# Flight Delay Trends at Chicago O'Hare airport

**Group 4:** Daniela Matinho, Hanna Kerr, Yuling Gu, Wei Yin

# Agenda

Flight delays trends in Chicago O'Hare airport

- ◆ Executive Summary & Business Use Case
- ◆ Visualization
- ◆ Data Source & Data preparation
- ◆ MongoDB & Neo4J Insights
- ◆ Modeling
- ◆ Summary & Future Work

# Executive Summary & Business Use Case

# Overview

- Analyze flight disruptions that occur at the busiest airports in the US – Atlanta, LA, O'Hare, Denver, Dallas, JFK

- In depth analysis of the performance at O'Hare airport in terms of delays

- Data captures domestic flights from Jan to Dec 2018

**7,213,446**
2018
Total Number of Flights in U.S

**1,304,214**
2018
Total Departure Delays in U.S.

**67,647 out of 332,953**
2018
Total Departure Delays per Number of Flights at O'Hare airport

# Business Use Case

Can we define a pattern in the
flight delays at O'Hare airport?

Define critical periods of the year in terms of delays to prepare and optimize resources

Identify the main causes of delays at O'Hare – air traffic, weather, etc

Create an on-time arrival tool that describes airlines' performance

Pinpoint origins & destinations that are most likely to have delays

# Data Source & Data preparation

# Data Sources

Data from January to December 2018

**Transit Data**
**(Bureau of transportation Statistics)**

**Weather Data**
**(Daily Summary of weather conditions)**

# Platform considerations

**01** **Cleaning & Analysis**

Python
- Data from different formats can easily be pulled in
- Automation of the cleaning process to save time for repeated tasks

**02** **Storing**

MySQL & Google Cloud
- Security and control
- Cheap storage
- Group work options

**03** **Visualization**

Tableau
- Connected to GCS
- Self-serving
- Easy to use

# Data Preparation

## Data Store

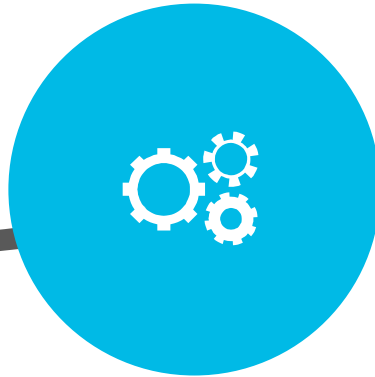- Store data in google cloud using buckets

## Data Source

- Monthly files with flight data
- Daily files with weather information

## Dimension Schema

- Pull data into a dimensional model
- Connect Mysql with tableau to visualize data

## Data Ingestion & Cleaning

- Format different data sources
- Clean columns

# Data Preparation – Step by Step

**01. Python**

**02. MySQL**

**03. Google Cloud**

**Iterate over 12 files with data**

```python
df3 = pd.DataFrame()
for file_name in glob.glob('2018_[0-9][0-9].csv')
    table = pd.read_csv(file_name)
    df3=df3.append(table)
```

```python
df3.head()
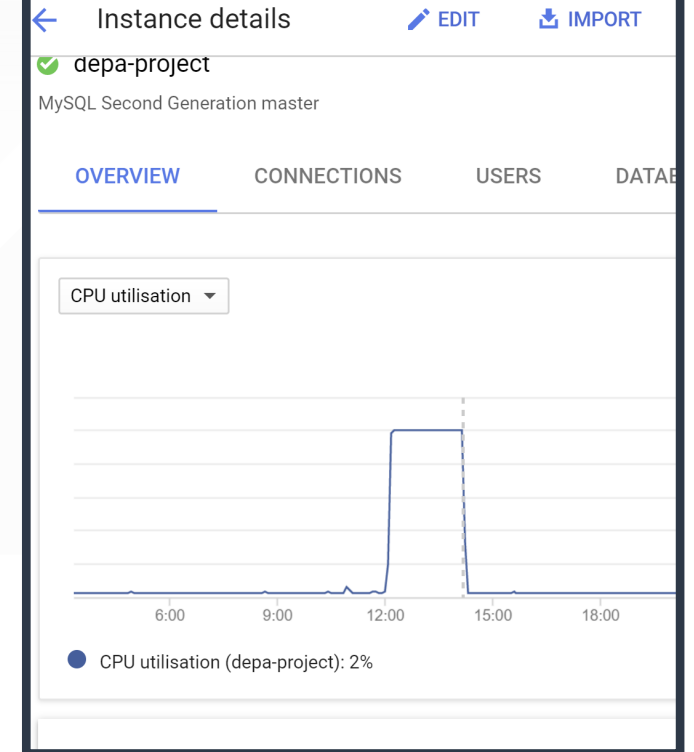```

| | ACTUAL_ELAPSED_TIME | AIR_TIME | ARR_DELAY | ARR_T |
|---|---|---|---|---|
| 0 | 250.0 | 225.0 | -23.0 | 17 |
| 1 | 83.0 | 65.0 | -24.0 | 12 |
| 2 | 126.0 | 106.0 | -13.0 | 16 |
| 3 | 182.0 | 157.0 | -2.0 | 17 |
| 4 | 106.0 | 83.0 | 14.0 | 9 |

