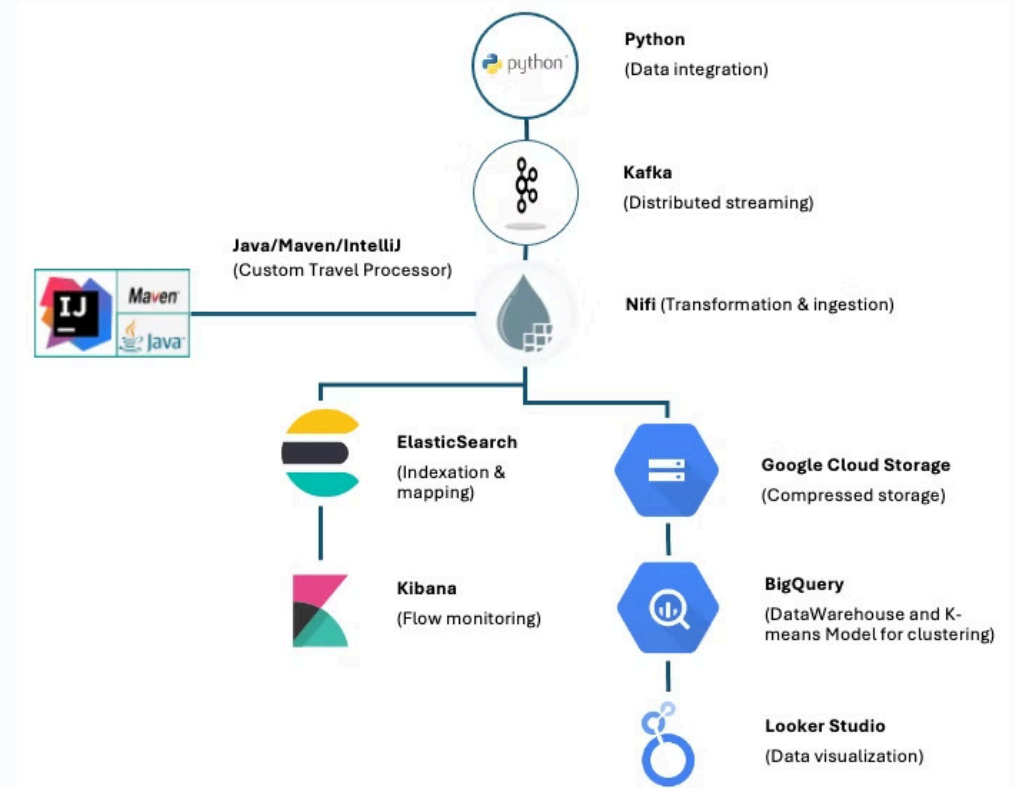# DATASTREAM

Back-end development for a real-time platform calculating taxi fare prices based on a selected comfort level.

# 1. Introduction

- **Objective**

Calculate the distance between a driver and a customer to determine the price of a trip based on the selected comfort type.

