

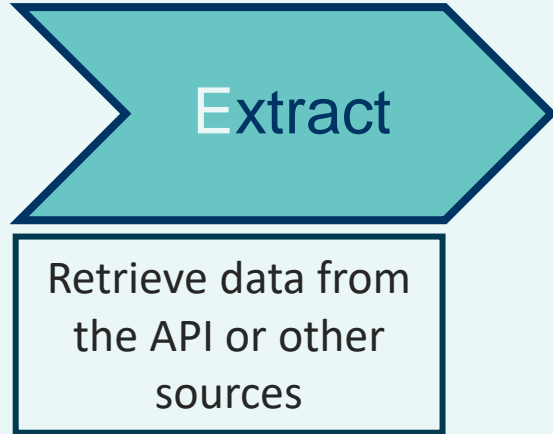
05 – Data Ingestion in FROST via Apache NiFi

Beatrice Olivari

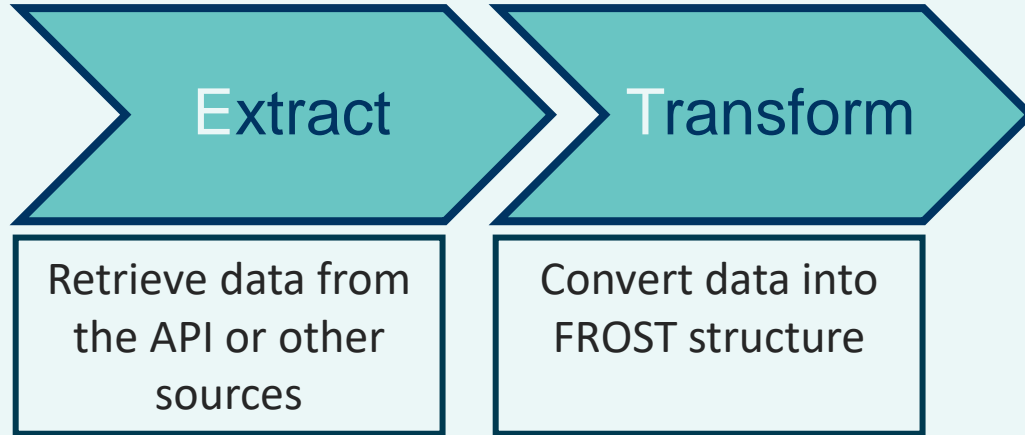


This project has received funding from the European Union's HE research and innovation programme under the grant agreement No. 101057497

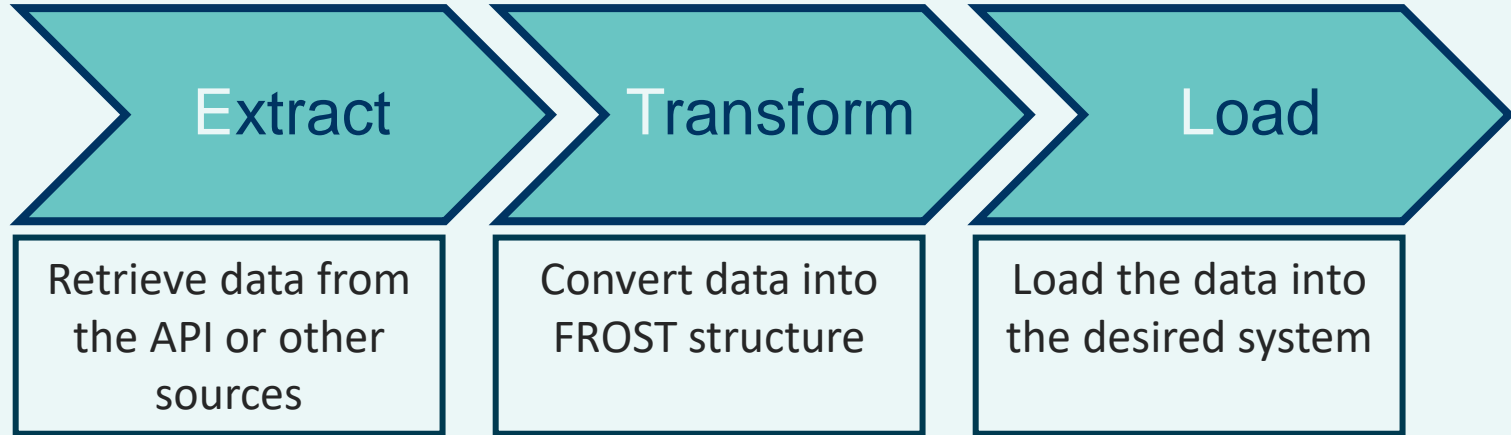
Go with the (ETL) FLOW



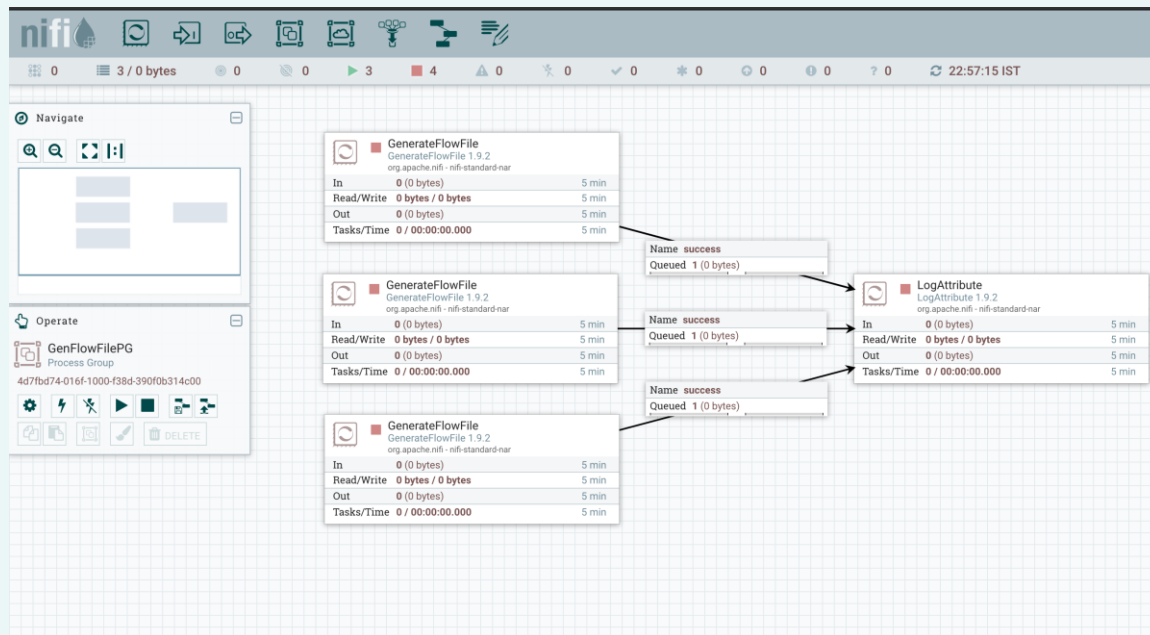
Go with the (ETL) FLOW



Go with the (ETL) FLOW



The tool: Apache NiFi



Example flow: Meteorological data



- MISTRAL is an **EU financed** project **started in 2018**
- The **goal** of the MISTRAL portal is to facilitate and foster the re-use of the datasets by the weather community, as well as by its cross-area communities

Mistral: User Interface

17 out of 21

Regions and Autonomous Provinces

12.177

Weather Stations

9 million +

Observations every day

Filter

Variable

Temperature/dry-bulb Temperature

Date

20/07/2023

11:00 - 13:59

Level

2.000m above ground

Time range

Analysis or observation, istantaneous value

Network

Any

Group of Licenses

CCBY COMPLIANT

Quality Control Filter

ON

Update Map



Data Stations Meteograms



Mistral: User Interface

17 out of 21

Regions and Autonomous Provinces

12.177

Weather Stations

9 million +

Observations every day

Filter

Variable

Temperature/dry-bulb Temperature

Date

20/07/2023

11:00 - 13:59

Level

2.000m above ground

Time range

Analysis or observation, istantaneous value

Network

Any

Group of Licenses

CCBY COMPLIANT

Quality Control Filter

ON

Update Map

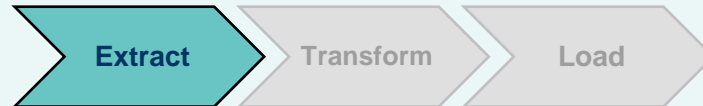


Data Stations Meteograms



Mistral: Data Extraction

How do we extract data in a continued, up to date and automated way?

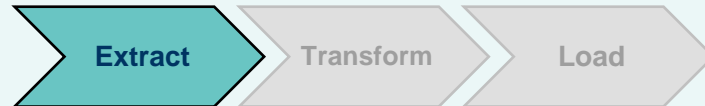


Mistral: Data Extraction

How do we extract data in a continued and automated way?

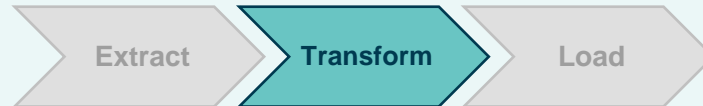


Mistral API



Mistral: Data Transformation

How do we make data STA compliant?

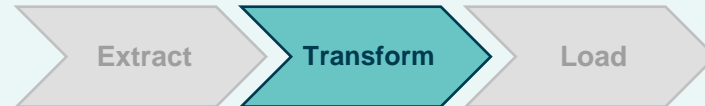


Mistral: Data Transformation

How do we make data STA compliant?