```sql
-- Table `flights_snowflake`.`fact_delay`
-- ---------------------------------------------
DROP TABLE IF EXISTS `flights_snowflake`.`fact_delay` ;

CREATE TABLE IF NOT EXISTS `flights_snowflake`.`fact_delay` (
  `Flight_delay_ID` INT(11) NOT NULL,
  `cancel_key` INT(11) NOT NULL,
  `airline_key` INT(11) NOT NULL,
  `origin_airport_key` INT(11) NOT NULL,
  `dep_delay` INT(11) NULL DEFAULT NULL,
  `destination_airport_key` INT(11) NOT NULL,
  `arr_dealy` INT(11) NULL DEFAULT NULL,
  `carrier_delay` INT(11) NULL DEFAULT NULL,
  `weather_delay` INT(11) NULL DEFAULT NULL,
  `NAS_delay` INT(11) NULL DEFAULT NULL,
  `security_delay` INT(11) NULL DEFAULT NULL,
  `late_aircraft_delay` INT(11) NULL DEFAULT NULL,
  `weather_key` INT(11) NOT NULL,
  `flight_key` INT NOT NULL,
  PRIMARY KEY (`Flight_delay_ID`),
  INDEX `fk_fact_flight_dim_cancellation1_idx` (`cancel_key`
```
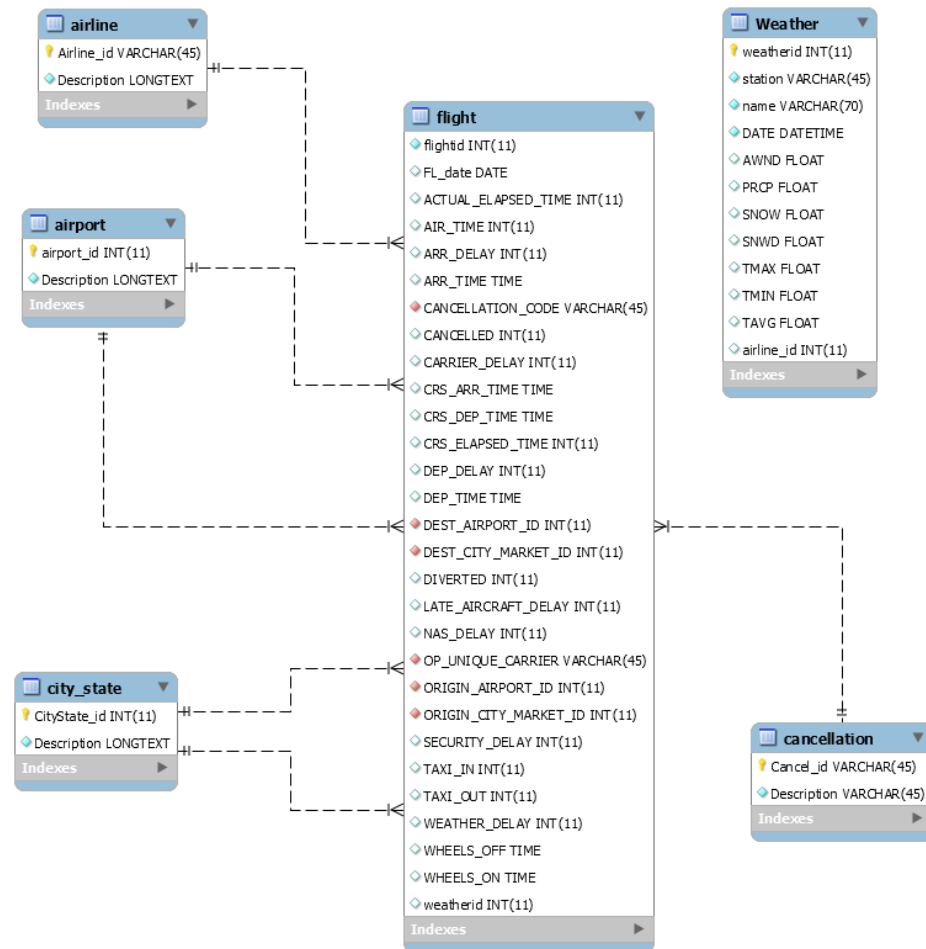
Instance details    ✏ EDIT    ⬇ IMPORT

✅ depa-project

MySQL Second Generation master

