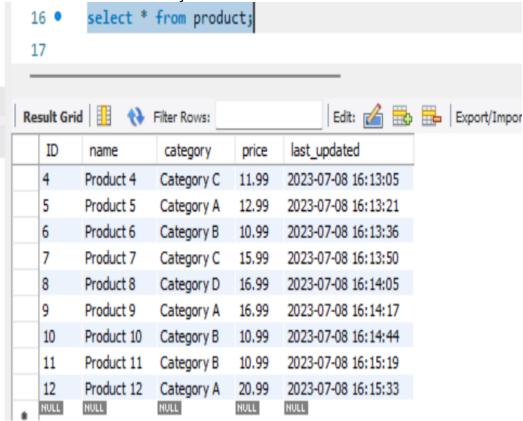
Orders Data Pub/Sub Using Kafka

- In this project, I builded a Kafka producer and a consumer group that worked with a MySQL database, Avro serialization, and multi-partition Kafka topics. The producer was fetching incremental data from a MySQL table and writing Avro serialized data into a Kafka topic. The consumers then deserializing this data and appending it to separate JSON files.
- 2. MySQL database storing product information such as product ID, name, category, price, and updated timestamp. Scenario is like I was updating the database frequently with new products and changing product information. I wanted to build a real-time system to stream these updates incrementally to a downstream system for real-time analytics.
- 3. Created a table named 'product' in MySQL database & initially loaded 12 records in the table 'Product' in MySQL database.



4. Wrote a Kafka producer in Python that used a MySQL connector to fetch data from the MySQL table. In my producer code, maintained a record of the last read timestamp. Each time I fetched data, used a SQL query to get records where the

last_updated timestamp is greater than the last read timestamp. Hence, records processed to producer successfully.

```
File Edit View Insert Cell Kernel Widgets Help

# Update the value in the config.json file
config_data['last_read_timestamp'] = max_date_str

with open('config.json', 'w') as file:
    json.dump(config_data, file)

# Close the cursor and database connection
cursor.close()
connection.close()
print("Data successfully produced to product_updates [3] at offset 153
Record b'1' successfully produced to product_updates [7] at offset 153
Record b'3' successfully produced to product_updates [1] at offset 153
Record b'3' successfully produced to product_updates [8] at offset 158
Record b'6' successfully produced to product_updates [9] at offset 158
Record b'6' successfully produced to product_updates [1] at offset 184
Record b'6' successfully produced to product_updates [6] at offset 184
Record b'6' successfully produced to product_updates [1] at offset 185
Record b'8' successfully produced to product_updates [1] at offset 185
Record b'9' successfully produced to product_updates [1] at offset 177
Record b'10' successfully produced to product_updates [1] at offset 177
Record b'10' successfully produced to product_updates [1] at offset 178
Record b'11' successfully produced to product_updates [1] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'12' successfully produced to product_updates [5] at offset 178
Record b'16' successfully produced to product_updates [5] at offset 178
Record b'16' successfully produced t
```

5. Serialized the data into Avro format and published the data to a Kafka topic named "product_updates". Configured this topic with 10 partitions. Used the product ID as the key when producing messages. This ensured that all updates for the same product end up in the same partition.

```
### Create Avro Serializer for the value

key_serializer = StringSerializer(vitf_B')

avro_serializer = AvroSerializer(schema_registry_client, schema_str)

### Define the SerializingProducer

producer = SerializingProducer()

| "bootstrap_servers': kafka_config('security.protocol'),
| "sasl.mechanisms': kafka
```

6. Wrote a Kafka consumer in Python and set it up as a consumer group of 5 consumers. Each consumer read data from the "product_updates" topic. Deserialized the Avro data back into a Python object & implemented data transformation logic.

```
consumer = DeserializingConsumer({
     'bootstrap.servers': kafka_config['bootstrap.servers'],
     'security.protocol': kafka_config['security.protocol'],
    'sasl.mechanisms': kafka_config['sasl.mechanisms'],
    'sasl.username': kafka_config['sasl.username'],
'sasl.password': kafka_config['sasl.password'],
     'key_deserializer': key_deserializer,
     'value.deserializer': avro_deserializer,
    'group.id': kafka_config['group.id'],
     'auto.offset.reset': kafka_config['auto.offset.reset'],
    #'enable.auto.commit': True
    if isinstance(obj, datetime):
return obj.isoformat()
# Path to the separate JSON file for each consumer
file path = 'consumer1.json'
# Python function to load append the json string data into json file.
def write_to_json_file(json_string, file_path):
   with open(file_path, 'a') as file:
        file.write(json_string + '\n')
consumer.subscribe(['product_updates'])
        msg = consumer.poll(1.0)
        if msg is None:
             continue
         if msg.error():
            print('Consumer error: {}'.format(msg.error()))
         msg.value()['category'] = msg.value()['category'].lower()
         if msg.value()['category'] == 'category a':
             msg.value()['price'] = msg.value()['price'] * 0.5
msg.value()['price'] = round(msg.value()['price'],2)
         print('Successfully consumed record with key {} and value {}'.format(msg.key(), msg.value()))
         json_string = json.dumps(msg.value(), default=datetime_encoder)
```

