

Introduction to Apache Hive



Transactional Systems vs. Analytical Systems

Feature	Transactional Systems	Analytical Systems
Type of Data Handled	Day-to-day transactional data	Historical data
Operations Performed	Insert, Delete, Update	Majorly Read operations to analyze large volumes of data
Example	ATM transactions, e-commerce transactions	Analyzing data of a sales campaign
Best Suited Systems	RDBMS (Databases: Ex - Oracle, MySQL, etc.), Monolithic Systems	Data Warehouses (Ex - Teradata, etc.), Distributed Systems

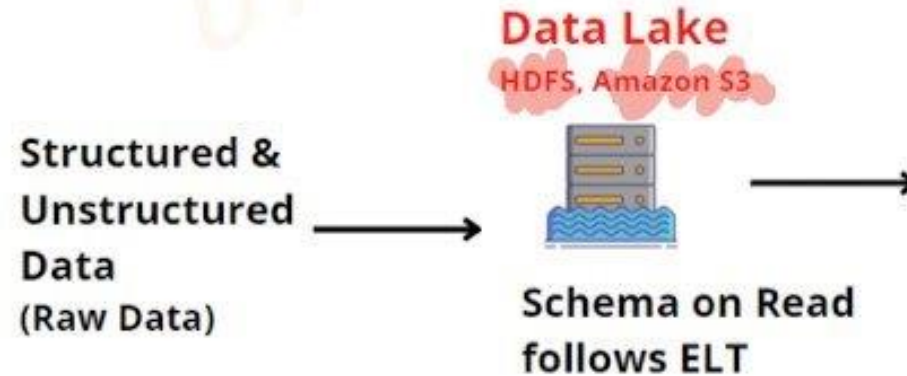
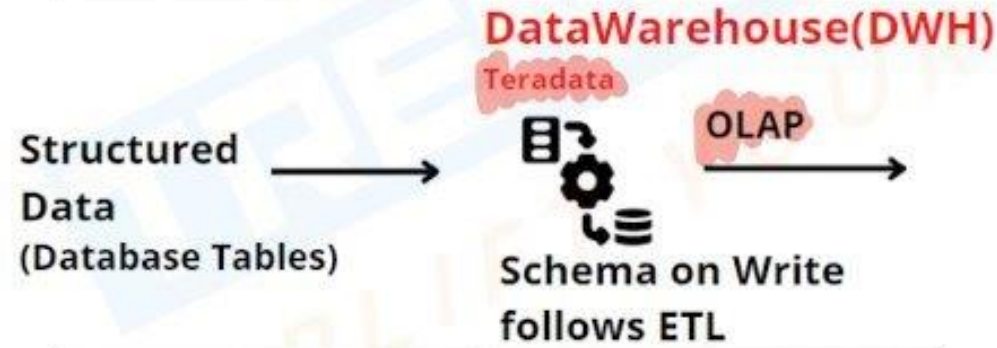
Key Takeaway:

- Transactional systems capture and process day 2 day transactions, while
- Analytical systems focus on analyzing historical data to gain insights.

Database Vs Data Warehouse Vs Data Lake

Feature	Database	Data Warehouse (DWH)	Data Lake
Purpose	OLTP (Transactional Processing)	OLAP (Analytical Processing)	Insights from large volumes of data
Data Structure	Structured (Rows/ Columns)	Structured	Raw (Structured & Unstructured)
Data Scope	Recent data for performance	Historical data	Both recent & historical data
Examples	Oracle, MySQL	Teradata	HDFS, Amazon S3
Schema Approach	Schema on Write	Schema on Write	Schema on Read
Process	-	ETL (Extract, Transform, Load)	ELT (Extract, Load, Transform)
Storage Cost	High	High, but less than Database	Cost-effective
Challenges	-	Complex transformations, rigidity	Flexibility, but requires management

