

Real-Time Accident Data Processing

Agenda

- Project Overview
- Architecture
- Batch Data
- ETL for batch data
- Real-Time Data
- Streaming Data from Kafka to Spark
- Processing Data with PySpark
- Analysis
- Huawei Cloud MRS

Project Overview

The project showcases a real-time and batch data pipeline using **Docker, PySpark, Kafka, MySQL, and Jupyter.**

Key components include:

- **Docker-Compose setup** for a PySpark Jupyter environment, Zookeeper, Kafka, and MySQL.
- Simulated streaming data source for real-time ingestion via **Kafka.**
- Batch data processed in a Jupyter notebook with **Pandas.**
- An ETL pipeline for transforming and loading data into **MySQL.**
- **Power BI** Dashboards for real-time data visualization.

The project's goal is to demonstrate building a scalable pipeline for processing both real-time and batch data.

Architecture



Kafka ingests real-time streaming data; Zookeeper manages its configuration.

KAFKA AND ZOOKEEPER



PySpark processes both the real-time streaming data from Kafka and batch data.

PYSPARK



The transformed data is loaded and stored in MySQL

MYSQL



Jupyter notebooks are Used for batch processing with Pandas and SQLAlchemy.

JUPYTER NOTEBOOKS



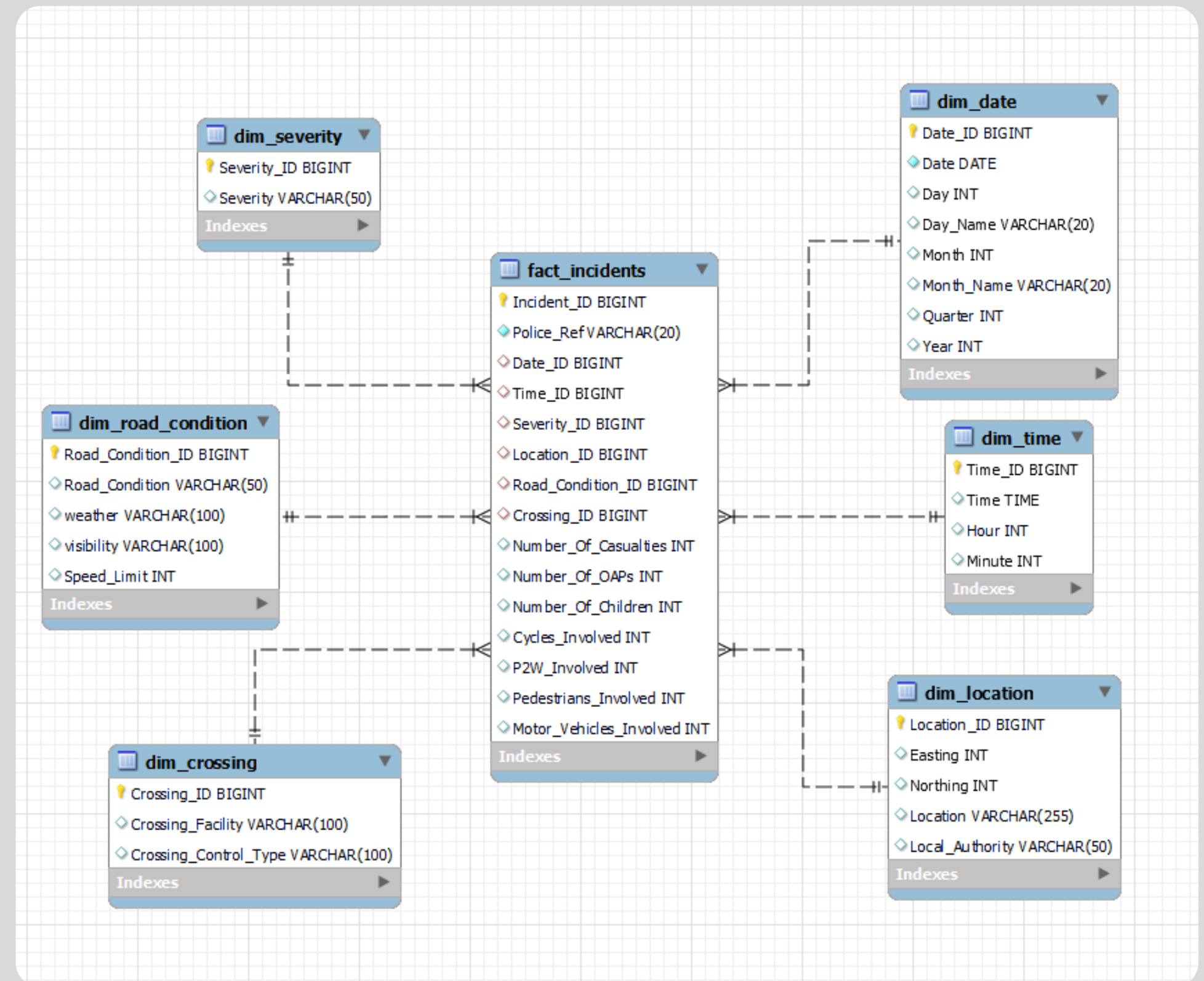
Visualizes data from MySQL, creating dashboards for insights into accident patterns and trends.

