# Data Engineering: API to AWS to SQL

Florian Bogner

June 17, 2022

This technical report describes the necessary steps to pipe a selection of data into a database. We will deploy various APIs to gather the desired data. These data will be pushed to a SQL database in the cloud. Dynamic data will be collected in the cloud and automatically pushed to the database.

## Introduction

The e-scooter company "Gans" seeks to expand into other European markets. E-scooters have to be strategically placed where they will be used a lot. Gans identified three key properties that will influence e-scooter demand: a city's features (population, elevation), weather and airport vicinity. Consequently, the characteristics of these three key aspects have to be investigated by data analysis. In turn, they represent the source of the necessary data.

This technical report documents the build-up of the heuristics to automatically collect data and populate a database with it. The documentation is structured as follows: First, the data collection by means of API calls will be described in Section 1. After that, Section 2 will be concerned with setting up a MySQL database. Finally, Section 3 will demonstrate the transition to the cloud.

## 1. Collecting Data

As Gans sought to expand into the German speaking markets, we decided to pick the top-ten German and the top-five Austrian cities by population. Initially, we performed web scraping on the Wikipedia pages of the respective cities in order to collect geo and population data. However, there are inconcistencies in the infobox in the Wikipedia entry of the cities which makes it impossible to collect uniform data. Therefore we have decided to use APIs for the entire data collection.

### 1.1. Get `cities` data

To collect data with a `get_cities(qcities)` function using the GeoDB Cities API from RapidAPI. This API takes the WikiData ID codes of the cities as input and provided

us with the city name as well as with information about its elevation, latitudinal and longitudinal geo data and population.

```
1  for qcity in qcities:
2    url = f"https://wft-geo-db.p.rapidapi.com/v1/geo/cities/{qcity}"
3    headers = {"X-RapidAPI-Host": "wft-geo-db.p.rapidapi.com", "X-RapidAPI-
       ↪ Key": "..."}
4    response = requests.request("GET", url, headers=headers)
5    time.sleep(2)
6    city_df = pd.json_normalize(response.json())
7    cities_list.append(city_df)
8  cities_df = pd.concat(cities_list, ignore_index = True)
```

Listing 1: `for` loop to iterate through all cities, call the API and concatenate the DataFrames.

We have programmed our requests such that all city data are returned to us in one concatenated DataFrame. We have dropped the redundant columns and renamed the rest with more intuitive names. This is necessary since some flattened column names returned from the API include dots which will pose problems while handling them in Python or querying them in MySQL.

## 1.2. Get `weather` data

To collect weather data with a `get_weather(cities)` function, we deployed the Open-Weather API. This API took the city names from the `get_cities(qcities)` function as input and provided us with information about precipitation probability, temperature, humidity, cloudiness as well as with wind speed and gust.

```
1  for city in cities:
2    url = f"http://api.openweathermap.org/data/2.5/forecast?q={city}&appid
       ↪ =...&units=metric"
3    weather = requests.get(url)
4    weather_df = pd.json_normalize(weather.json()["list"])
5    weather_df["city"] = city
6    df_list.append(weather_df)
7  weather_combined = pd.concat(df_list, ignore_index = True)
```

Listing 2: `for` loop to iterate through all cities, call the API and concatenate the DataFrames.

We have programmed our requests such that all weather data are returned to us in one concatenated DataFrame. We have dropped the redundant columns and renamed the rest with more intuitive names. This is necessary since some flattened column names returned from the API include dots which will pose problems while handling them in Python or querying them in MySQL.

## 1.3. Get `airports` data

To collect airports data with a `get_airports(lat, lon)` function, we deployed the AeroDataBox API. This API took the latitudinal and longitudinal geo data from the

`get_cities(qcities)` function as input and provided us with information about icao codes, airport names, municipality names, county codes as well as latitudinal and longitudinal geo data.

```python
for i in range(len(lat)):
  url = f"https://aerodatabox.p.rapidapi.com/airports/search/location/{
    ↪ lat[i]}/{lon[i]}/km/50/10"
  querystring = {"withFlightInfoOnly":"true"}
  headers = {"X-RapidAPI-Host": "aerodatabox.p.rapidapi.com", "X-RapidAPI
    ↪ -Key": "..."}
        response = requests.request("GET", url, headers=headers, params=
    ↪ querystring)
  airport_df = pd.json_normalize(response.json()["items"])
  airport_list.append(airport_df)
airports_df = pd.concat(airport_list, ignore_index = True)
```

Listing 3: `for` loop to iterate through all longitudinal and latitudinal geo data, call the API and concatenate the DataFrames.

