

## Real Time Orders Payments Data Ingestion

1. This project was about the ingestion of orders payments data in real time.
2. We needed 3 pub sub topics for this project. Topics are orders\_data, payments\_data and dlq\_payments\_data.
3. While creating topics, a default subscription gets created for the respective topics.
4. first we run producer of orders data so that some records might be able to ingest into cassandra table. Here attributes related to payment will be null when we check cassandra table.

orders data producer:

```
def generate_mock_data(order_id):
    items = ["Laptop", "Phone", "Book", "Tablet", "Monitor"]
    addresses = ["123 Main St, City A, Country", "456 Elm St, City B, Country", "789 Oak St, City C, Country"]
    statuses = ["Shipped", "Pending", "Delivered", "Cancelled"]

    return {
        "order_id": order_id,
        "customer_id": random.randint(100, 1000),
        "item": random.choice(items),
        "quantity": random.randint(1, 10),
        "price": random.uniform(100, 1500),
        "shipping_address": random.choice(addresses),
        "order_status": random.choice(statuses),
        "creation_date": "2024-01-21"
    }

order_id = 1
while True:
    data = generate_mock_data(order_id)
    json_data = json.dumps(data).encode('utf-8')

    try:
        future = publisher.publish(topic_path, data=json_data)
        future.add_done_callback(callback)
        future.result()
    except Exception as e:
        print(f"Exception encountered: {e}")

    time.sleep(2) # Wait for 2 seconds

    order_id += 1
    if order_id > 80:
        order_id = 1 # Reset the order_id back to 1 after reaching 500
```

orders data consumer:

```
34 cluster,session = cassandra_connection()
35
36 # Prepare the Cassandra insertion statement
37 insert_stmt = session.prepare("""
38     INSERT INTO orders_payments_facts (order_id, customer_id, item, quantity, price, shipping_address,
39     order_status, creation_date, payment_id, payment_method, card_last_four, payment_status, payment_datetime)
40     VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?)
41 """)
42
43 # Pull and process messages
44 def pull_messages():
45     while True:
46         response = subscriber.pull(request={"subscription": subscription_path, "max_messages": 10})
47         ack_ids = []
48
49         for received_message in response.received_messages:
50             # Extract JSON data
51             json_data = received_message.message.data.decode('utf-8')
52
53             # Deserialize the JSON data
54             deserialized_data = json.loads(json_data)
55
56             print(deserialized_data)
57
58             # Prepare data for Cassandra insertion
59             cassandra_data = (
60                 deserialized_data.get("order_id"),
61                 deserialized_data.get("customer_id"),
62                 deserialized_data.get("item"),
63                 deserialized_data.get("quantity"),
64                 deserialized_data.get("price"),
65                 deserialized_data.get("shipping_address"),
66                 deserialized_data.get("order_status"),
67                 deserialized_data.get("creation_date"),
68                 None,
69                 None,
70                 None,
71                 None,
72                 None
73             )
74
75             # Insert data into Cassandra
76             session.execute(insert_stmt, cassandra_data)
77
78             print("Data inserted in cassandra !!")
79
80             # Collect ack ID for acknowledgment
81             ack_ids.append(received_message.ack_id)
82
83             # Acknowledge the messages so they won't be sent again
84             if ack_ids:
85                 subscriber.acknowledge(request={"subscription": subscription_path, "ack_ids": ack_ids})
86
87
```

