

Swiss Transport ETL Project

This document serves as a short project overview to present some of it's procedures and final outcome.

1. Data source

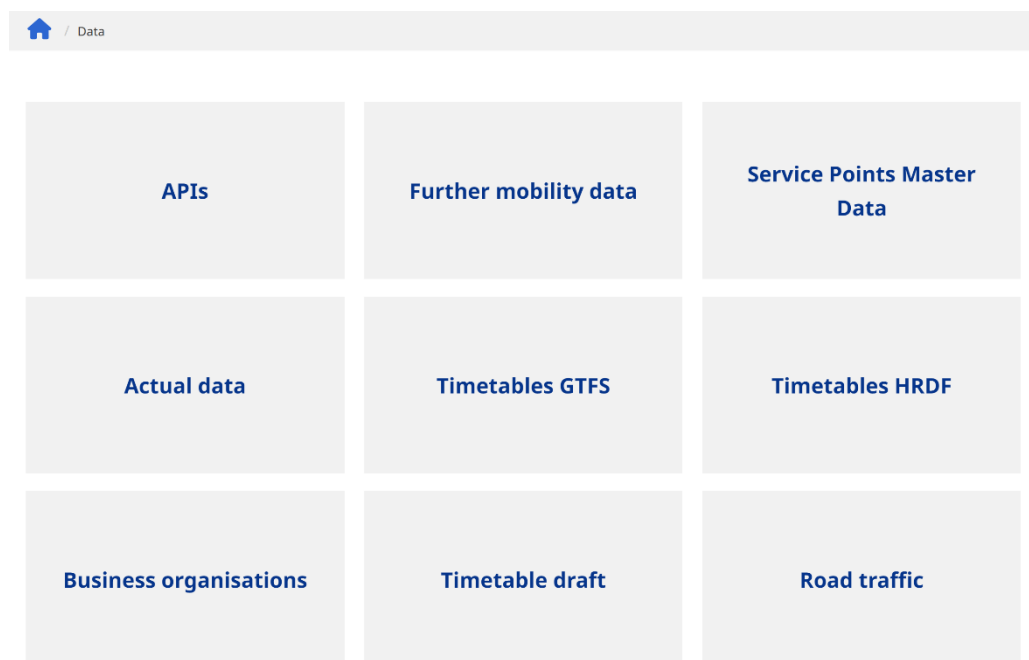
The data source for this project is the Swiss Transport Open Data Mobility Platform - <https://opentransportdata.swiss/en/>

The platform provides many files sorted in categories like Actual Data, Real-time data, Service Mobility Data, etc.

For this particular project, 12 files have been chosen to create a complex database for daily updates. The files are in different formats (.csv, .xlsx, .json) and some of them need additional processing before ingesting into the cloud.

The file downloads are fully automated with Selenium package running in headless mode on Docker container. A standard Chrome GUI mode is also available to run using files from *local_testing_files* folder.

Files are downloaded with different frequencies which are set intuitively and are fully customizable. However, the main file, which contains data on all trips that happened on the preceding day is updated daily.



Pic. 1.1 – Data platform categories

2. Data ingestion

Downloaded files are ingested into the Snowflake internal stage right after download – the python scripts obtains the file name and incorporates SnowSQL to do so.

Each download frequency has its own folder - both Selenium download directory & Snowflake internal stage. Furthermore, each file is being downloaded into separate subfolder. This way files can be downloaded in parallel mode without disrupting filename retrieval process.

An example of PUT command being successfully received in Snowflake is visible below.

SQL Text

```
PUT file:///opt/airflow/transport/daily/main_files/2024-05-30_IstDaten.csv @my_stg/daily auto_compress=true
```

Pic. 2.1 – SnowSQL PUT command sample

3. Data transformation & storage

Data from each file is copied into raw tables in the ‘raw’ schema right after ingestion process is done. In most cases only minor transformations are done during this step.

Files are deleted from the stage folders right after the COPY INTO command is successfully completed.

Further transformations are done while moving data to ‘curated’ & ‘consumption’ schemas – joins, string cleaning, formatting, removing useless data.

Copy processes and transformations are run directly from python scripts using Snowpark.

Query history for automated ‘consumption’ schema tables update is below.

