# Homework #5 Report: Airflow DAGs with Snowflake Integration

**Abstract**

This report documents the implementation of Homework #5. Apache Airflow orchestrates an ETL and a forecasting pipeline with Snowflake as the data warehouse. Screenshots demonstrate configuration and successful runs. Concise, relevant code snippets are embedded per the assignment requirement.

## I. INTRODUCTION

This submission shows:

- Airflow DAGs (`yfinance_etl`, `ml_forecast_tf`, and a simplified assignment DAG).
- Proper use of Airflow **Connections** and **Variables**.
- SQL **transaction** semantics for full refresh.
- Evidence of successful runs via the Airflow UI logs.

## II. AIRFLOW WEB UI: DAGS LOADED

Figure 1 confirms four DAGs are discovered by the scheduler and are active.
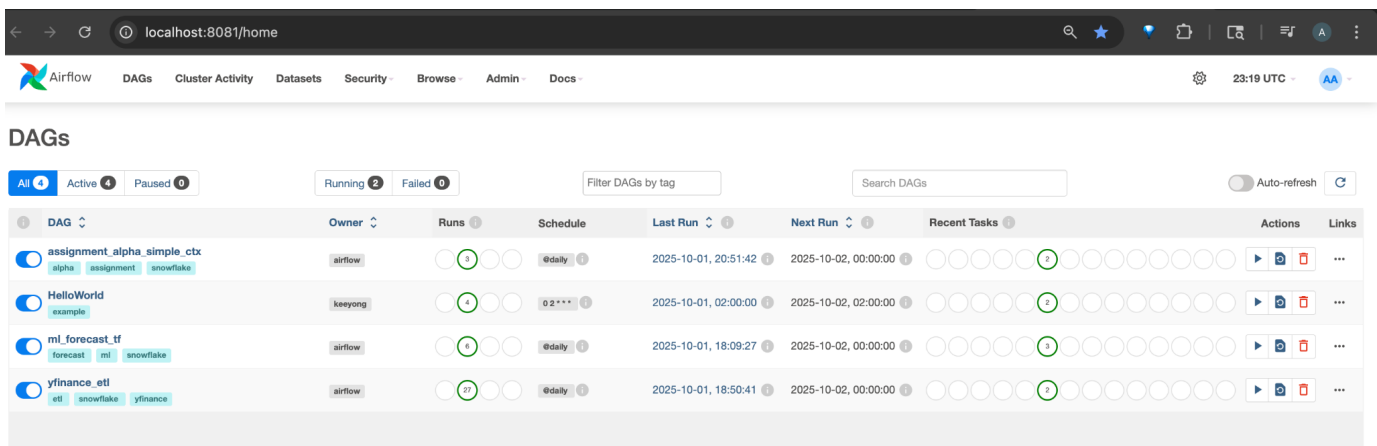


Fig. 1: Airflow Web UI showing all DAGs loaded.

## III. ADMIN → CONNECTION (SNOWFLAKE)

The connection `snowflake_catfish` was created with type *Snowflake*. Account, warehouse, database, schema, and role are set in the Extra/fields as required.

Fig. 2: Airflow Connection configuration for Snowflake.

*How to Configure the Airflow Snowflake Connection (UI)*

   **Path:** `Admin → Connections → + Add a new record`
1) **Connection Id:** `snowflake_catfish`   (must match the code)
2) **Connection Type:** `Snowflake`
3) **Login:** Snowflake user, e.g., `CATFISH`
4) **Password:** Snowflake password for the above user
5) **Schema:** Default schema for this connection, e.g., `RAW`
6) **Extra (JSON):** Provide account/warehouse/database/role (and optionally schema):

```
1  {
2    "account":   "xxxx-xxxx",
3    "warehouse": "COMPUTE_WH",
4    "database":  "YOUR_DB",
5    "schema":    "RAW",
6    "role":      "SYSADMIN"
7  }
```

   *Why this matters.*
   The DAG uses `BaseHook.get_connection("snowflake_catfish")` to read these values at runtime.
Keeping these secrets in the Airflow UI avoids hardcoding and makes deployments portable across environments.
   *Where it is used in code (exact snippet).*
*How the Snowflake Connection is Used in Assignment_5.py*

