**MINING BIG DATASETS PROJECT No2**

**DEPARTMENT:** DEPARTMENT OF MANAGEMENT SCIENCE AND TECHNOLOGY

**COURSE:** MINING BIG DATASETS

**PROFESSOR**: Y.KOTIDIS, I.FILIPPIDOU

**STUDENTS**: CHARILAOS PETRAKOGIANNIS (F2822112), I.DIMOS (F2822102)

**ACADEMIC PERIOD**: 2021-2022

## Table of Context

# INTRODUCTION

In this specific project we were given a part of open flights airports network, which contained airports, airlines and flights between airports. In particular, the data contained 7698 airports, 6161 airlines, 6956 cities, 237 countries and 65.935 flights between airports.

The aforementioned data was in three different ".csv" files, which we had to load in the neo4j program with which we would complete this project.

Below we are going to present you:

- A detailed description of our graph model using a chart and a verbal description of the elements.
- The commands we used in order to import the files to the database.
- The Cypher code for the required queries with their respective results.

# PART A: THE GRAPH MODEL

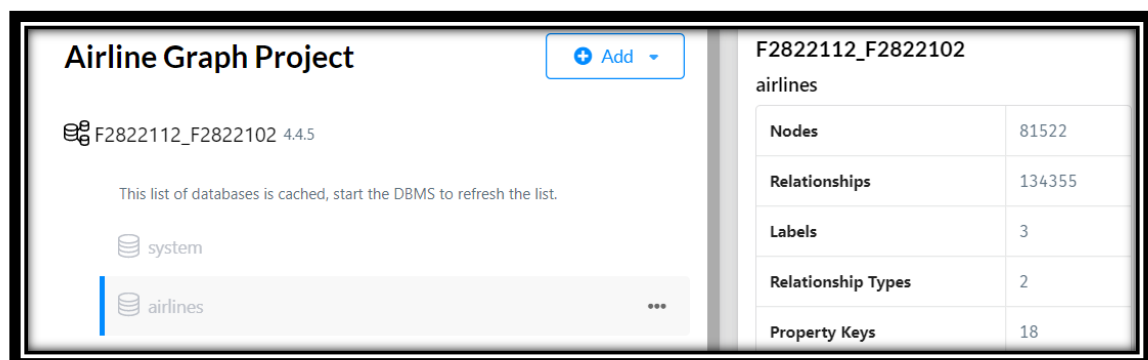For us to complete this project, we needed to create a new database within the neo4j desktop (Figure 1):



*Figure 1 : The neo4j graph database*

Our database consisted of three different nodes (Airports, Routes, and Airlines) and two types of directed relationships (From/To) between the vertices of Routes (flights) and Airports (Figures 2-6). In order for us to load the data in neo4j, we had to declare some constraints that would assure us that we would keep only the unique IDs and names of the existing airports:

Moreover, in order for us to conduct the project's queries, we had to equip our nodes and edges with the appropriate attributes:

**Nodes**

- Airlines (Node's ID, IATA, ICAO, Country of origin, Name of the company, Airline ID of each company)
- Routes (Node's ID, Airline's code, Airline ID of each company, Destination, Destination's ID, Equipment's code, Source's code, Source's ID, Number of stops)
- Airports (Node's ID, IATA, ICAO, City of origin, Country of origin, Airport ID, Latitude, Longitude, Name of the airport)

