

Final Project Report

1. Team members

Temirbayev Zhassulan

Kumbayeva Ingkar

Duisenbek Gulzat

2. API justification

We have selected API for binance, which is a trading platform, that provides real time data about cryptocurrency prices. Since the market is very volatile, prices are changing frequently, about once or twice a minute. It is a trustworthy source, so the data can be considered real. Data is sent in a JSON format which makes it easier to work with.

<https://api.binance.com/api/v3/ticker/24hr?symbol=BTCUSDT>

3. Kafka topic schema

raw_events - stores unprocessed API responses

```
{  
  "event_id": "uuid",  
  "source": "binance",  
  "ingested_at": "2025-12-19T14:26:46.552956+00:00",  
  "payload": {  
    "symbol": "BTCUSDT",  
    "lastPrice": "88006.09",  
    "volume": "27033.33",  
    "openTime": 1766068006013,  
    "closeTime": 1766154406013,  
  }  
}
```

4. Cleaning rules

- Drop rows missing symbol, price, or event_time
- Drop rows where price ≤ 0
- Remove duplicate event_id
- Remove duplicate (symbol, event_time) combinations
- Convert numeric columns to proper types (float/int)
- Standardize symbol format: uppercase, trimmed

5. SQLite Schema

```
events (
    event_id TEXT PRIMARY KEY,
    symbol TEXT NOT NULL,
    price REAL NOT NULL,
    event_time TEXT NOT NULL,
    source TEXT NOT NULL,
    ingested_at TEXT NOT NULL,

    priceChange REAL,
    priceChangePercent REAL,
    weightedAvgPrice REAL,
    prevClosePrice REAL,
    lastPrice REAL,
    lastQty REAL,
    bidPrice REAL,
    bidQty REAL,
    askPrice REAL,
    askQty REAL,
    openPrice REAL,
    highPrice REAL,
    lowPrice REAL,
    volume REAL,
    quoteVolume REAL,
    openTime INTEGER,
    closeTime INTEGER,
    firstId INTEGER,
    lastId INTEGER,
    count INTEGER
);
```

```
daily_summary (
    day TEXT NOT NULL,
    symbol TEXT NOT NULL,
    n_records INTEGER,
    first_event_time TEXT,
    last_event_time TEXT,
    open_price REAL,
    high_price REAL,
```

```

    low_price REAL,
    close_price REAL,
    avg_price REAL,
    min_price REAL,
    max_price REAL,
    avg_volume REAL,
    sum_volume REAL,
    avg_quote_volume REAL,
    sum_quote_volume REAL,
    avg_count REAL,
    sum_count REAL,
    avg_price_change REAL,
    avg_price_change_pct REAL,
    avg_spread REAL,
    avg_spread_pct REAL,
    daily_return REAL,
    log_return REAL,
    volatility REAL,
    computed_at TEXT NOT NULL,
    PRIMARY KEY (day, symbol)
);

