# Kafka-python Pipelines with Apache Airflow

## (*Proof of concept*)   *Last Edit 5/22/2025*

**Config settings:**
- acks='all',  # Wait for all replicas to acknowledge,
- batch_size=16384,  # 16KB batch size
- linger_ms=5,  # Wait up to 5ms to fill batches

**With 3 retries**
On *UnknownTopicOrPartitionError*

Monitoring
Topics, Messages

**UI for Apache Kafka**
*Python lib*
https://github.com/provectus/kafka-ui

**kafka**

**Producer job (DAG)**
*Creates multiple topics:*
- *Generates fake data with Kafka-faker*
- *Sends & flushes in batches*

**1**

**Airflow**

**Apache Airflow**

Periodically triggers
Kafka Producer and Kafka Consumer
jobs

**The Legend**
- A component written by me
- Open source
- Comments

**kafka**

**Consumer job ( DAG)**
- *Subscribed to topics*
- *Pulls messages*
- *Flushes messages in batches to the db*

**2**

Scheduled
jobs

**3**

**Microsoft SQL Server**

DLQ (dead-letter queues) python-based implementation
for failed messages using a secondary topic

**Topic-Pattern Subscription**
Subscribed to a List of Topics via *consumer.subscribe(pattern=r'test-.*')*

**EXTRACT**
(via python modules)
- *All of the integration logic is written here*

**5**

**Load**
(via stored procs)

**PROD**
*Final destination*

**STAGE**
1. Data Tables
2. 'Metadata Framework'
    -- Tbls & stored procs

**4**

**Transform**
(via stored procs)

**Config settings:**
```
auto_offset_reset='earliest',      # Start consuming from the topic's start
enable_auto_commit=False,          # Don't auto-commit to avoid skipping on retries
group_id= self.group_id,           # Consumer group ID for Kafka offset tracking
consumer_timeout_ms=5000,          # Stop after 5s if no new messages
```
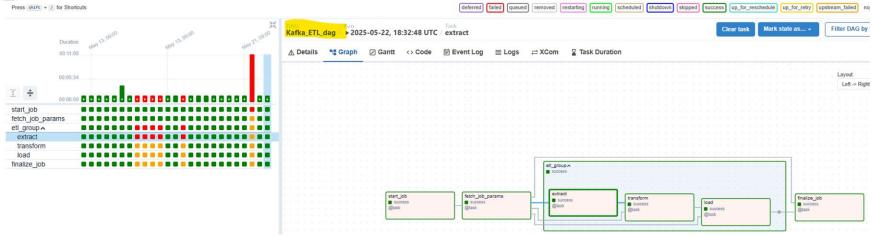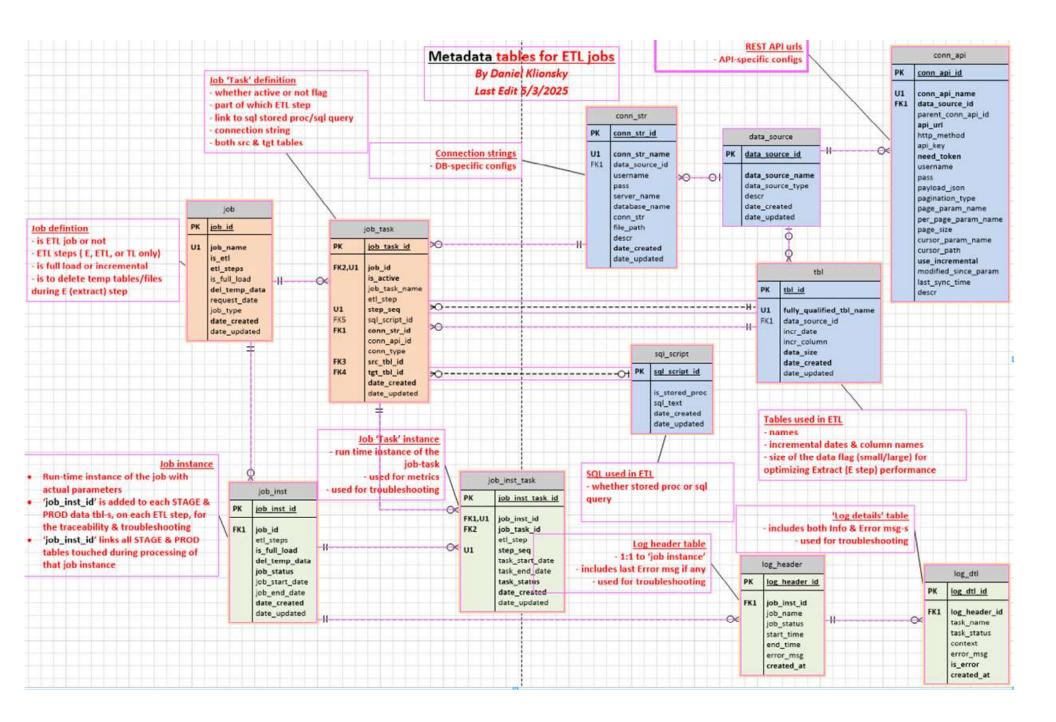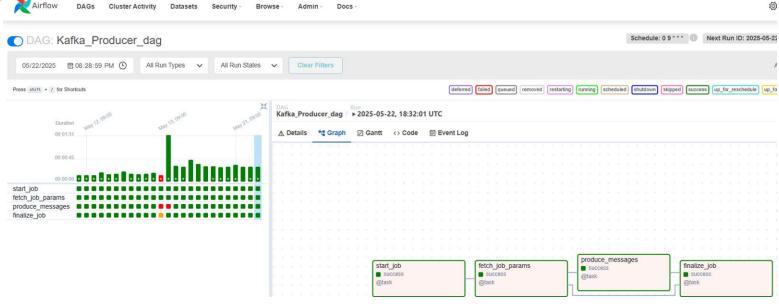
## 1. Kafka Consumer DAG