**Relationships**

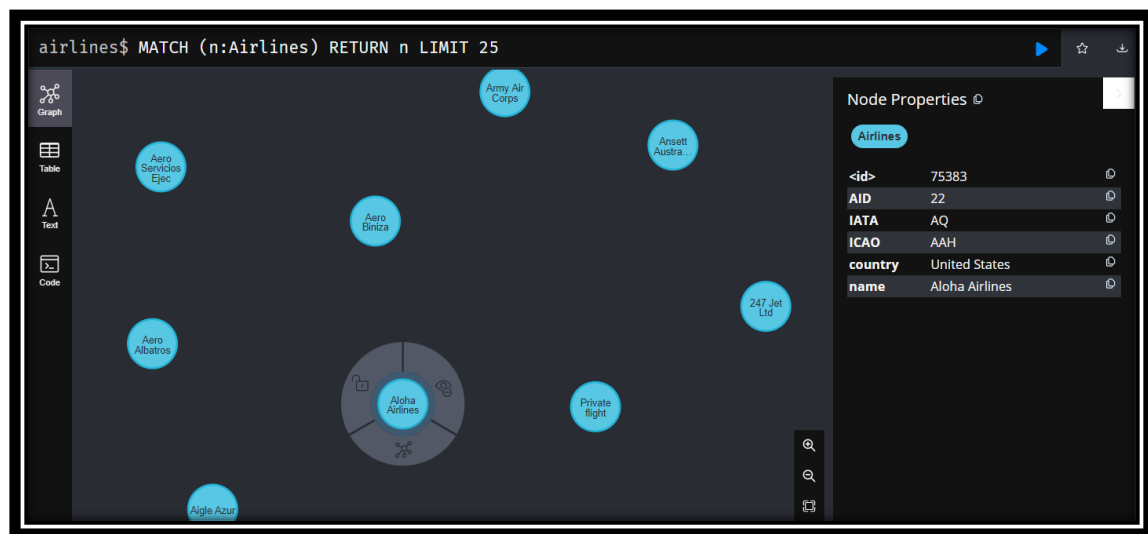- From/To (Relationship's ID, Airline's code, Airline's ID)
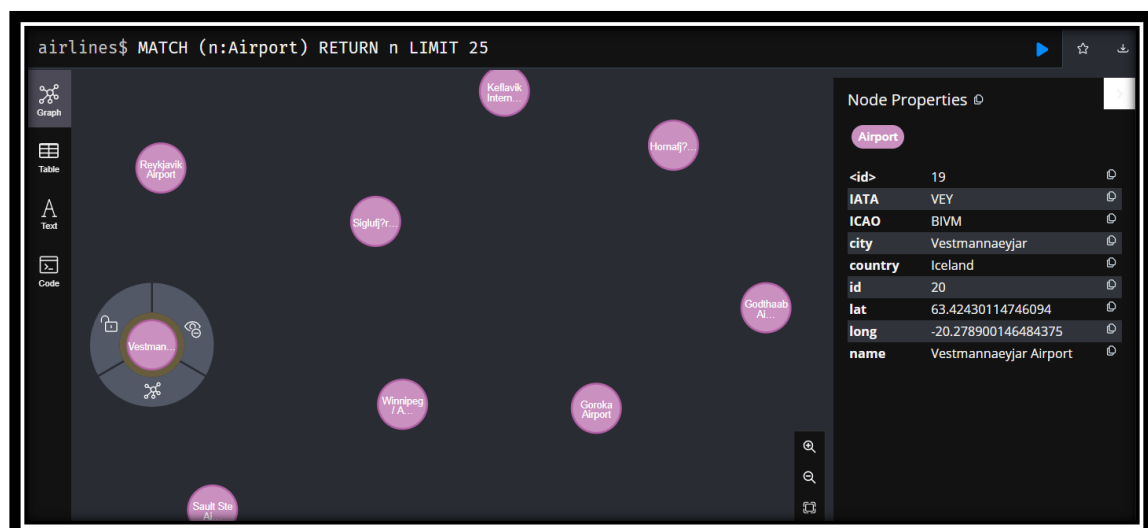


