# Massive Data Processing
## Taxi's data

Mariano Garralda Barrio
Oscar Ujaque Perez

# 1. Parsing Data

- **Initial data → CSV**
- **Final data → JSON**
- **Weight → 11MB and 1.9 GB**
- **How? → Using Open Refine**

# Data output in json

Trip ID : String

Call Type : Char

    Type A: Trip dispatched from the central

    Type B: Trip dispatched from a stand

    Type C: Trip dispatched randomly in a street

Origin Call: Integer if Call Type = 'A'.  Null otherwise.

Origin Stand: Integer if Call Type = 'B'.  Null otherwise

Taxi Id : Integer

TimeStamp: Integer

Day Type: Char

    Type A: normal day

    Type B: on holidays

    Type C: the day before holidays

Missing Data: Boolean. False if there is no missing data in polyline. True otherwise.

Polyline: String with GPS coordinates for each 15 seconds of the trip.

# 2. Cheking Attributes

- **All attributes checked**

# 3. Map Reduce Tasks

- **Statistics computed:**

  - **Max Distance**
  - **Distance average**
  - **Max velocity**
  - **Velocity average**
  - **Max Trip time**
  - **Trip time Average**
  - **Number of trips**

# 4. Structure



Input File: Taxi.json

**CONTROL AND JOBS DEPENDECY**

Read From File

**JOB1**
TAXITRIPINPUTFORMAT (CUSTOM JSON READER)

Key: Text
Value: TaxiTripWritable

FILTERSELECTIONMAPPER (HADOOP)

Key: Text
Value: TaxiTripWritable

Memory

CLEANUPANDFILTERINGMAPPER

Chain Mappers

TAXITRIPWRITABLE

COMPUTESTATISTICSREDUCER

Key: Text
Value: Text

TextFormat

**JOB2**
TOPTRENDINGTOPICMAPPER

TOPTRENDINGTOPICREDUCER

Key: TweetWritable
Value: Text

TextFormat

OutPut Results Files
Statistics and Top n distances

# 5. Advanced structures used

- **Custom Input File Format: For reading the json file.**

- **Custom Writables:**
    - **TaxiTripWritable: For storing the json data after reading**
    - **GpsPositionWritable: For storing the gps coordinates after reading.**
    - **ArrayWritable<GpsPositioWritable>:**

- **Chain Mappers: For chaining two mappers in first job.**

- **Mappers from hadoop: FieldSelectionMappers: For treating the input data.**

- **Job Control and dependencies: For running two jobs.**