**OVERVIEW**    CONNECTIONS    USERS    DATAB

CPU utilisation ▾

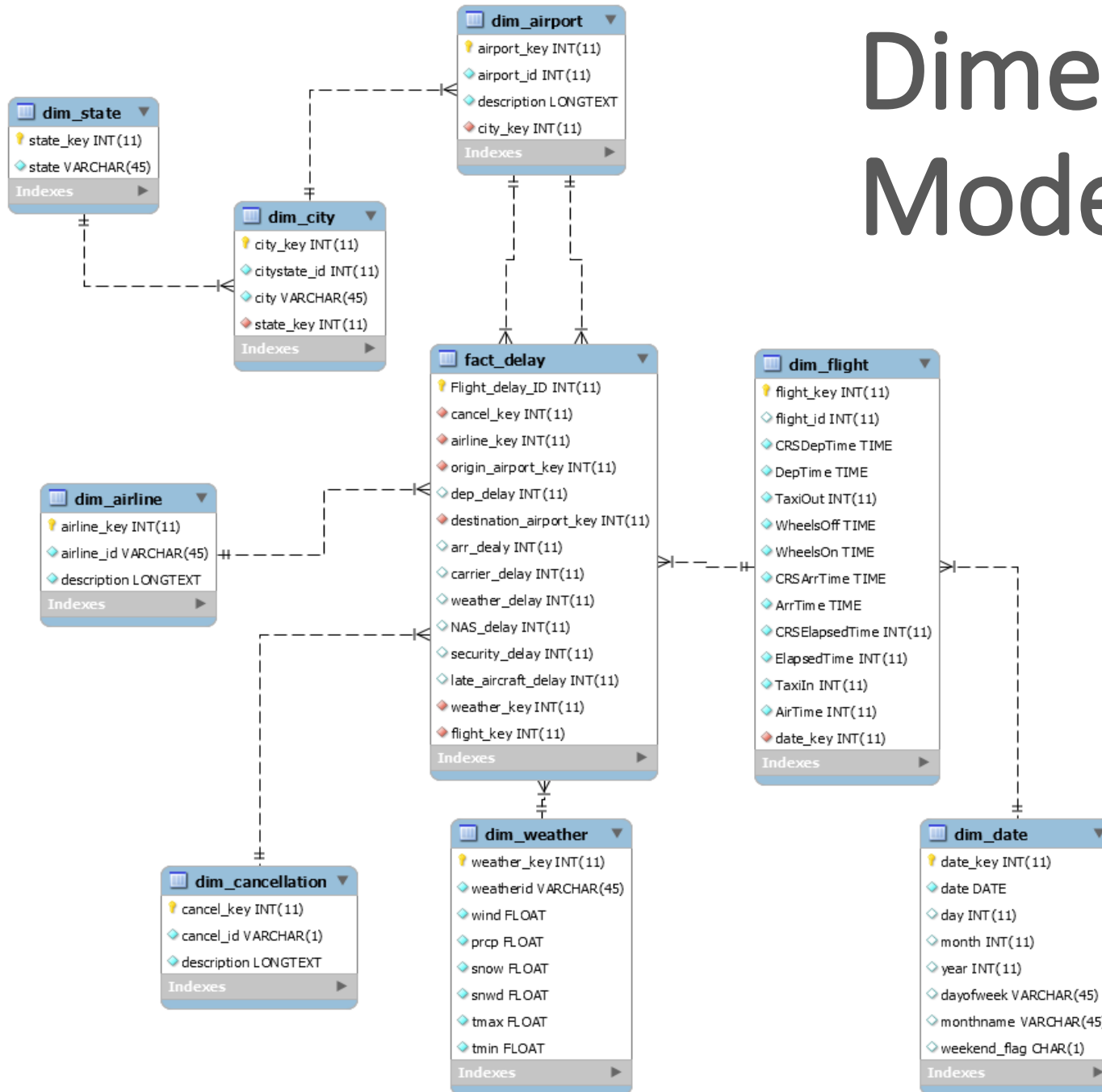6:00   9:00   12:00   15:00   18:00

● CPU utilisation (depa-project): 2%

# Modeling

# EER Diagram

# Dimensional Model

# Visualization

# Analysis

## O'Hare airport

- O'Hare has an average departure delay of 35 min

- The main type of delays are: Aircraft Delay & National Aircraft System & Carrier Delay

- Cancelation rate is higher during the severe weather periods – However rate decreases throughout the year

- Destinations that are more likely to have a flight delay: Wyoming & West Virgin

- O'Hare airport is performing better when the weather is better

- The cheapest airlines have the highest average departure delay minutes (Frontier: 54 min)

## About the 6 airport

- Atlanta is the busiest airport but is also the best in terms of delays

- Average departure delay for the 6 airports, range from 24 to 37 min (with JFK having the largest average)

# MongoDB & NEO 4J Insights

# MongoDB

## JSON Data Model