*Figure 2 : Airlines' nodes*
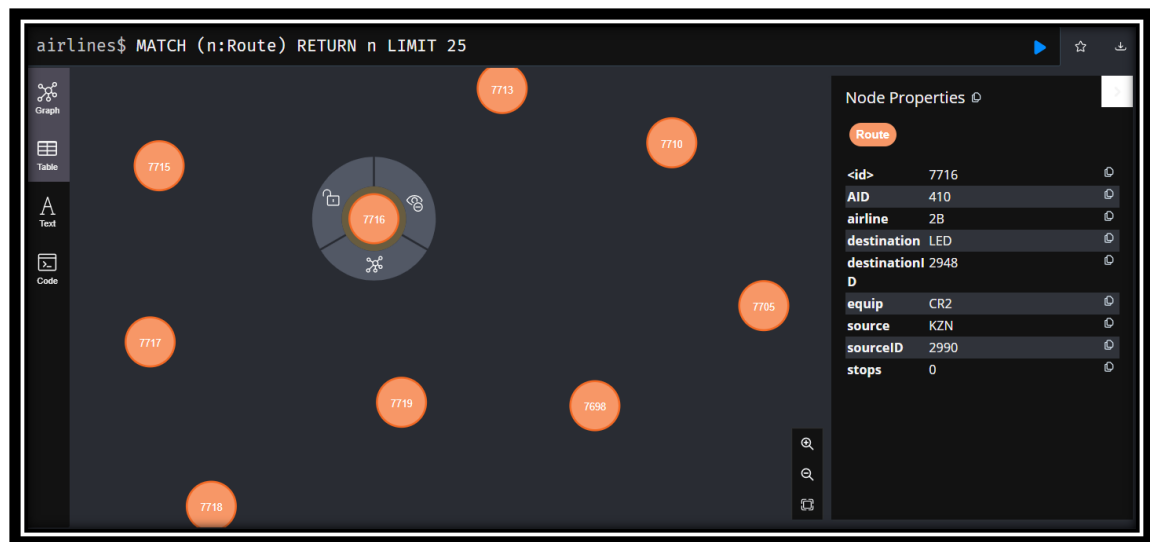


*Figure 3 : Airports' nodes*
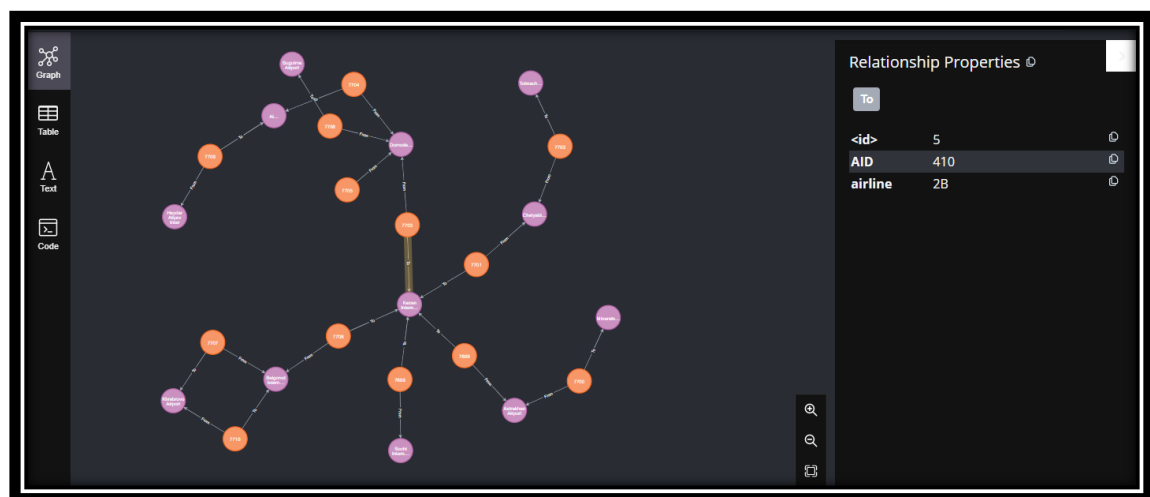
*Figure 4 : Routes' nodes*



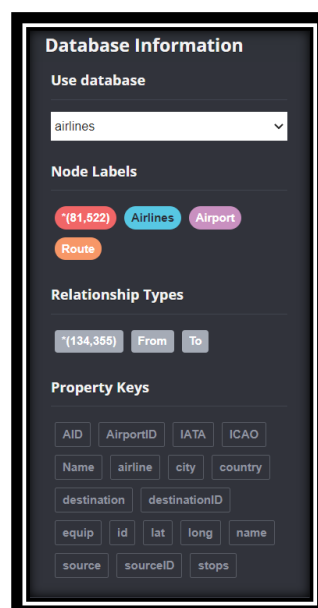*Figure 5 : Relationships between airports and routes (flights)*



*Figure 6 : Database's accumulated information*

It has to be mentioned that even though we did not use directly the "Airlines" nodes, their existence was crucial for our queries, due to the fact that they were necessary in order for us to find the name of each airline through the two-digit codes or the unique IDs we had in our possession.

## PART B: IMPORT OF THE FILES INTO THE DATABASE

In figure 7 below we can see the directory we placed the files so the neo4j be able to load them into our database.

It is very important to mention that we decided not to drop the null values from our csv files as the majority of the columns were containing information about the flights and only few of them were nulls.



*Figure 7 : Directory that contains the files for the database*

**Import of the file that contains information about the Airports**

```
1  LOAD CSV WITH HEADERS
2  FROM "file:///airports.csv" AS line
3  WITH DISTINCT line
4  CREATE (n:Airport {
5      id: line.AirportID,
6      name: line.Name,
7      country: line.Country,
8      city: line.City,
9      lat: line.Latitude,
10     long: line.Longitude,
11     IATA: line.IATA,
12     ICAO: line.ICAO })
```

*Figure 8 : Importation of the Airports.csv file*

**Import of the file that contain information about the Routes**

```
1   LOAD CSV WITH HEADERS
2   FROM "file:///routes.csv" AS line
3   WITH DISTINCT line
4   CREATE (n:Route {
5       source: line.Source,
6       destination: line.Destination,
7       stops: line.Stops,
8       equip: line.Equipment,
9       airline: line.Airline,
10      sourceID: line.SourceID,
11      destinationID: line.DestinationID,
12      AID: line.AirlineID})
```

