



ELT – Batch Processing

Airflow

Overview

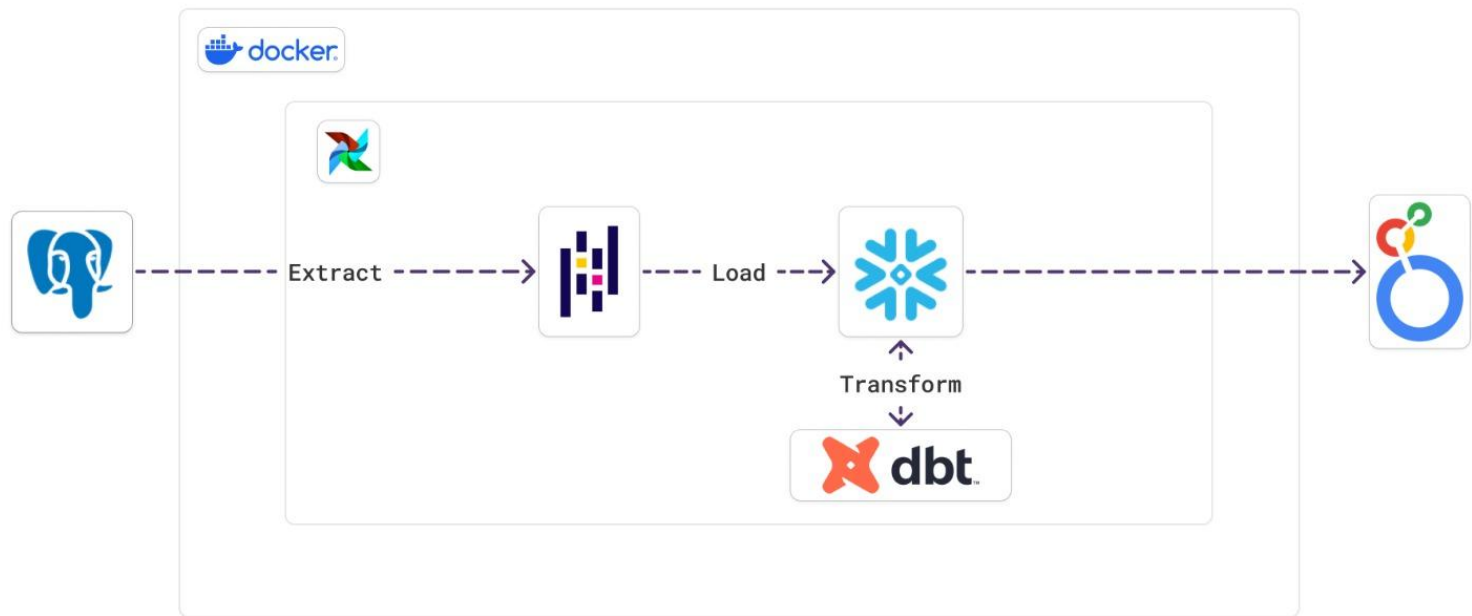
This project is an improvised result of project 2 where in project 3 we improvise to do scheduling with Airflow.

