

Lab 1: Stock Price Prediction Analytics using Snowflake & Airflow

Abhinita Sanabada (018320874)

October 1, 2025

Abstract

We implemented a secure, reproducible stock analytics workflow using Airflow, `yfinance`, and Snowflake. An ETL DAG ingests OHLCV data into `RAW.STOCK_PRICES`. A second DAG trains a Snowflake-native `SNOWFLAKE.ML.FORECAST` model and writes predictions to `MODEL.FORECASTS`. A final table, `ANALYTICS.FINAL_PRICES_FORECAST`, unions actuals and forecasts for downstream visualization. All Snowflake credentials (account, user, password, role, warehouse, database) are stored only in Airflow Connections, and pipeline parameters are managed via Airflow Variables.

1 Problem Statement

Build an end-to-end analytics pipeline that:

1. Extracts daily OHLCV for selected tickers via `yfinance`.
2. Forecasts daily close prices using Snowflake's built-in ML forecasting.
3. Unifies actuals and forecasts in a single analytics table.
4. Uses Airflow for orchestration, SQL/Python transactions for correctness, and Airflow Connections/Variables for secure configuration.
5. Produces reproducible runs, screenshots, and a public code repository.

Success criteria: both DAGs succeed; `RAW`, `MODEL`, `ANALYTICS` populated; final table supports plotting Actual vs. Forecast; screenshots and repo links are provided.

2 System Architecture

Overview

We use three schemas inside `USER_DB_CATFISH`: **RAW** (ingest), **MODEL** (predictions), **ANALYTICS** (consumption). Two DAGs orchestrate the flow:

- **DAG #1 `yfinance_etl`:** downloads OHLCV for `stock_symbols`, MERGEs into `RAW.STOCK_PRICES`

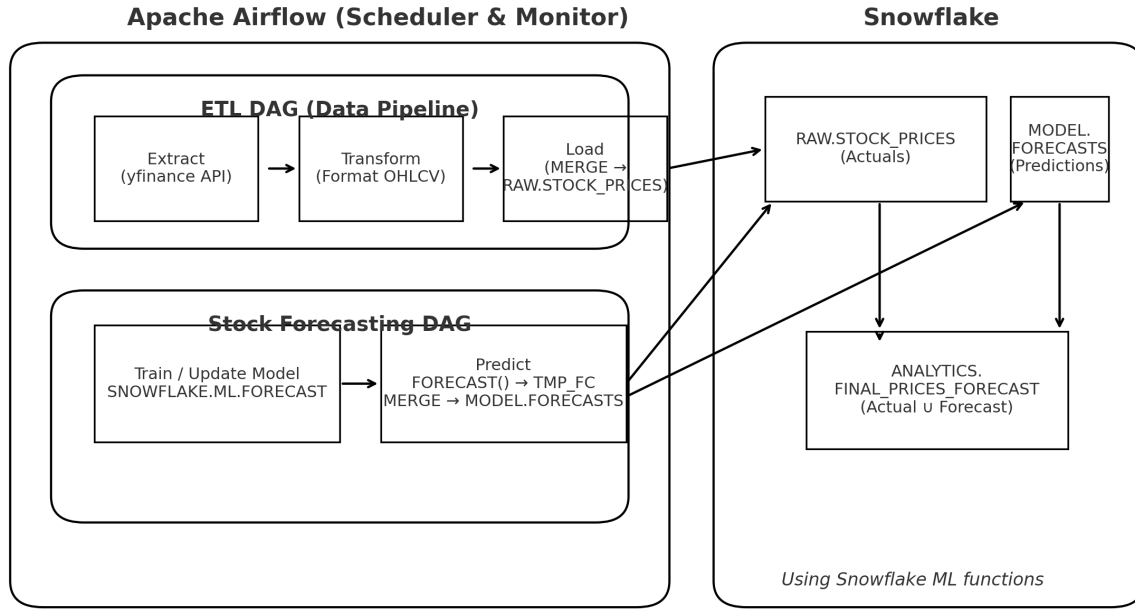


Figure 1: Architecture diagram

- **DAG #2 ml_forecast**: trains/updates `SNOWFLAKE.ML.FORECAST` on multi-series history; writes predictions to `MODEL.FORECASTS`; *unions* actuals + forecasts into `ANALYTICS.FINAL_PRICES_FORECAST`

Architecture Diagram

Screenshots (to be included)

- **Figure 1**: Airflow DAGs list showing both DAGs present.
- **Figure 2**: `yfinance_etl` Grid/Graph view with successful run.
- **Figure 3**: `ml_forecast` Grid/Graph view with successful run.