*Figure 9 : Importation of the routes.csv file*

**Import of the file that contain information about the Airlines**

```
1   LOAD CSV WITH HEADERS
2   FROM "file:///airlines.csv" AS line
3   WITH DISTINCT line
4   CREATE (n:Airlines {
5       name: line.Name,
6       IATA: line.IATA,
7       ICAO: line.ICAO,
8       country: line.Country,
9       AID:line.AirlineID})
```

*Figure 10 : Importation of the airlines.csv file*

In this step we have also created a query to delete one observation that has the name "Unknown".

*Figure 11 : Deletion of the Unknown observation*

# PART C: THE CYPHER CODE

1) Which are the top 5 airports with the most flights. Return airport name and number of flights.

**Solution**



```
1  MATCH (n:Airport)
2  WITH n, (SIZE((n)-[:To]-())) as NUMBER_OF_INCOMING_FLIGHTS, (SIZE((n)-[:From]-())) as NUMBER_OF_OUTCOMING_FLIGHTS
3  RETURN n.name as AIRPORT, NUMBER_OF_INCOMING_FLIGHTS, NUMBER_OF_OUTCOMING_FLIGHTS,
   (NUMBER_OF_INCOMING_FLIGHTS+NUMBER_OF_OUTCOMING_FLIGHTS) as TOTAL_NUMBER_OF_FLIGHTS
4  ORDER BY TOTAL_NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| AIRPORT | NUMBER_OF_INCOMING_FLIGHTS | NUMBER_OF_OUTCOMING_FLIGHTS | TOTAL_NUMBER_OF_FLIGHTS |
|---|---|---|---|
| "Hartsfield Jackson Atlanta International Airport" | 911 | 915 | 1826 |
| "Chicago O'Hare International Airport" | 550 | 558 | 1108 |
| "Beijing Capital International Airport" | 534 | 535 | 1069 |
| "London Heathrow Airport" | 524 | 527 | 1051 |
| "Charles de Gaulle International Airport" | 517 | 524 | 1041 |

*Figure 12 : Query of question No1*

2) Which are the top 5 countries with the most airports. Return country name and number of airports.

**Solution**

```
1 MATCH (n:Airport)
2 RETURN n.country AS COUNTRY, COUNT(n.country) AS NUMBER_OF_AIRPORTS
3 ORDER BY NUMBER_OF_AIRPORTS DESC LIMIT 5
```
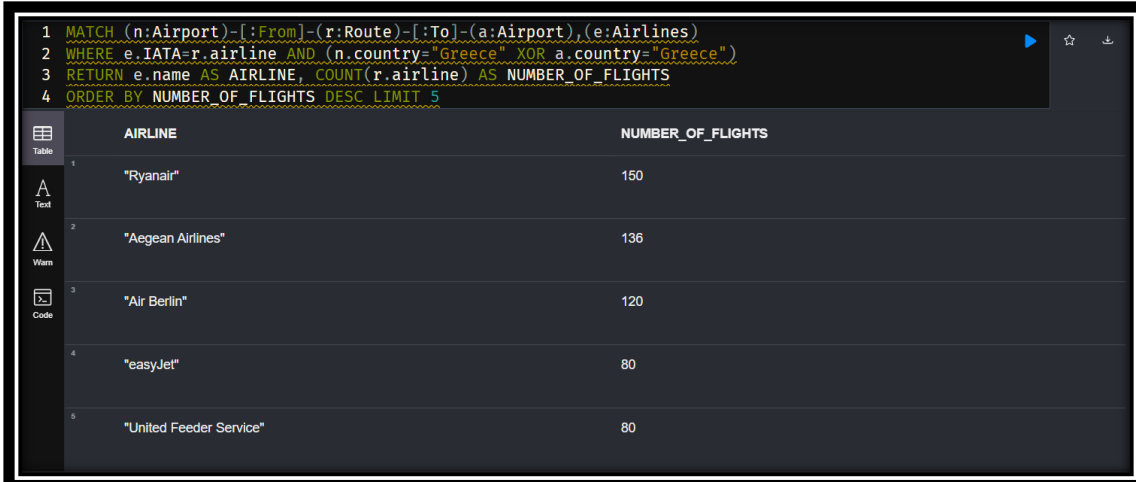
| COUNTRY | NUMBER_OF_AIRPORTS |
|---|---|
| "United States" | 1512 |
| "Canada" | 430 |
| "Australia" | 334 |
| "Brazil" | 264 |
| "Russia" | 264 |

*Figure 13 : Query of question No2*

3) Which are the top 5 airlines with international flights from/to 'Greece'. Return airline name and number of flights.

**Note:** In this task we've worked with two different approaches. In the first case we worked using IATA codes, which are a unique identifier for each airline. Their difference from Airline ID is that they take into account as a separate entity each subsidiary company that belongs to a larger group. This results in some minor differences in the results produced one way (IATA) and the other (Airline ID). We present both ways because, being non-experts in air data, we want to cover every possible possibility that can produce correct results.

**Solution 1**