Query History

Status	All	User	USER_01	Last 14 days	Filters	1000+ Queries	Columns	
SQL TEXT	QUERY ID	STATUS	USER	WAREHOUSE	DURATI...	STARTED		
INSERT OVERWRITE INTO consumption.t...	01b4b436-0000-b983-0000-a3a10...	Succe...	USER...	TRANSPORT...	5.8s	5/31/2024, 11:10:27 PM		
INSERT OVERWRITE INTO consumption.v...	01b4b436-0000-b983-0000-a3a10...	Succe...	USER...	TRANSPORT...	644ms	5/31/2024, 11:10:28 P...		
INSERT OVERWRITE INTO consumption.o...	01b4b436-0000-b983-0000-a3a10...	Succe...	USER...	TRANSPORT...	513ms	5/31/2024, 11:10:28 P...		
INSERT OVERWRITE INTO consumption.l...	01b4b436-0000-b98d-0000-a3a10...	Succe...	USER...	TRANSPORT...	591ms	5/31/2024, 11:10:27 PM		
INSERT OVERWRITE INTO consumption.a...	01b4b436-0000-b983-0000-a3a10...	Succe...	USER...	TRANSPORT...	636ms	5/31/2024, 11:10:27 PM		
INSERT OVERWRITE INTO consumption.s...	01b4b436-0000-b98d-0000-a3a10...	Succe...	USER...	TRANSPORT...	589ms	5/31/2024, 11:10:27 PM		
INSERT OVERWRITE INTO consumption.t...	01b4b436-0000-b983-0000-a3a10...	Succe...	USER...	TRANSPORT...	415ms	5/31/2024, 11:10:27 PM		
INSERT OVERWRITE INTO consumption.b...	01b4b436-0000-b98d-0000-a3a10...	Succe...	USER...	TRANSPORT...	413ms	5/31/2024, 11:10:27 PM		

Pic. 3.1 – Query history for consumption schema tables update

The final database structure is a snowflake data model made of 1 fact table and 10 dimension tables, which names, sample code and links you can see below.