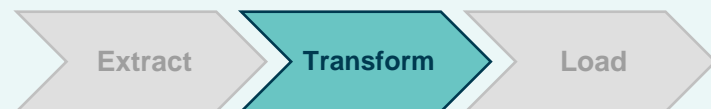
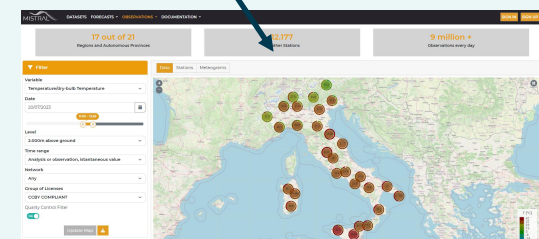
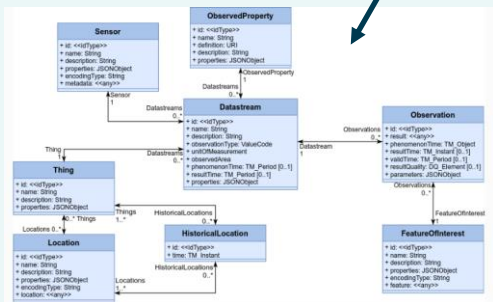
Mapping



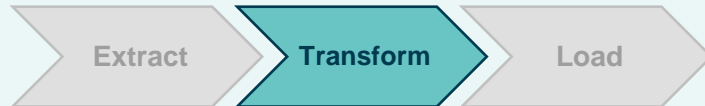
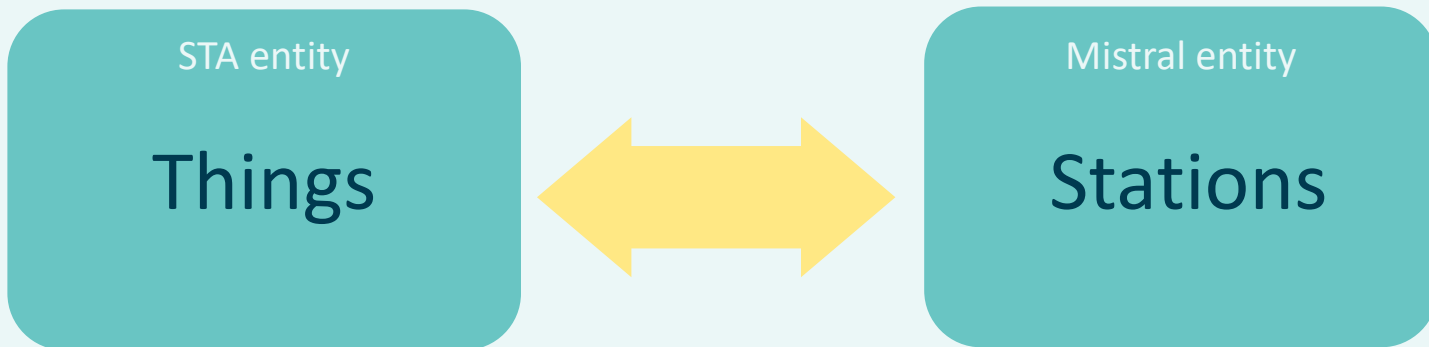
Mistral: Mapping

STA entity

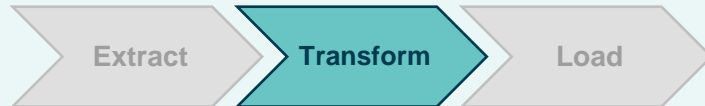
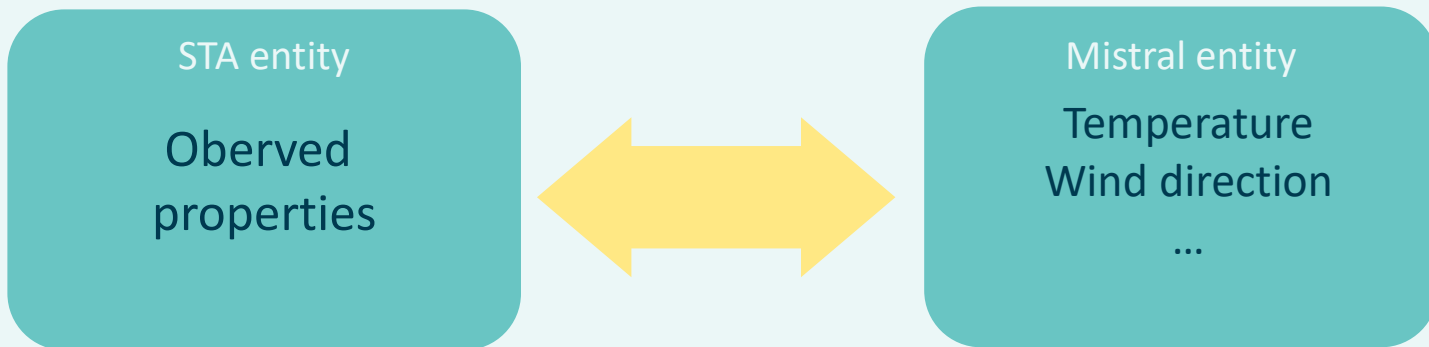
Mistral entity



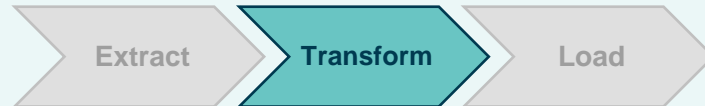
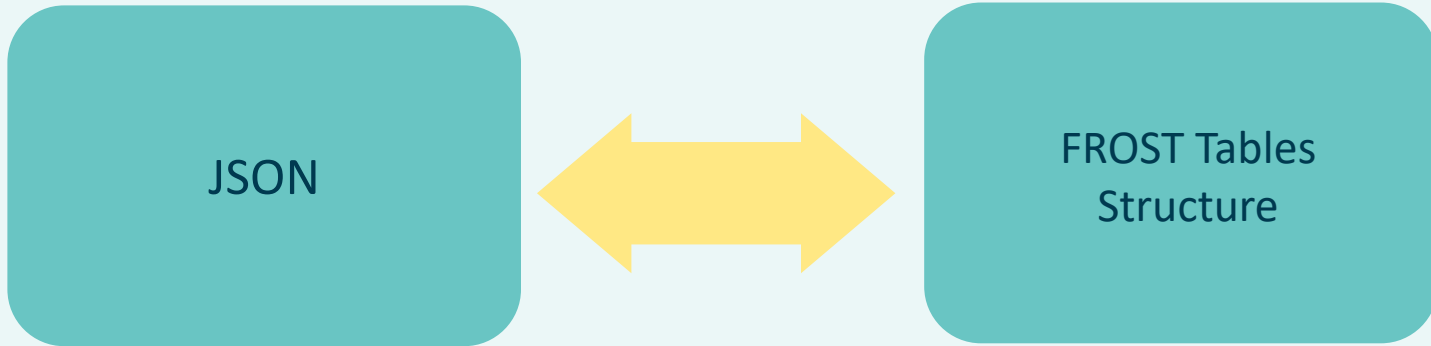
Mistral: Mapping



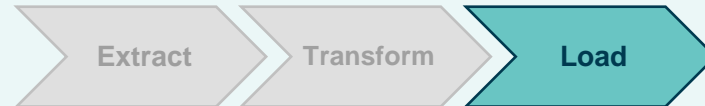
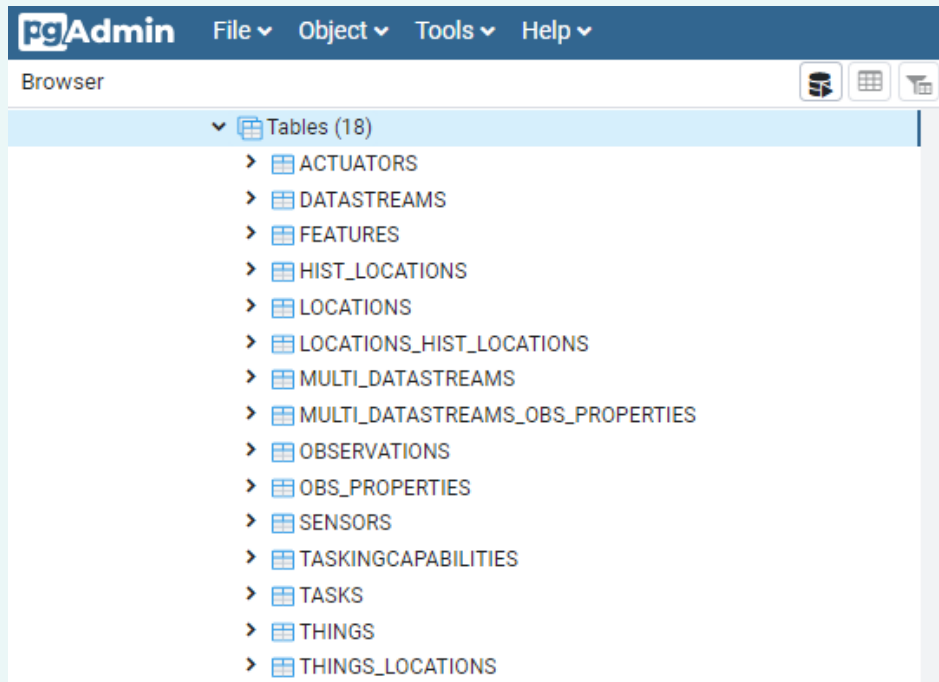
Mistral: Mapping