We have programmed our requests such that all airports data are returned to us in one concatenated DataFrame. We have dropped the redundant columns and renamed the rest with more intuitive names. This is necessary since some flattened column names returned from the API include dots which will pose problems while handling them in Python or querying them in MySQL.

## 1.4. Get `arrivals` data

To collect arrivals data with a `get_arrivals(icao)` function, we deployed the Aero-DataBox API. This API took the icao codes from the `get_airports(lat, lon)` function and returned flight numbers, ICAO codes of the departure airports, scheduled arrival times, airline names and aircraft models.

```python
for code in icao:
  url = f"https://aerodatabox.p.rapidapi.com/flights/airports/icao/{code
    ↪ }/{tomorrow}T10:00/{tomorrow}T22:00"
  querystring = {"withLeg":"false","direction":"Arrival","withCancelled":
    ↪ "false","withCodeshared":"true",
                  "withCargo":"false","withPrivate":"false","
    ↪ withLocation":"false"}
  headers = {"X-RapidAPI-Key": "...", "X-RapidAPI-Host": "aerodatabox.p.
    ↪ rapidapi.com"}
  response = requests.request("GET", url, headers = headers, params =
    ↪ querystring)
  arrival_df = pd.json_normalize(response.json()["arrivals"])
  arrival_df["arrival_icao"] = code
  arrival_list.append(arrival_df)
arrivals_df = pd.concat(arrival_list, ignore_index = True)
```

Listing 4: `for` loop to iterate through all icao airport codes, call the API and concatenate the DataFrames.

We have programmed our requests such that all flight arrivals data are returned to us in one concatenated DataFrame. We have dropped the redundant columns and renamed

the rest with more intuitive names. This is necessary since some flattened column names returned from the API include dots which will pose problems while handling them in Python or querying them in MySQL.

## 2. SQL

The gathered data for `cities`, `weather`, `airports` and `arrivals` have to be stored in a MySQL database. This allows for data consistency, reduction of redundancy and for querying and analyzing the data. A new database will be created with the following code:

```
1  CREATE DATABASE gans;
2  USE gans;
```
Listing 5: A new database.

To connect to a MySQL database, we used the Python library SQLAlchemy. The code in Listing 6 establishes a connection to a MySQL database to which the DataFrames from Section 1 can be pushed.

```
1  schema="gans"
2  host='...rds.amazonaws.com'
3  user="admin"
4  password='...'
5  port=3306
6  con = f'mysql+pymysql://{user}:{password}@{host}:{port}/{schema}'
```
Listing 6: Connection to MySQL database.

SQLAlchemy also handles the push of the DataFrames to the MySQL database with the following code in Listing 7. This will be done with `cities`, `weather`, `airports` and `arrivals` DataFrames.

```
1  DataFrame.to_sql('tablename', if_exists='append', con=con, index=False)
```
Listing 7: Push DataFrame to connected MySQL database.

This will populate the MySQL database with the respective DataFrames. However, the tables are still not connected and primary and foreign keys are missing. Altering the tables to have primary and foreign keys as well as changing the table's column's datatypes will be described in the following section.

As the DataFrames were pushed to an empty database, keys and datatypes had to be set afterwards. For the example of `weather`, a new column was added to act as a primary key:

```
1  ALTER TABLE weather
2  ADD weather_id INT NOT NULL AUTO_INCREMENT PRIMARY KEY FIRST;
```
Listing 8: Add a new column for the primary key of the `weather` table.