   In the DAG file (`Assignment_5.py`), the Snowflake connection is referenced through an Airflow Variable,
which points to the connection ID defined in the Admin → Connections page.
   *Code Snippet (from Assignment_5.py).*

```
1  from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
2  from airflow.models import Variable
3
4  # Connection Id is read from Airflow Variable, defaults to 'snowflake_catfish'
5  SNOWFLAKE_CONN_ID = Variable.get("SNOWFLAKE_CONN_ID", default_var="
     snowflake_catfish")
```

```
6
7   @task
8   def load_full_refresh():
9       hook = SnowflakeHook(snowflake_conn_id=SNOWFLAKE_CONN_ID)
10      conn = hook.get_conn()
11      cur = conn.cursor()
12
13      cur.execute("BEGIN")
14      # transactional load into RAW schema
15      cur.executemany(insert_sql, rows)
16      cur.execute("COMMIT")
```

- The Variable `SNOWFLAKE_CONN_ID` allows flexibility: in the UI you can point it to any valid Airflow Connection without changing code.
- The actual credentials and account/warehouse/database details come from the Airflow Connection with ID `snowflake_catfish`.
- The task wraps operations inside `BEGIN`/`COMMIT` ensuring that inserts are atomic. If an exception occurs, a `ROLLBACK` can be executed for consistency.

Thus, the connection part in your assignment is implemented via `SnowflakeHook`, pulling configuration securely from Airflow's Admin → Connections.

*Quick validation.*

After saving the connection, trigger a DAG run and check task logs. If the `conn_id` is wrong or fields are missing, Airflow will raise a clear error in the task log (e.g., "The conn_id `snowflake_catfish` isn't defined").

## IV. ADMIN → VARIABLES

An example variable (`ALPHAVANTAGE_API_KEY`) is configured and consumed in code.
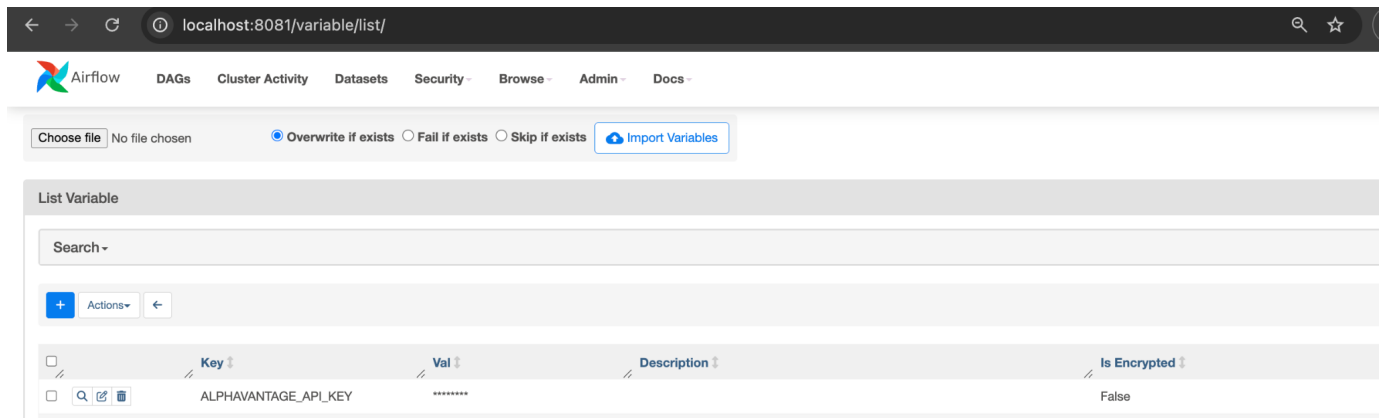


Fig. 3: Airflow Variables page.

*How to Add and Use Airflow Variables*

**Path:** `Admin → Variables → +`

1) Click **+** and set:
   - **Key:** `ALPHAVANTAGE_API_KEY`
   - **Val:** `<your-actual-api-key>`
2) Save and verify it appears in the list (Figure 4 shows it in use via successful runs).