POWER BI

Batch Data

Cambridgeshire County Council

Data Warehouse Schema



Why MySQL ?



- MySQL is the most widely adopted **open-source** database
- MySQL's **fast write performance**, especially for simple transactional operations, which is key for real-time data.
- MySQL tends to use **fewer system resources**, making it ideal for environments with limited resources

PostgreSQL can be an excellent choice too

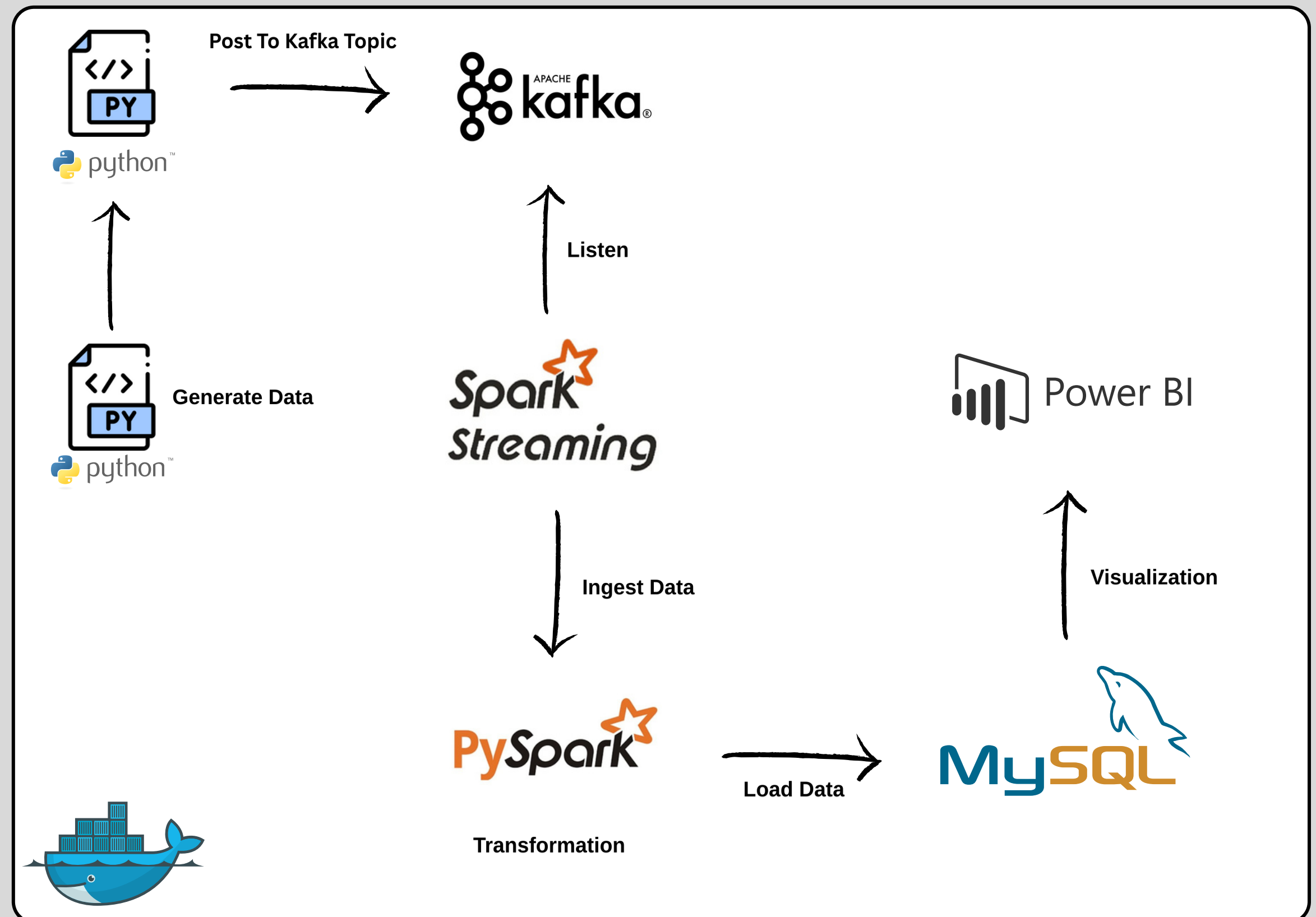


ETL For Batch Data

- **Extract:** Data is loaded from static sources using **Pandas**.
- **Transform:** Data is cleaned and preprocessed in the Jupyter notebook.
- **Load:** The batch data is inserted into MySQL using **SQLAlchemy**.

Real-Time Data

Pipeline



Streaming Data from Kafka to Spark

- Once the data is generated and posted to Kafka, our pipeline kicks into action.
- Spark Streaming seamlessly ingests real-time data from the Kafka topic.