```

6. Screenshots

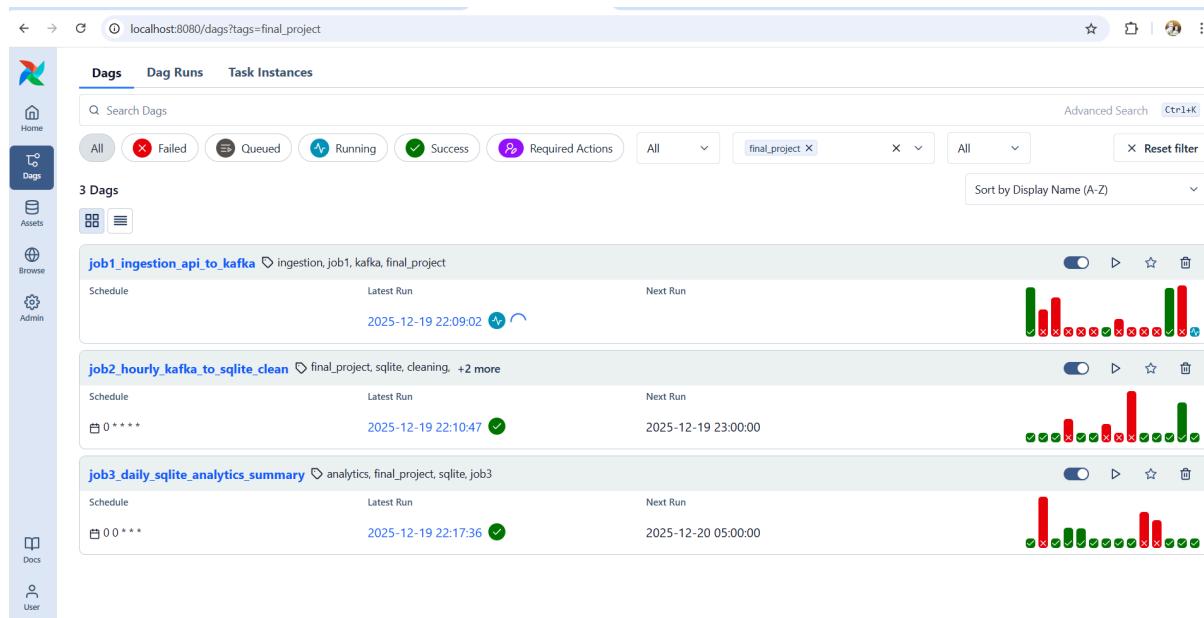


Figure 1. Overview of Airflow DAGs Execution

Figure 1 shows the Apache Airflow web interface displaying all DAGs related to the final_project pipeline.

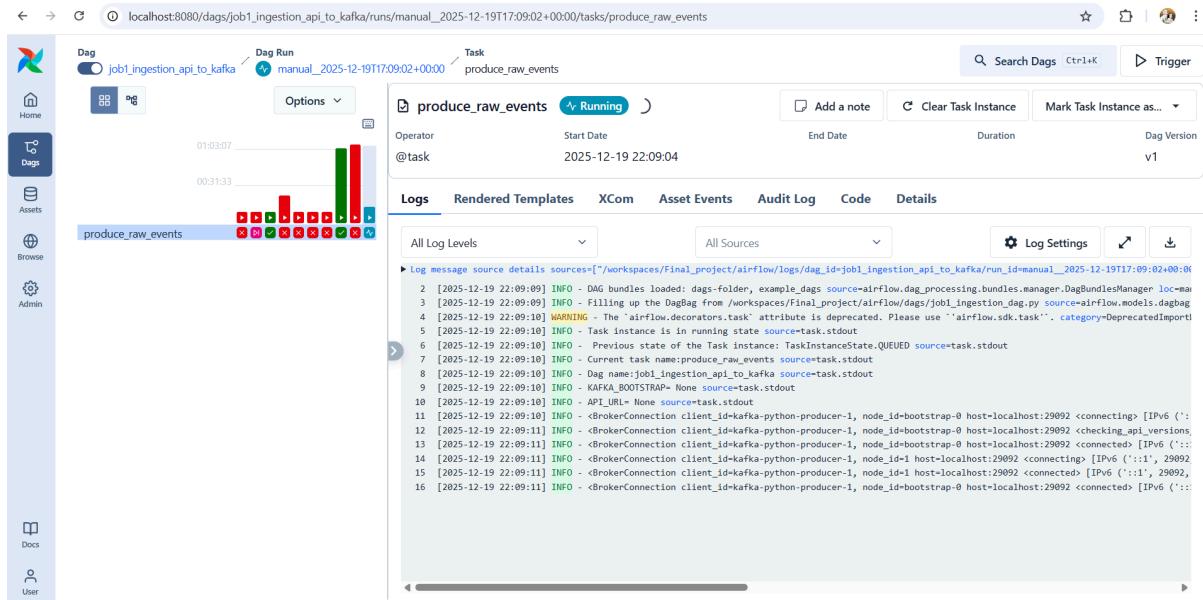


Figure 2. DAG 1 – API Ingestion to Kafka (Job 1)