# 2. Data Model

Incoming Data in Kafka via Python 'KafkaProducer'

```python
message = {
    "data": [
        {
            "confort": confort["confort"],
            "prix_base_per_km": confort["prix_base_per_km"],
            "properties-client":{
                "logitude": cluster_data["client"]["lon"],
                "latitude": cluster_data["client"]["lat"],
                "nomclient":"FALL",
                "telephoneClient":"060786575"
            },
            "properties-driver":{
                "logitude": cluster_data["driver"]["lon"],
                "latitude": cluster_data["driver"]["lat"],
                "nomDriver":"DIOP",
                "telephoneDriver":"0760786575"
            }
        }
    ]
}
```
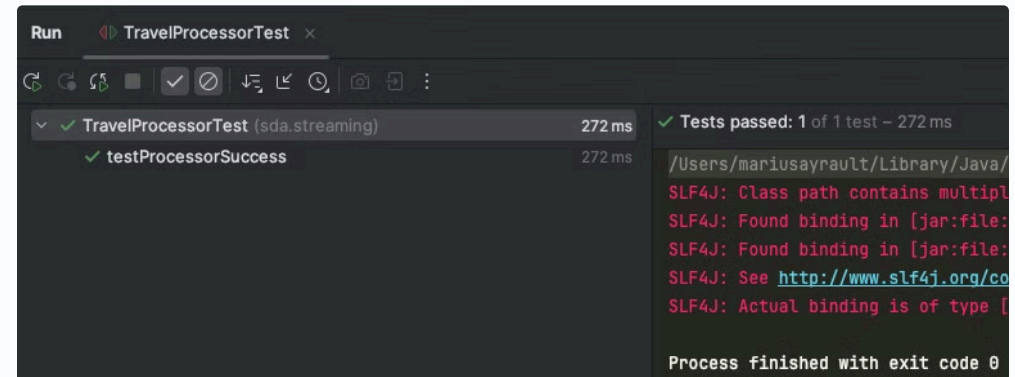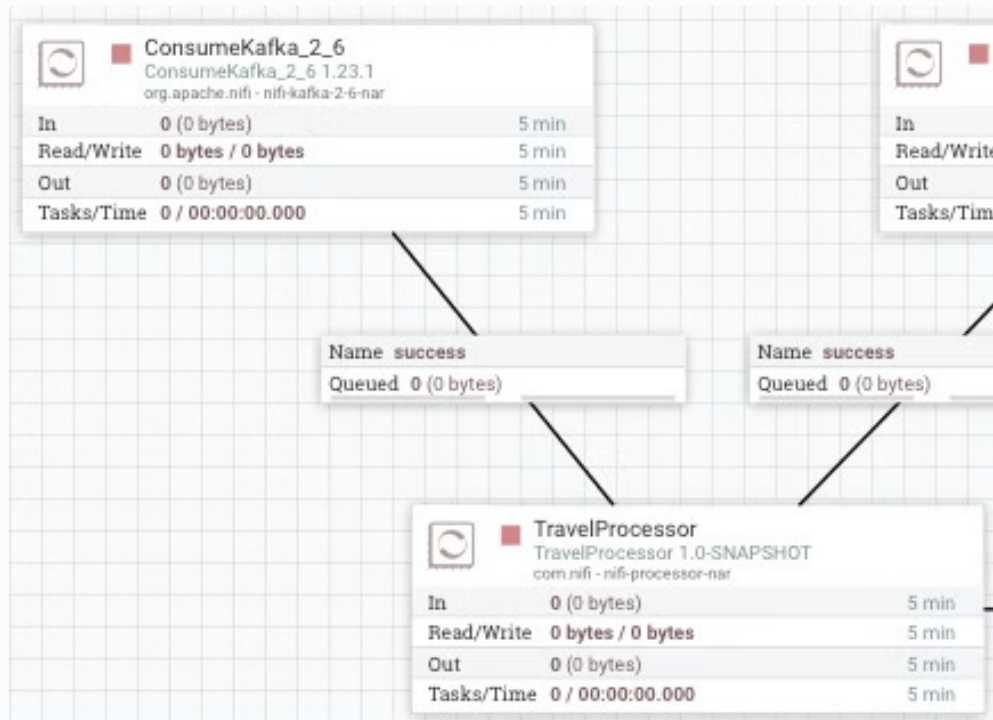
```python
conforts = [
    {"confort": "high", "prix_base_per_km": 0.5},
    {"confort": "medium", "prix_base_per_km": 0.3},
    {"confort": "low", "prix_base_per_km": 0.2}
]


clusters = [

    {"client": {"lat": 48.8566, "lon": 2.3522}, "driver": {"lat": 40.4168, "lo
    {"client": {"lat": 40.7128, "lon": -74.0060}, "driver": {"lat": 41.8781, "
    {"client": {"lat": 35.6895, "lon": 139.6917}, "driver": {"lat": 35.1814, "
    {"client": {"lat": 52.5200, "lon": 13.4050}, "driver": {"lat": 41.9028, "l
    {"client": {"lat": -12.0464, "lon": -77.0428}, "driver": {"lat": -33.4489,
    {"client": {"lat": 19.4326, "lon": -99.1332}, "driver": {"lat": 20.6597, "
    {"client": {"lat": -33.9249, "lon": 18.4241}, "driver": {"lat": -4.4419, "
    {"client": {"lat": 25.2048, "lon": 55.2708}, "driver": {"lat": 33.3152, "l
]
```

# 3. TravelProcessor Development

Development based on defined architecture and unit tests.