```
1 MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport),(e:Airlines)
2 WHERE e.IATA=r.airline AND (n.country="Greece" XOR a.country="Greece")
3 RETURN e.name AS AIRLINE, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4 ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| AIRLINE | NUMBER_OF_FLIGHTS |
|---|---|
| "Ryanair" | 150 |
| "Aegean Airlines" | 136 |
| "Air Berlin" | 120 |
| "easyJet" | 80 |
| "United Feeder Service" | 80 |

*Figure 14 : Query of question No3 (first way)*

**Solution 2**

```
1  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport),(e:Airlines)
2  WHERE e.AID=r.AID AND (n.country="Greece" XOR a.country="Greece")
3  RETURN e.name AS AIRLINE, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4  ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| AIRLINE | NUMBER_OF_FLIGHTS |
|---|---|
| "Ryanair" | 150 |
| "Aegean Airlines" | 136 |
| "Air Berlin" | 120 |
| "easyJet" | 80 |
| "TUIfly" | 64 |

*Figure 15 : Query of question No3 (second way)*

4) Which are the top 5 airlines with local flights inside 'Germany'. Return airline name and number of flights.

**Note:** As in the previous question, here too we present two different ways of solving that lead to similar, but not identical, results for the simple reason that the IATA codes consider subsidiaries as different entities, while the IDs do not.

**Solution 1**

```
1  MATCH (n:Airport{country:"Germany"})-[:From]-(r:Route)-[:To]-(a:Airport{country:"Germany"}),
   (e:Airlines)
2  WHERE e.IATA=r.airline
3  RETURN e.name AS AIRLINE, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4  ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| AIRLINE | NUMBER_OF_FLIGHTS |
|---|---|
| "Lufthansa Cargo" | 64 |
| "Lufthansa" | 64 |
| "Germanwings" | 54 |
| "Air Berlin" | 44 |
| "Hainan Airlines" | 16 |

*Figure 16 : Query of question No4 (first way)*

## Solution 2

```
1  MATCH (n:Airport{country:"Germany"})-[:From]-(r:Route)-[:To]-(a:Airport{country:"Germany"}),
   (e:Airlines)
2  WHERE e.AID=r.AID
3  RETURN e.name AS AIRLINE, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4  ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| AIRLINE | NUMBER_OF_FLIGHTS |
|---------|-------------------|
| "Lufthansa" | 64 |
| "Germanwings" | 54 |
| "Air Berlin" | 44 |
| "Hainan Airlines" | 16 |
| "Ethiopian Airlines" | 6 |

*Figure 17 : Query of question No4 (second way)*

5) Which are the top 10 countries with flights to Greece. Return country name and number of flights.

## Solution

```
1  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport{country:"Greece"}),(e:Airlines)
2  WHERE e.IATA=r.airline and e.country <> 'Greece'
3  RETURN e.country AS COUNTRY, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4  ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 10
```

| COUNTRY | NUMBER_OF_FLIGHTS |
|---------|-------------------|
| "Germany" | 174 |
| "United States" | 115 |
| "Ireland" | 85 |
| "United Kingdom" | 61 |
| "France" | 36 |
| "Cambodia" | 32 |
| "Netherlands" | 29 |
| "Austria" | 20 |
| "Russia" | 19 |
| "Canada" | 19 |

*Figure 18 : Query of question No5*

6)Find the percentage of air traffic (inbound and outbound) for every city in Greece. Return city name and the corresponding traffic percentage in descending order.

**Note:** In order to implement the query No 6, we created two extra queries to compute the total flights in Greece. The total flights in Greece will help us to compute the air traffic percentage , as their sum is the denominator of our fraction. The total is **1601** flights in Greece. The queries can be seen below in the figures 19 ,20 and the main query in figure 21 **:**

```
1  MATCH ((n:Airport{country:"Greece"})-[a:From]-())
2  return count(a) AS TOTAL_FLIGHTS_FROM_GREECE
```

| TOTAL_FLIGHTS_FROM_GREECE |
| --- |
| 787 |

*Figure 19 : Query to compute the total flights from Greece*

