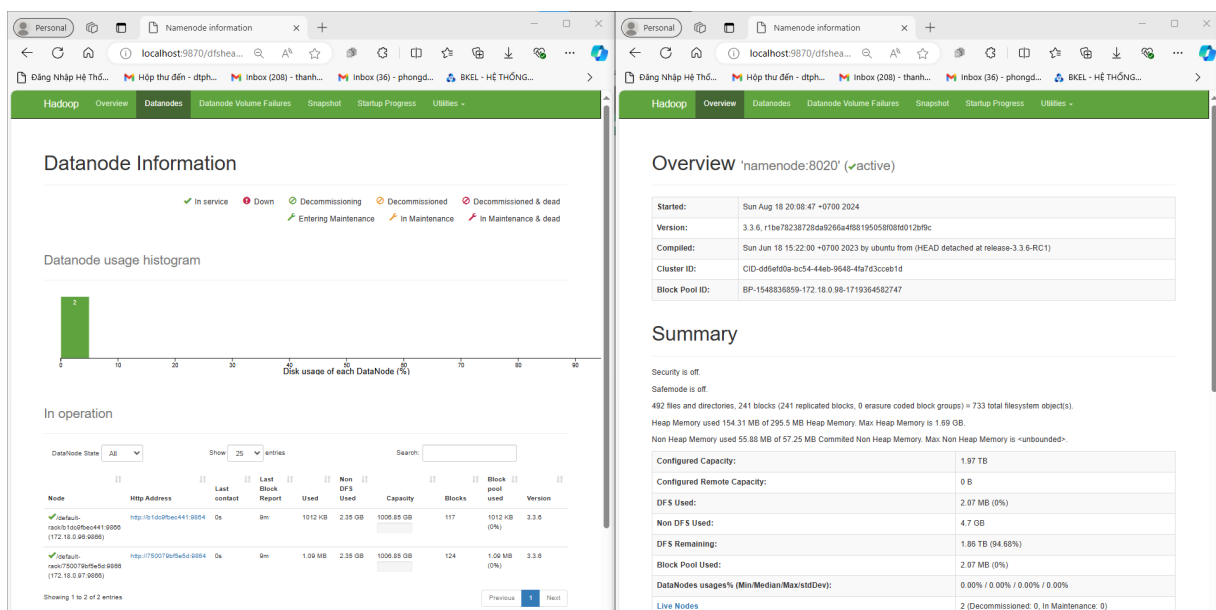


Figure 8: Checking Hadoop information from browser

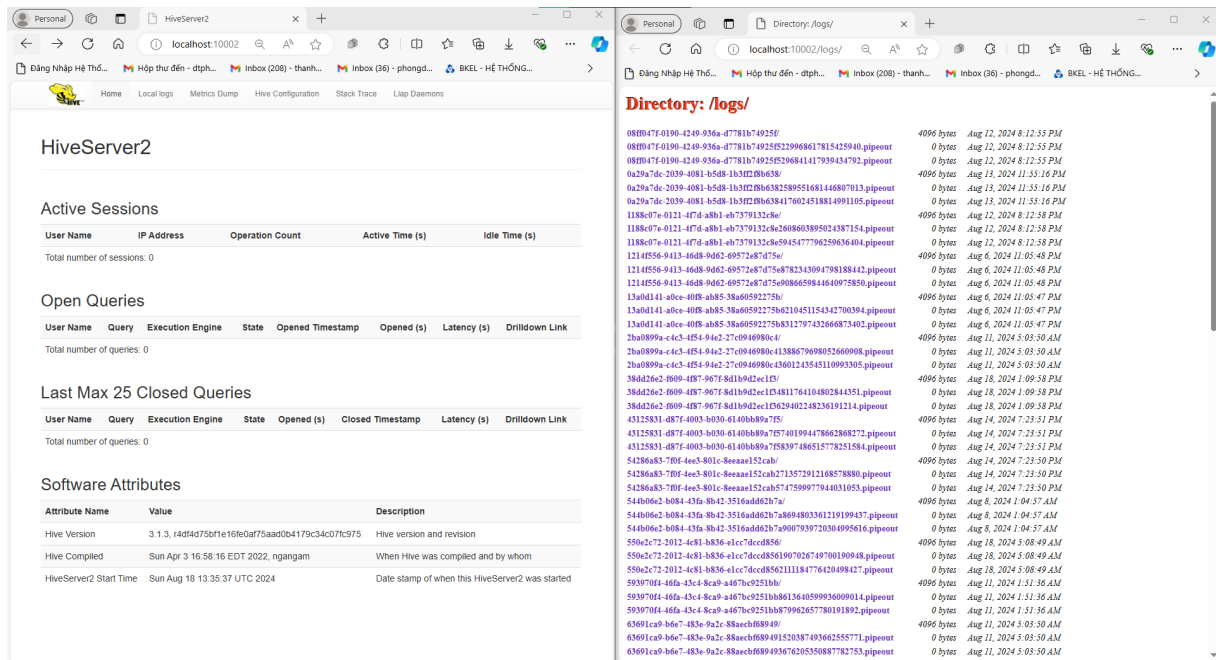


Figure 9: Checking HiveServer2 information from browser