SWISS_TRANSPORT

COMMON

CONSUMPTION

Tables

ACCESSIBILITY_DIM

BUSINESS_TYPES_DIM

LINES_DIM

MUNICIPALITY_DIM

OCCUPANCY_DIM

OPERATORS_DIM

PARKING_DIM

STOPS_DIM

TRANSPORT_FACT

TRANSPORT_TYPES_DIM

VEHICLES_DIM

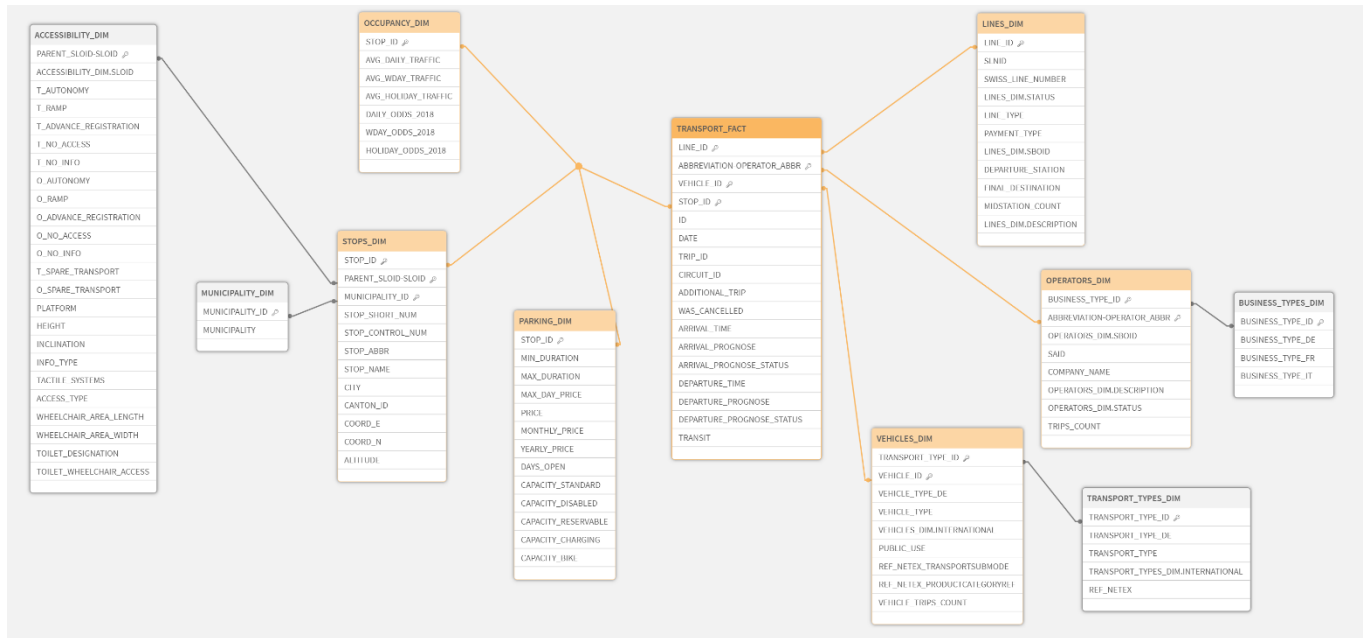
Table definition

```
1 create or replace TABLE SWISS_TRANSPORT.CONSUMPTION.TRANSPORT_FACT (  
2   ID NUMBER(38,0) NOT NULL autoincrement start 1 increment 1 noorder,  
3   DATE DATE,  
4   TRIP_ID VARCHAR(16777216),  
5   OPERATOR_ABBR VARCHAR(16777216),  
6   LINE_ID VARCHAR(16777216),  
7   CIRCUIT_ID VARCHAR(16777216),  
8   VEHICLE_ID VARCHAR(16777216),  
9   ADDITIONAL_TRIP BOOLEAN,  
10  WAS_CANCELLED BOOLEAN,  
11  STOP_ID NUMBER(38,0),  
12  ARRIVAL_TIME TIMESTAMP_NTZ(9),  
13  ARRIVAL_PROGNOSE TIMESTAMP_NTZ(9),  
14  ARRIVAL_PROGNOSE_STATUS VARCHAR(16777216),  
15  DEPARTURE_TIME TIMESTAMP_NTZ(9),  
16  DEPARTURE_PROGNOSE TIMESTAMP_NTZ(9),  
17  DEPARTURE_PROGNOSE_STATUS VARCHAR(16777216),  
18  TRANSIT BOOLEAN,  
19  primary key (ID)  
20 );
```

Show less

Pic. 3.2 – List of final tables

Pic. 3.3 – Sample SQL code – TRANSPORT_FACT table



Pic. 3.4 – ERD diagram of SWISS_TRANSPORT database

4. Orchestration

All above processes are orchestrated using Airflow. To make this project easy to run and more versatile, Airflow does not run locally – Docker container is used instead.

<input type="checkbox"/>	<div><div></div><div></div></div>	<div><div></div><div>swisstransport</div></div>	Running (3/3)	3.59%		0 seconds ago	<div><div></div><div></div><div></div></div>
<input type="checkbox"/>	<div><div></div><div>5bae63bbb2a4</div><div></div></div>	<div><div></div><div>webserver-1</div><div>swisstransport-webserver</div></div>	Running	0.08%	8080:8080	1 second ago	<div><div></div><div></div><div></div></div>
<input type="checkbox"/>	<div><div></div><div>e75fc34df131</div><div></div></div>	<div><div></div><div>scheduler-1</div><div>swisstransport-scheduler</div></div>	Running	2.4%		0 seconds ago	<div><div></div><div></div><div></div></div>
<input type="checkbox"/>	<div><div></div><div>87ccbf44c888</div><div></div></div>	<div><div></div><div>postgres-1</div><div>postgres:13</div></div>	Running	1.11%		1 second ago	<div><div></div><div></div><div></div></div>

Pic. 4.1 – Docker container – Airflow instances

The orchestration is done using 5 DAGs – which of 1 is externally triggered by the successful run of main, daily DAG – which tasks are visible in pic. 4.3.