Data



Reports



Machine
Learning

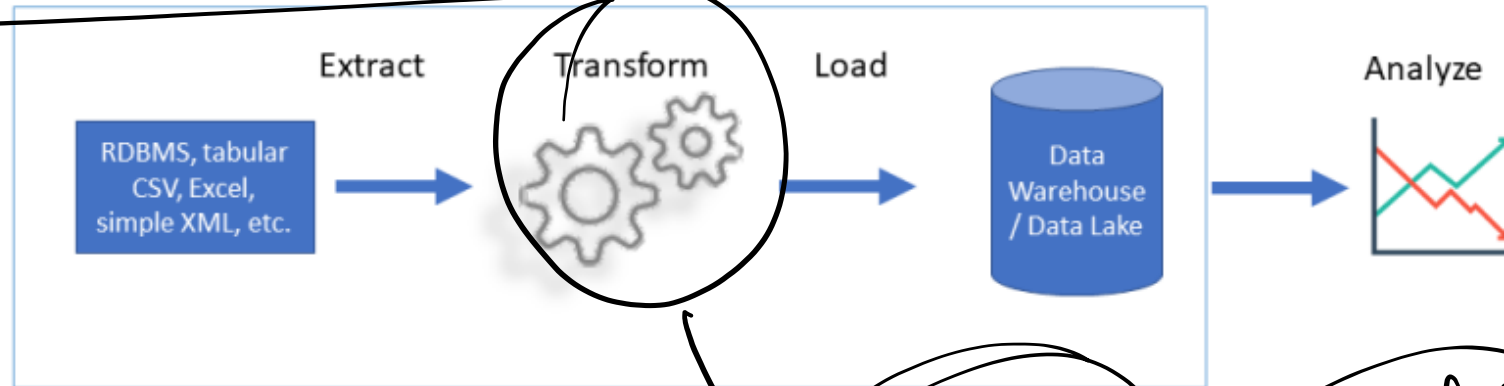
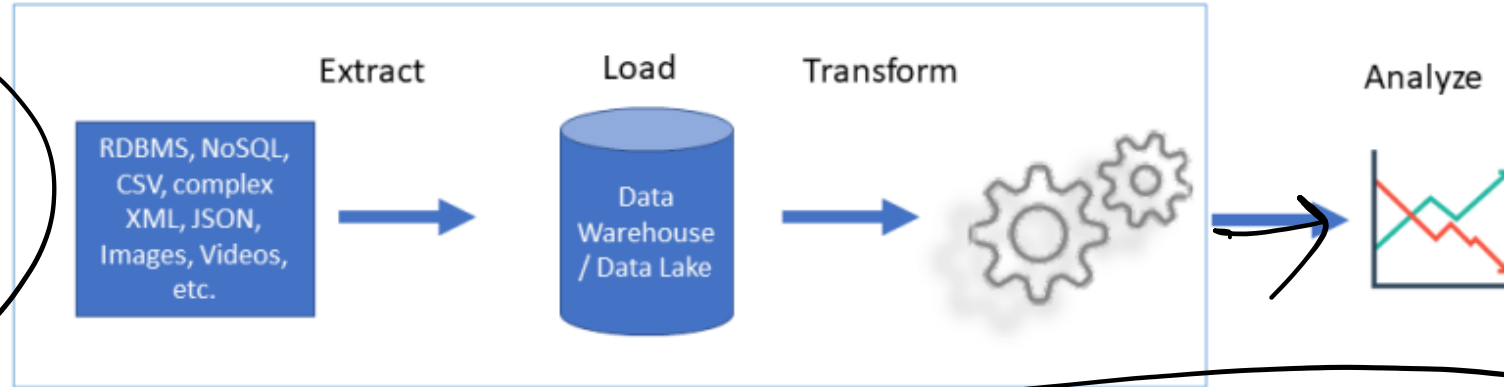