```
1  MATCH ((n:Airport{country:"Greece"})-[a:To]-())
2  return count(a) AS TOTAL_FLIGHTS_TO_GREECE
```

| TOTAL_FLIGHTS_TO_GREECE |
| --- |
| 814 |

*Figure 20 : Query to compute the total flights to Greece*

## Main Query

```
1  MATCH (n:Airport)
2  WITH n, (SIZE((n)-[:To]-())) as INBOUND_FLIGHTS,
3  (SIZE((n)-[:From]-())) as OUTBOUND_FLIGHTS
4  WHERE n.country='Greece' and INBOUND_FLIGHTS >= 1 and OUTBOUND_FLIGHTS >=1
5  RETURN n.city as CITY,
6  round(((INBOUND_FLIGHTS + OUTBOUND_FLIGHTS)*1.0/ 1601 )*100,2) as AIR_TRAFFIC_PERCENTAGE
7  order by AIR_TRAFFIC_PERCENTAGE DESC
```

| CITY | AIR_TRAFFIC_PERCENTAGE |
| --- | --- |
| "Athens" | 25.17 |
| "Heraklion" | 13.68 |
| "Thessaloniki" | 10.87 |
| "Rhodos" | 10.49 |
| "Kerkyra/corfu" | 6.43 |
| "Kos" | 5.62 |

*Figure 21 : Query of question No6*

| | CITY | AIR_TRAFFIC_PERCENTAGE | | | |
|---|---|---|---|---|---|
| 1 | CITY | AIR_TRAFFIC_PERCENTAGE | 25 | Kithira | 0.5 |
| 2 | Athens | 25.17 | 26 | Kasos | 0.5 |
| 3 | Heraklion | 13.68 | 27 | Astypalaia | 0.5 |
| 4 | Thessaloniki | 10.87 | 28 | Ikaria | 0.44 |
| 5 | Rhodos | 10.49 | 29 | Alexandroupolis | 0.37 |
| 6 | Kerkyra/corfu | 6.43 | 30 | Patras | 0.37 |
| 7 | Kos | 5.62 | 31 | Skiros | 0.37 |
| 8 | Chania | 5.62 | 32 | Ioannina | 0.25 |
| 9 | Zakynthos | 2.06 | 33 | Skiathos | 0.25 |
| 10 | Thira | 1.94 | 34 | Milos | 0.25 |
| 11 | Mykonos | 1.62 | 35 | Cyclades Islands | 0.25 |
| 12 | Preveza | 1.37 | 36 | Paros | 0.25 |
| 13 | Kalamata | 1.25 | 37 | Kastelorizo | 0.25 |
| 14 | Mytilini | 1.25 | 38 | Syros Island | 0.25 |
| 15 | Samos | 1.12 | 39 | Kastoria | 0.12 |
| 16 | Keffallinia | 0.81 | | | |
| 17 | Karpathos | 0.81 | | | |
| 18 | Chios | 0.75 | | | |
| 19 | Kavala | 0.75 | | | |
| 20 | Leros | 0.75 | | | |
| 21 | Kalymnos | 0.75 | | | |
| 22 | Nea Anghialos | 0.62 | | | |
| 23 | Limnos | 0.62 | | | |
| 24 | Sitia | 0.62 | | | |

*Figure 22 : Complete table that contains every Greek city and its corresponding traffic percentage*

7) Find the number of international flights to Greece with plane types '738' and '320'. Return for each plane type the number of flights.
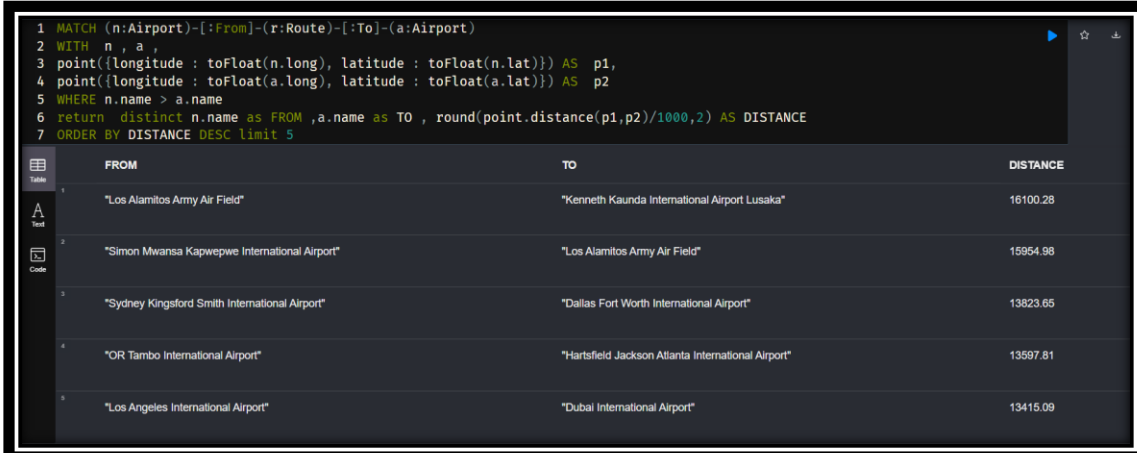
**Solution**