```json
{
  "_id": "5cf8a6e1733d928192870a56",
  "OP_CARRIER_FL_NUM": 1,
  "ORIGIN_AIRPORT_ID": 12478,
  "ORIGIN": "JFK",
  "ORIGIN_CITY_NAME": "New York, NY",
  "ORIGIN_STATE_NM": "New York",
  "DEST_AIRPORT_ID": 11697,
  "DEST": "FLL",
  "DEST_CITY_NAME": "Fort Lauderdale, FL",
  "DEST_STATE_NM": "Florida",
  "DEP_TIME": 1000,
  "DEP_DELAY": 0,
  "ARR_TIME": 1319,
  "ARR_DELAY": 4
     + {... }
}
```

## Pros

⚙ Easy to create queries, optimize, & maintain Work with data in a natural, intuitive way

⚙ Capacity to adapt & make changes quickly

⚙ Great performance with less code

⚙ Freedom to run anywhere

## Cons

⚙ High Memory Usage (due to no functionality of joins, there is data redundancy)

⚙ Limited Data Size

⚙ Limited Nesting (cannot perform nesting of documents for more than 100 levels)

⚙ Less Secure

# MongoDB – In action

# Summary & Future Work

# Lessons learned

## Data Preparation & Storing

- Cleaning the data is the most important step in the project – the better the quality of the data, the easier it is to treat, analyze & visualize data
- Keeping consistency in the cleaning process is crucial (e.g. format of variables)
- Creating indexes is necessary for querying large sets of data
- Data preparation & loading accounts for most of the overall process

## Visualization

- Keeping the business use case in mind while performing visualization

AND MUCH MORE… ▶

# Future Work

- Build a **model that predicts the probability** of delay & length of delay for upcoming flight
- Compare **ticket prices** for airlines with rate of delay

- Consider data from **international flights**
- **Analyze data for several years** to get a better understanding on O'Hare performance
- Analyze **costs related** to delays & find ways to reduce them
- **Automation of the process,** to constantly update new data

# References

## Articles

Perkins, E, 2019, 'The 10 Worst Airports for Flight Delays, Ranked', SmarterTravel, March 4 – LINK

Bai, Y 2006, Analysis of aircraft arrival delay and airport on-time performance, University of Central Florida Orlando – LINK

## Data

TRANSIT DATA  (BUREAU OF TRANSPORTATION STATISTICS) – LINK

WEATHER DATA  (DAILY SUMMARY OF WEATHER CONDITIONS) – LINK

# THANK YOU

# Q&A

**Group 4:** Daniela Matinho,
Hanna Kerr, Yuling Gu, Wei Yin