Mistral: Transformation



Mistral: Load



Mistral: NiFi Flow

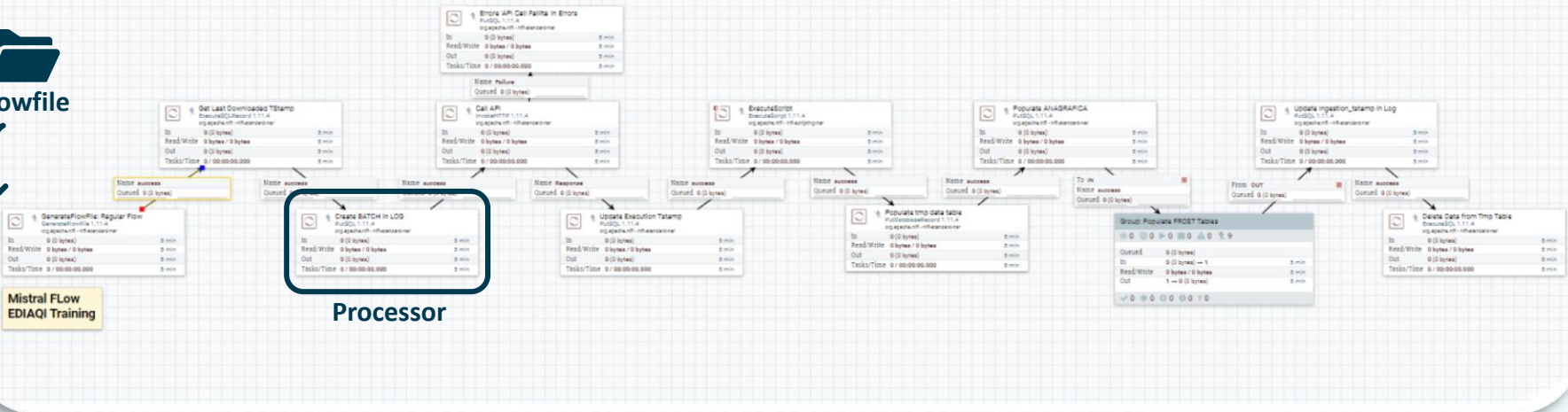


Flowfile

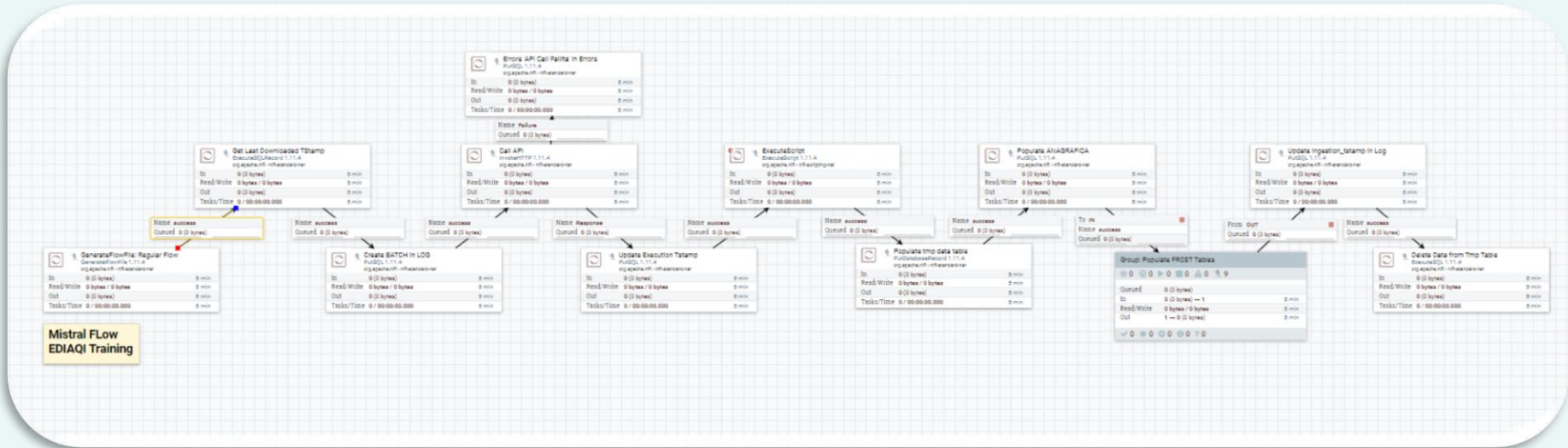


Mistral FLOW
EDIAQI Training

Processor

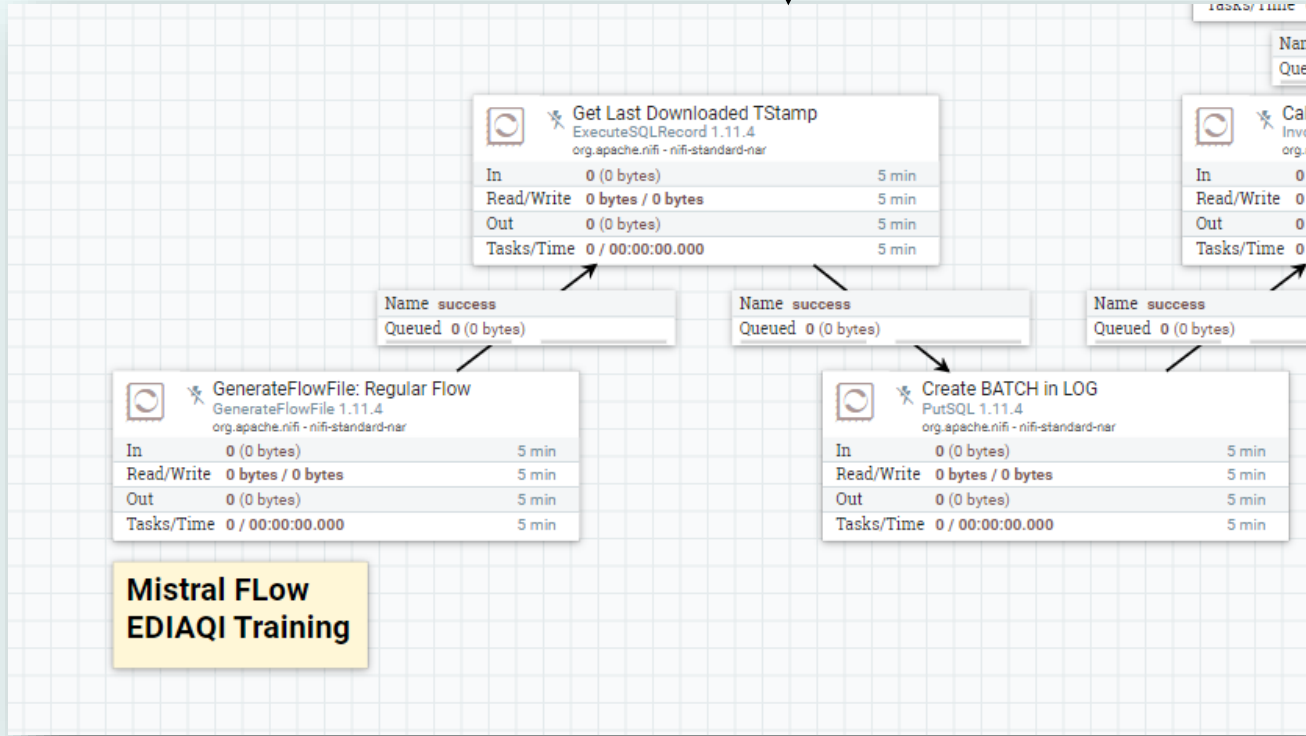
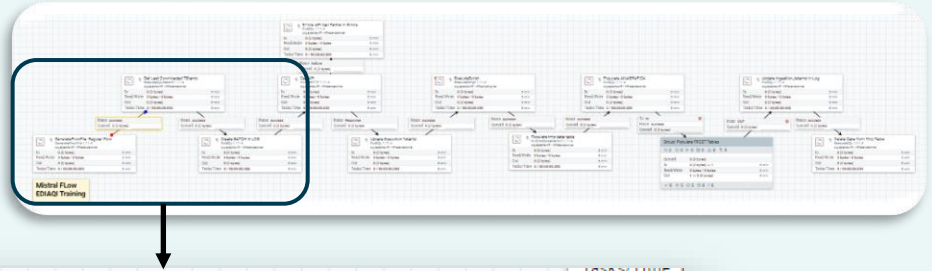


Mistral: NiFi Flow



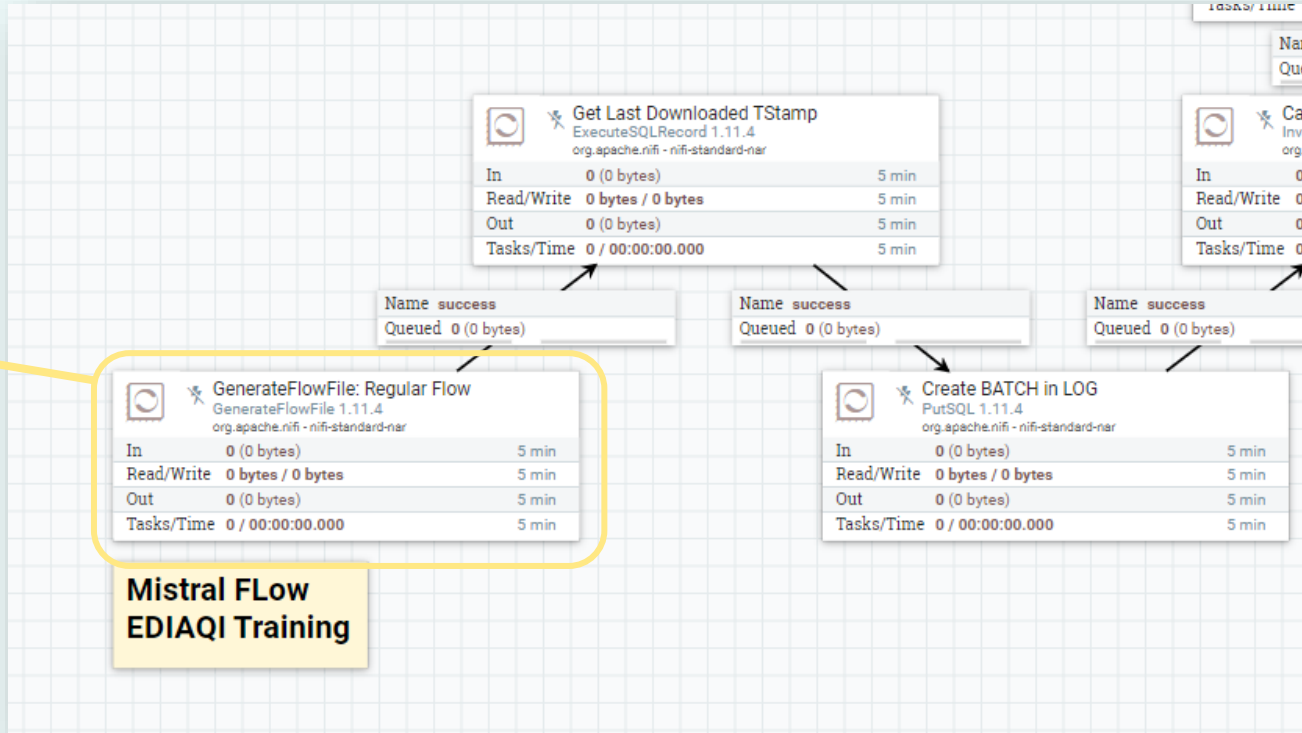
↪ 1 flow trip = 1 day of data ↻

Mistral: NiFi Flow



Mistral FLOW
EDIAQI Training

Mistral: NiFi Flow



Initial processor



1st of January

Mistral FLOW
EDIAQI Training

Mistral: NiFi Flow

In the next rounds,
how do we **keep track** of what day of data
has already been downloaded?