3 Data Model

All objects live in database `USER_DB_CATFISH` (warehouse: `CATFISH_QUERY_WH`, both configured via Airflow Connection). Schemas: `RAW`, `MODEL`, `ANALYTICS`.

RAW.STOCK_PRICES

```

SYMBOL STRING NOT NULL; TS TIMESTAMP_NTZ NOT NULL; OPEN FLOAT; HIGH
FLOAT; LOW FLOAT; CLOSE FLOAT; ADJ_CLOSE FLOAT; VOLUME NUMBER(38,0);
LOAD_TS TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP();

```

Primary key: (SYMBOL, TS).

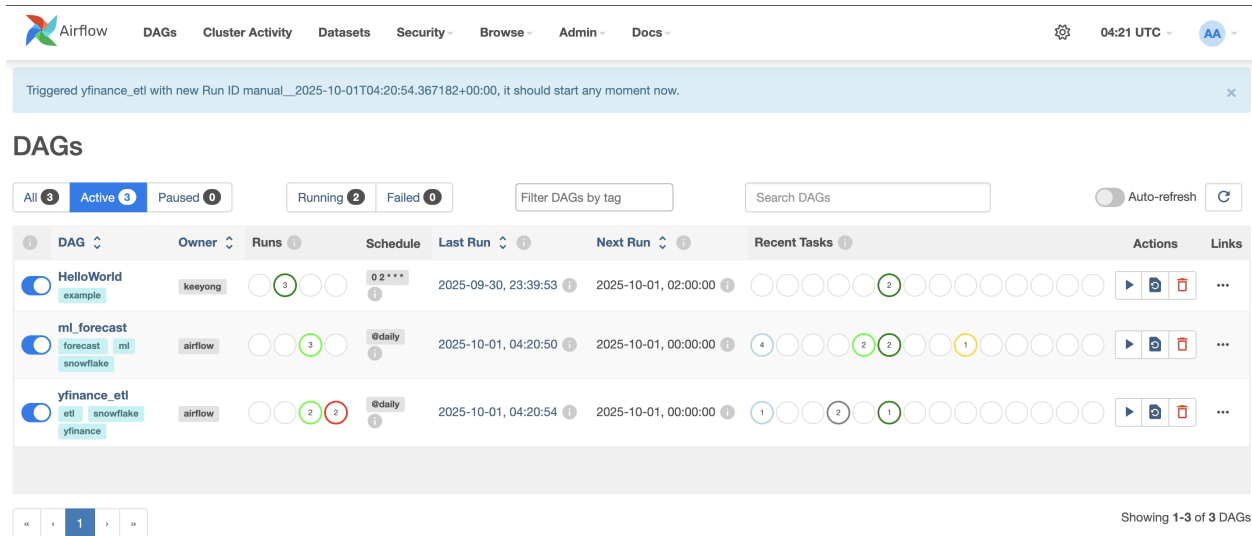


Figure 2: Airflow dags screen shot

MODEL.FORECASTS

SYMBOL STRING NOT NULL; TS DATE NOT NULL; PREDICTED_CLOSE FLOAT NOT NULL; MODEL_NAME STRING NOT NULL; TRAINED_AT TIMESTAMP_NTZ NOT NULL; HORIZON_D NUMBER(5,0) NOT NULL; LOAD_TS TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP
Primary key: (SYMBOL, TS, MODEL_NAME).

ANALYTICS.FINAL_PRICES_FORECAST

SYMBOL STRING NOT NULL; TS DATE NOT NULL; CLOSE FLOAT; SOURCE STRING NOT NULL; MODEL_NAME STRING; LOAD_TS TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP
Primary key: (SYMBOL, TS, SOURCE); SOURCE ∈ {ACTUAL, FORECAST}.