*Where it is used in code (exact snippet).*

```python
1  from airflow.models import Variable
2
3  API_KEY = Variable.get("ALPHAVANTAGE_API_KEY")  # raises if missing
4
5  # Optionally make it non-fatal with a default:
6  # API_KEY = Variable.get("ALPHAVANTAGE_API_KEY", default_var=None)
```

*JSON and typed Variables (optional).*

When storing structured config (symbols, thresholds, etc.), keep it as JSON and parse it in code:

```python
1  import json
2
3  symbols_json = Variable.get("stock_symbols", default_var='["AAPL","MSFT","TSLA"]')
4  SYMBOLS = json.loads(symbols_json)  # -> ["AAPL","MSFT","TSLA"]
5
6  lookback_days = int(Variable.get("lookback_days", default_var="365"))
```

*Best practices.*

- Use Variables for *non-secret* runtime knobs (symbols, lookback windows, flags).
- Keep *secrets* (passwords, tokens) in *Connections* or a secrets backend (Vault, AWS SM, etc.).
- Validate presence of critical Variables early in the DAG and fail fast with a clear message.

## V. RELEVANT CODE SNIPPETS

The following excerpts from `Assignment_5.py` capture the important mechanics: connection retrieval, ETL/ML orchestration, and transactional load into Snowflake. The full file is provided in the project and may be included as an appendix if required.

### A. Core DAG Logic

```python
1  # dags/assignment_alpha_simple_ctx.py
2  from airflow import DAG
3  from airflow.decorators import task
4  from airflow.models import Variable
5  from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
6
7  from datetime import datetime
8  import requests
9
10 # Read the Snowflake connection id from a Variable
11 SNOWFLAKE_CONN_ID = Variable.get("SNOWFLAKE_CONN_ID", default_var="
       snowflake_catfish")
12
13 with DAG(
14     dag_id="assignment_alpha_simple_ctx",
15     start_date=datetime(2025, 9, 1),
16     schedule="@daily",
17     catchup=False,
18     tags=["assignment", "alpha", "snowflake"],
19 ) as dag:
20
21     @task
22     def extract() -> list[dict]:
23         """Pull compact daily OHLCV from Alpha Vantage for a small symbol list."""
24         api_key = Variable.get("ALPHAVANTAGE_API_KEY")  # must exist
25         symbols_csv = Variable.get("ALPHAVANTAGE_SYMBOLS", default_var="AAPL,MSFT"
             )
```

```python
26          symbols = [s.strip().upper() for s in symbols_csv.split(",") if s.strip()]
27
28          base = "https://www.alphavantage.co/query"
29          rows: list[dict] = []
30
31          for sym in symbols:
32              r = requests.get(
33                  base,
34                  params={
35                      "function": "TIME_SERIES_DAILY",
36                      "symbol": sym,
37                      "outputsize": "compact",
38                      "datatype": "json",
39                      "apikey": api_key,
40                  },
41                  timeout=30,
42              )
43              r.raise_for_status()
44              data = r.json()
45              ts = data.get("Time Series (Daily)")
46              if not ts:
47                  continue
48              for ds, f in ts.items():
49                  try:
50                      rows.append(
51                          {
52                              "SYMBOL": sym,
53                              "PRICE_DATE": ds,
54                              "OPEN": float(f.get("1. open", 0) or 0),
55                              "HIGH": float(f.get("2. high", 0) or 0),
56                              "LOW":  float(f.get("3. low",  0) or 0),
57                              "CLOSE": float(f.get("4. close", 0) or 0),
58                              "VOLUME": int(float(f.get("5. volume", 0) or 0)),
59                              "NOTE": None,
60                          }
61                      )
62                  except Exception:
63                      # skip malformed row
64                      pass
65
66          # simple sanity filter
67          return [r for r in rows if r["PRICE_DATE"] and r["CLOSE"] and r["CLOSE"] >
              0]
68
69      @task
70      def load_full_refresh(rows: list[dict]) -> str:
71          """
72          Full refresh into RAW.STOCK_PRICES_AV:
73            - CREATE SCHEMA/TABLE IF NOT EXISTS
74            - TRUNCATE
75            - INSERT (executemany)
76            - COMMIT / ROLLBACK
77          """
78          schema = Variable.get("SNOWFLAKE_SCHEMA", default_var="RAW")
79          table  = Variable.get("SNOWFLAKE_TABLE",  default_var="STOCK_PRICES_AV")
80
81          hook = SnowflakeHook(snowflake_conn_id=SNOWFLAKE_CONN_ID)
```

```
82        conn = hook.get_conn()
83        cur  = conn.cursor()
84
85        try:
86            cur.execute("BEGIN")
87            cur.execute(f"CREATE SCHEMA IF NOT EXISTS {schema}")
88            cur.execute(f"""
89                CREATE TABLE IF NOT EXISTS {schema}.{table} (
90                    SYMBOL     STRING,
91                    PRICE_DATE DATE,
92                    OPEN       FLOAT,
93                    HIGH       FLOAT,
94                    LOW        FLOAT,
95                    CLOSE      FLOAT,
96                    VOLUME     NUMBER,
97                    NOTE       STRING
98                )
99            """)
100           cur.execute(f"TRUNCATE TABLE {schema}.{table}")
101
102           if rows:
103               insert_sql = f"""
104                   INSERT INTO {schema}.{table}
105                   (SYMBOL, PRICE_DATE, OPEN, HIGH, LOW, CLOSE, VOLUME, NOTE)
106                   VALUES (%(SYMBOL)s, %(PRICE_DATE)s, %(OPEN)s, %(HIGH)s, %(LOW)
                        s, %(CLOSE)s, %(VOLUME)s, %(NOTE)s)
107               """
108               cur.executemany(insert_sql, rows)
109
110           cur.execute("COMMIT")
111           return f"Loaded {len(rows)} rows into {schema}.{table}"
112       except Exception:
113           cur.execute("ROLLBACK")
114           raise
115       finally:
116           cur.close()
117           conn.close()
118
119   # wiring
120   load_full_refresh(extract())
```