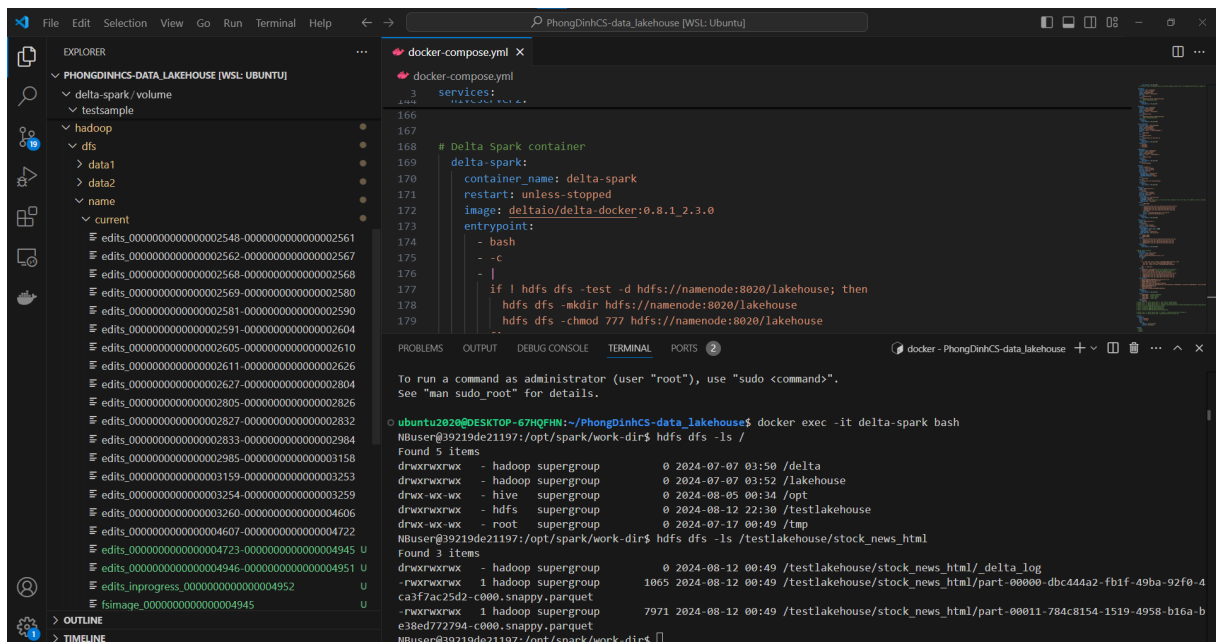


Figure 10: Data in the Delta table is stored in Parquet format along with the delta_log file and stored in HDFS