Mistral: NiFi Flow

In the next rounds,
how do we **keep track** of what day of data
has already been downloaded?



Log table

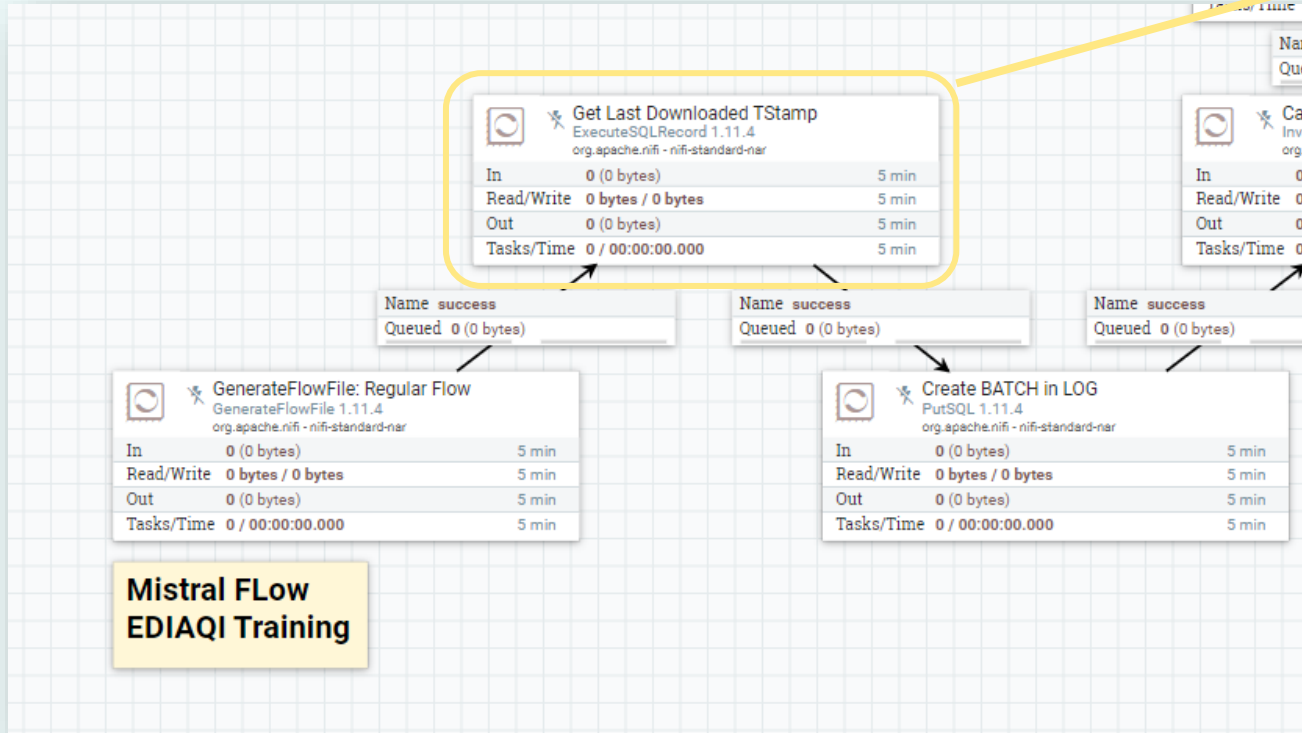
Batch_ID	Start_Date	Stop_Date	API_Excecuton_Tstamp	FROST_Ingestion_Tstamp

Mistral: NiFi Flow

Keeping track of ingested data



Log table



Mistral FLOW
EDIAQI Training

Mistral: NiFi Flow

In the next rounds,
how do we **keep track** of what day of data
has already been downloaded?



Log table

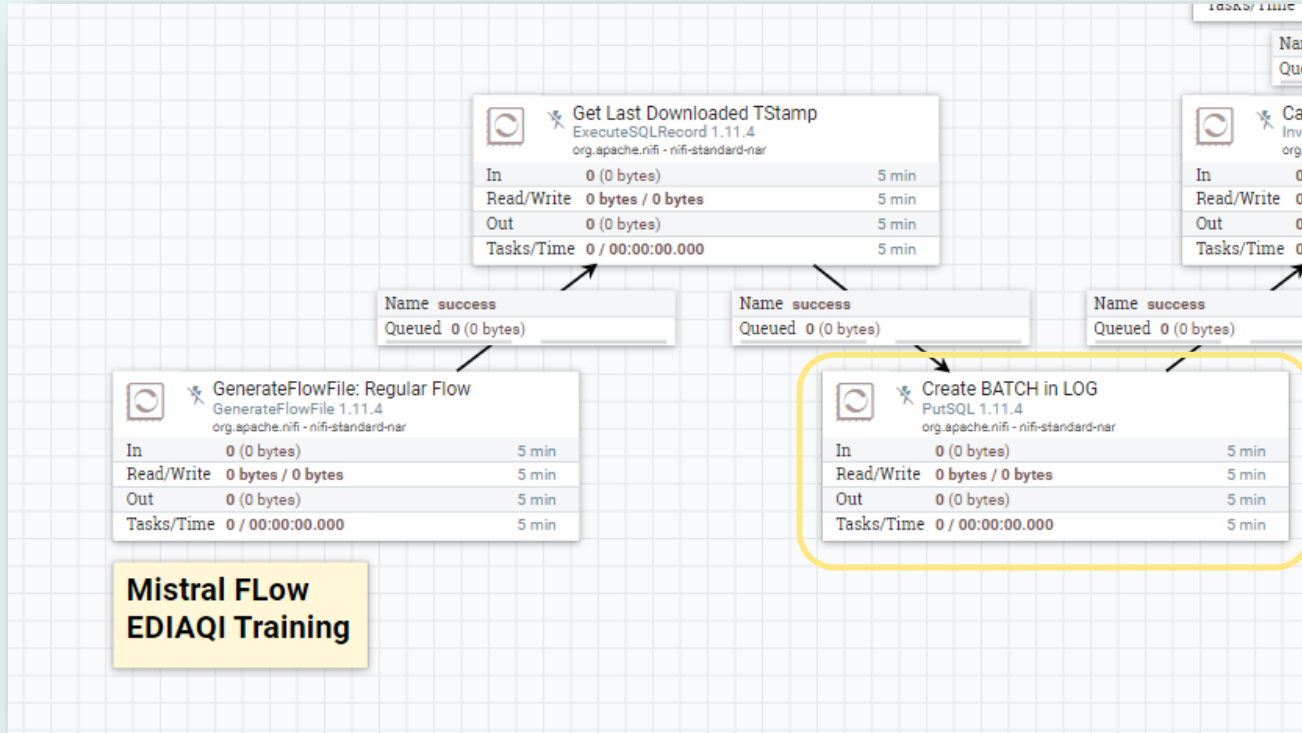
Batch_ID	Start_Date	Stop Date	API_Excecuton_Tstamp	FROST_Ingestion_Tstamp
1	01/01/2023 00:00	01/01/2023 23:59	07/21/2023 10:00	07/21/2023 10:01

Mistral: NiFi Flow

How do we keep track of what data have already been downloaded?