### B. SQL Transaction for Full Refresh (If Implemented in Python Task)

If the DAG uses Python to run SQL, ensure BEGIN / COMMIT / ROLLBACK are used appropriately, and MERGE/DELETE+INSERT are atomic.

```
1     def load_full_refresh(rows: list[dict]) -> str:
2         """
3         Full refresh into RAW.STOCK_PRICES_AV:
4           - CREATE SCHEMA/TABLE IF NOT EXISTS
5           - TRUNCATE
6           - INSERT (executemany)
7           - COMMIT / ROLLBACK
8         """
9         schema = Variable.get("SNOWFLAKE_SCHEMA", default_var="RAW")
10        table  = Variable.get("SNOWFLAKE_TABLE",  default_var="STOCK_PRICES_AV")
11
```

```python
12          hook = SnowflakeHook(snowflake_conn_id=SNOWFLAKE_CONN_ID)
13          conn = hook.get_conn()
14          cur  = conn.cursor()
15
16          try:
17              cur.execute("BEGIN")
18              cur.execute(f"CREATE SCHEMA IF NOT EXISTS {schema}")
19              cur.execute(f"""
20                  CREATE TABLE IF NOT EXISTS {schema}.{table} (
21                      SYMBOL      STRING,
22                      PRICE_DATE DATE,
23                      OPEN        FLOAT,
24                      HIGH        FLOAT,
25                      LOW         FLOAT,
26                      CLOSE       FLOAT,
27                      VOLUME      NUMBER,
28                      NOTE        STRING
29                  )
30              """)
31              cur.execute(f"TRUNCATE TABLE {schema}.{table}")
32
33              if rows:
34                  insert_sql = f"""
35                      INSERT INTO {schema}.{table}
36                      (SYMBOL, PRICE_DATE, OPEN, HIGH, LOW, CLOSE, VOLUME, NOTE)
37                      VALUES (%(SYMBOL)s, %(PRICE_DATE)s, %(OPEN)s, %(HIGH)s, %(LOW)
                            s, %(CLOSE)s, %(VOLUME)s, %(NOTE)s)
38                  """
39                  cur.executemany(insert_sql, rows)
40
41              cur.execute("COMMIT")
42              return f"Loaded {len(rows)} rows into {schema}.{table}"
43          except Exception:
44              cur.execute("ROLLBACK")
45              raise
46          finally:
47              cur.close()
48              conn.close()
49
50      # wiring
51      load_full_refresh(extract())
```