One-time Bootstrap DDL

Listing 1: Bootstrap DDL (run once).

```

1 BEGIN;
2 CREATE SCHEMA IF NOT EXISTS RAW;
3 CREATE SCHEMA IF NOT EXISTS MODEL;
4 CREATE SCHEMA IF NOT EXISTS ANALYTICS;
5
6 CREATE TABLE IF NOT EXISTS RAW.STOCK_PRICES (
7     SYMBOL STRING NOT NULL, TS TIMESTAMP_NTZ NOT NULL,
8     OPEN FLOAT, HIGH FLOAT, LOW FLOAT, CLOSE FLOAT, ADJ_CLOSE FLOAT,
9     VOLUME NUMBER(38,0),
10    LOAD_TS TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP(),
11    CONSTRAINT PK_STOCK_PRICES PRIMARY KEY (SYMBOL, TS)
12 );

```

```

12
13 CREATE TABLE IF NOT EXISTS MODEL.FORECASTS (
14     SYMBOL STRING NOT NULL, TS DATE NOT NULL, PREDICTED_CLOSE FLOAT
        NOT NULL,
15     MODEL_NAME STRING NOT NULL, TRAINED_AT TIMESTAMP_NTZ NOT NULL,
        HORIZON_D NUMBER(5,0) NOT NULL,
16     LOAD_TS TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP(),
17     CONSTRAINT PK_FORECASTS PRIMARY KEY (SYMBOL, TS, MODEL_NAME)
18 );
19
20 CREATE TABLE IF NOT EXISTS ANALYTICS.FINAL_PRICES_FORECAST (
21     SYMBOL STRING NOT NULL, TS DATE NOT NULL, CLOSE FLOAT,
22     SOURCE STRING NOT NULL, MODEL_NAME STRING, LOAD_TS TIMESTAMP_NTZ
        DEFAULT CURRENT_TIMESTAMP(),
23     CONSTRAINT PK_FINAL PRIMARY KEY (SYMBOL, TS, SOURCE)
24 );
25 COMMIT;

```

4 Implementation

4.1 Airflow Connections & Variables

We created a Snowflake Connection **snowflake_catfish** with account, user, password, role, default warehouse CATFISH_QUERY_WH, and database USER_DB_CATFISH. **No secrets in code.** Variables:

- stock_symbols: JSON list, e.g., ["AAPL", "MSFT", "TSLA"].
- lookback_days: e.g., 365.
- forecast_horizon_days: e.g., 14.
- target_schema_raw=RAW, target_schema_model=MODEL, target_schema_analytics=ANALYTICS

4.2 DAG #1: yfinance_etl (ETL)

Python downloads OHLCV for stock_symbols over the last lookback_days and MERGES into RAW.STOCK_PRICES. Credentials/DB/WH/role come from the Airflow Connection at run-time (via snowflake.connector). We use a transactional pattern (commit/rollback).

Listing 2: ETL MERGE (executed each run).
Runtime DML (executed by the DAG after staging rows):

```

1 BEGIN;
2 CREATE TEMP TABLE TMP_LOAD (
3     SYMBOL STRING, TS TIMESTAMP_NTZ, OPEN FLOAT, HIGH FLOAT, LOW FLOAT
        ,
4     CLOSE FLOAT, ADJ_CLOSE FLOAT, VOLUME NUMBER(38,0)

```

```

5 );
6 -- Python inserts many rows into TMP_LOAD via executemany(...)
7
8 MERGE INTO RAW.STOCK_PRICES AS t
9 USING TMP_LOAD AS s
10 ON t.SYMBOL = s.SYMBOL
11 AND t.TS = s.TS
12 WHEN MATCHED THEN UPDATE SET
13 OPEN=s.OPEN, HIGH=s.HIGH, LOW=s.LOW, CLOSE=s.CLOSE,
14 ADJ_CLOSE=s.ADJ_CLOSE, VOLUME=s.VOLUME, LOAD_TS=CURRENT_TIMESTAMP
15 ()
16 WHEN NOT MATCHED THEN INSERT (
17 SYMBOL, TS, OPEN, HIGH, LOW, CLOSE, ADJ_CLOSE, VOLUME
18 ) VALUES (
19 s.SYMBOL, s.TS, s.OPEN, s.HIGH, s.LOW, s.CLOSE, s.ADJ_CLOSE, s.
20 VOLUME
21 );
22 COMMIT;

```

4.3 DAG #2: ml_forecast (Snowflake ML) + Final Union

This DAG uses only `SnowflakeOperator`. It trains/updates a multi-series model via `SNOWFLAKE.ML.FORECAST`, stages horizon-wide predictions, MERGEs them into `MODEL.FORECASTS`, and rebuilds `ANALYTICS.FINAL_RESULTS` by unioning ACTUAL and FORECAST rows.