```
1  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport{country:"Greece"})
2  WHERE n.country<>"Greece" AND (r.equip="738" OR r.equip="320")
3  RETURN r.equip AS PLAIN_TYPE, COUNT(r.airline) AS NUMBER_OF_FLIGHTS
4  ORDER BY NUMBER_OF_FLIGHTS DESC LIMIT 5
```

| PLAIN_TYPE | NUMBER_OF_FLIGHTS |
|---|---|
| "320" | 147 |
| "738" | 110 |

*Figure 23 : Query of question No7*

8) Which are the top 5 flights that cover the biggest distance between two airports (use function point({ longitude: s1.longitude, latitude: s1.latitude }) and function distance(point1, point2)). Return From (airport) To (airport) and distance in km.

**Solution**

```
1  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport)
2  WITH  n , a ,
3  point({longitude : toFloat(n.long), latitude : toFloat(n.lat)}) AS  p1,
4  point({longitude : toFloat(a.long), latitude : toFloat(a.lat)}) AS  p2
5  WHERE n.name > a.name
6  return  distinct n.name as FROM ,a.name as TO , round(point.distance(p1,p2)/1000,2) AS DISTANCE
7  ORDER BY DISTANCE DESC limit 5
```
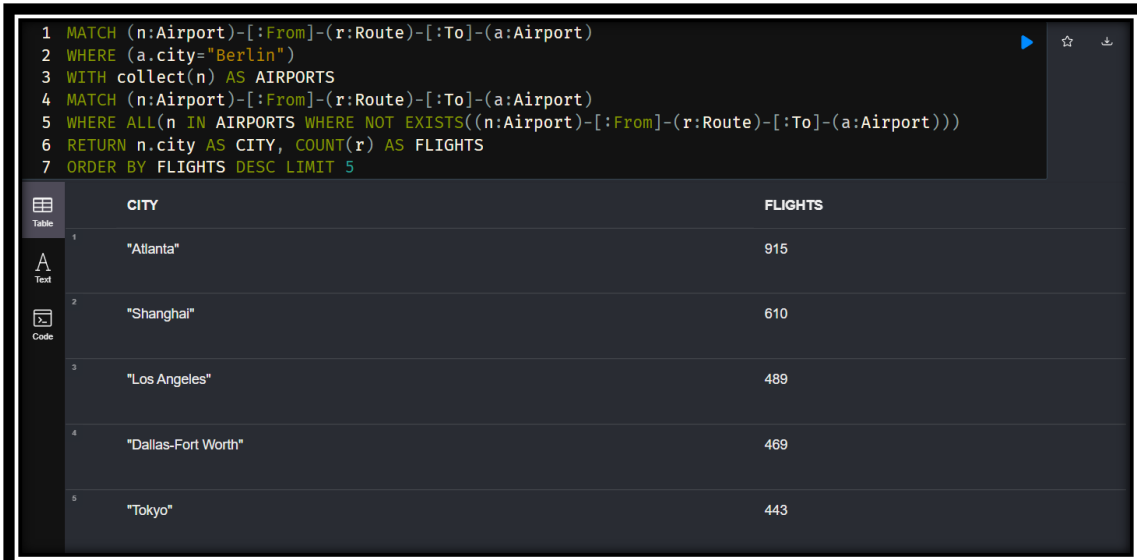
| FROM | TO | DISTANCE |
|------|----|----|
| "Los Alamitos Army Air Field" | "Kenneth Kaunda International Airport Lusaka" | 16100.28 |
| "Simon Mwansa Kapwepwe International Airport" | "Los Alamitos Army Air Field" | 15954.98 |
| "Sydney Kingsford Smith International Airport" | "Dallas Fort Worth International Airport" | 13823.65 |
| "OR Tambo International Airport" | "Hartsfield Jackson Atlanta International Airport" | 13597.81 |
| "Los Angeles International Airport" | "Dubai International Airport" | 13415.09 |

*Figure 24 : Query of question No8*

9) Find 5 cities that are not connected with direct flights to 'Berlin'. Score the cities in descending order with the total number of flights to other destinations. Return city name and score.

**Solution**