This figure shows the execution of DAG 1, which is responsible for ingesting data from the external API and publishing it to a Kafka topic.

```

• @InkarKum →/workspaces/Final_project (main) $ docker exec -it final_project-kafka-1 bash -lc '
kafka-console-consumer \
--bootstrap-server kafka:9092 \
--topic raw_events \
--from-beginning \
--max-messages 3 \
--timeout-ms 15000
'

{"event_id": "3c0a0f88-e000-48bf-8255-c02f3368a36e", "source": "binance", "ingested_at": "2025-12-19T14:26:16.324017+00:00", "payload": {"symbol": "BTCUSDT", "priceChange": "-478.1700000", "priceChangePercent": "-0.540", "weightedAvgPrice": "86753.39950733", "prevClosePrice": "88485.0700000", "lastPrice": "88006.9000000", "lastQty": "0.00044000", "bidPrice": "88006.9000000", "bidQty": "1.57164000", "askPrice": "88006.9100000", "askQty": "1.64070000", "openPrice": "88485.0700000", "highPrice": "89362.7900000", "lowPrice": "84450.0100000", "volume": "27040.07056000", "quoteVolume": "2345818043.99810160", "openTime": 1766067976013, "closeTime": 1766154376013, "firstId": 5677474367, "lastId": 5685238592, "count": 7764226}}
{"event_id": "2f4266cf-769b-45f2-b6e8-988d754b187b", "source": "binance", "ingested_at": "2025-12-19T14:26:46.552956+00:00", "payload": {"symbol": "BTCUSDT", "priceChange": "-500.1200000", "priceChangePercent": "-0.565", "weightedAvgPrice": "86752.90331844", "prevClosePrice": "88506.2100000", "lastPrice": "88006.0900000", "lastQty": "0.00226000", "bidPrice": "88006.0900000", "bidQty": "1.45080000", "askPrice": "88006.1000000", "askQty": "3.24192000", "openPrice": "88506.2100000", "highPrice": "89362.7900000", "lowPrice": "84450.0100000", "volume": "27033.33077000", "quoteVolume": "2345219930.66519880", "openTime": 1766068006013, "closeTime": 1766154406013, "firstId": 5677477025, "lastId": 5685240605, "count": 7763581}}
{"event_id": "c3a526da-02d0-4689-9241-53e7ae0ddb62", "source": "binance", "ingested_at": "2025-12-19T14:27:16.752865+00:00", "payload": {"symbol": "BTCUSDT", "priceChange": "-497.4500000", "priceChangePercent": "-0.562", "weightedAvgPrice": "86752.09715245", "prevClosePrice": "88490.8100000", "lastPrice": "87993.3500000", "lastQty": "0.81248000", "bidPrice": "87993.3500000", "bidQty": "0.13243000", "askPrice": "87993.3600000", "askQty": "5.96701000", "openPrice": "88490.8000000", "highPrice": "89362.7900000", "lowPrice": "84450.0100000", "volume": "27022.32434000", "quoteVolume": "2344243306.42881230", "openTime": 1766068836007, "closeTime": 1766154436007, "firstId": 5677481564, "lastId": 5685241971, "count": 7760408}}
Processed a total of 3 messages
○ @InkarKum →/workspaces/Final_project (main) $

```

Figure 3. Kafka topic inspection (raw_events)

This figure demonstrates direct inspection of the Kafka topic using the Kafka console consumer. The output shows multiple raw JSON messages produced by Job 1, each containing market data fields such as symbol, prices, volumes, timestamps, and trade statistics.

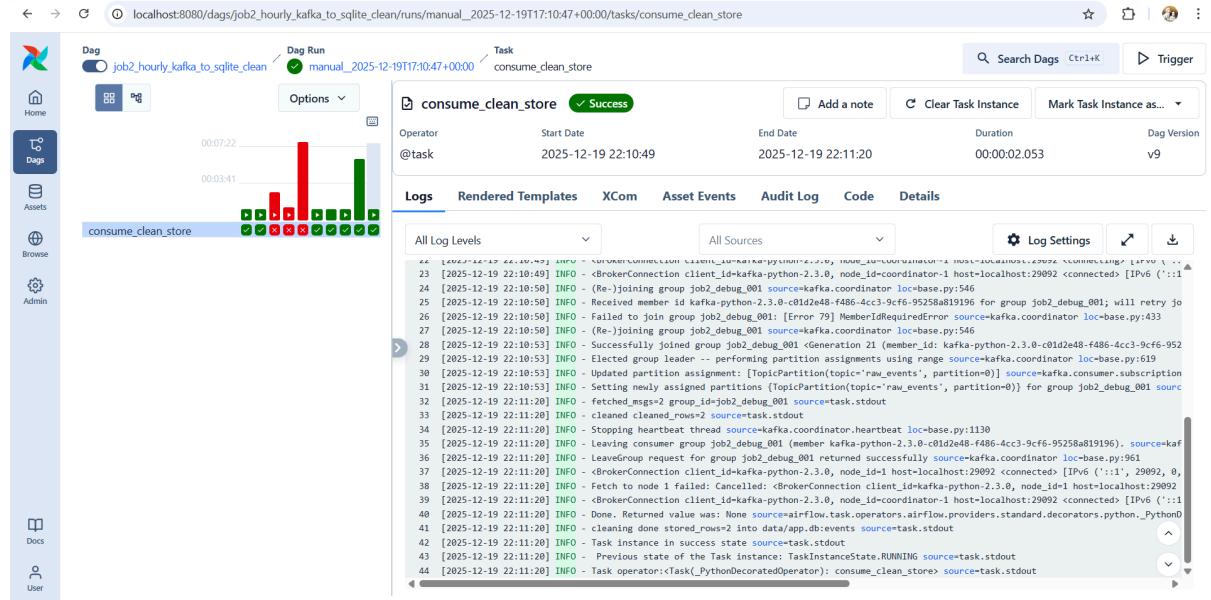


Figure 4. DAG 2 – Kafka to SQLite Cleaning and Storage (Job 2)

