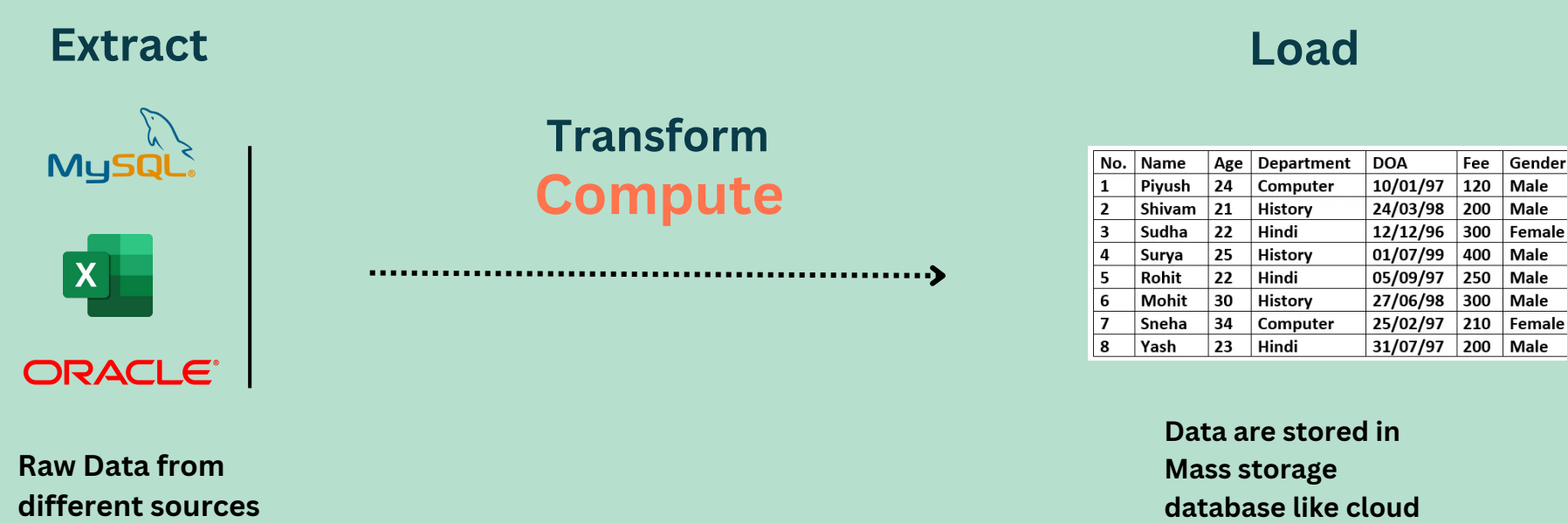# Setting Enviroment

## for Data Pipeline

To build a local ETL(Extract,Transform,Load) pipeline that reads data from:

-Mysql table(online sales)
-An CSV file (offline sales)

**Extract**

**Load**

**Transform**
**Compute**

| No. | Name | Age | Department | DOA | Fee | Gender |
|---|---|---|---|---|---|---|
| 1 | Piyush | 24 | Computer | 10/01/97 | 120 | Male |
| 2 | Shivam | 21 | History | 24/03/98 | 200 | Male |
| 3 | Sudha | 22 | Hindi | 12/12/96 | 300 | Female |
| 4 | Surya | 25 | History | 01/07/99 | 400 | Male |
| 5 | Rohit | 22 | Hindi | 05/09/97 | 250 | Male |
| 6 | Mohit | 30 | History | 27/06/98 | 300 | Male |
| 7 | Sneha | 34 | Computer | 25/02/97 | 210 | Female |
| 8 | Yash | 23 | Hindi | 31/07/97 | 200 | Male |

**Raw Data from different sources**

**Data are stored in Mass storage database like cloud**

# Tools Required for

## Data Pipeline Process (ETL)

### ➡ 1. APACHE SPARK

A distributed computing system designed for fast data processing, analytics, and machine learning, ideal for handling large-scale data.

**Prerequisites for Installing Apache Spark on Windows:**

● **Java Development Kit (JDK) – Install JDK and set the JAVA_HOME environment variable.**

● **Python – Install Python 3.x (Anaconda includes Python).**

● **Anaconda (Optional) – Recommended for managing Python dependencies.**

● **Hadoop winutils – Required for running Spark on Windows.**

● **Set Environment Variables – Add Spark, Hadoop, and Java paths to the system environment variables.**

## Installation Steps:

### Create Folders:

C:\spark\ (for Spark)
C:\hadoop\bin\ (for winutils)

### Download & Extract Spark:

Download Spark from
https://spark.apache.org/downloads.html (Pre-built for
Hadoop 3.3).