Objective

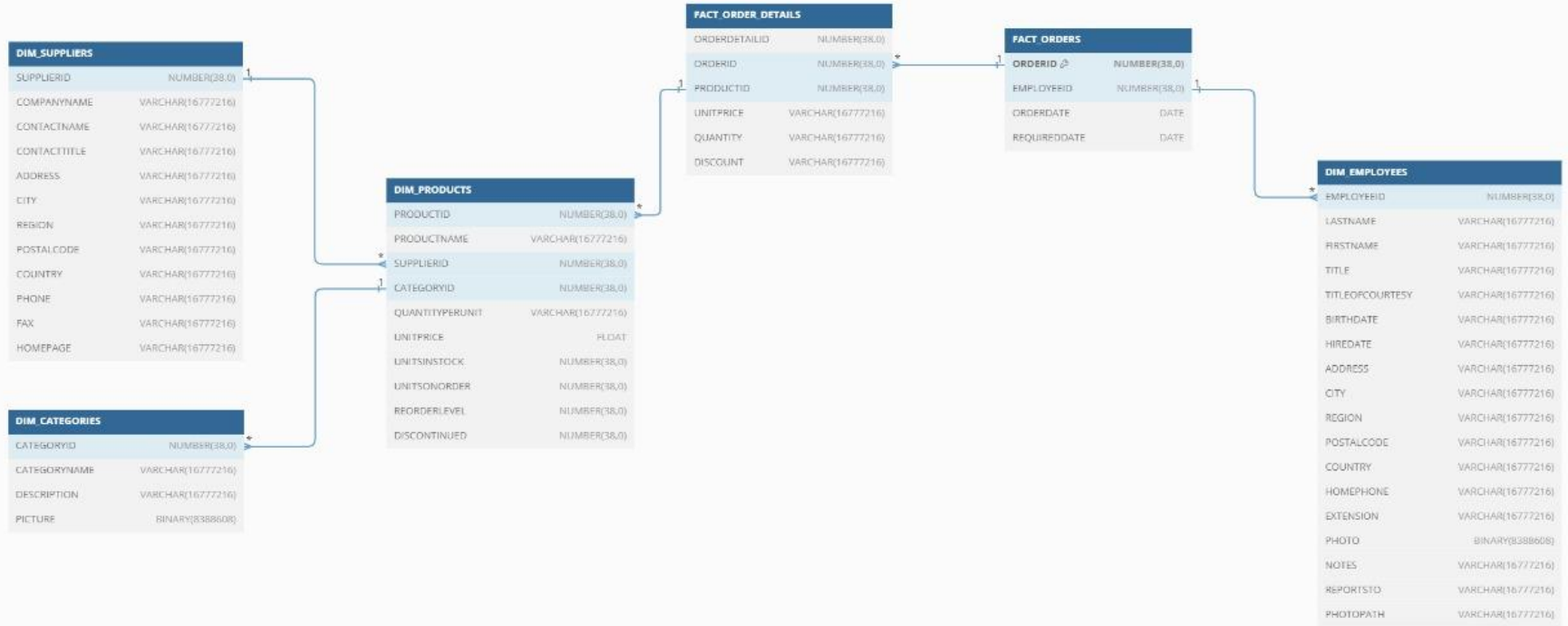
The objective of this project is to perform scheduling that runs at 12pm every day.



ARCHITECTURE



Data Warehouse Modeling



DAG Scripts


```
from datetime import datetime, timedelta
from airflow.decorators import dag, task, task_group
from airflow.operators.dummy_operator import DummyOperator
from airflow.sensors.time_delta import TimeDeltaSensor
from cosmos import DbtTaskGroup
from dbt.dbt_project.dbt_config import (
    profile_config,
    project_config,
    execution_config,
    render_config
)
from el_module.extractors import extract_func
from el_module.loaders import load_func
from el_module.filenames import names
import os

# Define Date & Run Schedule
startDate = datetime(2024, 6, 27)
schedule = '@ 0 * * * *'

# DAG Configuration
default_args = {
    "owner": "Data-Ninja",
    "depends_on_past": False,
    "start_date": startDate,
    "retries": 1,
    "retry_delay": timedelta(minutes=5),
}

@dag(
    dag_id="ELT-PROJECT",
    default_args=default_args,
    schedule_interval=schedule,
    catchup=True,
    fail_stop=True,
    description="ELT Orchestrator",
    tags=['PostgresHooks', 'Snowflake', 'DBT'],
)
def elt_dag():
```

Scheduling – Airflow

 Airflow

DAGsCluster ActivityDatasetsSecurityBrowseAdminDocs

04:48 UTC → Log In

DAGs

All 1Active 1Paused 0Running 12Failed 0

Filter DAGs by tag

Search DAGs

☒ Auto-refresh

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Acti
<div><div><input checked="" type="checkbox"/></div><div>ELT-PROJECT</div><div><div>DBT</div><div>PostgresHooks</div><div>Snowflake</div></div></div>	Data-Ninja	<div><div></div><div>14</div><div>6</div><div>5</div><div>00</div></div>	00 *	2024-07-14, 03:47:05	2024-07-14, 00:00:00	<div><div>12</div><div></div><div>30</div><div>2</div><div>14</div><div>302</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>	<div></div>

« < 1 > »

Showing 1-1 of 1 DAGs

Activate Windows
Go to Settings to activate Windows.