Log table



Mistral: NiFi Flow

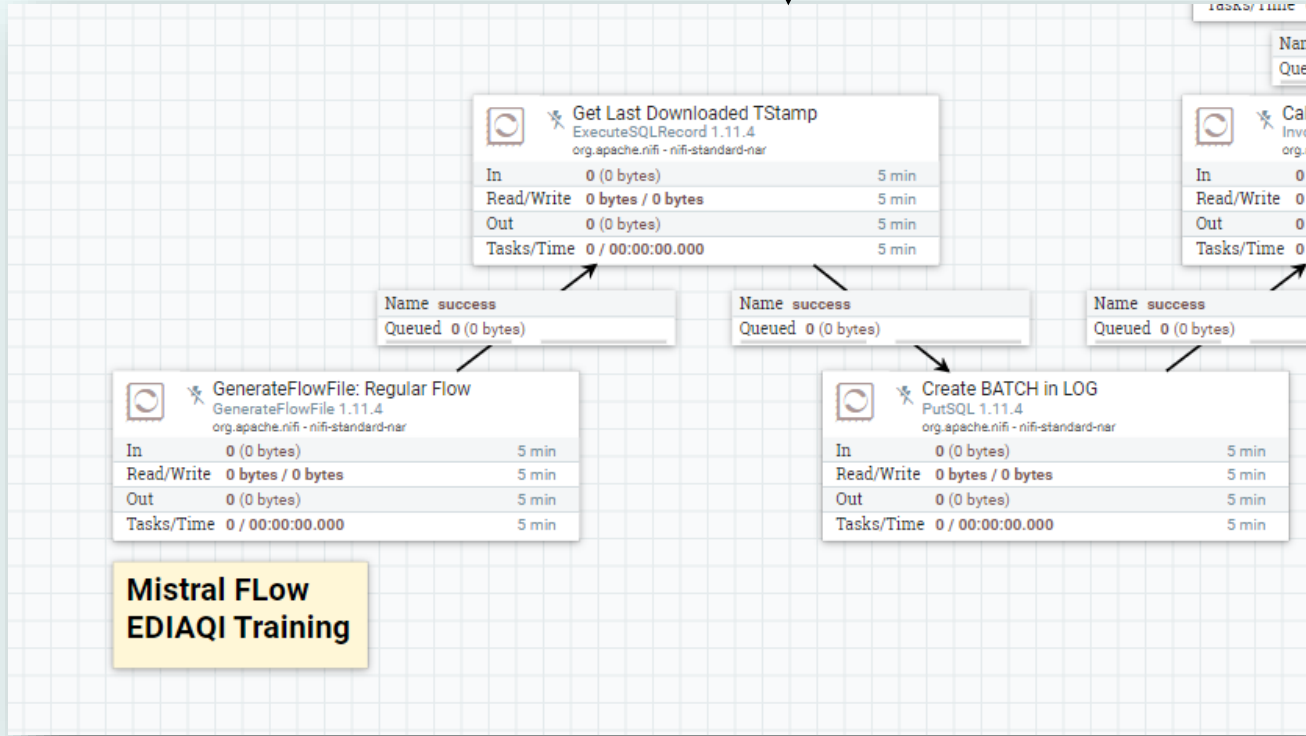
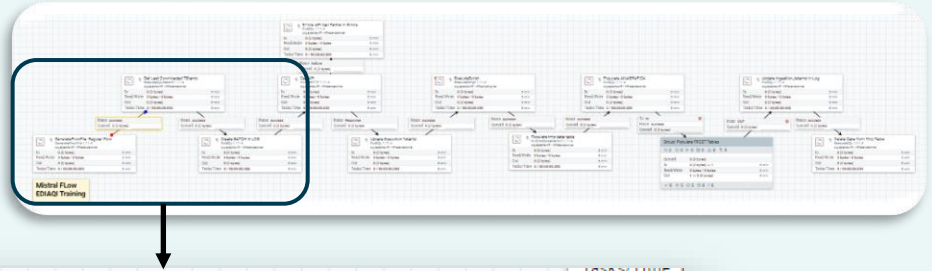
In the next rounds,
how do we **keep track** of what day of data
has already been downloaded?



Log table

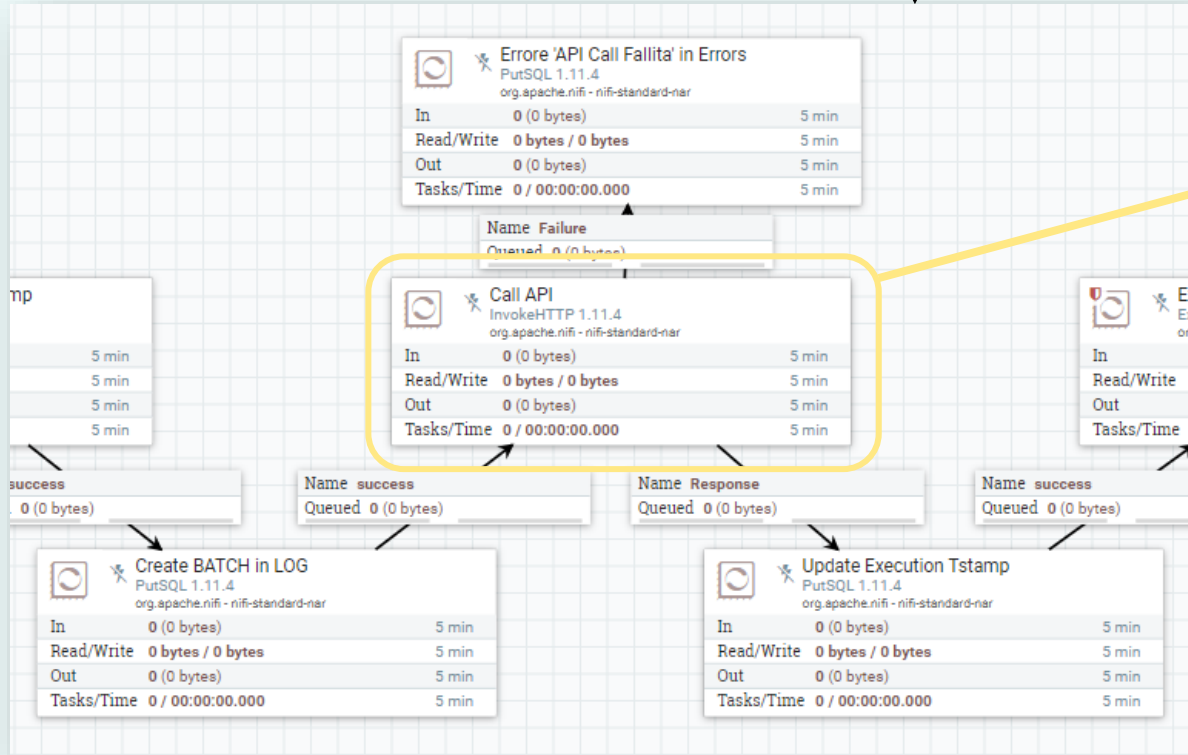
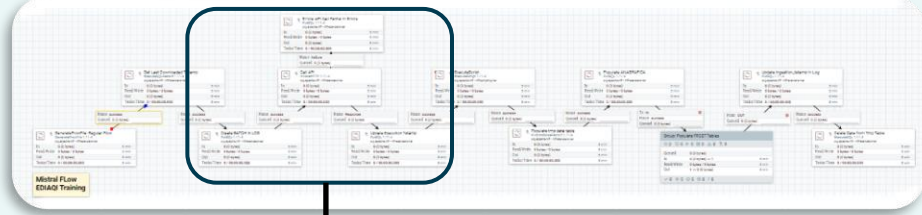
Batch_ID	Start_Date	Stop_Date	API_Excecuion_Tstamp	FROST_Ingestion_Tstamp
1	01/01/2023 00:00	01/01/2023 23:59	07/21/2023 10:00	07/21/2023 10:01
2	02/01/2023 00:00	02/01/2023 23:59		

Mistral: NiFi Flow



Mistral FLOW
EDIAQI Training

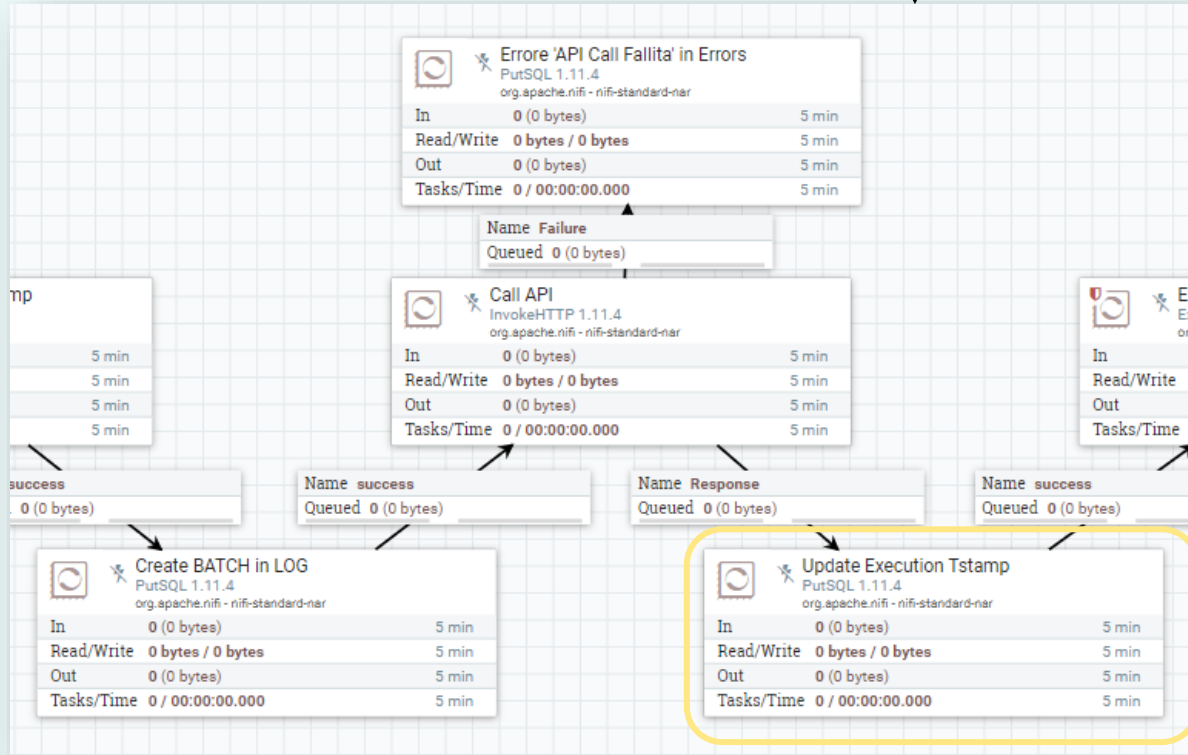
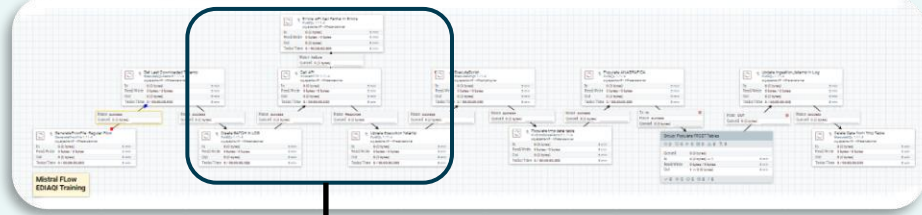
Mistral: NiFi Flow