# 4. Data Model

## Transformation Examples

### Custom TravelProcessor
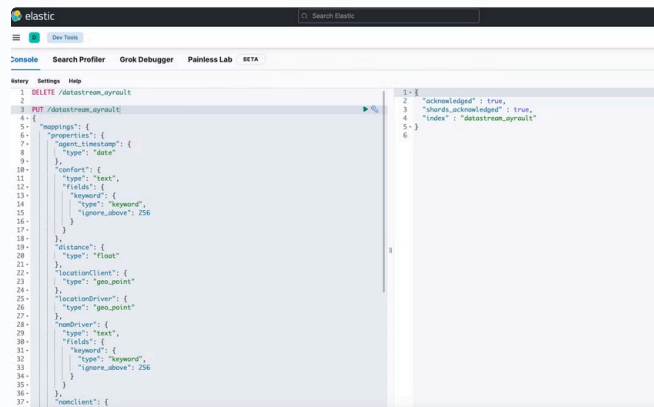


```json
{"data": [{
    "properties-client": {
        "nomclient": "FALL",
        "telephoneClient": "060786575",
        "location": "2.3522,48.8566"
    },
    "distance": 944.494,
    "properties-driver": {
        "nomDriver": "DIOP",
        "location": "3.7038,40.4168",
        "telephoneDriver": "0760786575"
    },
    "prix_base_per_km": 2,
    "confort": "standard",
    "prix_travel": 1888.99
}]}
```
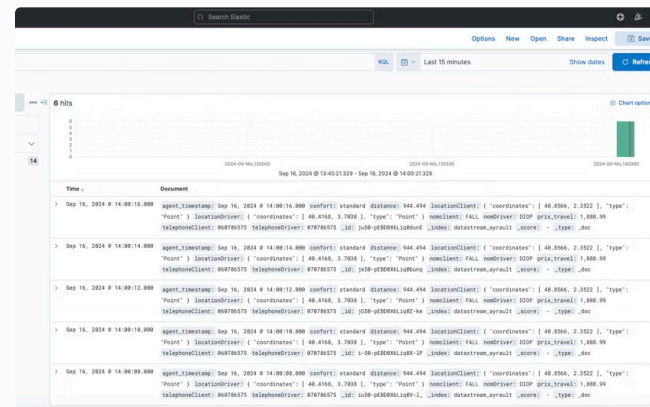
### JoltTransformJSON



```json
{
    "nomclient" : "FALL",
    "telephoneClient" : "060786575",
    "locationClient" : "2.3522,48.8566",
    "distance" : 944.494,
    "confort" : "standard",
    "prix_travel" : 1888.99,
    "nomDriver" : "DIOP",
    "locationDriver" : "3.7038,40.4168",
    "telephoneDriver" : "0760786575",
    "agent_timestamp" : "2024-09-17T18:06:39Z"
}
```