Data
Science



BI

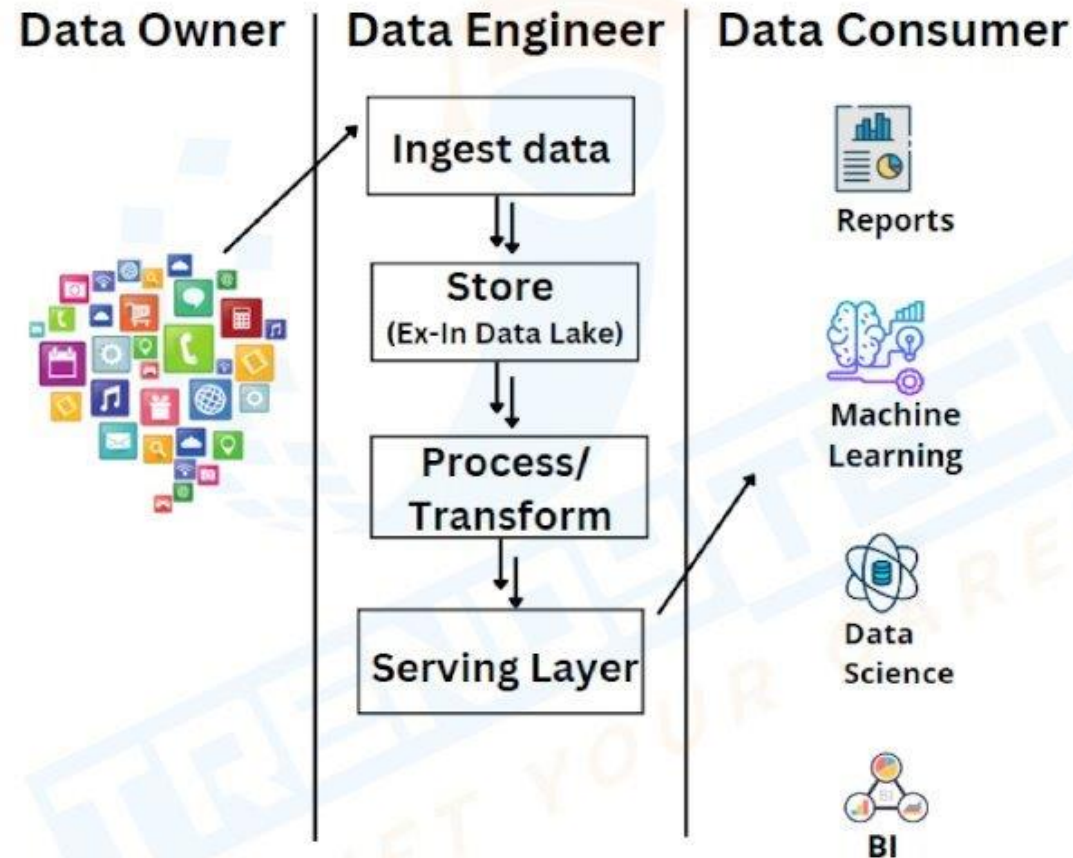
ELT vs ETL



processing

MapReduce

Role of Data Engineers



History

- Facebook started generating large volumes of data at a very high pace
- These data had to be analyzed and reports had to be generated to get some meaningful inferences.
- Used a Datawarehouse on a commercial RDBMS.
- **Problems:**
 - Handling large volumes
 - scalability

- **Problems with using hadoop:**

- Required to write map reduce program
- SQL developers who were earlier getting the reports had to be reskilled/upskilled in writing java/python programs
- Writing map reduce programs took time

History

- **Solution**

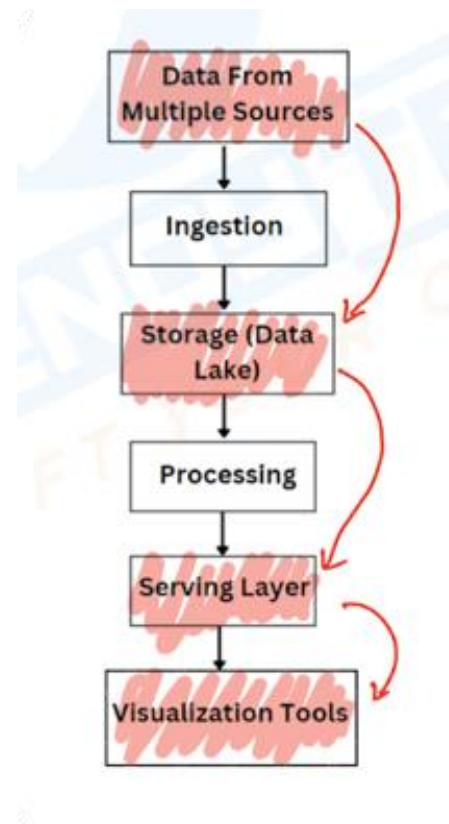
- Started using hadoop
- Data retrieval Jobs gave faster results
- Face book decided to improve the capabilities of hadoop, so that their developers were allowed to use SQL on top of hadoop
- Hive was born in 2007 and was open sourced in 2008