5. Then we start the payments producer so that it can push the data related to the respective order ids.

```
# Initialize the Pub/Sub publisher client
publisher = pubsub_v1.PublisherClient()

# Project and Topic details
project_id = "big-data-projects-411817"
topic_name = "payments_data"
topic_path = publisher.topic_path(project_id, topic_name)

# Payment methods mock data
payment_methods = ["Credit Card", "Debit Card", "PayPal", "Google Pay", "Apple Pay"]

def generate_mock_payment(order_id):
    return {
        # Starting from 1001 as in the example
        "payment_id": order_id + 1000,
        "order_id": order_id,
        "payment_method": choice(payment_methods),
        # Just using order_id to make last 4 digits
        "card_last_four": str(order_id).zfill(4)[-4:],
        "payment_status": "Completed",
        # Using order_id to vary the hour for variety
        "payment_datetime": f"2024-01-21T{str(order_id).zfill(2)}:01:30Z"
    }

def callback(future):
    try:
        # Get the message_id after publishing.
        message_id = future.result()
        print(f"Published message with ID: {message_id}")
    except Exception as e:
        print(f"Error publishing message: {e}")

# Produce mock payment data for each order_id and publish to the topic
for order_id in range(1, 501): # 500 order_ids
    mock_payment = generate_mock_payment(order_id)
    json_data = json.dumps(mock_payment).encode('utf-8')
    try:
        future = publisher.publish(topic_path, data=json_data)
        future.add_done_callback(callback)
        time.sleep(1) # Sleep for 2 seconds before producing the next message
    except Exception as e:
        print(f"Exception encountered: {e}")
```

6. Next executed ingestion in fact table using payments consumer which starts reading data from payments topic so that it can check order id existence in cassandra table and update the payment data for that order id there at table in real time.

```
# Setup Cassandra connection
cluster, session = cassandra_connection()

# Pull and process messages
def pull_messages():
    while True:
        response = subscriber.pull(request={"subscription": subscription_path, "max_messages": 10})
        ack_ids = []

        for received_message in response.received_messages:
            # Extract JSON data
            json_data = received_message.message.data.decode('utf-8')

            # Deserialize the JSON data
            deserialized_data = json.loads(json_data)

            # Check if order_id exists in the Cassandra table
            query = f"SELECT order_id FROM orders_payments_facts WHERE order_id = {deserialized_data.get('order_id')}"
            rows = session.execute(query)
            if rows.one(): # If order_id found
                # Assuming you have a 'payments' column, you can update it like this:
                update_query = """
                UPDATE orders_payments_facts
                SET payment_id = %,
                    payment_method = %,
                    card_last_four = %,
                    payment_status = %,
                    payment_datetime = %
                WHERE order_id = %
                """
                values = (
                    deserialized_data.get('payment_id'),
                    deserialized_data.get('payment_method'),
                    deserialized_data.get('card_last_four'),
                    deserialized_data.get('payment_status'),
                    deserialized_data.get('payment_datetime'),
                    deserialized_data.get('order_id')
                )
                session.execute(update_query, values)
                print("data updated in Cassandra -> ", deserialized_data)
            else: # if order_id not found
                publisher.publish(dlq_topic_path, data=json_data.encode('utf-8'))
                print("Data thrown in dlq because order_id not found -> ", deserialized_data)

        # Collect ack ID for acknowledgment
        ack_ids.append(received_message.ack_id)

    # Acknowledge the messages so they won't be sent again
    if ack_ids:
        subscriber.acknowledge(request={"subscription": subscription_path, "ack_ids": ack_ids})
```

7. If order id not found, those data thrown to dlq payments topic.

Google Cloud | big-data-projects | pubsub

Pub/Sub | dlq\_payments\_data | EDIT | + TRIGGER CLOUD FUNCTION | IMPORT | DELETE

Topics | Subscriptions | Snapshots | Schemas

Lite | Lite Reservations | Lite Topics | Lite Subscriptions

Release Notes

Export options have moved to the Create subscription dropdown menu under the Subscriptions tab below. GOT IT