# 5. Indexing & monitoring
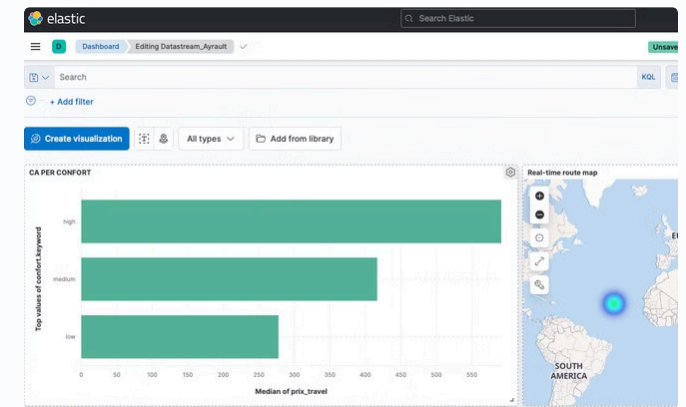






## Indexing and mapping via *ElasticSearch*

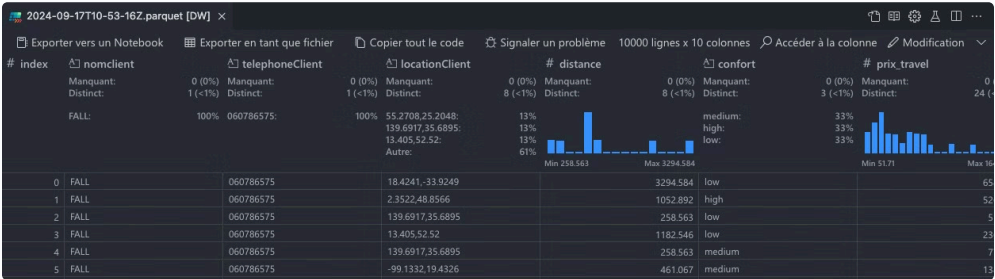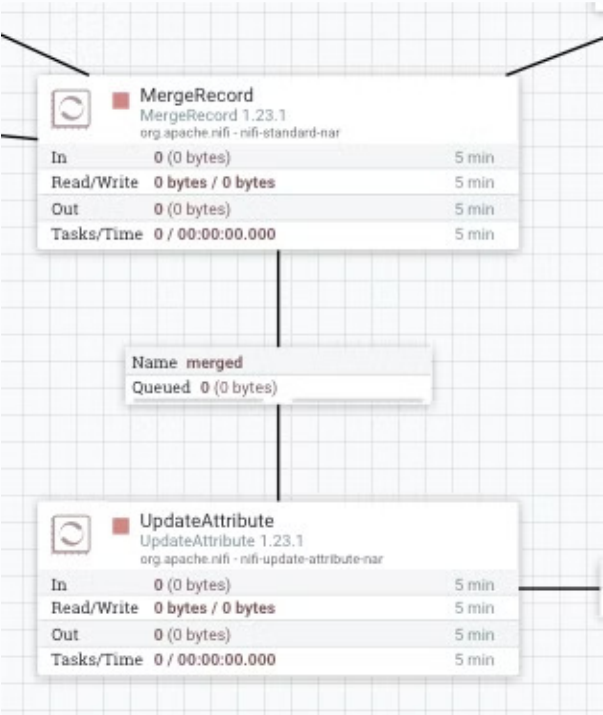Definition of data indexing and mapping.

## Performance analysis

Monitoring of data streams.

## Data analysis via *Kibana*

Detailed real-time visualization.

# 6. Datawarehouse and BigQuery

Merge all 10,000 records and compress data (.parquet*) with timestamp agent.

# 7. DataWarehouse Configuration



| OBJECTS | CONFIGURATION | PERMISSIONS | PROTECTION | LIFECYCLE | OBSERVABILITY | INVENTORY REPORTS | OPERATIONS |
|---|---|---|---|---|---|---|---|

**Folder browser** ⟨

🛍 sda-datastreaming-bucket                               ⋮

Buckets > sda-datastreaming-bucket ⧉

CREATE FOLDER     UPLOAD ▾     TRANSFER DATA ▾     OTHER SERVICES ▾

Filter by name prefix only ▾     ≡ Filter  Filter objects and folders          Show  Live objects onl

| ☐ | Name | Size | Type | Created ❓ |
|---|---|---|---|---|
| ☐ | 📄 2024-09-17T10:45:43Z.parquet | 23.6 KB | application/parquet | Sep 17, 2024, 12:45:43 PM |
| ☐ | 📄 2024-09-17T10:46:13Z.parquet | 23.7 KB | application/parquet | Sep 17, 2024, 12:46:14 PM |

🏠 ▾   ✕        @ *1) Create...ble ▾   ✕        ➕ ▾

🔍  1) CreateExternalTable        ▶ RUN      📤 SAVE QUERY ▾      ⬇ DOWNLOAD

```
1  CREATE EXTERNAL TABLE `sda-datastreaming.the_dataset.the_external_table`
2  OPTIONS (
3    format = 'PARQUET',
4    uris = ['gs://sda-datastreaming-bucket/*.parquet']
5  );
```

# 8. Revenue Calculation by Cluster and Comfort Level

## Analysis process with BigQuery ML

Using BigQuery ML to create a K-Means model.