Task Group

```
@task_group(group_id="extract")
def extract_data():
    tg_extract = []

    # Generate Tasks
    for gen_name in names:
        @task(task_id=f'extract_{gen_name}')
        def extract(name):
            return extract_func(
                table_name=name
            )

        # Populate Task Group
        tg_extract.append(
            extract(name=gen_name)
        )

    return tg_extract

@task_group(group_id="load")
def load_data(dfs):
    tg_load = []

    # Generate Tasks
    for i, gen_name in enumerate(names):
        @task(task_id=f'load_{gen_name}')
        def load(df, table_name):
            load_func(
                df=df, filename=table_name
            )

        # Populate Task Group
        tg_load.append(
            load(df=dfs[i], table_name=gen_name)
        )

    return tg_load
```

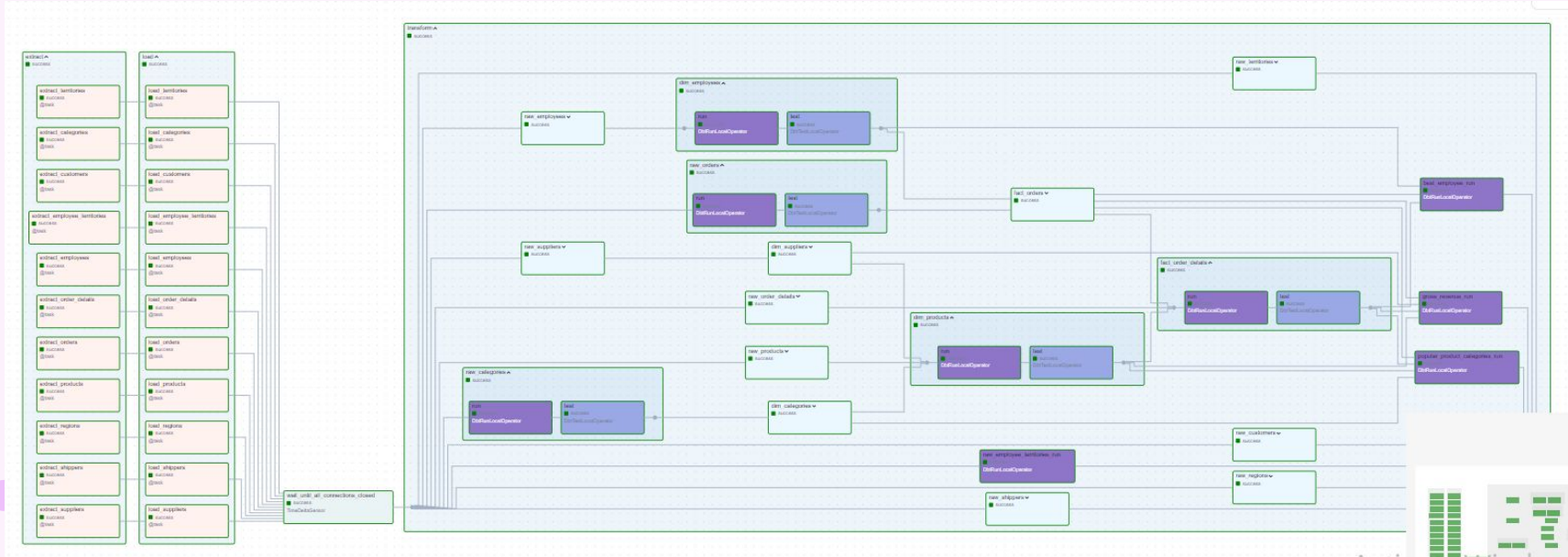
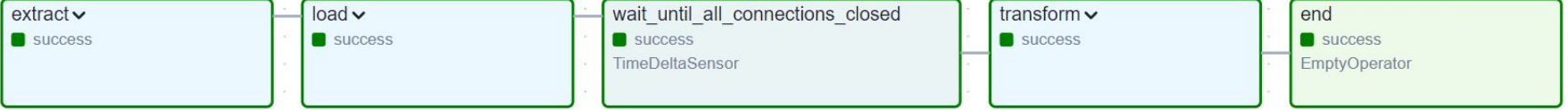
```
wait = TimeDeltaSensor(
    task_id="wait_until_all_connections_closed",
    delta=timedelta(seconds=10)
)

transform_data = DbtTaskGroup(
    group_id="transform",
    project_config=project_config,
    profile_config=profile_config,
    execution_config=execution_config,
    render_config=render_config,
    operator_args={
        "install_deps": True
    }
)

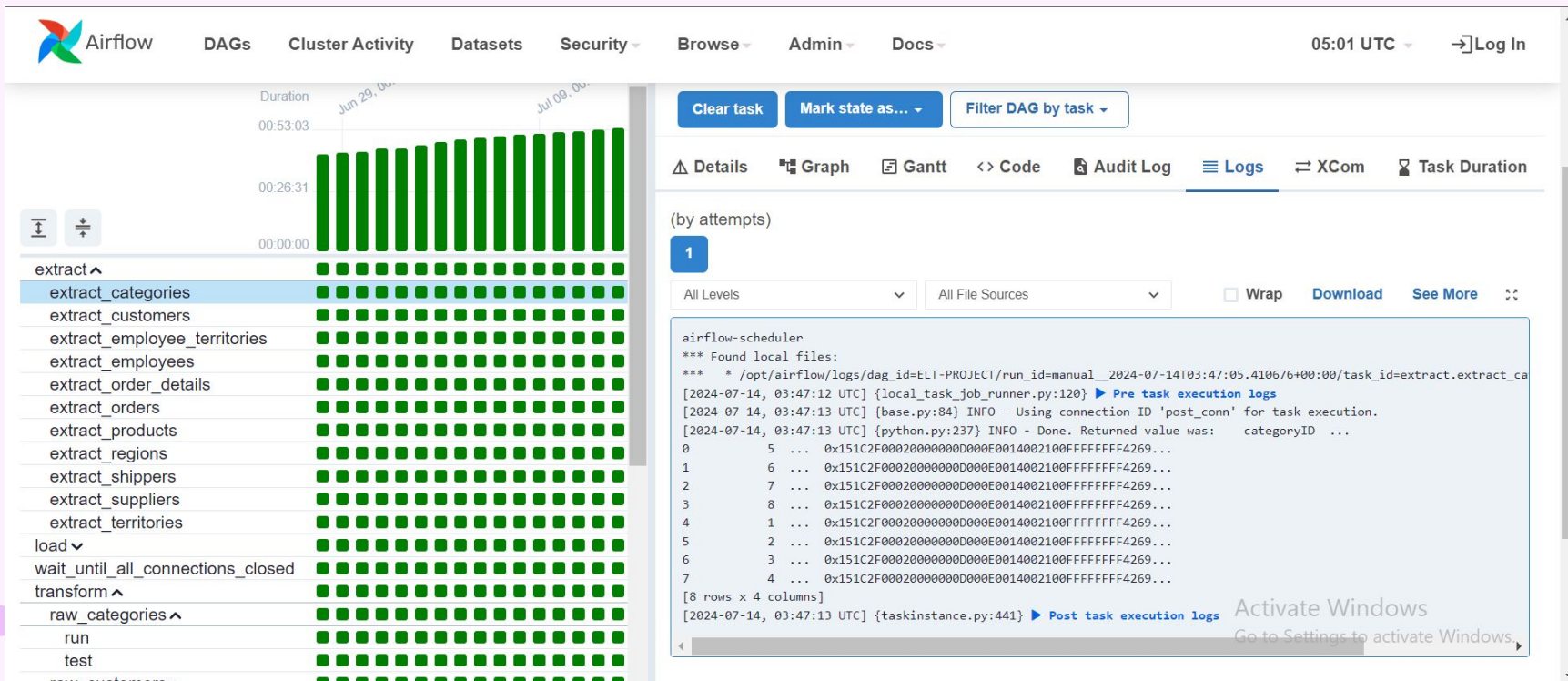
end = DummyOperator(
    task_id="end"
)

# Task Flow: Mixing Operator
extract = extract_data()
load_data(dfs=extract) >> wait
wait >> transform_data >> end
```