## VI. EXECUTION EVIDENCE: EVENT LOGS

Figures 4 and 5 show successful runs of the two primary tasks (extract and load_full_refresh).
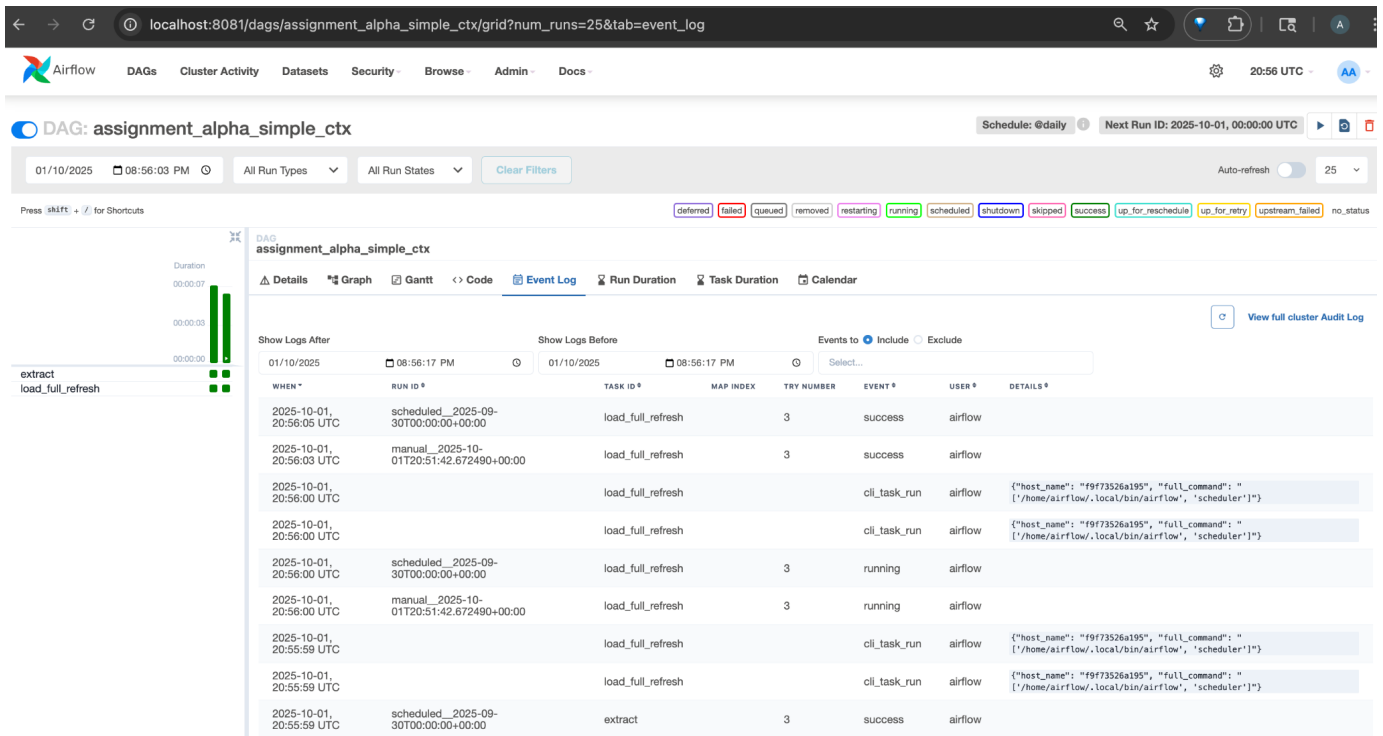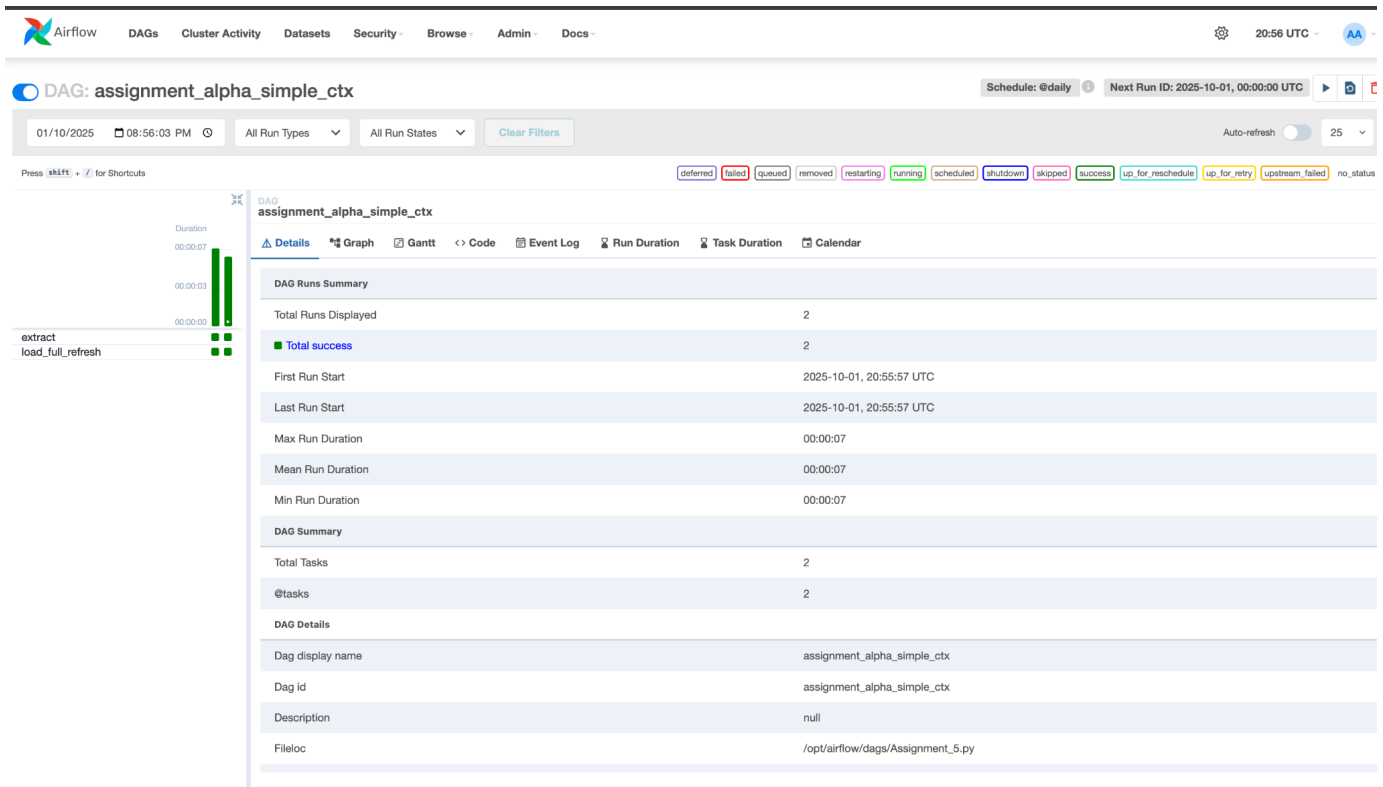
Fig. 4: Event log view with task successes.