## 2. Database log for Kafka Consumerr DAG

| ob_name | job_inst_id | task_name | task_status | context | error_msg | is_error | created_at |
|---|---|---|---|---|---|---|---|
| ClientD_Kafka_ETL | 531 | finalize_job | succeeded | finalize_job | *** FINISHED | 0 | 2025-05-22 11:33:28.700 |
| ClientD_Kafka_ETL | 531 | load | running | load | Skipping 'L' step | 0 | 2025-05-22 11:33:27.030 |
| ClientD_Kafka_ETL | 531 | transform | running | transform | Skipping 'T' step | 0 | 2025-05-22 11:33:25.090 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | JobTask.process() | Finished task ‖ extract-kafka_topic-step-1 | 0 | 2025-05-22 11:33:23.327 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Kafka consumer closed. Message count: {'test-top... | 0 | 2025-05-22 11:33:22.890 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Reached max_messages for all topics. Exiting loop. | 0 | 2025-05-22 11:33:22.763 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Topic ‖ test-topic2 ‖ Reached max_messages of 1... | 0 | 2025-05-22 11:33:22.640 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | _flush_to_db | Topic ‖ test-topic2 ‖ Flushed 1000 records to ‖ dm... | 0 | 2025-05-22 11:33:22.480 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Topic ‖ test-topic3 ‖ Reached max_messages of 1... | 0 | 2025-05-22 11:33:13.770 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | _flush_to_db | Topic ‖ test-topic3 ‖ Flushed 1000 records to ‖ dm... | 0 | 2025-05-22 11:33:13.647 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Topic ‖ test-topic1 ‖ Reached max_messages of 1... | 0 | 2025-05-22 11:33:05.463 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | _flush_to_db | Topic ‖ test-topic1 ‖ Flushed 1000 records to ‖ dm... | 0 | 2025-05-22 11:33:05.323 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | run | Starting Kafka consumer group_id ‖ airflow-consu... | 0 | 2025-05-22 11:32:53.847 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | __init__ | KafkaConsumer initialized successfully; topic_patt... | 0 | 2025-05-22 11:32:53.347 |
| ClientD_Kafka_ETL | 531 | extract-kafka_topic-step-1 | running | JobTask.process() | [#1] Started task ‖ extract-kafka_topic-step-1 ‖ kaf... | 0 | 2025-05-22 11:32:53.003 |
| ClientD_Kafka_ETL | 531 | fetch_job_params | running | get_job_inst_info | Starting Job: ClientD_Kafka_ETL ‖ ETL steps: E ‖ ... | 0 | 2025-05-22 11:32:51.533 |
| ClientD_Kafka_ETL | 531 | start_job | started | sp sp_crud_log_header | *** STARTED ‖ ClientD_Kafka_ETL with job_inst_... | 0 | 2025-05-22 11:32:49.647 |