Only after that, it was possible to define the datatypes as shown below:

```
1  ALTER TABLE `gans`.`weather`
2  CHANGE COLUMN `precip_prob` `precip_prob` FLOAT NULL DEFAULT NULL ,
3  CHANGE COLUMN `temperature` `temperature` FLOAT NULL DEFAULT NULL ,
4  CHANGE COLUMN `humidity` `humidity` INT NULL DEFAULT NULL ,
5  CHANGE COLUMN `cloudiness` `cloudiness` INT NULL DEFAULT NULL ,
6  CHANGE COLUMN `wind_speed` `wind_speed` FLOAT NULL DEFAULT NULL ,
7  CHANGE COLUMN `wind_gust` `wind_gust` FLOAT NULL DEFAULT NULL ,
8  CHANGE COLUMN `city` `city` VARCHAR(100) NULL DEFAULT NULL ,
9  CHANGE COLUMN `city_id` `city_id` VARCHAR(10) NULL DEFAULT NULL ;
```

Listing 9: Changing the datatypes of the `weather` table.

Lastly, the foreign key could be set:

```
1  ALTER TABLE `gans`.`weather`
2  ADD INDEX `weather_city_id_idx` (`city_id` ASC) VISIBLE;
3  ALTER TABLE `gans`.`weather`
4  ADD CONSTRAINT `weather_city_id`
5    FOREIGN KEY (`city_id`)
6    REFERENCES `gans`.`cities` (`city_id`)
7    ON DELETE NO ACTION
8    ON UPDATE NO ACTION;
```

Listing 10: Set foreign key for the `weather` table.

After changing all of the table's column's datatypes and defining primary and foreign keys, the database was set. Figure 1 shows the final schema with all necessary relations between the tables. The entire MySQL script that leads to this very schema is appended in the Listing 11 in Appendix A.1.
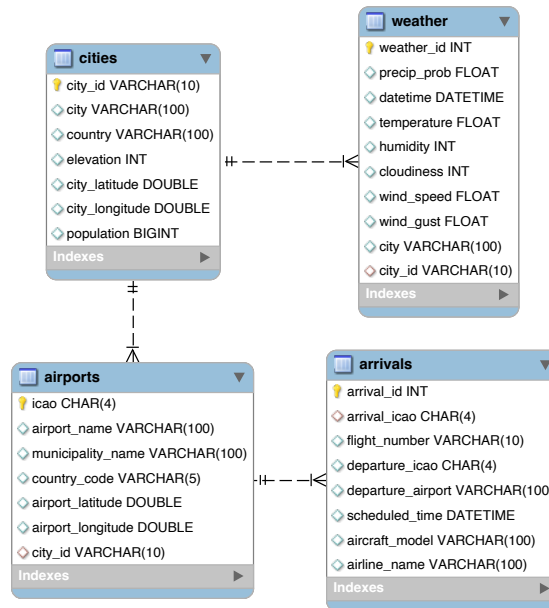


Figure 1: MySQL schema.

# 3. AWS

Gathering dynamic data can be facilitated by cloud computing and cloud storage. Both the prior Python code as well als the MySQL database can be transferred to the cloud. This enables API calls to be periodically triggered as well as their data to be automatically pushed to a cloud instance of a MySQL database. The service provider of choice will be Amazon Web Services (AWS)

In order to have a cloud instance of a MySQL database, we deployed one in AWS' Relational Database Service. Now, it would not be any local database that would be filled with data but the one in the AWS cloud servers.

Similarly, Python code would not be executed locally but also in the cloud. Therefore, we used AWS' Lambda functions. Their Lambda handler functions allowed us to execute the API calls, database connections and data pushes to this very database in one go for `cities`, `weather`, `airports` and `arrivals`. The contents for all four of them are appended in Appendix A.2.

As a final step, two of the Lambda functions had to be automated. While `cities` and `airports` are rather static in nature `weather` and flight `arrivals` data are dynamic and thus subject to rapid changes over time. Therefore, we set for the latter two Lambda functions a custom trigger, namely an EventBridge (CloudWatch Events) that pulled data from the OpenWeather and AeroDataBox APIs on a daily basis by deploying `59 21 ? * * *` as a Cron expression.