Topic name: projects/big-data-projects-411817/topics/dlq\_payments\_data

SUBSCRIPTIONS	SNAPSHOTS	METRICS	DETAILS	MESSAGES
'payment_id': 1039, 'order_id': 39, 'payment_method': 'Debit Card', 'card_last_four':				Deadline exceeded
'payment_id': 1041, 'order_id': 41, 'payment_method': 'Apple Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1044, 'order_id': 44, 'payment_method': 'PayPal', 'card_last_four': '0044',				Deadline exceeded
'payment_id': 1045, 'order_id': 45, 'payment_method': 'Debit Card', 'card_last_four':				Deadline exceeded
'payment_id': 1047, 'order_id': 47, 'payment_method': 'Google Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1048, 'order_id': 48, 'payment_method': 'PayPal', 'card_last_four': '0048',				Deadline exceeded
'payment_id': 1049, 'order_id': 49, 'payment_method': 'Credit Card', 'card_last_four':				Deadline exceeded
'payment_id': 1051, 'order_id': 51, 'payment_method': 'Apple Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1053, 'order_id': 53, 'payment_method': 'Google Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1054, 'order_id': 54, 'payment_method': 'Apple Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1055, 'order_id': 55, 'payment_method': 'Debit Card', 'card_last_four':				Deadline exceeded
'payment_id': 1056, 'order_id': 56, 'payment_method': 'PayPal', 'card_last_four': '0056',				Deadline exceeded
'payment_id': 1058, 'order_id': 58, 'payment_method': 'PayPal', 'card_last_four': '0058',				Deadline exceeded
'payment_id': 1059, 'order_id': 59, 'payment_method': 'Google Pay', 'card_last_four':				Deadline exceeded
'payment_id': 1060, 'order_id': 60, 'payment_method': 'Google Pay', 'card_last_four':				Deadline exceeded

dlq\_payments\_data | PERMISSIONS | LABELS

Edit or delete permissions below, or select "Add Principal" to grant new access.

Show inherited permissions

Filter: Enter property name or value

Role / Principal

- Artifact Registry Service Agent (1)
- Cloud Build Service Account (1)
- Cloud Build Service Agent (1)
- Cloud Composer API Service Agent
- Container Registry Service Agent
- Editor (2) sb260219 +9190075
- Kubernetes Engine Service Agent
- Owner (1)
- Pub/Sub Publisher (1)
- Pub/Sub Subscriber (1)

8. We wrote another streaming process on dlq topic to reprocess dlq data in cassandra. That's how we can handle failure data which was somehow didn't ingest into table.

```
, 'card_last_four': '0046', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T46:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1047, 'order_id': 47, 'payment_method': 'Google Pay',
, 'card_last_four': '0047', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T47:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1048, 'order_id': 48, 'payment_method': 'PayPal', 'c
ard_last_four': '0048', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T48:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1049, 'order_id': 49, 'payment_method': 'Credit Card
', 'card_last_four': '0049', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T49:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1050, 'order_id': 50, 'payment_method': 'PayPal', 'c
ard_last_four': '0050', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T50:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1051, 'order_id': 51, 'payment_method': 'Apple Pay',
, 'card_last_four': '0051', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T51:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1052, 'order_id': 52, 'payment_method': 'Credit Card
', 'card_last_four': '0052', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T52:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1053, 'order_id': 53, 'payment_method': 'Google Pay',
, 'card_last_four': '0053', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T53:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1054, 'order_id': 54, 'payment_method': 'Apple Pay',
, 'card_last_four': '0054', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T54:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1055, 'order_id': 55, 'payment_method': 'Debit Card',
, 'card_last_four': '0055', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T55:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1056, 'order_id': 56, 'payment_method': 'PayPal', 'c
ard_last_four': '0056', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T56:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1057, 'order_id': 57, 'payment_method': 'Google Pay',
, 'card_last_four': '0057', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T57:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1058, 'order_id': 58, 'payment_method': 'PayPal', 'c
ard_last_four': '0058', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T58:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1059, 'order_id': 59, 'payment_method': 'Google Pay',
, 'card_last_four': '0059', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T59:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1060, 'order_id': 60, 'payment_method': 'Google Pay',
, 'card_last_four': '0060', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T60:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1061, 'order_id': 61, 'payment_method': 'Apple Pay',
, 'card_last_four': '0061', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T61:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1062, 'order_id': 62, 'payment_method': 'Credit Card
', 'card_last_four': '0062', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T62:01:30Z'}
Data thrown in dlq because order_id not found -> {'payment_id': 1063, 'order_id': 63, 'payment_method': 'Debit Card',
, 'card_last_four': '0063', 'payment_status': 'Completed', 'payment_datetime': '2024-01-21T63:01:30Z'}
```