Listing 3: Snowflake ML model + forecast + upsert.
Model Training, Forecasting & Upsert

```

1 BEGIN;
2 USE SCHEMA MODEL;
3
4 WITH symbols AS (
5     SELECT value::string AS symbol
6     FROM TABLE(FLATTEN(input => PARSE_JSON('{{ var.value.stock_symbols
7     }'))))
8 ),
9 training_data AS (
10     SELECT
11         TO_VARIANT(sp.SYMBOL) AS SERIES,
12         sp.TS,
13         sp.CLOSE
14     FROM RAW.STOCK_PRICES sp
15     JOIN symbols s ON s.symbol = sp.SYMBOL
16     WHERE sp.TS >= DATEADD('day', -{{ var.value.lookback_days |
17         default('365', true) }}, CURRENT_TIMESTAMP())

```

```

18 CREATE OR REPLACE SNOWFLAKE.ML.FORECAST PRICE_FORECASTER (
19   INPUT_DATA          => SYSTEM$QUERY_REFERENCE($$ SELECT SERIES, TS,
      CLOSE FROM training_data $$),
20   SERIES_COLNAME      => 'SERIES',
21   TIMESTAMP_COLNAME  => 'TS',
22   TARGET_COLNAME     => 'CLOSE',
23   CONFIG_OBJECT       => {{ '{{$' }} 'method':'fast','on_error':'skip'
      {{ '}}' }}
24 );
25
26 CREATE OR REPLACE TEMP TABLE TMP_FC AS
27 SELECT
28   SERIES::STRING          AS SYMBOL,
29   CAST(TS AS DATE)        AS TS,
30   FORECAST                AS PREDICTED_CLOSE,
31   'SNOWFLAKE_ML'         AS MODEL_NAME,
32   CURRENT_TIMESTAMP()     AS TRAINED_AT,
33   {{ var.value.forecast_horizon_days | default('14', true) }}::
      NUMBER AS HORIZON_D
34 FROM TABLE(PRICE_FORECASTER!FORECAST(
35   FORECASTING_PERIODS => {{ var.value.forecast_horizon_days |
      default('14', true) }}
36 ));
37
38 MERGE INTO MODEL.FORECASTS AS t
39 USING TMP_FC AS s
40   ON t.SYMBOL      = s.SYMBOL
41  AND t.TS          = s.TS
42  AND t.MODEL_NAME  = s.MODEL_NAME
43 WHEN MATCHED THEN UPDATE SET
44   PREDICTED_CLOSE = s.PREDICTED_CLOSE,
45   TRAINED_AT      = s.TRAINED_AT,
46   HORIZON_D       = s.HORIZON_D,
47   LOAD_TS         = CURRENT_TIMESTAMP()
48 WHEN NOT MATCHED THEN INSERT (
49   SYMBOL, TS, PREDICTED_CLOSE, MODEL_NAME, TRAINED_AT, HORIZON_D
50 ) VALUES (
51   s.SYMBOL, s.TS, s.PREDICTED_CLOSE, s.MODEL_NAME, s.TRAINED_AT, s.
      HORIZON_D
52 );
53 COMMIT;

```

Listing 4: Rebuild ANALYTICS final table.

Final Union Build (ACTUAL UNION FORECAST)

```

1 BEGIN;
2 USE SCHEMA ANALYTICS;
3

```

```

4 WITH symbols AS (
5     SELECT value::string AS symbol
6     FROM TABLE(FLATTEN(input => PARSE_JSON('{{ var.value.stock_symbols
7         }}}'))
8 )
9 TRUNCATE TABLE ANALYTICS.FINAL_PRICES_FORECAST;
10
11 -- ACTUALS from RAW
12 INSERT INTO ANALYTICS.FINAL_PRICES_FORECAST (SYMBOL, TS, CLOSE,
13     SOURCE, MODEL_NAME)
14 SELECT
15     sp.SYMBOL,
16     CAST(sp.TS AS DATE) AS TS,
17     sp.CLOSE,
18     'ACTUAL' AS SOURCE,
19     NULL AS MODEL_NAME
20 FROM RAW.STOCK_PRICES sp
21 JOIN symbols s ON s.symbol = sp.SYMBOL;
22
23 -- FORECASTS from MODEL
24 INSERT INTO ANALYTICS.FINAL_PRICES_FORECAST (SYMBOL, TS, CLOSE,
25     SOURCE, MODEL_NAME)
26 SELECT
27     f.SYMBOL,
28     f.TS,
29     f.PREDICTED_CLOSE AS CLOSE,
30     'FORECAST' AS SOURCE,
31     f.MODEL_NAME
32 FROM MODEL.FORECASTS f
33 JOIN symbols s ON s.symbol = f.SYMBOL;
34 COMMIT;

```

4.4 Transactions & Error Handling

- **DAG #1:** `snowflake.connector` uses `conn.autocommit(False)` and wraps DDL/DML in `try/except` with `commit()` on success and `rollback()` on failure.
- **DAG #2:** Each `SnowflakeOperator` task uses `BEGIN; ...COMMIT;` so the entire step is atomic.