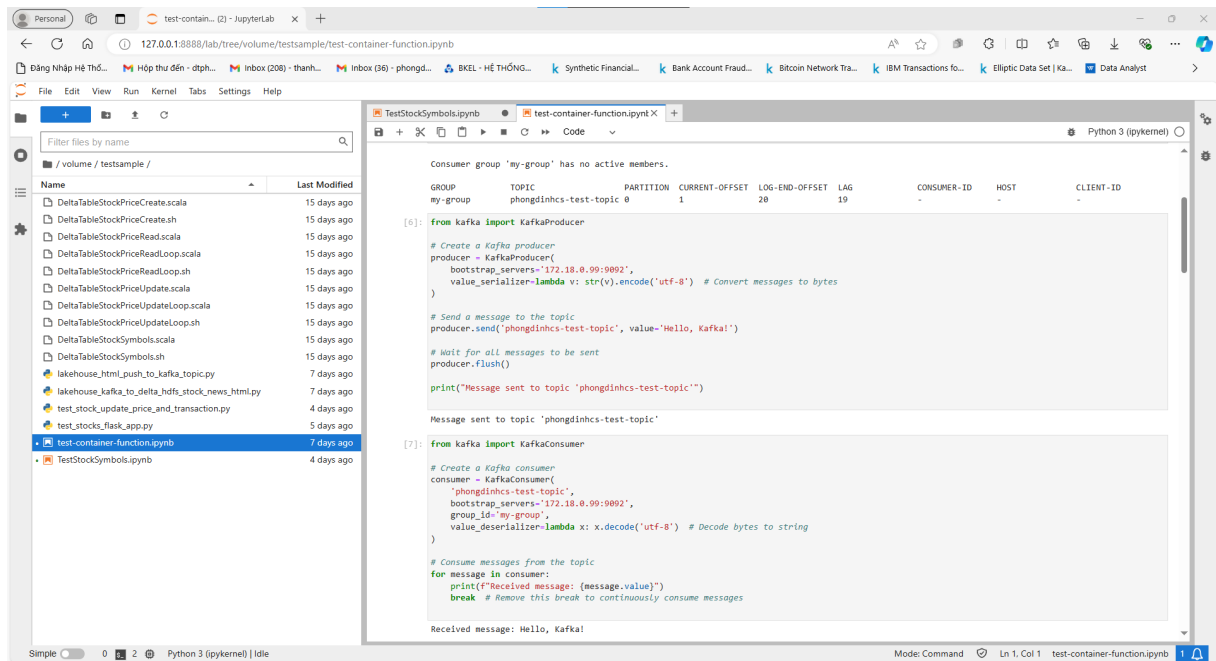


Figure 11: Running Jupyter Notebook from the Delta-Spark container to test Kafka functionality