- This allows us to process large volumes of data as they arrive, ensuring efficient data flow without latency.



Processing Data with PySpark

- With the data ingested, **PySpark** takes over, transforming the raw streams into structured, actionable data.
- We apply business rules, clean the dataset, and prepare it for storage, enabling the data to meet all analysis and reporting requirements.



Loading Data into MySQL

- After transformation, the next step is to persist the data. **PySpark** loads the processed data into a **MySQL** database.
- This provides a reliable and scalable storage solution, where the data can be queried and accessed efficiently for further operations.

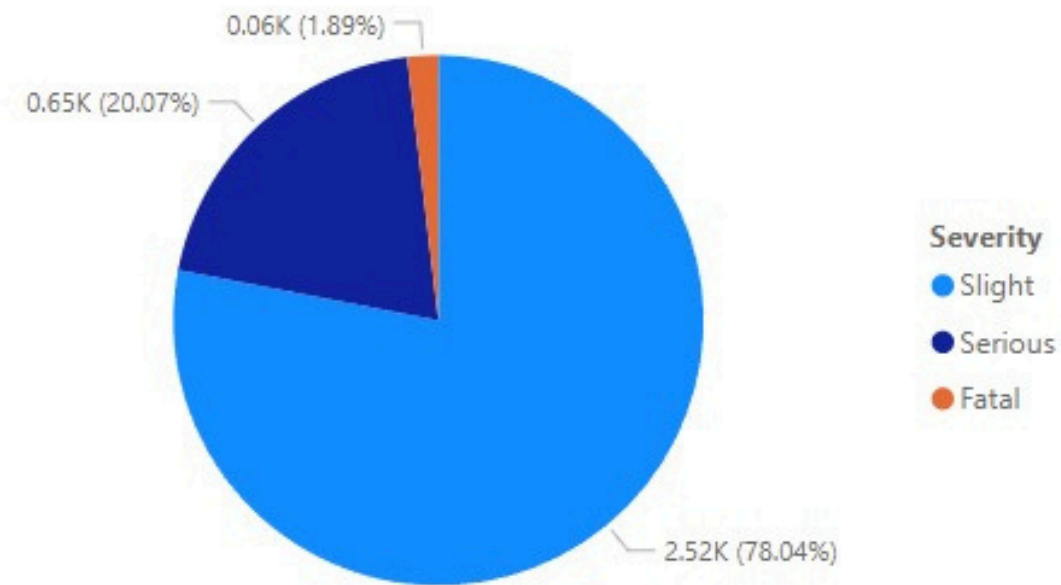


Ensuring Data Integrity and Availability

- Throughout the process, **Spark** ensures that data integrity is maintained.
- Real-time error handling and monitoring are incorporated, guaranteeing that only high-quality, clean data reaches the database.

