# DATA 226: DATA WAREHOUSE AND PIPELINE

# INSTRUCTOR: PROF KEEYONG HAN

# ASSIGNMENT 5

## STEP 1:

Create tasks using @task decorator
- You can use as many tasks as you want
- Schedule the tasks properly (task dependency)

```python
@task
def ensure_objects():
    ddl_schema = f"CREATE SCHEMA IF NOT EXISTS {TARGET_SCHEMA}"
    ddl_table = f"""
    CREATE TABLE IF NOT EXISTS {TARGET_SCHEMA}.{TARGET_TABLE} (
      SYMBOL VARCHAR(10) NOT NULL,
      DATE    DATE        NOT NULL,
      OPEN    FLOAT       NOT NULL,
      CLOSE   FLOAT       NOT NULL,
      HIGH    FLOAT       NOT NULL,
      LOW     FLOAT       NOT NULL,
      VOLUME INTEGER      NOT NULL,
      PRIMARY KEY (SYMBOL, DATE)
    )
    """

    conn, cur = _snowflake_cursor()
    try:
        # If your role cannot create schemas, pre-create RAW once and comment the next line.
        cur.execute(ddl_schema)
        cur.execute(ddl_table)
        conn.commit()
        print(f"Ensured {TARGET_SCHEMA}.{TARGET_TABLE} exists.")
    except Exception as e:
        conn.rollback()
        # Optional: richer error output
        try:
            from snowflake.connector.errors import Error as SFError
            if isinstance(e, SFError):
                print(f"Snowflake error: sqlstate={e.sqlstate}, errno={e.errno}, msg={e.msg}")
        except Exception:
            pass
        raise
    finally:
        cur.close()
        conn.close()
```

```python
@task
def extract_prices(symbol: str, days: int = 90) -> List[Dict]:
    api_key = Variable.get("vantage_api_key")
    if not api_key:
        raise ValueError("Missing Airflow Variable: 'vantage_api_key'")

    url = (
        "https://www.alphavantage.co/query"
        f"?function=TIME_SERIES_DAILY&symbol={symbol}"
        f"&outputsize=compact&datatype=json&apikey={api_key}"
    )
    r = requests.get(url, timeout=30)
    r.raise_for_status()
    data = r.json()

    # Alpha Vantage throttle note or error
    if "Time Series (Daily)" not in data:
        raise RuntimeError(f"Alpha Vantage response error: {data}")

    series = data["Time Series (Daily)"]
    cutoff = (datetime.now(timezone.utc) - timedelta(days=days)).date()

    rows: List[Dict] = []
    for d_str, vals in series.items():
        d = datetime.strptime(d_str, "%Y-%m-%d").date()
        if d >= cutoff:
            rows.append({
                "symbol": symbol.upper(),
                "date":   d,    # Python date (binds cleanly to Snowflake DATE)
                "open":   float(vals["1. open"]),
                "high":   float(vals["2. high"]),
                "low":    float(vals["3. low"]),
                "close":  float(vals["4. close"]),
                "volume": int(vals["5. volume"]),
            })

    rows.sort(key=lambda r: r["date"])  # oldest -> newest
    return rows
```

```python
@task
def load_prices_idempotent(rows: List[Dict]):
    """
    Load via a session TEMP table (unqualified) then MERGE into RAW.ASSIGNMENTS.
    Using an unqualified TEMP table avoids needing CREATE privilege on RAW for the stage.
    """
    if not rows:
        print("No rows to load.")
        return

    conn, cur = _snowflake_cursor()
    try:
        cur.execute("BEGIN")

        # IMPORTANT: no schema qualifier for TEMP table
        cur.execute("""
            CREATE TEMP TABLE IF NOT EXISTS ASSIGNMENTS_STAGE (
                SYMBOL VARCHAR(10),
                DATE    DATE,
                OPEN   FLOAT,
                CLOSE  FLOAT,
                HIGH   FLOAT,
                LOW    FLOAT,
                VOLUME INTEGER
            )
        """)
        cur.execute("DELETE FROM ASSIGNMENTS_STAGE")

        insert_sql = """
            INSERT INTO ASSIGNMENTS_STAGE
            (SYMBOL, DATE, OPEN, CLOSE, HIGH, LOW, VOLUME)
            VALUES (%(symbol)s, %(date)s, %(open)s, %(close)s, %(high)s, %(low)s, %(volume)s)
        """
        cur.executemany(insert_sql, rows)

        merge_sql = f"""
            MERGE INTO {TARGET_SCHEMA}.{TARGET_TABLE} t
            USING ASSIGNMENTS_STAGE s
            ON t.SYMBOL = s.SYMBOL AND t.DATE = s.DATE
            WHEN MATCHED THEN UPDATE SET
                OPEN = s.OPEN,
                CLOSE = s.CLOSE,
                HIGH = s.HIGH,
```

```
cur.executemany(insert_sql, rows)

merge_sql = f"""
    MERGE INTO {TARGET_SCHEMA}.{TARGET_TABLE} t
    USING ASSIGNMENTS_STAGE s
    ON t.SYMBOL = s.SYMBOL AND t.DATE = s.DATE
    WHEN MATCHED THEN UPDATE SET
        OPEN = s.OPEN,
        CLOSE = s.CLOSE,
        HIGH = s.HIGH,
        LOW = s.LOW,
        VOLUME = s.VOLUME
    WHEN NOT MATCHED THEN INSERT (SYMBOL, DATE, OPEN, CLOSE, HIGH, LOW, VOLUME)
    VALUES (s.SYMBOL, s.DATE, s.OPEN, s.CLOSE, s.HIGH, s.LOW, s.VOLUME)
"""
        cur.execute(merge_sql)
        conn.commit()
    except Exception as e:
        conn.rollback()
        try:
            from snowflake.connector.errors import Error as SFError
            if isinstance(e, SFError):
                print(f"Snowflake error: sqlstate={e.sqlstate}, errno={e.errno}, msg={e.msg}")
        except Exception:
            pass
        raise
    finally:
        cur.close()
        conn.close()

# wiring
ensure = ensure_objects()
data = extract_prices(SYMBOL)
ensure >> load_prices_idempotent(data)
```
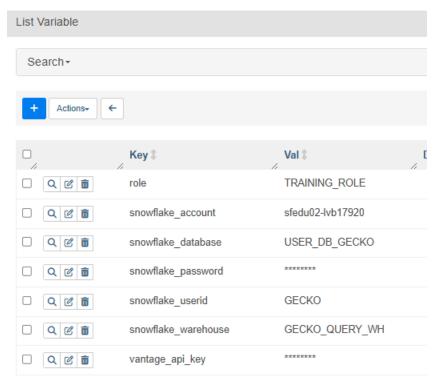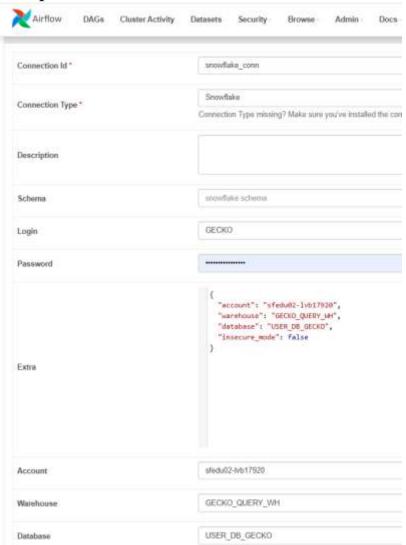
## Step 2:
Set up a variable for Alpha Vantage API key