Extract and move it to C:\spark\.

### Download Winutils:

Download winutils.exe from GitHub.

Copy it to C:\hadoop\bin\.

### Set Environment Variables:

SPARK_HOME = C:\spark\spark-3.x.x-bin-hadoop3.3
HADOOP_HOME = C:\hadoop
JAVA_HOME = C:\Program Files\Java\jdk-XX.X.X
Add %SPARK_HOME%\bin and %HADOOP_HOME%\bin to Path.

### Verify Installation in cmd:

java -version,  python --version,  spark-shell,  pyspark

```
Microsoft Windows [Version 10.0.22631.5039]
(c) Microsoft Corporation. All rights reserved.

C:\Users\sarav>spark-shell
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
25/04/03 20:18:20 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform..
Spark context Web UI available at http://DHARSHINI:4040
Spark context available as 'sc' (master = local[*], app id = local-1743691702250).
Spark session available as 'spark'.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 3.5.5
      /_/

Using Scala version 2.12.18 (Java HotSpot(TM) 64-Bit Server VM, Java 17.0.12)
Type in expressions to have them evaluated.
Type :help for more information.
```

## ➡ 2. MySQL

MySQL is used in a data pipeline to store, manage, and retrieve structured data efficiently with SQL queries. It integrates with Spark via JDBC, enabling seamless data processing, transformation, and analytics.
.

# Installation Steps:

## Download & Install MySQL (MySQL Installer)

**Select MySQL Server & Workbench
→ Set root password → Complete installation.**

## Download MySQL JDBC Driver (Connector/J)

**MySQL Connector/J (JDBC) is used to establish a connection
between Apache Spark (or any Java-based application) and
MySQL, enabling seamless data transfer.**

**Extract & copy mysql-connector-java-9.1.0.jar
to C:\spark\jars\.**

## Load Data:

```sql
CREATE DATABASE sales_data;
USE sales_data;

CREATE TABLE online_sales ( order_id INT,
 customer VARCHAR(50),
 amount FLOAT,
order_date DATE );

INSERT INTO online_sales VALUES
(101, 'Alice', 300.50, '2024-03-01'),
 (102, 'Charlie', 90.00, '2024-03-03');
```

# ➡️3. pySpark code (MySql)

- MySQL Connector/J (JDBC) in PySpark enables seamless reading and writing of MySQL data using Spark DataFrames. It allows efficient SQL-based data processing, making it essential for ETL and analytics pipelines

- **Jupyter Notebook** is commonly used with PySpark for interactive development, debugging, and visualization. It allows executing PySpark code step by step, making it easier to analyze MySQL data in a local data pipeline.



```python
from pyspark.sql import SparkSession

# Initialize Spark Session with MySQL Connector/J 9.1.0
spark = SparkSession.builder \
    .appName("Test Spark") \
    .config("spark.jars", r"C:\mysql-connector\mysql-connector-java-9.1.0.jar") \
    .getOrCreate()  # ✅ No indentation error here

print(spark)  # Print Spark session info

# MySQL Connection Details
mysql_url = "jdbc:mysql://localhost:3306/sales_data"
mysql_properties = {
    "user": "root",
    "password": "root12345",
    "driver": "com.mysql.cj.jdbc.Driver"
}

# Read MySQL Table into PySpark DataFrame
table_name = "online_sales"
df = spark.read.jdbc(url=mysql_url, table=table_name, properties=mysql_properties)

# Show Data
df.show()
```

```
<pyspark.sql.session.SparkSession object at 0x000002185BB73560>
+--------+--------+------+----------+
|order_id|customer|amount|order_date|
+--------+--------+------+----------+
|     101|   Alice| 500.0|2024-01-01|
|     102|   priya| 600.0|2024-01-01|
+--------+--------+------+----------+
```

# ➡️ 4. pySpark code (CSV file)

## Loading a CSV File in PySpark
### Why Load CSV in PySpark?

**Handles Large Datasets –** PySpark efficiently processes large CSV files in a distributed manner.

**Schema Inference & Customization –** Supports automatic and manual schema definitions.

**Flexible Data Processing –** Allows filtering, transformation, and SQL-like operations on CSV data.