Analysis

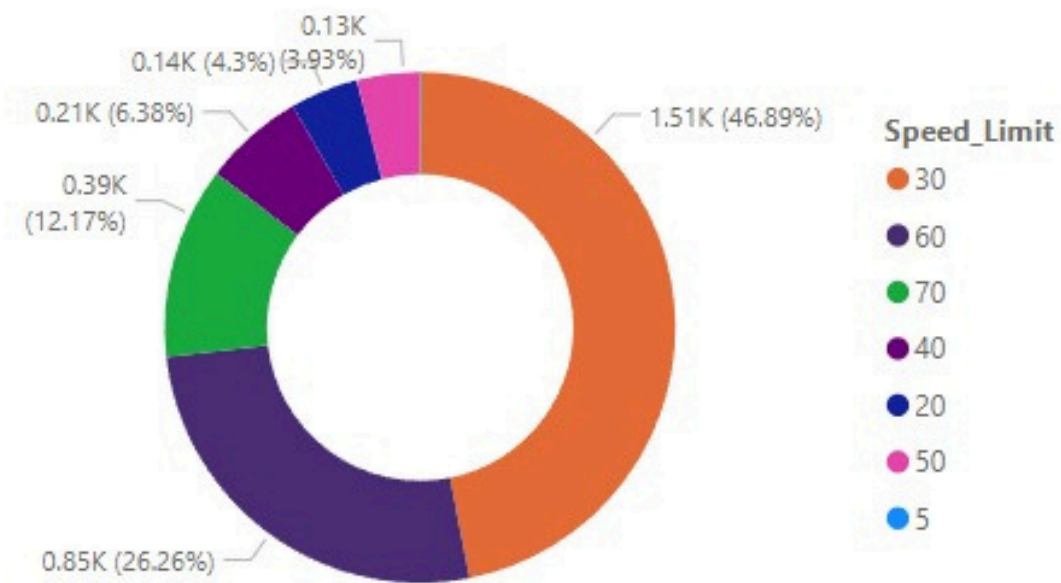
Count of Severity_ID by Severity



visibility



Count of Road_Condition by Speed_Limit

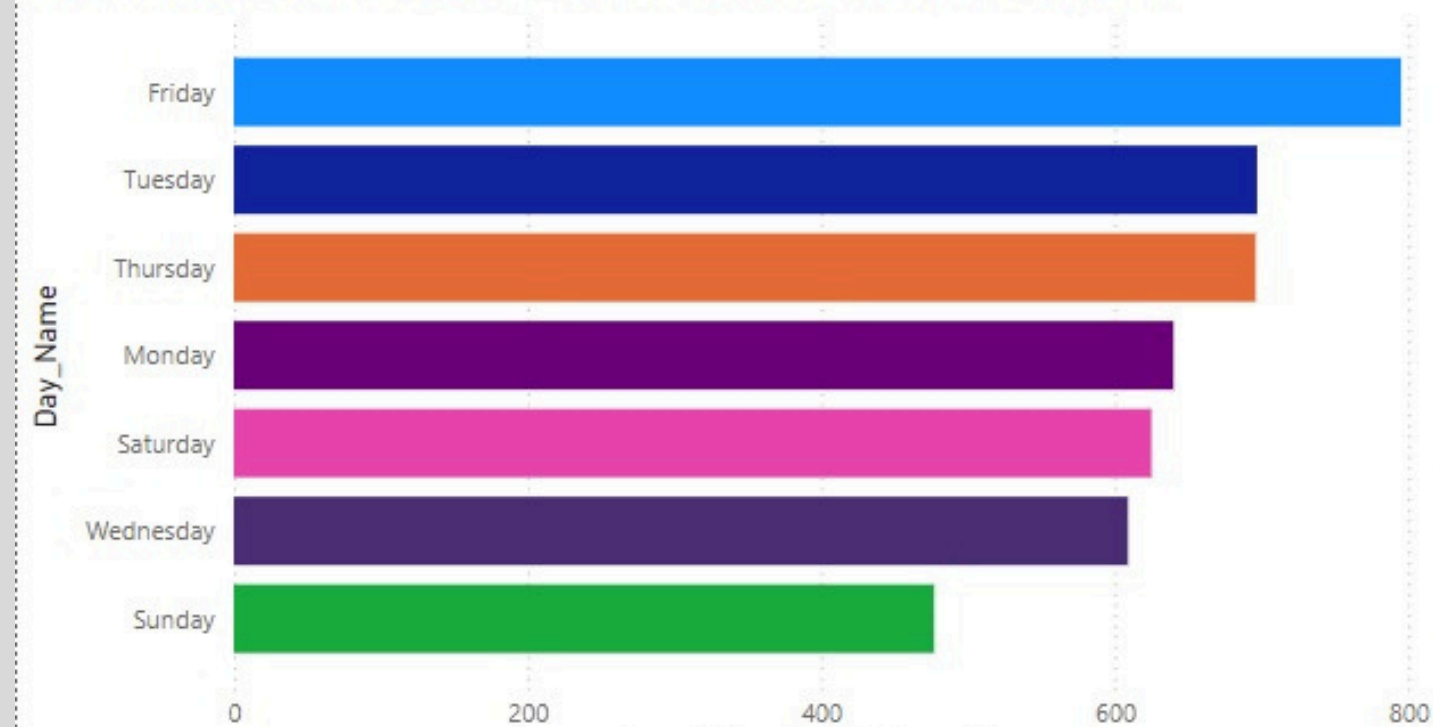


Sum of Speed_Limit by Speed_Limit

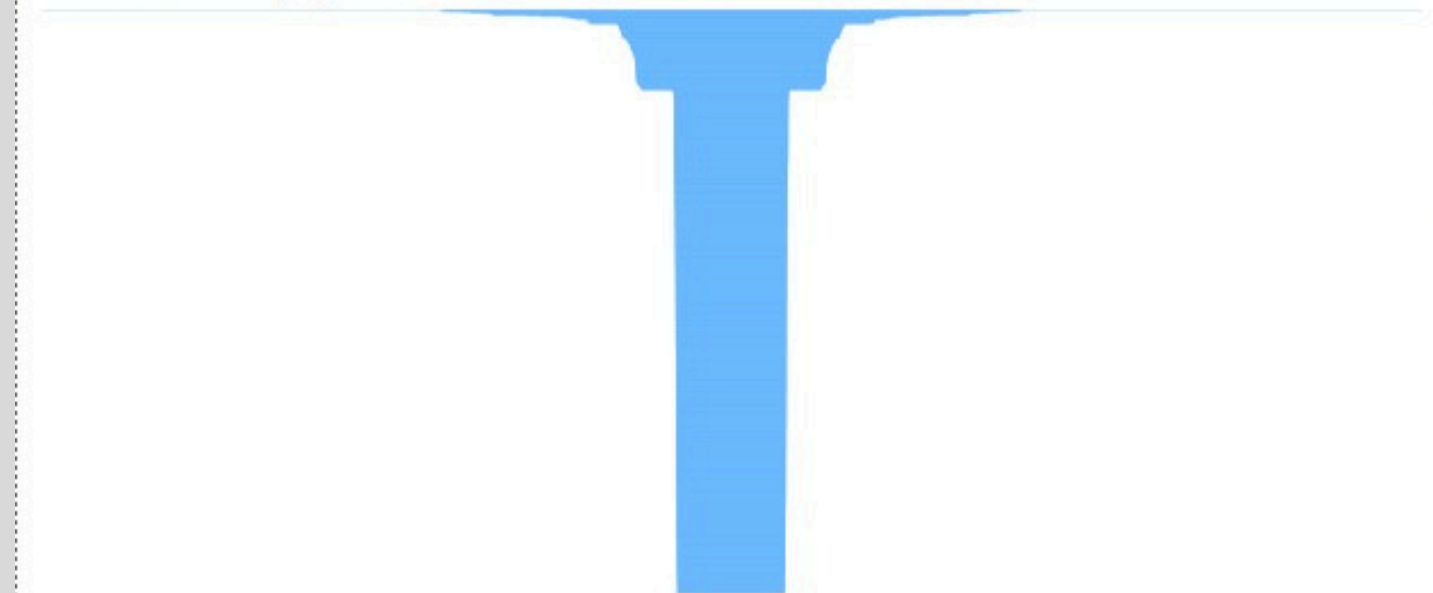


Sum of Number_Of_Casualties by Day_Name and Day_Name

Day_Name ● Friday ● Tuesday ● Thursday ● Monday ● Saturday ● Wednesday ● Sunday



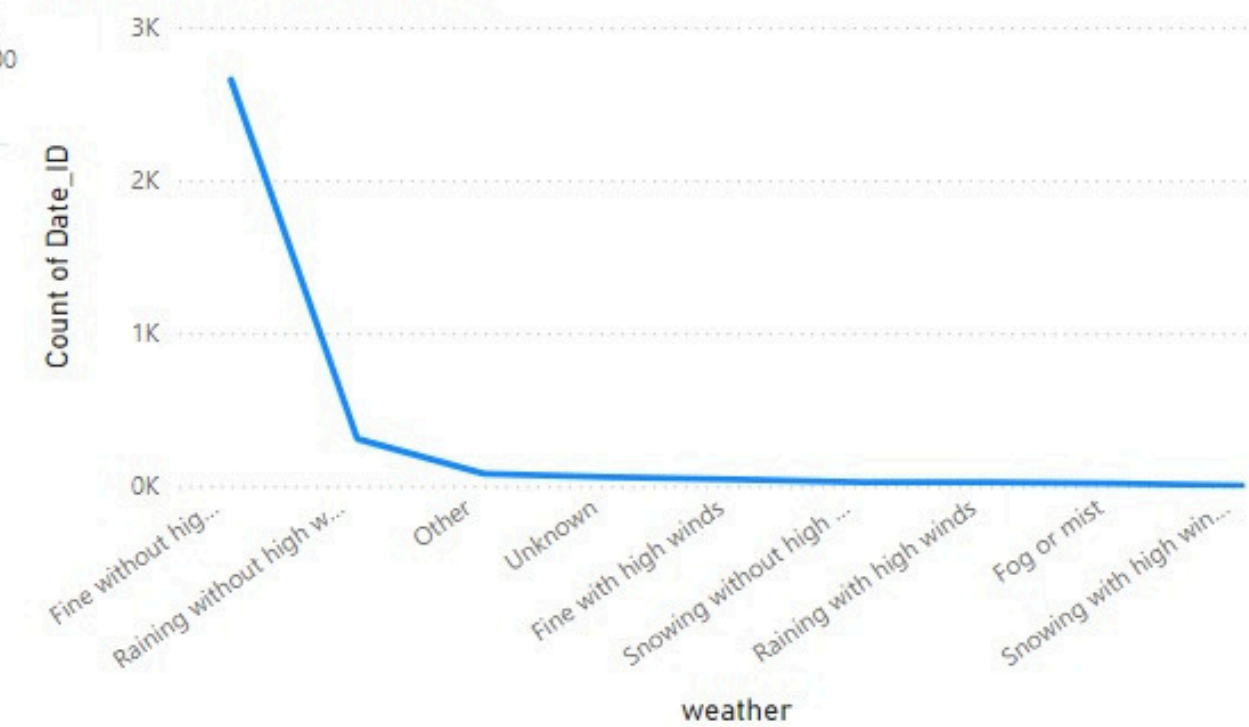