Endpoint
+ start date
+ stop date
+ variable to
measure



Mistral: NiFi Flow



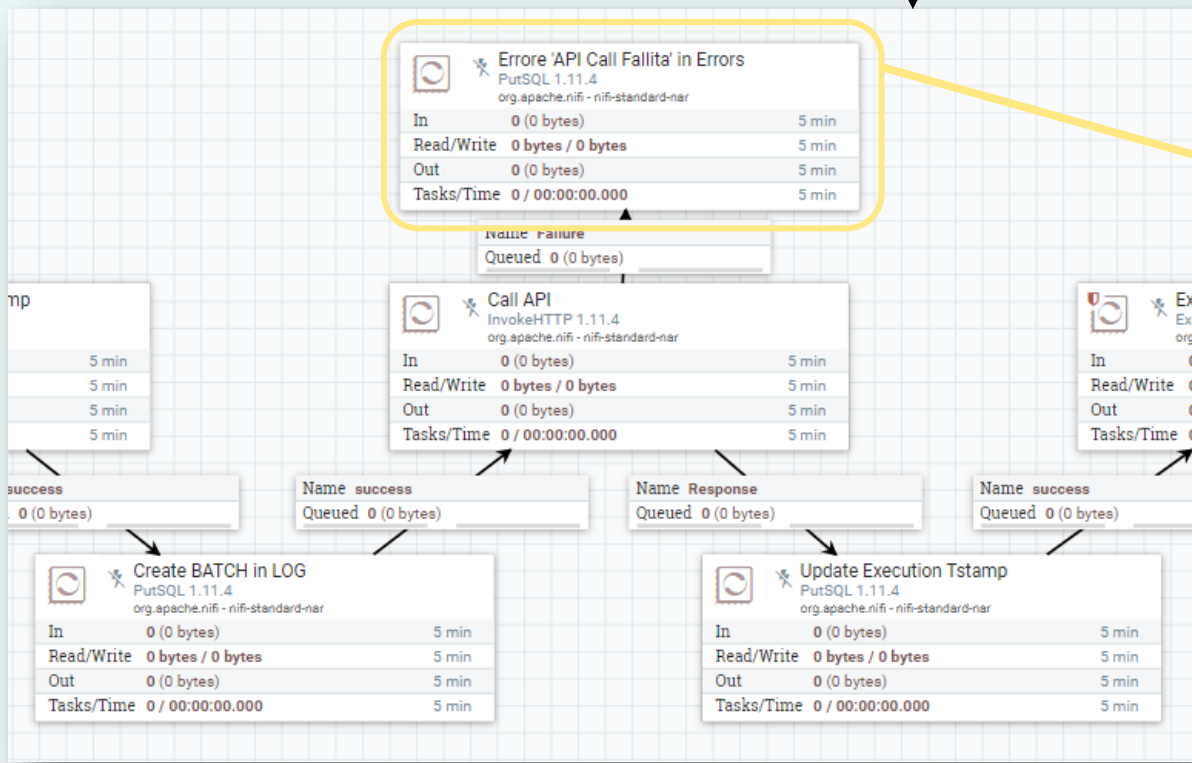
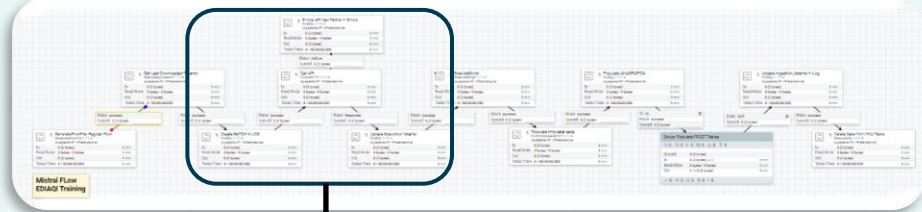
Log table

Mistral: NiFi Flow

Log table

Batch_ID	Start_Date	Stop_Date	API_Execucion_Tstamp	FROST_Ingestion_Tstamp
1	01/01/2023 00:00	01/01/2023 23:59	07/21/2023 10:00	07/21/2023 10:01
2	02/01/2023 00:00	02/01/2023 23:59	07/21/2023 10:30	

Mistral: NiFi Flow



How do we deal with errors?

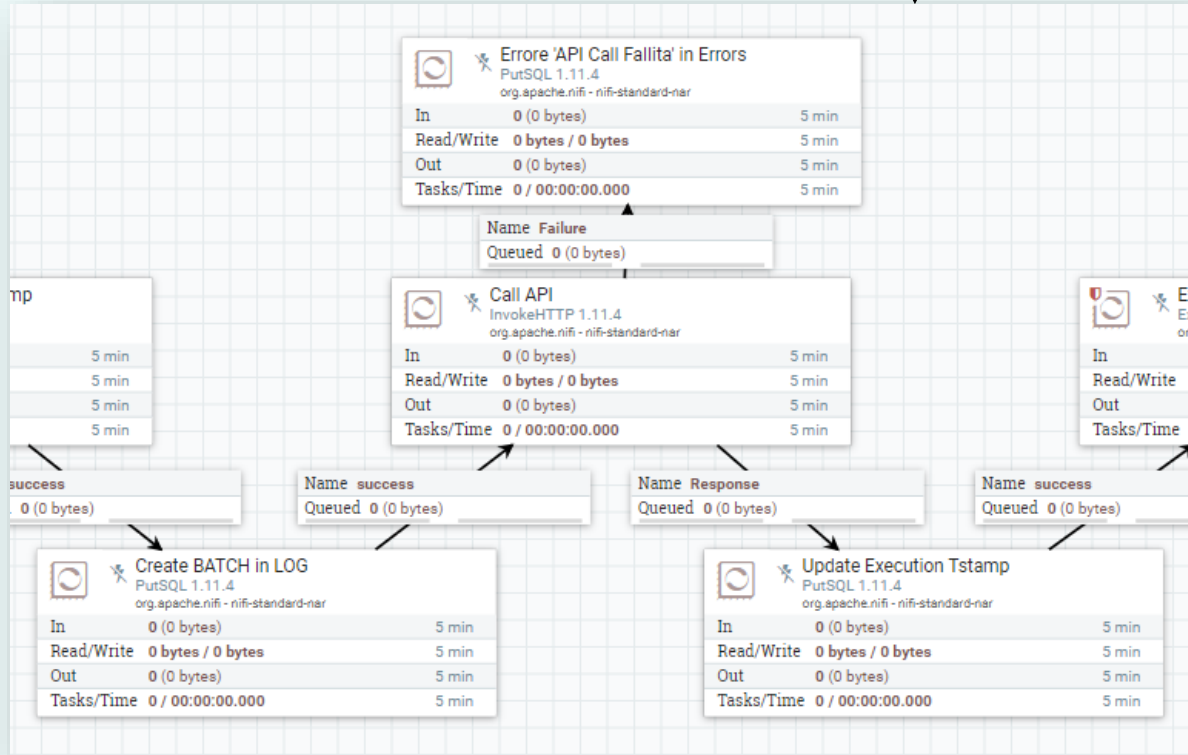
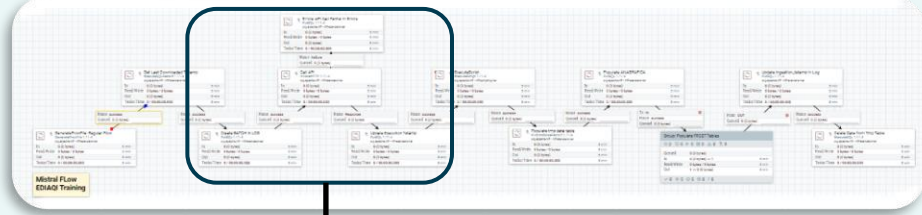


Error table

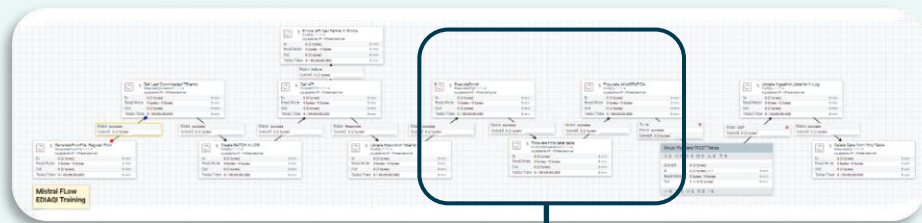
Mistral: NiFi Flow

How do you actually go from JSON data to FROST data?