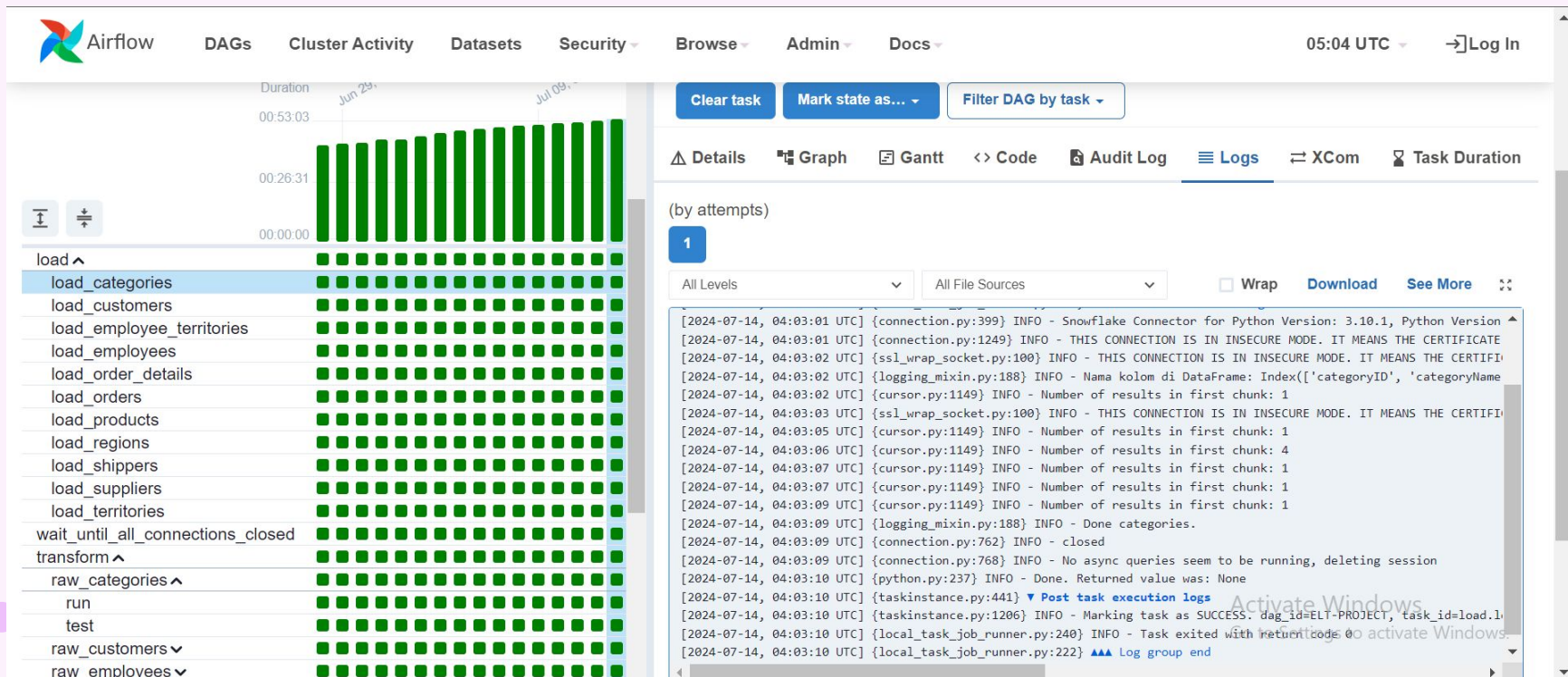

Scheduling – Airflow




Extract Task



Load Task



Wait Task

 Airflow

DAGs

Cluster Activity

Datasets

Security

Browse

Admin

Docs

05:06 UTC

Log In

wait_until_all_connections_closed

transform ^

raw_categories ^

run

test

raw_customers v

raw_employees v

raw_employee_territories_run

raw_orders ^

run

test

raw_order_details v

raw_products v

raw_regions v

raw_shippers v

raw_suppliers v

raw_territories v

dim_employees ^

run

Duration

00:53:03

00:26:31

00:00:00

Jun 29

Jul 09

Clear task

Mark state as...

Filter DAG by task

Details

Graph

Gantt

Code

Audit Log

Logs

XCom

Task Duration

(by attempts)

1

All Levels

All File Sources

☐ Wrap

Download

See More

airflow-scheduler

*** Found local files:

*** /opt/airflow/logs/dag_id=ELT-PROJECT/run_id=scheduled__2024-07-12T00:00:00+00:00/task_id=wait_until_all_connect

[2024-07-14, 04:03:18 UTC] {local_task_job_runner.py:120} Pre task execution logs

[2024-07-14, 04:03:20 UTC] {time_delta.py:51} INFO - Checking if the time (2024-07-13 00:00:10+00:00) has come

[2024-07-14, 04:03:20 UTC] {base.py:294} INFO - Success criteria met. Exiting.

[2024-07-14, 04:03:20 UTC] {taskinstance.py:441} Post task execution logs

[2024-07-14, 04:03:20 UTC] {taskinstance.py:1206} INFO - Marking task as SUCCESS. dag_id=ELT-PROJECT, task_id=wait_unti

[2024-07-14, 04:03:20 UTC] {local_task_job_runner.py:240} INFO - Task exited with return code 0

[2024-07-14, 04:03:20 UTC] {local_task_job_runner.py:222} Log group end

Activate Windows

Go to Settings to activate Windows.

Transform Task

The screenshot displays the Apache Airflow web interface. At the top, there are navigation tabs: Airflow, DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The current time is 05:09 UTC, and there is a 'Log In' button.