List Variable

Search ▾

| + | Actions▾ | ← |

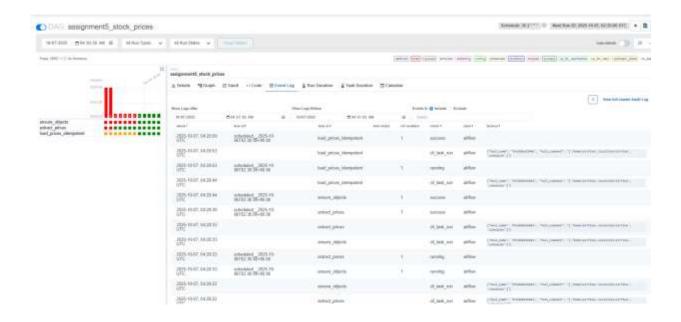| | | Key ⇕ | Val ⇕ | |
|---|---|---|---|---|
| ☐ | Q ✎ 🗑 | role | TRAINING_ROLE | |
| ☐ | Q ✎ 🗑 | snowflake_account | sfedu02-lvb17920 | |
| ☐ | Q ✎ 🗑 | snowflake_database | USER_DB_GECKO | |
| ☐ | Q ✎ 🗑 | snowflake_password | ******** | |
| ☐ | Q ✎ 🗑 | snowflake_userid | GECKO | |
| ☐ | Q ✎ 🗑 | snowflake_warehouse | GECKO_QUERY_WH | |
| ☐ | Q ✎ 🗑 | vantage_api_key | ******** | |

Step 3:
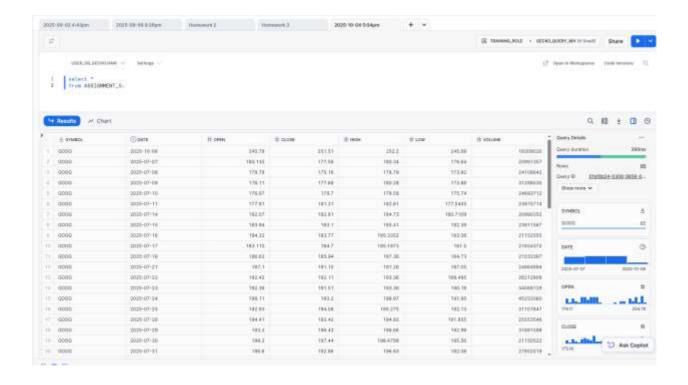Set up Snowflake Connection



Step 4: Ensure the overall DAG is implemented properly and runs successfully



Step 5: Capture two screenshot of your Airflow Web UI

## Snowflake Answer:



The Github Link:
https://github.com/NAMAN-CHHEDA/DATA-226-ASSIGNEMNT-5-