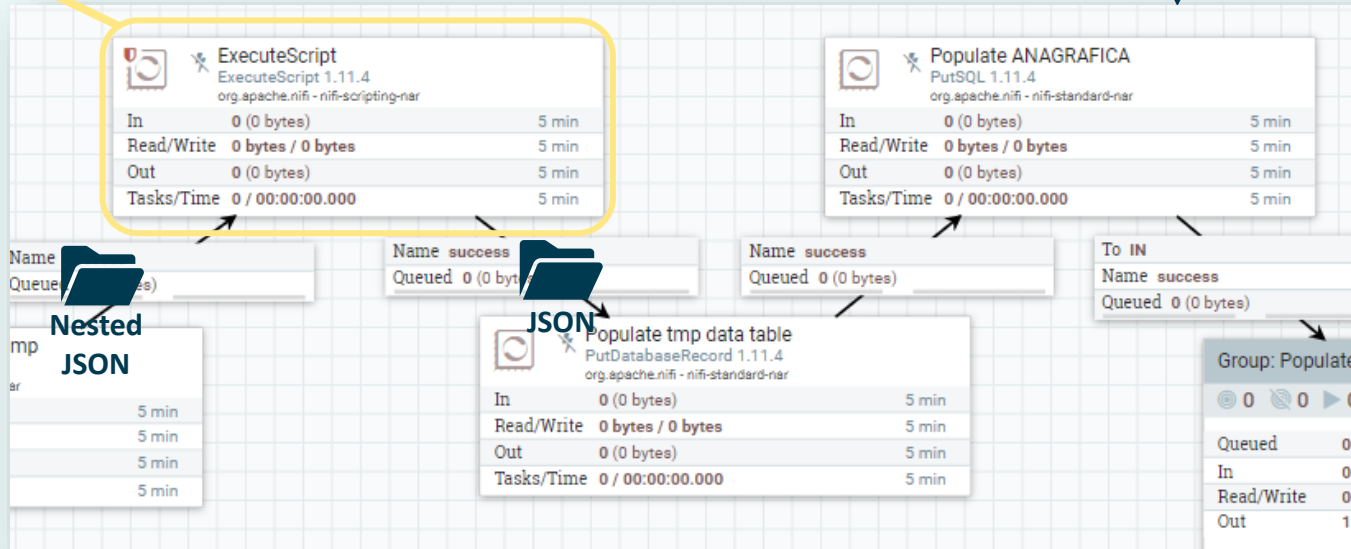
Mistral: NiFi Flow



Mistral: NiFi Flow

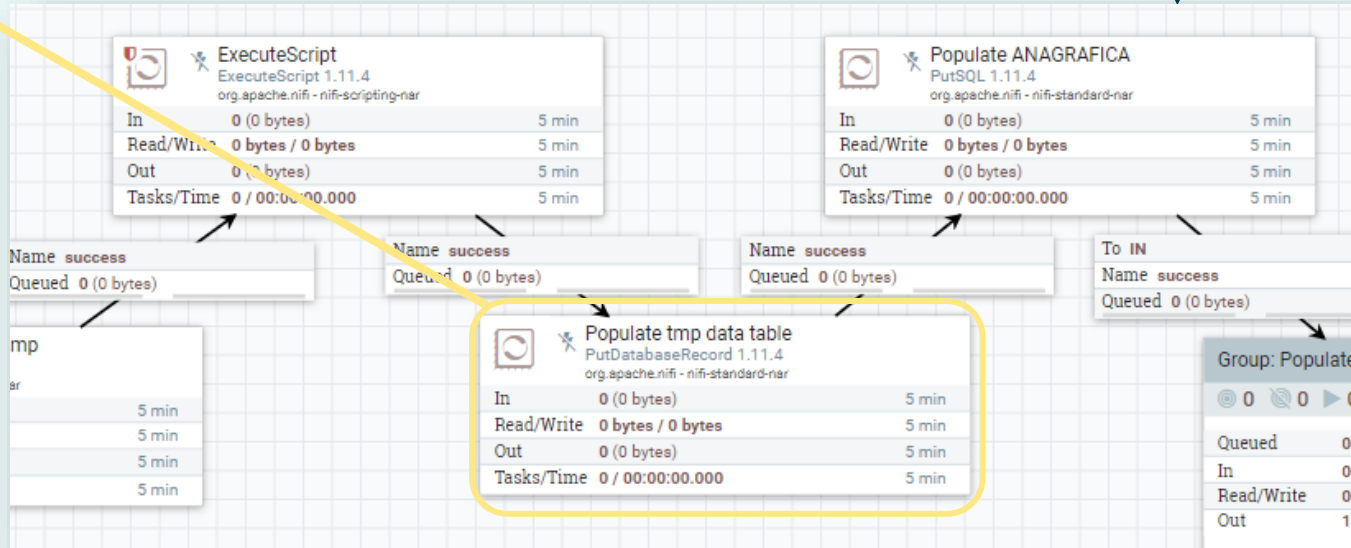
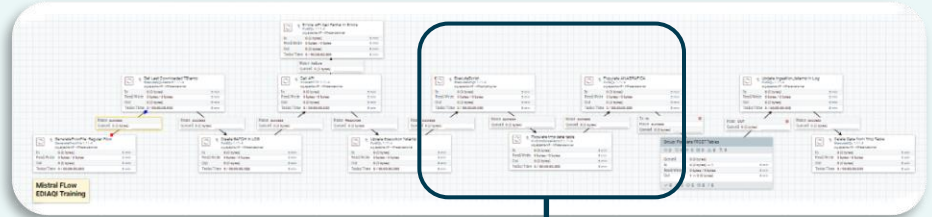


Jython scrip



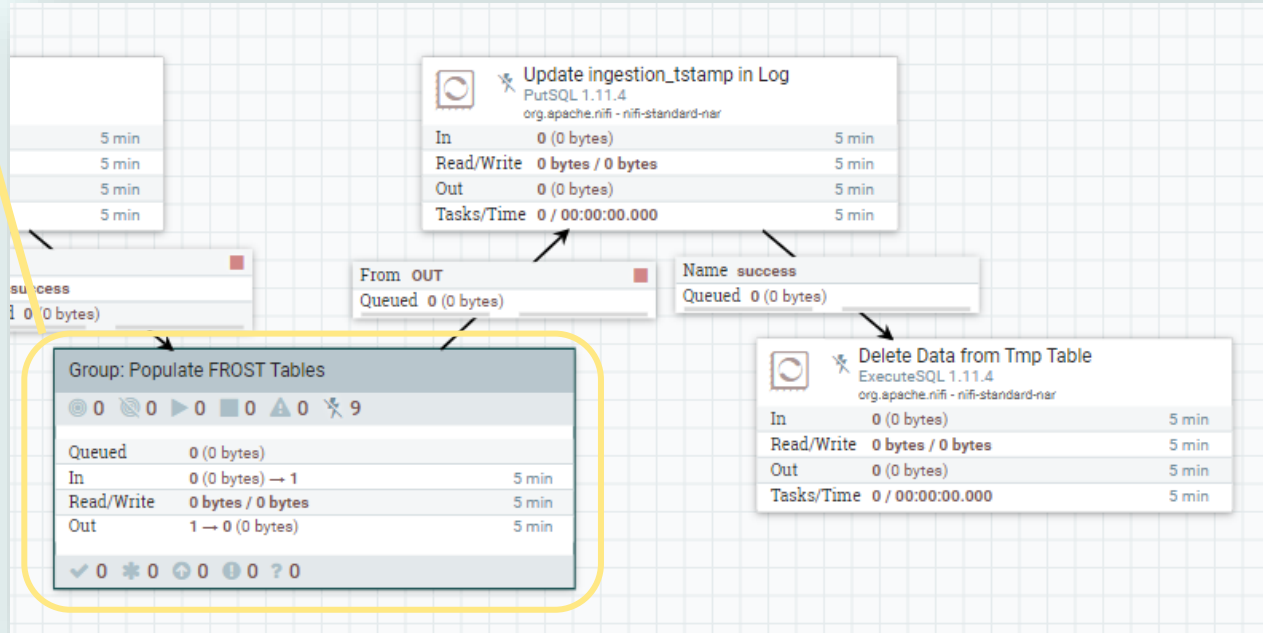
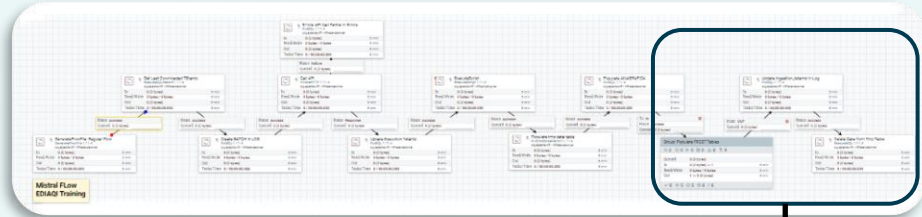
Mistral: NiFi Flow