5 Results

After triggering `yfinance_etl` and then `ml_forecast`, we validated:

- `RAW.STOCK_PRICES` contains daily OHLCV for each ticker.

- `MODEL.FORECASTS` contains `forecast_horizon_days` predictions per ticker with `MODEL_NAME='SNOWFLAKE_ML'`.
- `ANALYTICS.FINAL_PRICES_FORECAST` holds both `ACTUAL` and `FORECAST` rows.

Validation Queries

```

1 SELECT 'RAW' AS t, COUNT(*) c FROM RAW.STOCK_PRICES
2 UNION ALL SELECT 'MODEL', COUNT(*) FROM MODEL.FORECASTS
3 UNION ALL SELECT 'ANALYTICS', COUNT(*) FROM ANALYTICS.
  FINAL_PRICES_FORECAST;
4
5 SELECT SYMBOL, MIN(TS) AS min_ts, MAX(TS) AS max_ts, COUNT(*) AS n
6 FROM ANALYTICS.FINAL_PRICES_FORECAST
7 GROUP BY SYMBOL
8 ORDER BY SYMBOL, min_ts;

```

6 Discussion

Pros. Snowflake-side training avoids heavy Python deps and simplifies daily orchestration. Airflow manages schedule, retries, and secret handling via Connections/Variables. Transactions ensure atomic loads and clear failure modes.

Limitations. Calendar effects (weekends/holidays) and lack of exogenous regressors.

Future Work. Switch `CONFIG_OBJECT` method from `'fast'` to `'best'` for accuracy, add exchange calendars/holidays, or incorporate regressors (macro/news) via Snowflake features.

7 Conclusion

The pipeline meets the lab goals: ETL → Snowflake ML Forecast → union table; clean Airflow orchestration with transactions; secure configuration using Connections/Variables; and results suitable for analytics and visualization.

Appendix A: Colab (Read-Only) Snippet

Use Colab only to *read* results for plotting; credentials are prompted at runtime (not saved).

Listing 5: Colab: query final table and plot.

```

1 !pip install snowflake-connector-python pandas matplotlib --quiet
2
3 import snowflake.connector, pandas as pd
4 from getpass import getpass
5 import matplotlib.pyplot as plt
6
7 conn = snowflake.connector.connect(
8     account=input("Account: "),

```



```

9     user=input("User: "),
10    password=getpass("Password: "),
11    warehouse="CATFISH_QUERY_WH",
12    database="USER_DB_CATFISH",
13    schema="ANALYTICS",
14 )
15 df = pd.read_sql("""
16     SELECT SYMBOL, TS, CLOSE, SOURCE
17     FROM FINAL_PRICES_FORECAST
18     WHERE SYMBOL IN ('AAPL','MSFT','TSLA')
19     ORDER BY SYMBOL, TS
20 """, conn)
21 conn.close()
22
23 for sym, g in df.groupby("SYMBOL"):
24     g = g.sort_values("TS")
25     plt.figure()
26     plt.plot(g[g["SOURCE"]=="ACTUAL"]["TS"], g[g["SOURCE"]=="ACTUAL"]
27             ["CLOSE"], label=f"{sym} ACTUAL")
28     plt.plot(g[g["SOURCE"]=="FORECAST"]["TS"], g[g["SOURCE"]=="
29             FORECAST"]["CLOSE"], linestyle="--", label=f"{sym} FORECAST")
30     plt.title(f"{sym}: Actual vs Forecast")
31     plt.legend(); plt.xticks(rotation=45); plt.tight_layout(); plt.
32     show()

```

References

- yfinance (PyPI): <https://pypi.org/project/yfinance/>
- Apache Airflow: <https://airflow.apache.org/>
- Snowflake Documentation: <https://docs.snowflake.com/>