## Conclusion

This technical report documented the data flow from API to AWS to SQL. Figure 2 illustrates the pipeline. Data is pulled from APIs for `cities`, `weather`, `airports` and `arrivals` data. The API for `cities` takes WikiData IDs and returns the respecitve city name as well as information about its elevation, latitudinal and longitudinal geo data and population. The city name is then piped into the API for the `weather` call. This provides us with information about precipitation probability, temperature, humidity, cloudiness as well as with wind speed and gust. Apart from that, latitudinal and longitudinal data from `cities` is also used by the API for `airports` which returns information about icao codes, airport names, municipality names, county codes as well as latitudinal and longitudinal geo data. Finally, the ICAO airport codes from `airports` are used to gather flights `arrivals` data about flight numbers, ICAO codes of the departure airports, scheduled arrival times, airline names and aircraft models.

Dynamic `weather` and flight `arrivals` data were pulled on a daily basis. All data was pushed to a MySQL database instance on AWS' RDS.

This pipeline allows the database to be queried, analyze and visualized by Python pandas.
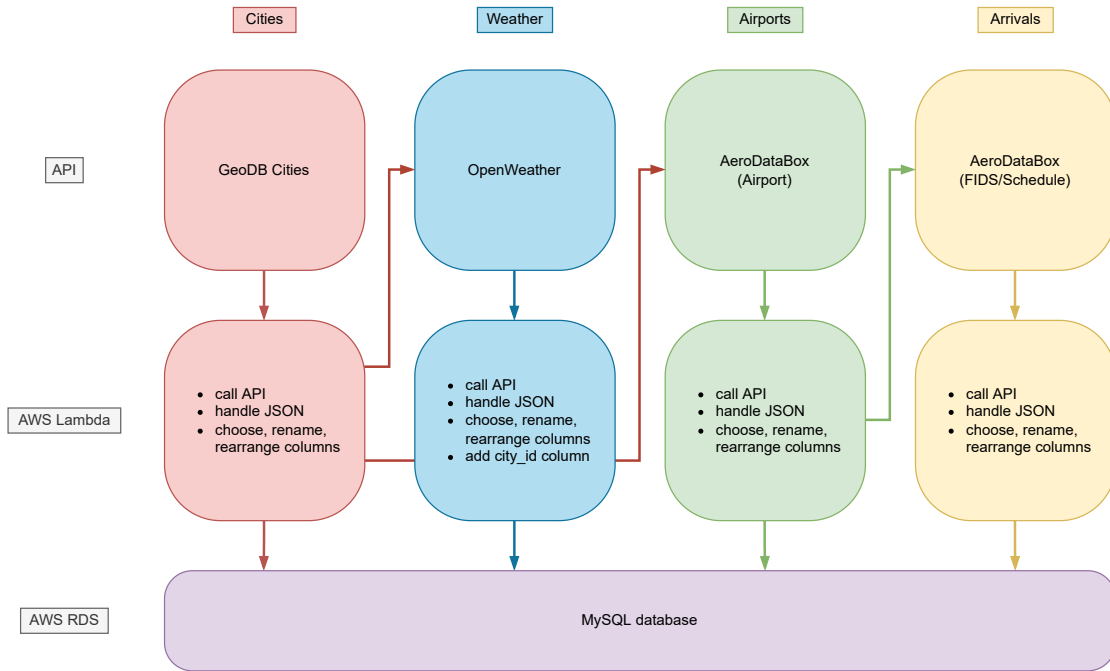
Figure 2: Pipeline from API to database.

# A. Appendix

## A.1. SQL

```
1  CREATE DATABASE gans;
2  USE gans;
3
4  ---------------------------------------------------------------------------
5
6  -- SETTING UP CITIES:
7  -- Alter table datatypes of the columns and set primary key:
8  ALTER TABLE `gans`.`cities`
9  CHANGE COLUMN `city_id` `city_id` VARCHAR(10) NOT NULL ,
10 CHANGE COLUMN `city` `city` VARCHAR(100) NULL DEFAULT NULL ,
11 CHANGE COLUMN `country` `country` VARCHAR(100) NULL DEFAULT NULL ,
12 CHANGE COLUMN `elevation` `elevation` INT NULL DEFAULT NULL ,
13 ADD PRIMARY KEY (`city_id`);
14
15 ---------------------------------------------------------------------------
16
17 -- SETTING UP WEATHER:
18 -- BEFORE altering datatypes and setting keys: Weather has duplicate
      ↪ datetime values; therefore a new column with unique values is
      ↪ necessary:
19 ALTER TABLE weather
20 ADD weather_id INT NOT NULL AUTO_INCREMENT PRIMARY KEY FIRST;
```