## Clustering Methodology

Identification of 8 clusters based on geographical coordinates.

## Revenue Calculation

Calculation of revenue per cluster and comfort type.
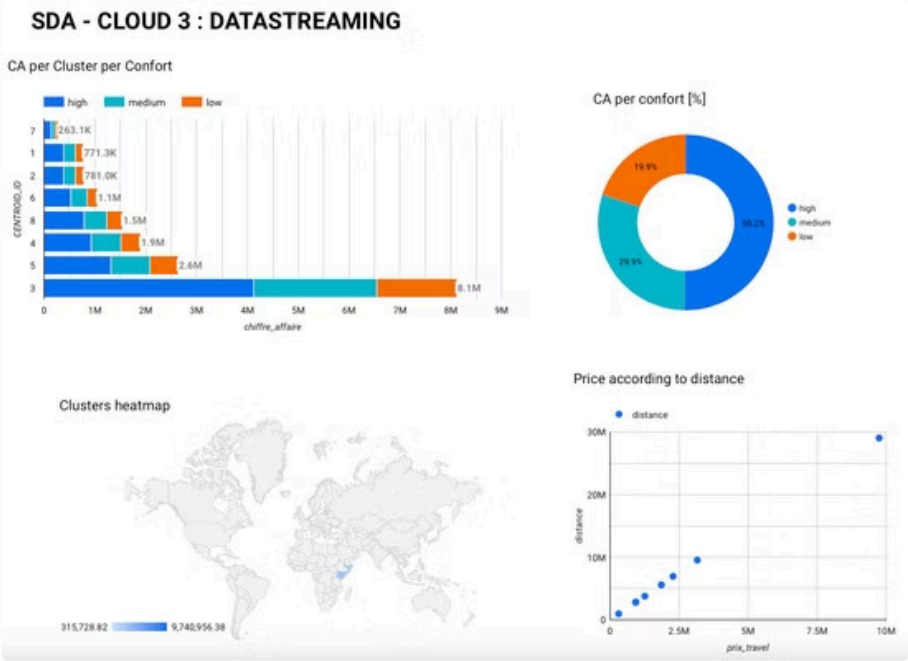


```
2) ProjectModel

1  CREATE OR REPLACE MODEL `sda-datastreaming.the_dataset.mymodel`
2  OPTIONS
3  (MODEL_TYPE='KMEANS',
4  NUM_CLUSTERS=8,
5  standardize_features=true)AS
6  SELECT
7    CAST(REGEXP_EXTRACT(locationClient, '(.*)?,') AS FLOAT64) AS lon,
8    CAST(REGEXP_EXTRACT(locationClient, ',(.*)?') AS FLOAT64) AS lat
9  FROM `sda-datastreaming.the_dataset.the_external_table`
```

```
3) CA per Confort per Cluster

1   SELECT
2     CENTROID_ID,
3     confort,
4     SUM(prix_travel) AS chiffre_affaire
5   FROM (
6     SELECT
7       prix_travel,
8       confort,
9       CENTROID_ID -- cluster
10    FROM ML.PREDICT(MODEL `sda-datastreaming.the_dataset.mymodel`,
11      (
12        SELECT
13          prix_travel,
14          confort,
15          CAST(REGEXP_EXTRACT(locationClient, '(.*)?,') AS FLOAT64) AS lon,
16          CAST(REGEXP_EXTRACT(locationClient, ',(.*)?') AS FLOAT64) AS lat
17        FROM `sda-datastreaming.the_dataset.the_external_table`
18      )
19    )
20  )
21  GROUP BY CENTROID_ID, confort
```

# 9. Data visualization

Visualizing the model results via Looker Studio