DAGs

All 5		Active 5	Paused 0	Running 2	Failed 0	Filter DAGs by tag	Search DAGs	Auto-refresh
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	
transport_data_etl_daily	airflow		@daily	2024-05-31, 13:56:52	2024-06-01, 00:00:00			
transport_data_etl_monthly	airflow		@monthly	2024-05-01, 00:00:00	2024-06-01, 00:00:00			
transport_data_etl_weekly	airflow		@weekly	2024-05-19, 00:00:00	2024-05-26, 00:00:00			
transport_data_etl_yearly	airflow		@yearly	2024-05-31, 21:06:14	2024-01-01, 00:00:00			
update_curated_consumption	airflow		None	2024-05-31, 21:10:17				

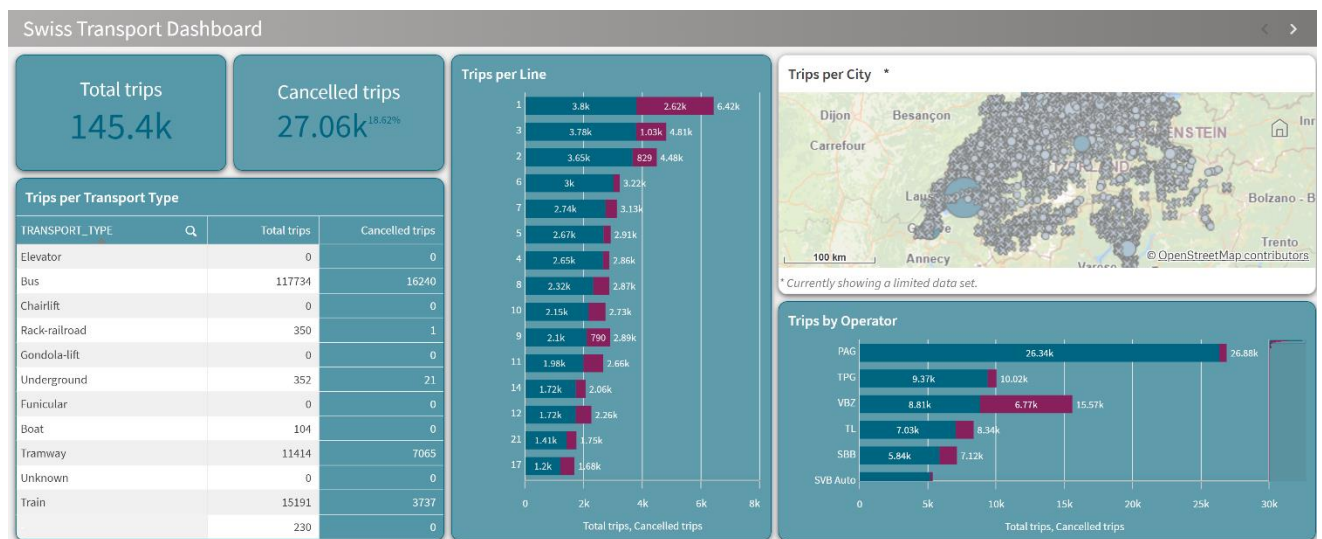
Pic. 4.2 – Airflow web UI – DAG list



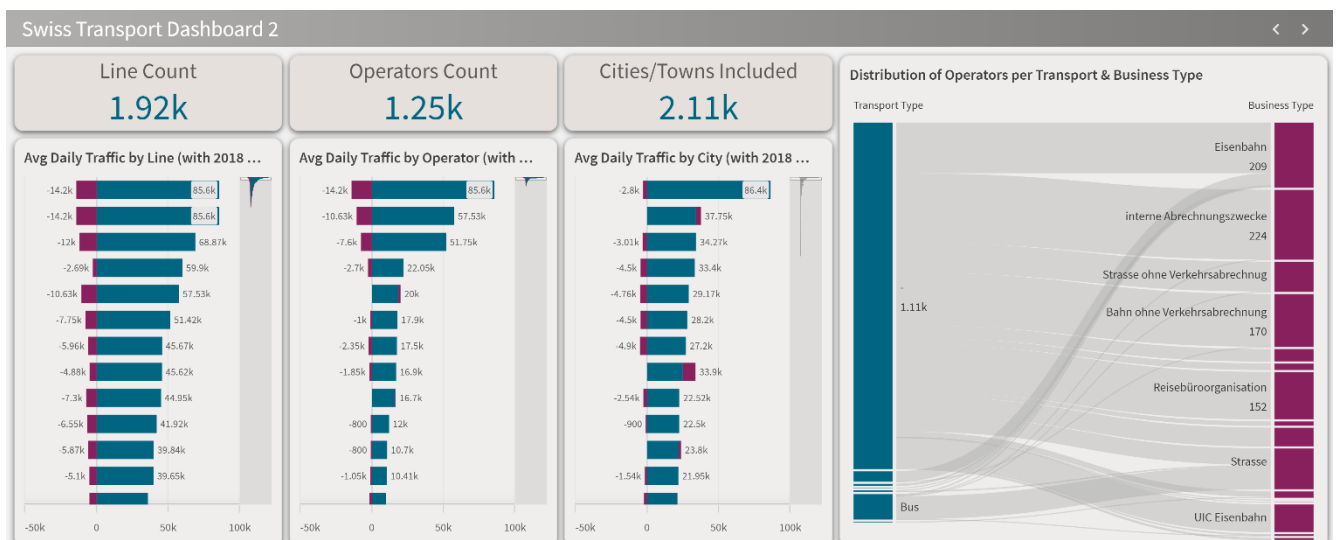