7

```sql
21
22 -- Alter table datatypes of the columns and set primary key:
23 ALTER TABLE `gans`.`weather`
24 CHANGE COLUMN `precip_prob` `precip_prob` FLOAT NULL DEFAULT NULL ,
25 CHANGE COLUMN `temperature` `temperature` FLOAT NULL DEFAULT NULL ,
26 CHANGE COLUMN `humidity` `humidity` INT NULL DEFAULT NULL ,
27 CHANGE COLUMN `cloudiness` `cloudiness` INT NULL DEFAULT NULL ,
28 CHANGE COLUMN `wind_speed` `wind_speed` FLOAT NULL DEFAULT NULL ,
29 CHANGE COLUMN `wind_gust` `wind_gust` FLOAT NULL DEFAULT NULL ,
30 CHANGE COLUMN `city` `city` VARCHAR(100) NULL DEFAULT NULL ,
31 CHANGE COLUMN `city_id` `city_id` VARCHAR(10) NULL DEFAULT NULL ;
32
33 -- Setting foreign key:
34 ALTER TABLE `gans`.`weather`
35 ADD INDEX `weather_city_id_idx` (`city_id` ASC) VISIBLE;
36 ALTER TABLE `gans`.`weather`
37 ADD CONSTRAINT `weather_city_id`
38   FOREIGN KEY (`city_id`)
39   REFERENCES `gans`.`cities` (`city_id`)
40   ON DELETE NO ACTION
41   ON UPDATE NO ACTION;
42
43 -------------------------------------------------------------------------
44
45 -- SETTING UP ARRIVALS:
46 -- BEFORE altering datatypes and setting keys: Arrivals has duplicate
      ↪ datetime values; therefore a new column with unique values is
      ↪ necessary:
47 ALTER TABLE arrivals
48 ADD arrival_id INT NOT NULL AUTO_INCREMENT PRIMARY KEY FIRST;
49
50 -- Alter table datatypes of the columns and set primary key:
51 ALTER TABLE `gans`.`arrivals`
52 CHANGE COLUMN `arrival_icao` `arrival_icao` CHAR(4) NULL DEFAULT NULL ,
53 CHANGE COLUMN `flight_number` `flight_number` VARCHAR(10) NULL DEFAULT
      ↪ NULL ,
54 CHANGE COLUMN `departure_icao` `departure_icao` CHAR(4) NULL DEFAULT NULL
      ↪   ,
55 CHANGE COLUMN `departure_airport` `departure_airport` VARCHAR(100) NULL
      ↪ DEFAULT NULL ,
56 CHANGE COLUMN `scheduled_time` `scheduled_time` DATETIME NULL DEFAULT
      ↪ NULL ,
57 CHANGE COLUMN `aircraft_model` `aircraft_model` VARCHAR(100) NULL DEFAULT
      ↪   NULL ,
58 CHANGE COLUMN `airline_name` `airline_name` VARCHAR(100) NULL DEFAULT
      ↪ NULL ;
59
60 -- Setting foreign key:
61 ALTER TABLE `gans`.`arrivals`
62 ADD INDEX `arrivals_icao_idx` (`arrival_icao` ASC) VISIBLE;
63 ALTER TABLE `gans`.`arrivals`
64 ADD CONSTRAINT `arrivals_icao`
65   FOREIGN KEY (`arrival_icao`)
66   REFERENCES `gans`.`airports` (`icao`)
```

```
67    ON DELETE NO ACTION
68    ON UPDATE NO ACTION ;
69
70  --------------------------------------------------------------------------
71
72  -- SETTING UP AIRPORTS :
73  -- Alter table datatypes of the columns and set primary key :
74  ALTER TABLE `gans `.` airports `
75  CHANGE COLUMN `icao ` `icao ` CHAR (4) NOT NULL ,
76  CHANGE COLUMN `airport_name ` `airport_name ` VARCHAR (100) NULL DEFAULT
         ↪ NULL ,
77  CHANGE COLUMN `municipality_name ` `municipality_name ` VARCHAR (100) NULL
         ↪ DEFAULT NULL ,
78  CHANGE COLUMN `country_code ` `country_code ` VARCHAR (5) NULL DEFAULT NULL
         ↪ ,
79  CHANGE COLUMN `city_id ` `city_id ` VARCHAR (10) NULL DEFAULT NULL ,
80  ADD PRIMARY KEY (`icao `);
81
82  -- Drop localCode column from airports :
83  ALTER TABLE `gans `.` airports `
84  DROP COLUMN `localCode `;
85
86  -- Setting foreign key :
87  ALTER TABLE `gans `.` airports `
88  ADD INDEX `airports_city_id_idx ` (`city_id ` ASC) VISIBLE ;
89  ALTER TABLE `gans `.` airports `
90  ADD CONSTRAINT `airports_city_id `
91    FOREIGN KEY (`city_id `)
92    REFERENCES `gans `.` cities ` (`city_id `)
93    ON DELETE NO ACTION
94    ON UPDATE NO ACTION ;
```