- It is a **data warehouse** in Hadoop with fault tolerant
- Facilitates **querying and managing** Massive petabytes of data residing on distributed storage
- Provides an SQL like interface called **HiveQL** which translates a query into java map reduce program and runs the same on a Hadoop cluster
- Developed by **Facebook** and open sourced to **apache**

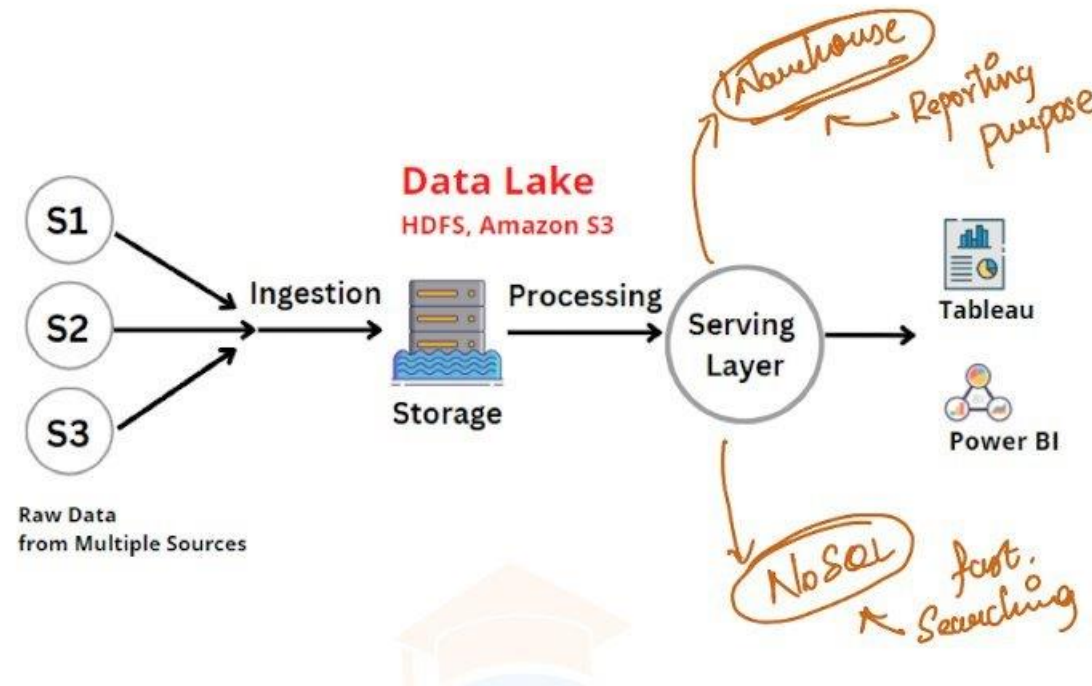
What is Hive?

Data engineering Flow

- **Data Collection:** Gather data from various sources into a **Data Lake** for centralized storage.
- **Data Ingestion:** Use an **Ingestion Framework** to move data from different sources into the Data Lake.
- **Data Processing:** Follow the **ELT process**—load data into the Data Lake, then transform it (e.g., cleaning, aggregation, joins) as needed.
- **Data Serving:** Store processed data in the **Serving Layer** for visualization tools like Tableau and Power BI to display results graphically.



Data engineering Flow



Data engineering Flow

ELT

Phases / Stages of Data Pipeline	Example Data Pipeline Workflow tools and technologies for On-Premise and Cloud			
		On-Premise (Hadoop)	Azure Cloud	AWS Cloud
	Source	MySQL Database table	Multiple sources	Multiple sources
	Ingestion	Sqoop	ADF	AWS GLUE
	Storage	HDFS	ADLS Gen2	Amazon S3
	Processing	MapReduce / SPARK	Azure Databricks / Synapse	Athena / Redshift
	Serving layer	HBase	Azure SQL / Cosmos DB	AWS RDS / Dynamo DB

Why Hive?