Pic. 4.3 – Sample DAG task dependencies – transport_data_etl_daily DAG

5. Visualization

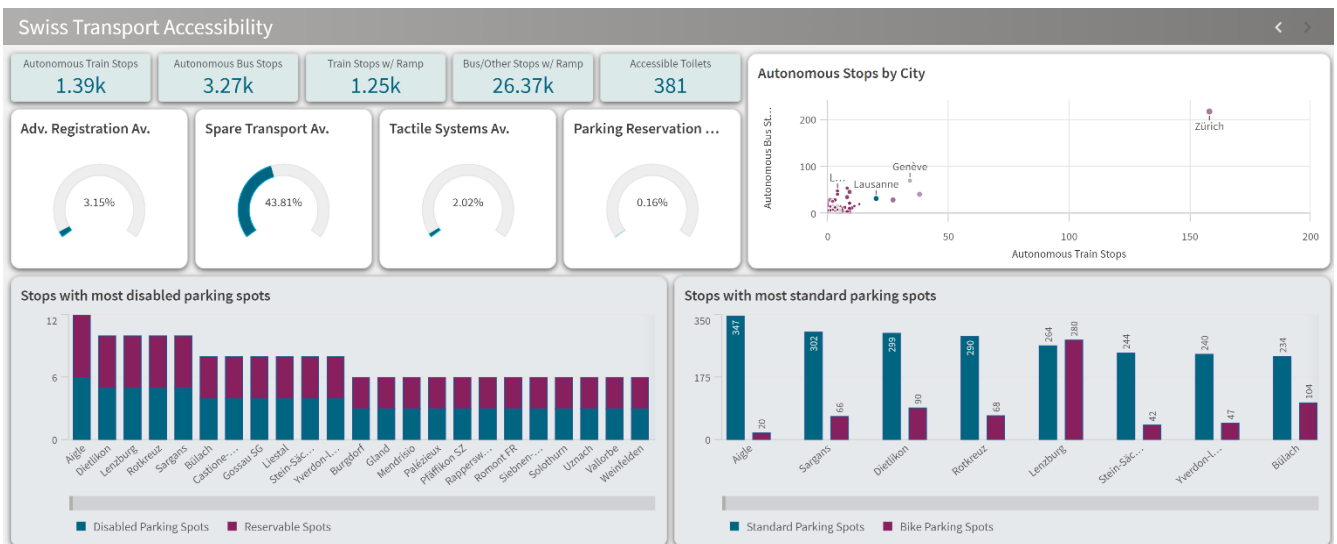
The final data is linked to Qlik Sense, in order to create interactive dashboards. Three different dashboards are created, to show some interesting dependencies, KPIs and more.



Pic. 5.1 – QlikSense dashboard – No 1



Pic. 5.2 – QlikSense dashboard – No 2



Pic. 5.3 – QlikSense dashboard – No 3

6. Summary

A lot of new things were learned during realization of this project in all it's aspects – automation, Snowflake tools, BI visualization.

The final outcome could obviously be expanded by more and more data, possibly including real-time traffic as well.