Listing 11: MySQL script to set up database schema.

## A.2. AWS Lambda functions

### A.2.1. Cities

```
1   import requests
2   import pymysql
3   import sqlalchemy
4   import pandas as pd
5   import json
6   import time
7   from datetime import datetime , date , timedelta
8   from pytz import timezone
9
10  def lambda_handler (event , context ):
11      # Define cities
12      qcities = [# 10 German cities by population :
13      'Q64 ',      # Berlin
14      #'Q1055 ',   # Hamburg
15      #'Q1726 ',   # München
```

```python
     'Q365',     # Köln
#'Q1794',   # Frankfurt am Main
#'Q1022',   # Stuttgart
     'Q1718',    # Düsseldorf
#'Q1295',   # Dortmund
#'Q2066',   # Essen
#'Q2079',   # Leipzig

            # 5 Austrian cities by population:
#'Q1741',   # Wien
#'Q13298', # Graz
#'Q41329', # Linz
#'Q34713', # Salzburg
#'Q1735'   # Innsbruck
]

# Get cities:
def get_cities(qcities):

    # Connect to database:
    schema="gans"
    host='...rds.amazonaws.com'
    user="admin"
    password='...'
    port=3306
    con = f'mysql+pymysql://{user}:{password}@{host}:{port}/{schema}'

    # Make cities_df:
    cities_list = []
    for qcity in qcities:
        url = f"https://wft-geo-db.p.rapidapi.com/v1/geo/cities/{
    qcity}"
        headers = {"X-RapidAPI-Host": "wft-geo-db.p.rapidapi.com", "X
    -RapidAPI-Key": "..."}
        response = requests.request("GET", url, headers=headers)
        time.sleep(2)
        city_df = pd.json_normalize(response.json())
        cities_list.append(city_df)
    cities_df = pd.concat(cities_list, ignore_index = True)
    cities_df = cities_df[["data.wikiDataId",
                    "data.city",
                    "data.country",
                    "data.elevationMeters",
                    "data.latitude",
                    "data.longitude",
                    "data.population"]]
    cities_df.rename(columns = {'data.wikiDataId': 'city_id',
                       'data.city': 'city' ,
                       'data.country': 'country',
                       'data.elevationMeters': 'elevation',
                       'data.latitude': 'city_latitude',
                       'data.longitude': 'city_longitude',
                       'data.population': 'population'},
                   inplace = True)
```

```
68
69          # Push to MySQL:
70          cities_df.to_sql('cities', if_exists='append', con=con, index=
    ↪ False)
71
72      get_cities(qcities)
```

Listing 12: Lambda function pulls cities data from API, applies post-hoc processing to
dataframe, connects to the RDS cloud instance MySQL and pushes the data
to the database.