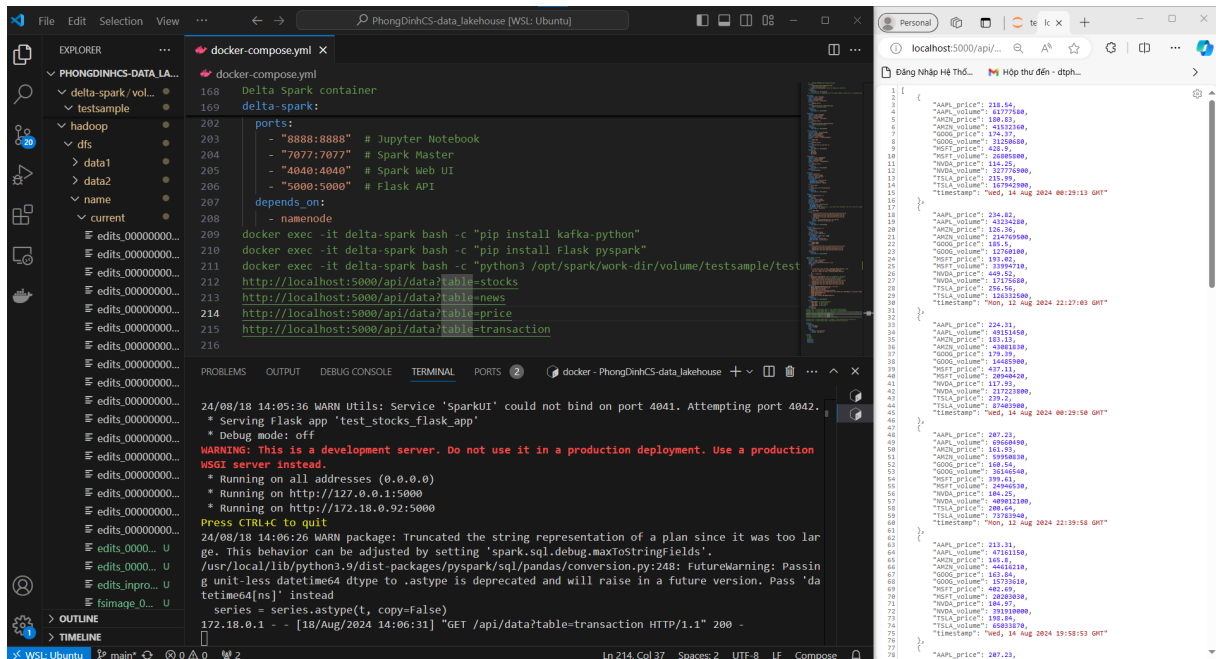


Figure 12: Using Flask in the Delta-Spark container to implement a Web API for easy data sharing

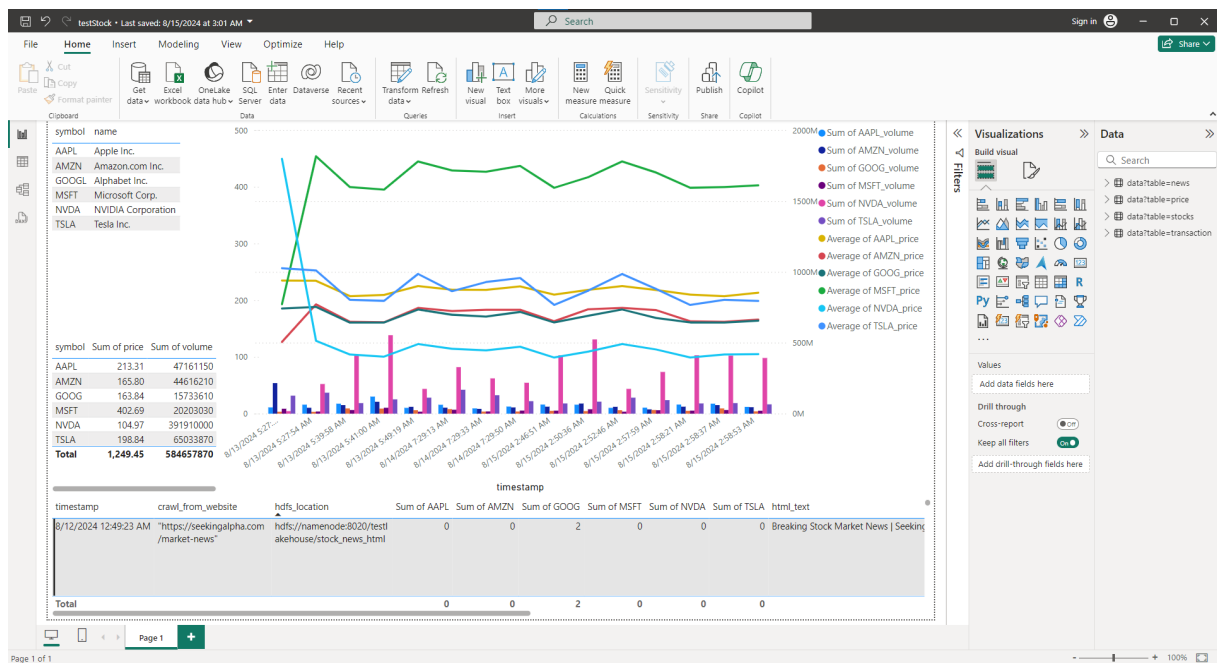


Figure 13: Using Power BI to retrieve data from the Web API and generate reports

6 Next Steps for Completing the On-Premises Data Lakehouse Architecture

In the current complete-version-1.0 repository, I have completed the setup and operation of the ingestion layer, storage layer, catalog layer, processing layer, and serving layer on the Docker and Docker Compose platform.

The next goal is to build a Monitoring and Management Layer and finalize the architecture on Kubernetes with Prometheus and Grafana to facilitate system activity control and resource monitoring.

Additionally, it is necessary to use Apache Airflow to manage automated pipelines, ensuring accurate and stable operational processes.

Finally, the On-Premises Data Lakehouse Architecture also needs a solution to integrate with a Public Cloud Data Lakehouse, such as Databricks on Microsoft Azure, to create a hybrid platform. This will ensure that the data lakehouse maintains security for personal information while providing cost-effective solutions for less sensitive data.