Sum of Northing by Location



Average of Northing and Average of Easting by Location



Count of Date_ID by weather



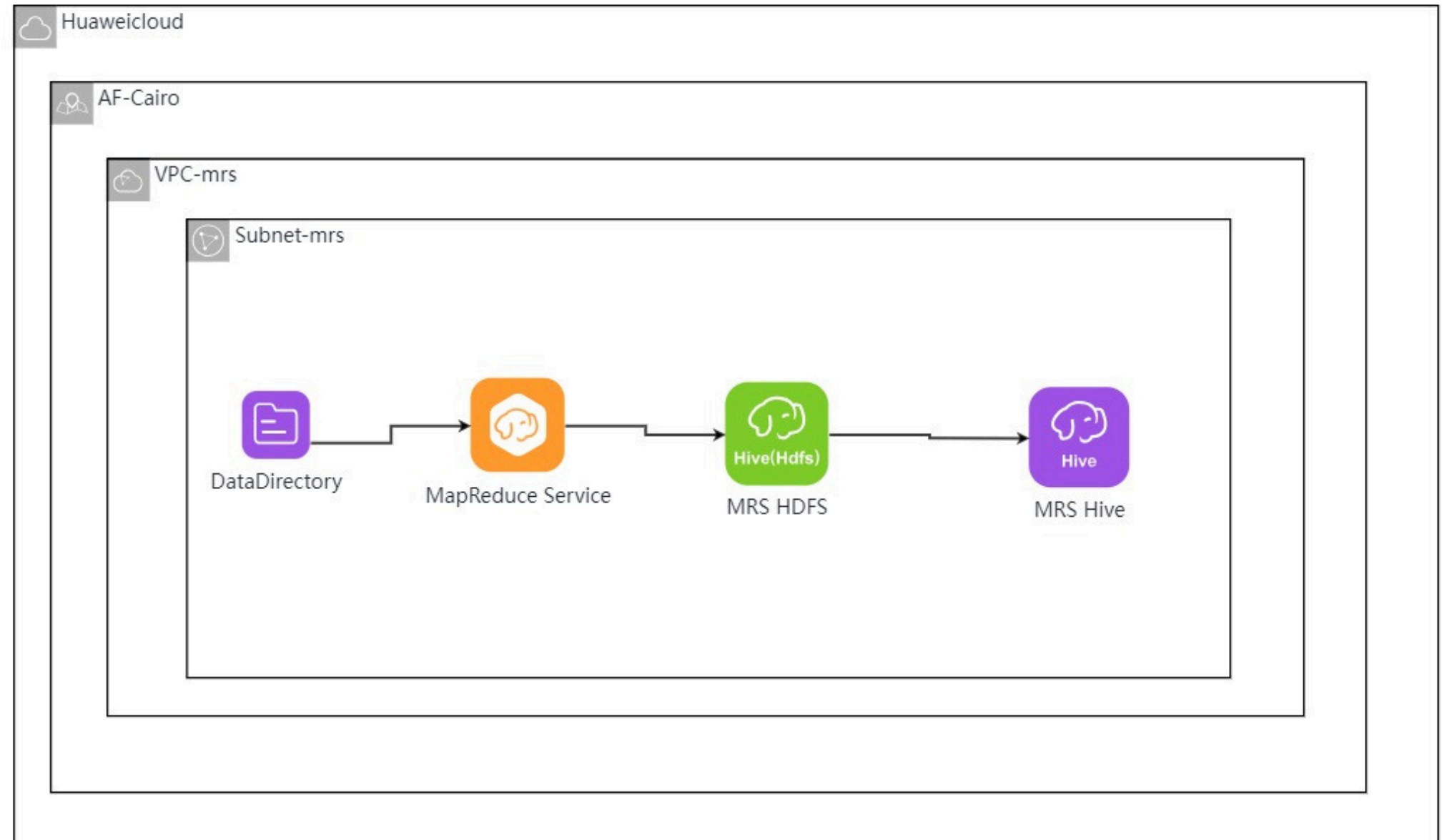
Huawei Cloud MRS

Huawei Cloud MRS

- Huawei Cloud MRS stands for **MapReduce Service**, which provides a cloud-based big data platform.
- It's built on top of **Hadoop** and supports tools like HDFS, Hive, Spark, Kafka, and more.
- It simplifies data storage and processing by managing everything in the cloud, which is especially useful for large-scale projects.
- With MRS, we don't need to manage servers directly. it handles configuration, monitoring, and scaling automatically.



Architecture



-
- **Ingest data into HDFS:** We uploaded the raw CSV files into HDFS for distributed storage.
 - **Create Hive tables:** We designed a Hive database and created multiple tables to store the raw data.
 - **Load data into Hive tables:** Using LOAD DATA commands, we imported the data from HDFS into the Hive tables.
 - **Build the analysis layer:** We wrote complex SQL queries to create summary tables for analysis, using the CTAS method (Create Table As Select).
 - **Export to RDS:** Finally, we exported the results from the analysis layer in Hive to an RDS database on Huawei Cloud for reporting and further analysis.

MEET THE TEAM

- Yousof Hatem
- Ahmed Tarek
- Martin Amgad
- Omar Akram

Thank you!