### A.2.2. Weather

```
1  import requests
2  import pymysql
3  import sqlalchemy
4  import pandas as pd
5  import json
6  import time
7  from datetime import datetime, date, timedelta
8  from pytz import timezone
9
10 def get_weather(cities):
11     df_list = []
12     for city in cities:
13         url = f"http://api.openweathermap.org/data/2.5/forecast?q={city}&
    ↪ appid=...&units=metric"
14         weather = requests.get(url)
15         weather_df = pd.json_normalize(weather.json()["list"])
16         weather_df["city"] = city
17         df_list.append(weather_df)
18     weather_combined = pd.concat(df_list, ignore_index = True)
19     weather_combined = weather_combined[["pop", "dt_txt", "main.temp", "
    ↪ main.humidity", "clouds.all", "wind.speed", "wind.gust", "city"]]
20     weather_combined.rename(columns = {'pop': 'precip_prob',
21                             'dt_txt': 'datetime',
22                             'main.temp': 'temperature',
23                             'main.humidity': 'humidity',
24                             'clouds.all': 'cloudiness',
25                             'wind.speed': 'wind_speed',
26                             'wind.gust': 'wind_gust'},
27                 inplace = True)
28     weather_combined['datetime'] = pd.to_datetime(weather_combined['
    ↪ datetime'])
29     return weather_combined
30
31 def lambda_handler(event, context):
32     # Connect to database:
33     schema="gans"
34     host='...rds.amazonaws.com'
35     user="admin"
36     password='...'
```

```
37     port=3306
38     con = f'mysql+pymysql://{user}:{password}@{host}:{port}/{schema}'
39
40     # Make new df for city_id and city columns from SQL query so that we
       ↪ eventually have the city_id as a foreign key in the weather df:
41     city_df = pd.read_sql('SELECT city_id, city FROM cities', con=con)
42
43     # Use the queried SQL df to have the city names as a later input for
       ↪ the get_weather function:
44     cities = city_df['city'].to_list()
45
46     # Calling the get_weather function:
47     weather_output = get_weather(cities)
48
49     # Merging the queried SQL df on 'city' onto the weather_output df
       ↪ which results in an additional column with 'city_id':
50     weather_d = weather_output.merge(city_df, how = 'left')
51
52     # Push to MySQL:
53     weather_d.to_sql('weather', if_exists='append', con=con, index=False)
       ↪
```

Listing 13: Lambda function pulls weather data from API, applies post-hoc processing to dataframe, connects to the RDS cloud instance MySQL and pushes the data to the database.

### A.2.3. Airports

```
1  import requests
2  import pymysql
3  import sqlalchemy
4  import pandas as pd
5  import json
6  import time
7  from datetime import datetime, date, timedelta
8  from pytz import timezone
9
10 def get_airports(lat, lon):
11     airport_list = []
12     for i in range(len(lat)):
13         url = f"https://aerodatabox.p.rapidapi.com/airports/search/
       ↪ location/{lat[i]}/{lon[i]}/km/50/10"
14         querystring = {"withFlightInfoOnly":"true"}
15         headers = {"X-RapidAPI-Host": "aerodatabox.p.rapidapi.com", "X-
       ↪ RapidAPI-Key": "..."}
16         response = requests.request("GET", url, headers=headers, params=
       ↪ querystring)
17         airport_df = pd.json_normalize(response.json()["items"])
18         airport_list.append(airport_df)
19     airports_df = pd.concat(airport_list, ignore_index = True)
20     airports_df.drop_duplicates(subset ='icao', inplace = True)
21     airports_df = airports_df[~airports_df.name.str.contains("Air Base",
       ↪ case = False)]
```

```python
22    airports_df.drop(columns = ['iata', 'shortName'], inplace = True)
23    airports_df.reset_index(drop = True, inplace = True)
24    airports_df['city_id'] = [
25                             'Q64',
26                             'Q64',
27                             #'Q1055',
28                             #'Q1726',
29                             'Q365',
30                             'Q1718',
31                             #'Q1794',
32                             #'Q1022',
33                             #'Q1295',
34                             #'Q2079',
35                             #'Q1741',
36                             #'Q13298',
37                             #'Q41329',
38                             #'Q34713',
39                             #'Q1735'
40                             ]
41    airports_df.rename(columns = {'name': 'airport_name',
42                                  'municipalityName': 'municipality_name',
43                                  'countryCode': 'country_code',
44                                  'location.lat': 'airport_latitude',
45                                  'location.lon': 'airport_longitude',
46                                  },
47                                  inplace = True)
48    return airports_df
49
50 def lambda_handler(event, context):
51    # Connect to database:
52    schema="gans"
53    host='...rds.amazonaws.com'
54    user="admin"
55    password='...'
56    port=3306
57    con = f'mysql+pymysql://{user}:{password}@{host}:{port}/{schema}'
58
59    # Make new df for city_id and city columns from SQL query so that we
   ↪ eventually have the city_id as a foreign key in the weather df:
60    city_df = pd.read_sql('SELECT city_latitude, city_longitude FROM
   ↪ cities', con=con)
61
62    # Calling the get_airports function:
63    airports_output = get_airports(city_df['city_latitude'].to_list(),
   ↪ city_df['city_longitude'])
64
65    # Push to MySQL:
66    airports_output.to_sql('airports', if_exists='append', con=con, index
   ↪ =False)
```