# Metadata tables for ETL jobs
### By Daniel Klionsky
### Last Edit 5/3/2025

**REST API urls**
- API-specific configs

**conn_api**

| | |
|---|---|
| PK | conn_api_id |
| U1 | conn_api_name |
| FK1 | data_source_id |
| | parent_conn_api_id |
| | api_url |
| | http_method |
| | api_key |
| | need_token |
| | username |
| | pass |
| | payload_json |
| | pagination_type |
| | page_param_name |
| | per_page_param_name |
| | page_size |
| | cursor_param_name |
| | cursor_path |
| | use_incremental |
| | modified_since_param |
| | last_sync_time |
| | descr |

**Job 'Task' definition**
- whether active or not flag
- part of which ETL step
- link to sql stored proc/sql query
- connection string
- both src & tgt tables

**conn_str**

| | |
|---|---|
| PK | conn_str_id |
| U1 | conn_str_name |
| FK1 | data_source_id |
| | username |
| | pass |
| | server_name |
| | database_name |
| | conn_str |
| | file_path |
| | descr |
| | date_created |
| | date_updated |

**data_source**

| | |
|---|---|
| PK | data_source_id |
| | data_source_name |
| | data_source_type |
| | descr |
| | date_created |
| | date_updated |

**Connection strings**
- DB-specific configs

**job**

| | |
|---|---|
| PK | job_id |
| U1 | job_name |
| | is_etl |
| | etl_steps |
| | is_full_load |
| | del_temp_data |
| | request_date |
| | job_type |
| | date_created |
| | date_updated |

**job_task**

| | |
|---|---|
| PK | job_task_id |
| FK2,U1 | job_id |
| | is_active |
| | job_task_name |
| | etl_step |
| U1 | step_seq |
| FK5 | sql_script_id |
| FK1 | conn_str_id |
| | conn_api_id |
| | conn_type |
| FK3 | src_tbl_id |
| FK4 | tgt_tbl_id |
| | date_created |
| | date_updated |

**Job defintion**
- is ETL job or not
- ETL steps ( E, ETL, or TL only)
- is full load or incremental
- is to delete temp tables/files during E (extract) step

**tbl**

| | |
|---|---|
| PK | tbl_id |
| U1 | fully_qualified_tbl_name |
| FK1 | data_source_id |
| | incr_date |
| | incr_column |
| | data_size |
| | date_created |
| | date_updated |

**sql_script**

| | |
|---|---|
| PK | sql_script_id |
| | is_stored_proc |
| | sql_text |
| | date_created |
| | date_updated |

**Tables used in ETL**
- names
- incremental dates & column names
- size of the data flag (small/large) for optimizing Extract (E step) performance

**SQL used in ETL**
- whether stored proc or sql query

**Job 'Task' instance**
- run time instance of the job-task
- used for metrics
- used for troubleshooting

**Job instance**
- Run-time instance of the job with actual parameters
- 'job_inst_id' is added to each STAGE & PROD data tbl-s, on each ETL step, for the traceability & troubleshooting
- 'job_inst_id' links all STAGE & PROD tables touched during processing of that job instance

**job_inst**

| | |
|---|---|
| PK | job_inst_id |
| FK1 | job_id |
| | etl_steps |
| | is_full_load |
| | del_temp_data |
| | job_status |
| | job_start_date |
| | job_end_date |
| | date_created |
| | date_updated |

**job_inst_task**

| | |
|---|---|
| PK | job_inst_task_id |
| FK1,U1 | job_inst_id |
| FK2 | job_task_id |
| | etl_step |
| U1 | step_seq |
| | task_start_date |
| | task_end_date |
| | task_status |
| | date_created |
| | date_updated |

**Log header table**
- 1:1 to 'job instance'
- includes last Error msg if any
- used for troubleshooting

**'Log details' table**
- includes both Info & Error msg-s
- used for troubleshooting

**log_header**

| | |
|---|---|
| PK | log_header_id |
| FK1 | job_inst_id |
| | job_name |
| | job_status |
| | start_time |
| | end_time |
| | error_msg |
| | created_at |

**log_dtl**

| | |
|---|---|
| PK | log_dtl_id |
| FK1 | log_header_id |
| | task_name |
| | task_status |
| | context |
| | error_msg |
| | is_error |
| | created_at |

**Kafka Producer**

## 1. AirFlow Graph for Kafka Producer DAG
*to generate sample messages*

## 2. Database log for Kafka Producer DAG

*Processing steps are logged here, in addition to AirFlow logging*