Support DB table
to temporarily
store the data

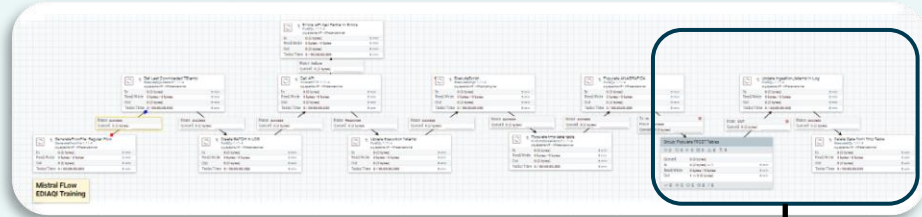


Mistral: NiFi Flow

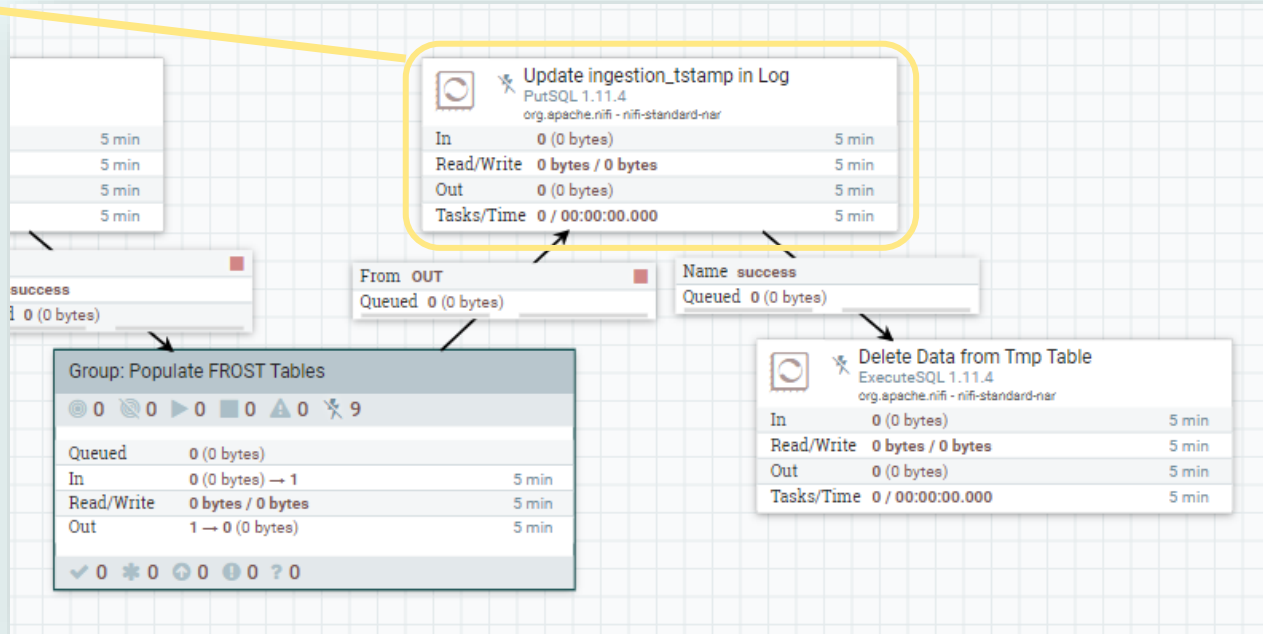
SQL queries to populate FROST tables



Mistral: NiFi Flow



Log table

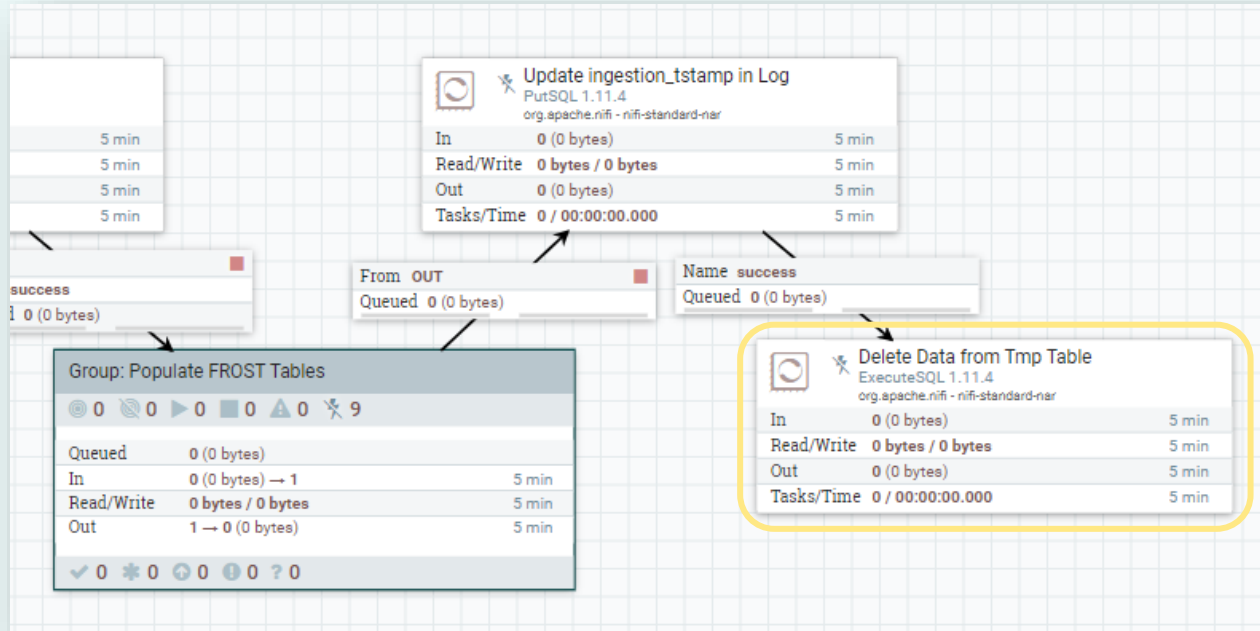
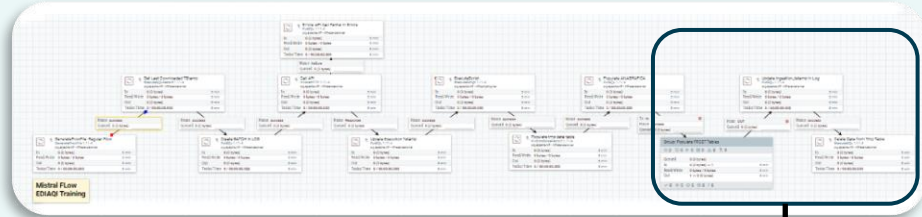


Mistral: NiFi Flow

Log table

Batch_ID	Start_Date	Stop_Date	API_Excecuton_Tstamp	FROST_Ingestion_Tstamp
1	01/01/2023 00:00	01/01/2023 23:59	07/21/2023 10:00	07/21/2023 10:01
2	02/01/2023 00:00	02/01/2023 23:59	07/21/2023 10:30	07/21/2023 10:32

Mistral: NiFi Flow

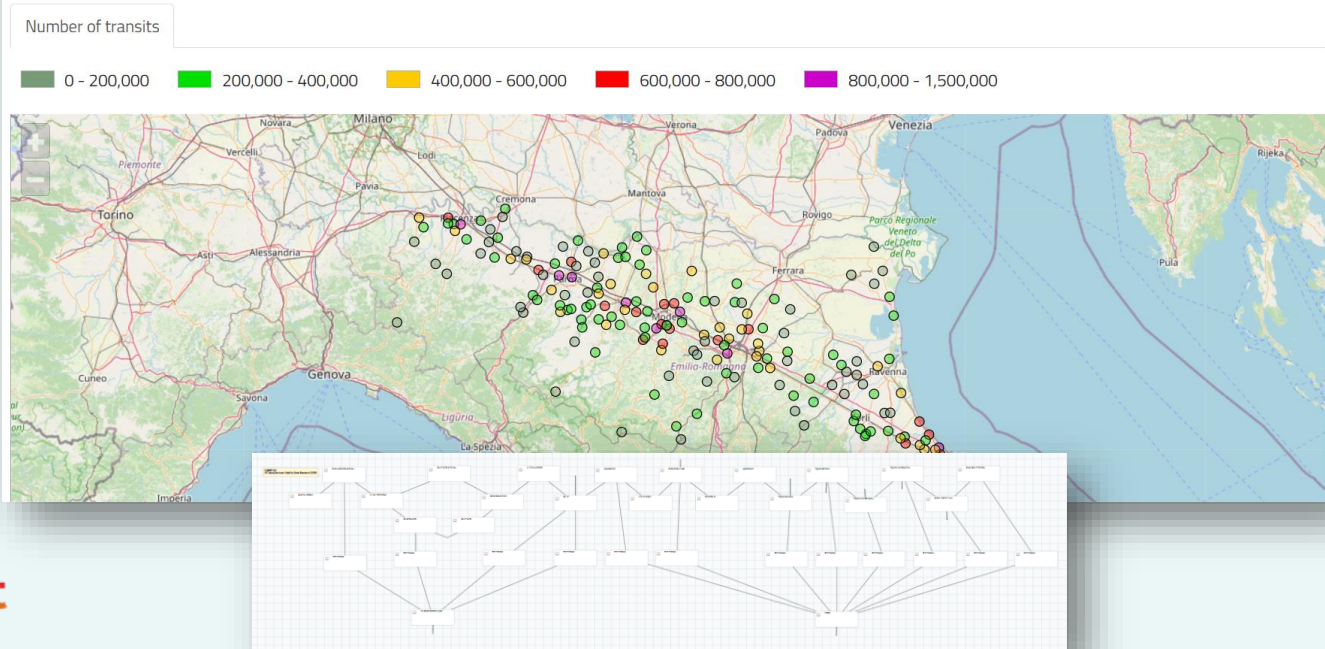


Other examples: road traffic data

Online streams

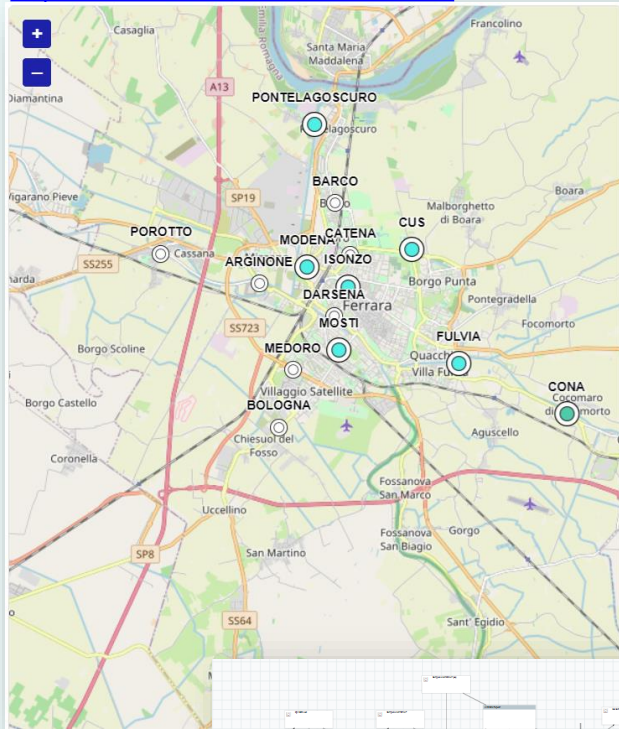
The Regional Road Traffic Flow Monitoring System (MTS) of Emilia-Romagna, created by the Region itself, the Provinces and Anas, is made up of 283 stations, operating 24 hours a day, mainly installed on state and provincial roads.

The consultation and download of the traffic flows detected allow the use of the data recorded by the MTS System, managed by the Traffic, Logistics, Waterways and Airports Area.



Other examples: outdoor AQ data

<https://airbreakferrara.net/che-aria-tira/>



- ✳ BARCO ✕
- ✳ MODENA ✕
- ✳ CATENA ✕
- ✳ ARGINONE ✕
- ✳ ISONZO ✕
- ✳ MOSTI ✕
- ✳ BOLOGNA ✕
- ✳ CONA ✕
- ✳ FULVIA ✕
- ✳ PONTELAGOSCURO ✕
- ✳ CUS ✕
- ✳ POROTTO ✕
- ✳ MEDORO ✕
- ✳ DARSENA ✕

Nascondi tutte le centraline

Indice AirBreak