- Data Transformation and ETL
- Data Warehouse Features
- Scalability
- SQL-Like Interface
- Batch Processing
- Schema on Read
- Support for Various Data Formats

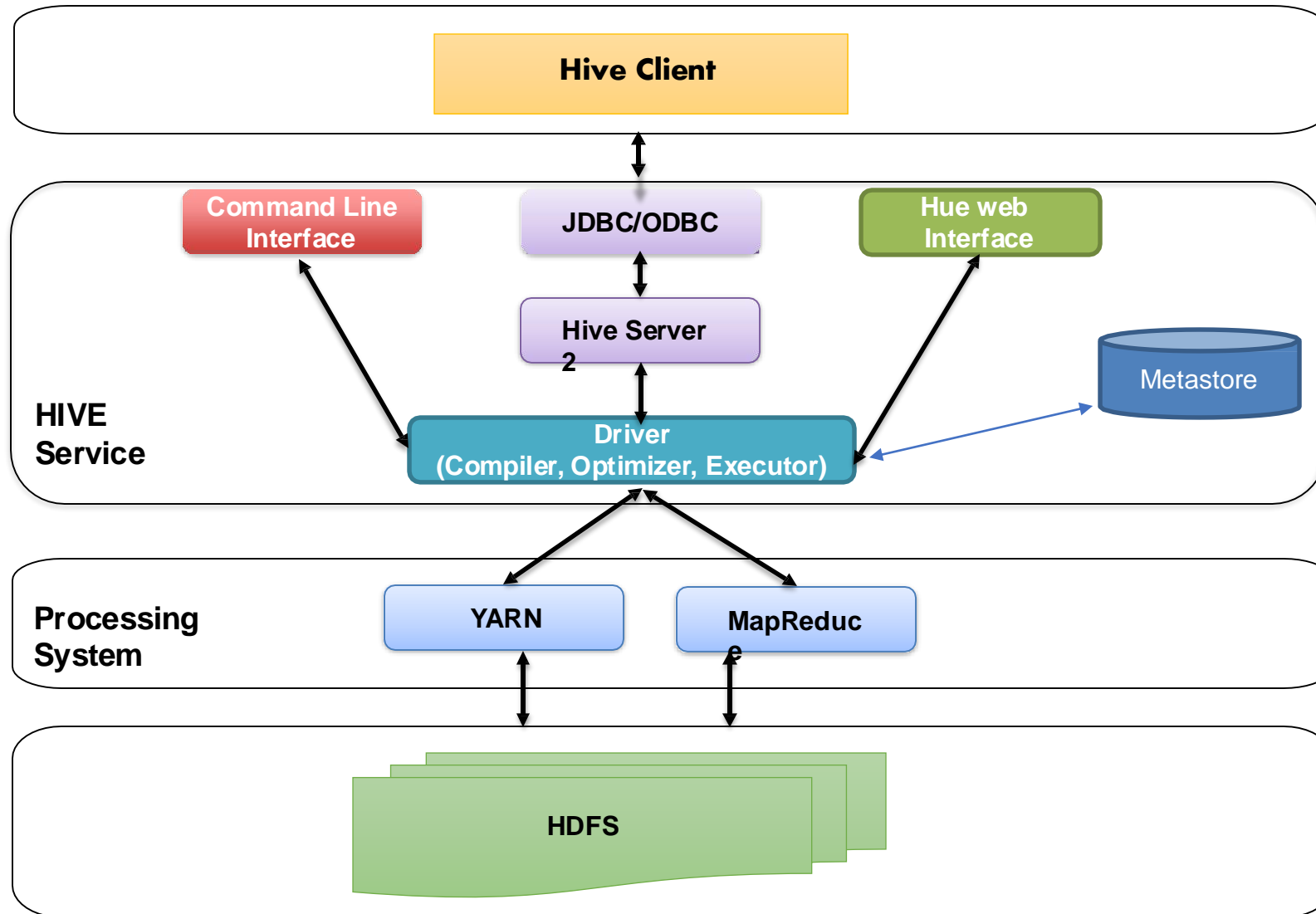
Is hive a database?