Fig. 5: DAG details summary confirming successful runs.

## VII. CONCLUSION

This report demonstrates a working Airflow setup with Snowflake integration, proper configuration via Connections/Variables, transactional loading, and verified successful DAG runs.

APPENDIX (OPTIONAL): FULL CODE

```python
1   # dags/assignment_alpha_simple_ctx.py
2   from airflow import DAG
3   from airflow.decorators import task
4   from airflow.models import Variable
5   from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
6
7   from datetime import datetime
8   import requests
9
10  # Read the Snowflake connection id from a Variable
11  SNOWFLAKE_CONN_ID = Variable.get("SNOWFLAKE_CONN_ID", default_var="
        snowflake_catfish")
12
13  with DAG(
14      dag_id="assignment_alpha_simple_ctx",
15      start_date=datetime(2025, 9, 1),
16      schedule="@daily",
17      catchup=False,
18      tags=["assignment", "alpha", "snowflake"],
19  ) as dag:
20
21      @task
22      def extract() -> list[dict]:
23          """Pull compact daily OHLCV from Alpha Vantage for a small symbol list."""
24          api_key = Variable.get("ALPHAVANTAGE_API_KEY")  # must exist
25          symbols_csv = Variable.get("ALPHAVANTAGE_SYMBOLS", default_var="AAPL,MSFT"
                )
26          symbols = [s.strip().upper() for s in symbols_csv.split(",") if s.strip()]
27
28          base = "https://www.alphavantage.co/query"
29          rows: list[dict] = []
30
31          for sym in symbols:
32              r = requests.get(
33                  base,
34                  params={
35                      "function": "TIME_SERIES_DAILY",
36                      "symbol": sym,
37                      "outputsize": "compact",
38                      "datatype": "json",
39                      "apikey": api_key,
40                  },
41                  timeout=30,
42              )
43              r.raise_for_status()
44              data = r.json()
45              ts = data.get("Time Series (Daily)")
46              if not ts:
47                  continue
48              for ds, f in ts.items():
49                  try:
50                      rows.append(
51                          {
52                              "SYMBOL": sym,
53                              "PRICE_DATE": ds,
54                              "OPEN": float(f.get("1. open", 0) or 0),
```

```python
                              "HIGH": float(f.get("2. high", 0) or 0),
                              "LOW":  float(f.get("3. low",  0) or 0),
                              "CLOSE": float(f.get("4. close", 0) or 0),
                              "VOLUME": int(float(f.get("5. volume", 0) or 0)),
                              "NOTE": None,
                        }
                    )
                except Exception:
                    # skip malformed row
                    pass

        # simple sanity filter
        return [r for r in rows if r["PRICE_DATE"] and r["CLOSE"] and r["CLOSE"] >
            0]

    @task
    def load_full_refresh(rows: list[dict]) -> str:
        """
        Full refresh into RAW.STOCK_PRICES_AV:
          - CREATE SCHEMA/TABLE IF NOT EXISTS
          - TRUNCATE
          - INSERT (executemany)
          - COMMIT / ROLLBACK
        """
        schema = Variable.get("SNOWFLAKE_SCHEMA", default_var="RAW")
        table  = Variable.get("SNOWFLAKE_TABLE",  default_var="STOCK_PRICES_AV")

        hook = SnowflakeHook(snowflake_conn_id=SNOWFLAKE_CONN_ID)
        conn = hook.get_conn()
        cur  = conn.cursor()

        try:
            cur.execute("BEGIN")
            cur.execute(f"CREATE SCHEMA IF NOT EXISTS {schema}")
            cur.execute(f"""
                CREATE TABLE IF NOT EXISTS {schema}.{table} (
                    SYMBOL      STRING,
                    PRICE_DATE  DATE,
                    OPEN        FLOAT,
                    HIGH        FLOAT,
                    LOW         FLOAT,
                    CLOSE       FLOAT,
                    VOLUME      NUMBER,
                    NOTE        STRING
                )
            """)
            cur.execute(f"TRUNCATE TABLE {schema}.{table}")

            if rows:
                insert_sql = f"""
                    INSERT INTO {schema}.{table}
                    (SYMBOL, PRICE_DATE, OPEN, HIGH, LOW, CLOSE, VOLUME, NOTE)
                    VALUES (%(SYMBOL)s, %(PRICE_DATE)s, %(OPEN)s, %(HIGH)s, %(LOW)
                        s, %(CLOSE)s, %(VOLUME)s, %(NOTE)s)
                """
                cur.executemany(insert_sql, rows)
```

```
110            cur.execute("COMMIT")
111            return f"Loaded {len(rows)} rows into {schema}.{table}"
112        except Exception:
113            cur.execute("ROLLBACK")
114            raise
115        finally:
116            cur.close()
117            conn.close()
118
119    # wiring
120    load_full_refresh(extract())
```