This figure presents DAG 2, which consumes raw Kafka messages, cleans and normalizes the data, and stores it in a SQLite database.

```

job1_producer.py 1  job3_analytics.py 312.py  docker-compose.yml  job2_cleaner.py 2  job2_clean_store_dag.py  job3_daily_summary_d...
src > __pycache__ >  job3_analytics.py-312.pyc
The file is not displayed in the text editor because it is either binary or
uses an unsupported text encoding.
Open Anyway

PROBLEMS 11 OUTPUT DEBUG CONSOLE TERMINAL 4
PORTS
Port Forward
8080 127.0.0.1: kafka-console-consumer \
8793 127.0.0.1: --bootstrap-server kafka:9092 \
8794 127.0.0.1: --topic raw_events \
29092 127.0.0.1: --from-beginning \
Add Port

TERMINE
Port Forward
8080 127.0.0.1: kafka-console-consumer \
8793 127.0.0.1: --bootstrap-server kafka:9092 \
8794 127.0.0.1: --topic raw_events \
29092 127.0.0.1: --from-beginning \
eighitedAvgPrice": "86752.9931844", "prevClosePrice": "88966.2100000", "lastPrice": "88006.0900000", "lastQty": "0.00226000", "bidPrice": "88006.0900000", "bidQty": "1.4500000", "askPrice": "88006.1000000", "askQty": "3.24192000", "openPrice": "88966.2100000", "highPrice": "89362.7900000", "lowPrice": "84458.0100000", "volume": "27033.3307000", "quoteVolume": "2345219930.66519880", "openTime": 17660680806013, "closeTime": 1766154406013, "firstId": 5677474020, "lastId": 5685246095, "count": 7763581}
{"event_id": "c3aeb41fe-4be3-4bb8-bf64-822c962fd6f", "symbol": "BTUSDT", "priceChange": "-497.4500000", "priceChangePercent": "-0.562", "weightedAvgPrice": "86752.09715245", "prevClosePrice": "88940.8100000", "lastPrice": "87993.3500000", "lastQty": "0.81242000", "bidPrice": "87993.3500000", "bidQty": "0.13243000", "askPrice": "87993.3600000", "askQty": "5.06701000", "openPrice": "88490.8000000", "highPrice": "89362.7900000", "lowPrice": "84458.0100000", "volume": "27022.32434000", "quoteVolume": "234423306.42881230", "openTime": 17660680806007, "closeTime": 1766154436007, "firstId": 5677481564, "lastId": 5685241971, "count": 7764088}
Processed a total of 3 messages
@linkarkum ~ /workspaces/Final_project (main) $ python -c "PY"
import sqlite3
con = sqlite3.connect("data/app.db")
cur = con.cursor()
print("events count:", cur.execute("SELECT COUNT(*) FROM events").fetchone()[0])
print("one row:", cur.execute("SELECT * FROM events ORDER BY event_time DESC LIMIT 1").fetchone())
PY
events_count: 103
one_row: ('3ae4b1fe-4be3-4bb8-bf64-822c962fd6f', 'BTUSDT', 88225.42, '2025-12-19T17:11:12.008000+00:00', 'binance', '2025-12-19T17:11:12.698925+00:00', None, 1621.42, 1.872, 86770.90368861, 86605.74, 88225.42, 0.0033, 88225.41, 2.16934, 88225.42, 0.86658, 86664.0, 89399.7, 84450.01, 27241.9627, 2363979710.682837, 17660680806008, 5679241363, 568678998, 7557636)
@linkarkum ~ /workspaces/Final_project (main) $

```

Figure 5. SQLite verification – events table

This figure shows a direct query of the SQLite database after Job 2 execution. The query confirms that cleaned records are stored in the events table, including price, symbol, timestamps, and additional numeric fields extracted from the API response.