- **Answer** : No. Many people consider hive as a database management system. But, the truth is different.
- Apache **hive itself is not a database**.
- Consider hive as **logical view** of the underlying data in HDFS.
- It **cannot store any data** of its own. It always uses **HDFS for storing** the processed data.
- The only thing it can do is **enforcing the structure** in which the data can be stored in HDFS.

Hive - Drawbacks

- Is Not an OLTP system
- Is Not a full fledged database – data is stored in **HDFS** as flat files.
- Does not support all the SQL queries especially complicated ones
- Has high latency – since it converts queries to Java MapReduce code and executes

Hive Architecture



Advantages of Hive Against Map Reduce

- 100 lines of code can be easily solved using 2-3 lines of SQL code
- Best suited for SQL developers, who are novice to Java/python programming

```

1 package org.myorg;
2
3 import java.io.IOException;
4 import java.util.*;
5
6 import org.apache.hadoop.fs.Path;
7 import org.apache.hadoop.conf.*;
8 import org.apache.hadoop.io.*;
9 import org.apache.hadoop.mapreduce.*;
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
11 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
14
15 public class WordCount {
16
17     public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
18         private final static IntWritable one = new IntWritable(1);
19         private Text word = new Text();
20
21         public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
22             String line = value.toString();
23             StringTokenizer tokenizer = new StringTokenizer(line);
24             while (tokenizer.hasMoreTokens()) {
25                 word.set(tokenizer.nextToken());
26                 context.write(word, one);
27             }
28         }
29     }
30
31     public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {
32
33         public void reduce(Text key, Iterable<IntWritable> values, Context context)
34             throws IOException, InterruptedException {
35             int sum = 0;
36             for (IntWritable val : values) {
37                 sum += val.get();
38             }
39             context.write(key, new IntWritable(sum));
40         }
41     }
42
43     public static void main(String[] args) throws Exception {
44         Configuration conf = new Configuration();
45
46         Job job = new Job(conf, "wordcount");
47
48         job.setOutputKeyClass(Text.class);
49         job.setOutputValueClass(IntWritable.class);
50
51         job.setMapperClass(Map.class);
52         job.setReducerClass(Reduce.class);
53
54         job.setInputFormatClass(TextInputFormat.class);
55         job.setOutputFormatClass(TextOutputFormat.class);
56

```

CREATE TABLE docs (line STRING);

LOAD DATA INPATH 'docs' OVERWRITE INTO TABLE docs;

CREATE TABLE word_counts AS
 SELECT word, count(1) AS count FROM
 (SELECT explode(split(line, '\s')) AS word FROM docs) w
 GROUP BY word
 ORDER BY word;

Recap : Apache Hive

- Hive is NOT an RDBMS
- Hive is a SQL for HDFS
- No constraints !!!
- Hive SQL syntax is similar to MySQL's SQL syntax
- Hive is flexible
 - Data can come first, Schema can come later or vice versa
 - Create schema definitions only for the datasets we wish to query
 - Hence its Warehouse.
- Hive's data will be on HDFS, and the schema will be in a metastore DB
- Hive is Schema-On-Read → The structure is validated against the data when we issue a select query



Assignment

- **Based on the extra reads and research, write a LinkedIn Blog(100-words) on following**
 1. DE workflow(on-premises & cloud) and DE role
 2. Data lake vs DW vs DB

DEMO

Hive

Some Extra Reads

Useful info

GUIDELINES FOR HiveQL STATEMENTS

- Can be Executed by invoking a hive shell or webUI
- Not case sensitive
- Keywords cannot be abbreviated or split across lines
- Clauses (like SELECT, FROM, WHERE etc.) are usually placed on separate lines
- Indents are used to enhance readability
 - In the command line client it is mandatory to terminate each SQL statement end with a semicolon(;)