Listing 14: Lambda function pulls cities data from API, applies post-hoc processing to dataframe, connects to the RDS cloud instance MySQL and pushes the data to the database.

### A.2.4. Arrivals

```python
import requests
import pymysql
import sqlalchemy
import pandas as pd
import json
import time
from datetime import datetime, date, timedelta
from pytz import timezone

def get_arrivals(icao):
  arrival_list = []
  today = datetime.now().astimezone(timezone('Europe/Berlin')).date()
  tomorrow = (today + timedelta(days = 1))
  for code in icao:
    url = f"https://aerodatabox.p.rapidapi.com/flights/airports/icao/{
    ↪ code}/{tomorrow}T10:00/{tomorrow}T22:00"
    querystring = {"withLeg":"false","direction":"Arrival","withCancelled
    ↪ ":"false","withCodeshared":"true",
                   "withCargo":"false","withPrivate":"false","
    ↪ withLocation":"false"}
    headers = {"X-RapidAPI-Key": "...", "X-RapidAPI-Host": "aerodatabox.p
    ↪ .rapidapi.com"}
    response = requests.request("GET", url, headers = headers, params =
    ↪ querystring)
    arrival_df = pd.json_normalize(response.json()["arrivals"])
    arrival_df["arrival_icao"] = code
    arrival_list.append(arrival_df)
  arrivals_df = pd.concat(arrival_list, ignore_index = True)
  arrivals_df.drop(columns = ['isCargo',
                              'status',
                              'callSign',
                              'codeshareStatus',
                              'movement.airport.iata',
                              #'movement.actualTimeLocal',
                              'movement.quality',
                              'aircraft.reg',
                              'aircraft.modeS',
                              'movement.terminal',
                              'movement.scheduledTimeUtc',
                              #'movement.actualTimeUtc',
                              #'movement.baggageBelt',
                              #'movement.gate'
                              ],
                   inplace = True)
  arrivals_df.rename(columns = {'number': 'flight_number',
                                'movement.airport.icao': '
    ↪ departure_icao',
                                'movement.airport.name': '
    ↪ departure_airport',
                                'movement.scheduledTimeLocal': '
    ↪ scheduled_time',
                                'aircraft.model': 'aircraft_model',
```

```
45                                     'airline.name': 'airline_name'
46
47   },
48                     inplace = True)
49   arrivals_df['scheduled_time'] = pd.to_datetime(arrivals_df['
       ↪ scheduled_time'])
50   return arrivals_df
51
52 def lambda_handler(event, context):
53     # Connect to database:
54     schema="gans"
55     host='...rds.amazonaws.com'
56     user="admin"
57     password='...'
58     port=3306
59     con = f'mysql+pymysql://{user}:{password}@{host}:{port}/{schema}'
60
61     # Make new df for city_id and city columns from SQL query so that we
       ↪ eventually have the city_id as a foreign key in the weather df:
62     icao_df = pd.read_sql('SELECT icao FROM airports', con=con)
63
64     icao = icao_df['icao'].to_list()
65
66     # Call the get_arrivals:
67     arrivals_output = get_arrivals(icao)
68
69      # Push to MySQL:
70     arrivals_output.to_sql('arrivals', if_exists='append', con=con, index
       ↪ =False)
```

Listing 15: Lambda function pulls flight arrivals data from API, applies post-hoc processing to dataframe, connects to the RDS cloud instance MySQL and pushes the data to the database.