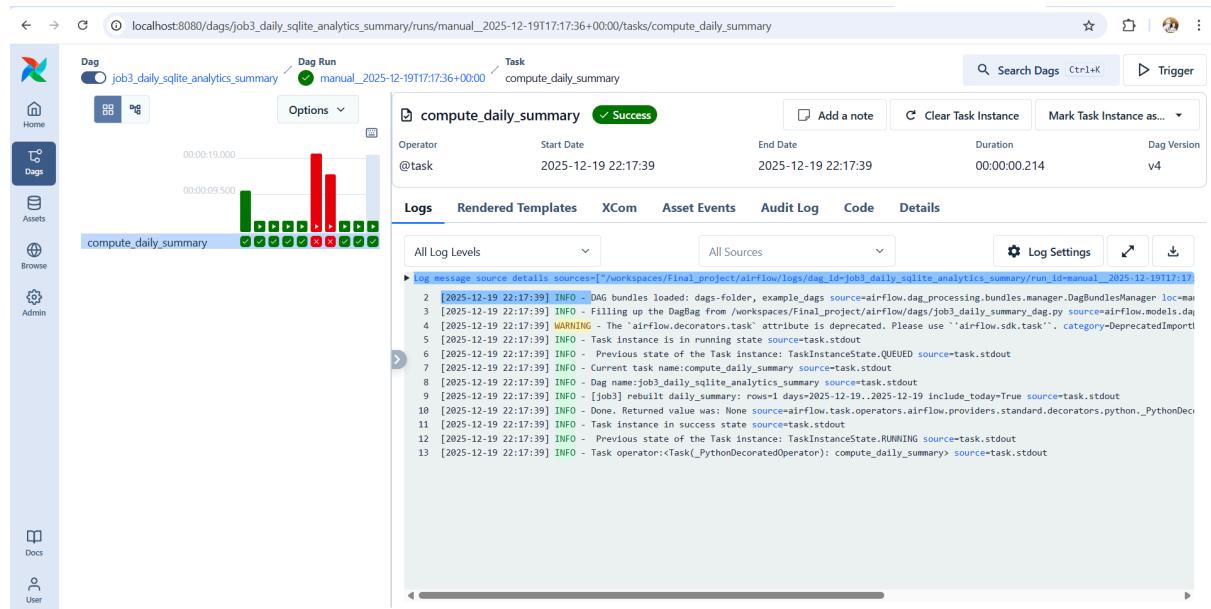
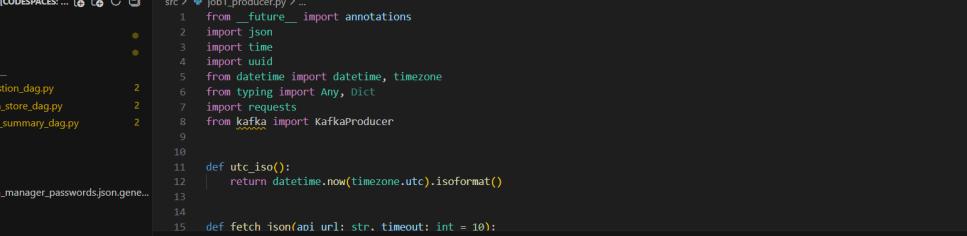


Figure 5. DAG 3 – Daily Analytics and Aggregation (Job 3)

This figure shows DAG 3, which performs analytical processing on the cleaned data stored in SQLite.



The screenshot shows a Jupyter Notebook interface with the following details:

- File Explorer:** The left sidebar shows a project structure for "FINAL_PROJECT [CODESPACES: potential journey]". It includes a ".env" file, an "airflow" directory containing "airflow.cfg", "airflow.db", and "simple_auth_manager_passwords.json", a "data" directory, a "src" directory containing "job1_producer.py", "job2_cleaner.py", "job3_analytics.py", and ".gitignore", and files like "README.md" and "requirements.txt".
- Code Cell:** The main area contains a code cell for "job1_producer.py". The code imports annotations from future, json, time, uuid, datetime, timezone, Any, Dict, requests, and KafkaProducer. It defines two functions: `def utc_iso()` which returns the current UTC timestamp in ISO format, and `def fetch_json(api_url: str, timeout: int = 10)`.
- Terminal:** Below the code cell is a terminal tab labeled "TERMINAL". It shows a session with the user "@InkarKum" running a Python command to execute the code.
- Ports:** A "PORTS" section is visible, showing forwarded ports 8080, 8793, 8794, and 29092 to local ports 127.0.0.1:12345, 12346, 12347, and 12348 respectively. An "Add Port" button is present.
- Bottom Status Bar:** The status bar at the bottom indicates the file is "CODESPACES: potential journey", the cell index is "4", and the Python version is "3.12.1".

Figure 6. Job 3 – Daily analytics and aggregation

This figure shows Job 3 execution, which reads cleaned event data from SQLite and computes daily analytical metrics such as minimum, maximum, average price, and record counts per symbol. The task successfully rebuilds the daily summary table on each run.