The main area shows a DAG with a task named 'transform'. The task is highlighted in blue. To the right of the task name, there is a green bar indicating the task's status. Below the task name, there is a list of tasks, including 'load_territories', 'wait_until_all_connections_closed', 'raw_categories', 'run', 'test', 'raw_customers', 'raw_employees', 'raw_employee_territories_run', 'raw_orders', 'run', 'test', 'raw_order_details', 'raw_products', 'raw_regions', 'raw_shippers', 'raw_suppliers', 'raw_territories', 'dim_employees', 'run', 'test', 'dim_categories', and 'dim_suppliers'.

The 'transform' task is expanded, showing its execution logs. The logs are filtered by attempts, showing attempt 1. The logs include the following information:

- Running command: `['/opt/airflow/dbt_venv/bin/dbt', 'run', '--mo`
- Command output:
- Running with dbt=1.8.3
- Registered adapter: snowflake=1.8.3
- Found 20 models, 27 data tests, 11 sources, 684 macro
- Concurrency: 1 threads (target='dev')
- 1 of 1 START sql table model raw.raw_categories
- 1 of 1 OK created sql table model raw.raw_categories
- Finished running 1 table model in 0 hours 0 minutes a
- Completed successfully
- Done. PASS=1 WARN=0 ERROR=0 SKIP=0 TOTAL=1
- Command exited with return code 0
- Done. PASS=1 WARN=0 ERROR=0 SKIP=0 TOTAL=1
- Artifact schema version: <https://schemas.getdbt.com/dbt/manifest>
- Artifact schema version: <https://schemas.getdbt.com/dbt/run-resu>
- Inlets: [Dataset(uri='snowflake://[redacted]/ELT_BATCH
- Outlets: [Dataset(uri='snowflake://[redacted]/ELT_BATCH
- Sync 1 DAGs
- Post task execution logs

[illegible]

Data Mart

The screenshot displays the Snowflake Data Mart interface. On the left, a sidebar shows the database hierarchy: **ELT_BATCH** > **DM** > **Tables**. The **BEST_EMPLOYEE** table is selected and highlighted in blue. Other tables listed under **DM** are **GROSS_REVENUE** and **POPULAR_PRODUCT_CATEG...**. Below this, other schemas like **DWH**, **INFORMATION_SCHEMA**, **PUBLIC**, **RAW**, **SNOWFLAKE**, and **SNOWFLAKE_SAMPLE_DATA** are visible.

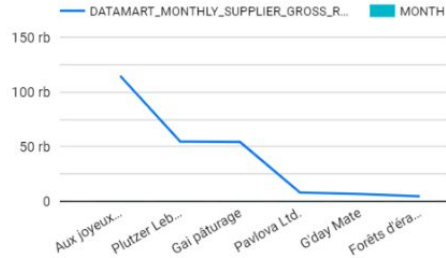
The main panel shows the details for the **ELT_BATCH / DM / BEST_EMPLOYEE** table. It includes a **Load Data** button, a table icon, the user **SYSDADMIN**, a timestamp of **32 minutes ago**, a row count of **23**, and a size of **2.0KB**. Below this, tabs for **Table Details**, **Columns**, **Data Preview** (which is active), and **Copy History** are present.

The **Data Preview** tab shows a table with 23 rows and 4 columns: **MONTH**, **EMPLOYEE_NAME**, **GROSS_REVENUE**, and **OVERALL_RANKING**. The first 8 rows are visible, showing data for various employees and their rankings. A **Select Warehouse** dropdown is located above the table, and a **23 Rows • Updated 2 minutes ago** status is shown. A **Refresh** button is also present.

	MONTH	EMPLOYEE_NAME	GROSS_REVENUE	OVERALL_RANKING
1	8-04-01	Andrew Fuller	29152.28	1
2	8-01-01	Janet Leverling	25705.0075	2
3	8-03-01	Nancy Davolio	24827.45	3
4	7-01-01	Margaret Peacock	23736.465	4
5	8-02-01	Andrew Fuller	23127.55	5
6	7-07-01	Nancy Davolio	19530.93	6
7	7-05-01	Janet Leverling	18049.6	7
8	7-12-01	Janet Leverling	17636.659	8

At the bottom of the interface, a Windows taskbar is visible with the text "Activate Windows" and "Go to Settings to activate Windows."

NORTHWIND ANALYTIC



MONTH

Masukkan nilai

EMPLOYEE_NAME

Sama den...

Masukkan nilai

Pilih rentang tanggal

DATAMART_MONTHLY_BEST_E...

6.517,47

30.990,28