```
1  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport)
2  WHERE (a.city="Berlin")
3  WITH collect(n) AS AIRPORTS
4  MATCH (n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport)
5  WHERE ALL(n IN AIRPORTS WHERE NOT EXISTS((n:Airport)-[:From]-(r:Route)-[:To]-(a:Airport)))
6  RETURN n.city AS CITY, COUNT(r) AS FLIGHTS
7  ORDER BY FLIGHTS DESC LIMIT 5
```
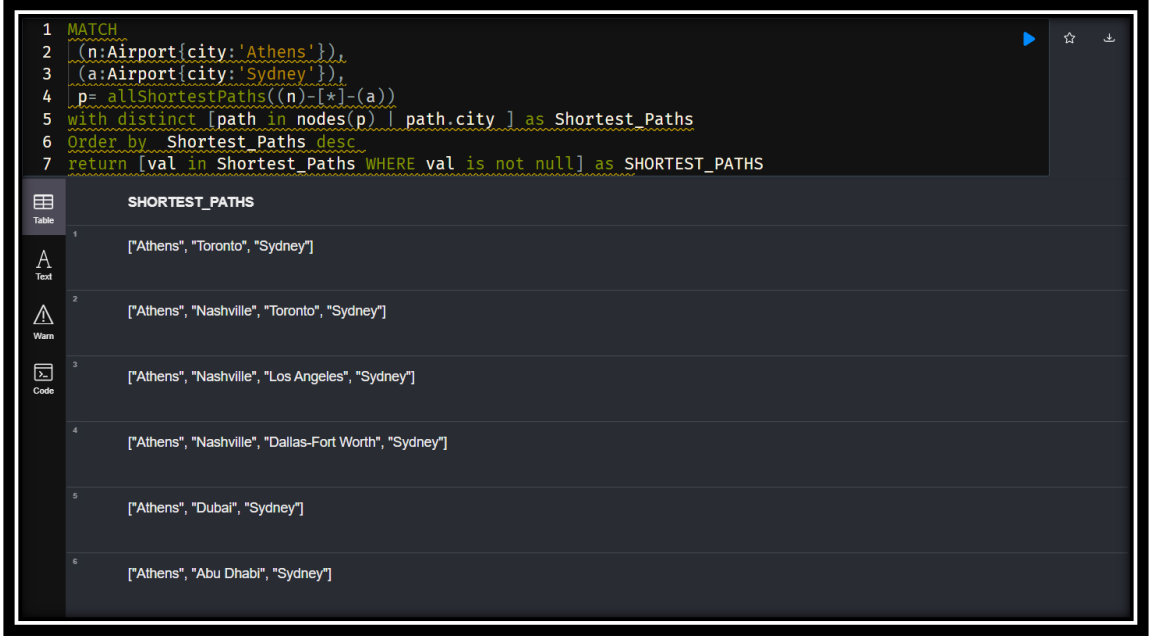
| CITY | FLIGHTS |
|------|---------|
| "Atlanta" | 915 |
| "Shanghai" | 610 |
| "Los Angeles" | 489 |
| "Dallas-Fort Worth" | 469 |
| "Tokyo" | 443 |

*Figure 25 : Query of question No9*

10) Find all shortest paths from 'Athens' to 'Sydney'. Use only relations between flights and city airports.

**Solution**

```
1  MATCH
2  (n:Airport{city:'Athens'}),
3  (a:Airport{city:'Sydney'}),
4  p= allShortestPaths((n)-[*]-(a))
5  with distinct [path in nodes(p) | path.city ] as Shortest_Paths
6  Order by Shortest_Paths desc
7  return [val in Shortest_Paths WHERE val is not null] as SHORTEST_PATHS
```

| SHORTEST_PATHS |
| --- |
| 1 ["Athens", "Toronto", "Sydney"] |
| 2 ["Athens", "Nashville", "Toronto", "Sydney"] |
| 3 ["Athens", "Nashville", "Los Angeles", "Sydney"] |
| 4 ["Athens", "Nashville", "Dallas-Fort Worth", "Sydney"] |
| 5 ["Athens", "Dubai", "Sydney"] |
| 6 ["Athens", "Abu Dhabi", "Sydney"] |

*Figure 26 : Query of question No10*