7. Consumers start consuming the 12 records(screenshots attached for consumer 1 to 5 consuming all 12 records in the table initially).

Each consumer converting the transformed Python object into a JSON string and appending the JSON string to a separate JSON file. Opened the file in append mode

avro data consumer 1 Last Checkpoint: 15 minutes ago (unsaved changes)

```
View
         Insert
                 Cell
                       Kernel
                                Widgets
                                          Help
2
   ↑ ↓
               # Check if the file exists
         if not os.path.isfile(file_path):
             # Create the file and write the initial data
             with open(file_path, 'w') as file:
                 file.write(json_string + '\n')
         else:
             # Append the data to the existing file
             write_to_json_file(json_string, file_path)
             print("json_string data is added to the JSON file.")
         file.close()
  except KeyboardInterrupt:
     pass
  finally:
     consumer.close()
  Successfully consumed record with key 2 and value {'ID': 2, 'name': 'Product 2', 'catego
  1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 10, 36, tzinfo=datetime.timez
  ison string data is added to the JSON file.
  Successfully consumed record with key 5 and value {'ID': 5, 'name': 'Product 5', 'catego
  updated': datetime.datetime(2023, 7, 8, 16, 13, 21, tzinfo=datetime.timezone.utc)
  json_string data is added to the JSON file.
  Successfully consumed record with key 7 and value {'ID': 7, 'name': 'Product 7', 'catego
  1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 13, 50, tzinfo=datetime.timez
  json_string data is added to the JSON file.
```

avro data consumer 2 Last Checkpoint: 16 minutes ago (autosaved)

```
View
         Insert
                 Cell
                        Kernel
                                 Widgets
                                            Help
                        ■ C → Code
42 ■
                ► Run
              msg.value()['price'] = round(msg.value()['price'],2)
          print('Successfully consumed record with key {} and value {}'.for
          json_string = json.dumps(msg.value(), default=datetime_encoder)
          def write_to_json_file(json_string, file_path):
              with open(file_path, 'a') as file:
                  file.write(json_string + '\n')
          # Check if the file exists
          if not os.path.isfile(file_path):
              # Create the file and write the initial data
              with open(file_path, 'w') as file:
                  file.write(json_string + '\n')
          else:
              # Append the data to the existing file
              write_to_json_file(json_string, file_path)
              print("json_string data is added to the JSON file.")
          file.close()
  except KeyboardInterrupt:
      pass
  finally:
      consumer.close()
  Successfully consumed record with key 12 and value {'ID': 12, 'name': 'Pr
  'last_updated': datetime.datetime(2023, 7, 8, 16, 15, 33, tzinfo=datetime
```

avro_data_consumer_3 Last Checkpoint: 16 minutes ago (autosaved)

```
View
                Cell
                                Widgets
       Insert
                       Kernel
                                          Help
 4
               ► Run ■ C
      4
                            >>
                                                 -----
                                 Code
            with open(file path, 'a') as file:
                file.write(json_string + '\n')
        # Check if the file exists
        if not os.path.isfile(file_path):
            # Create the file and write the initial data
            with open(file_path, 'w') as file:
                file.write(json_string + '\n')
        else:
            # Append the data to the existing file
            write_to_json_file(json_string, file_path)
            print("json_string data is added to the JSON file.")
        file.close()
except KeyboardInterrupt:
   pass
finally:
   consumer.close()
```

Successfully consumed record with key 4 and value {'ID': 4, 'name': 'Product 4' 1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 13, 5, tzinfo=dateti Successfully consumed record with key 9 and value {'ID': 9, 'name': 'Product 9' _updated': datetime.datetime(2023, 7, 8, 16, 14, 17, tzinfo=datetime.timezone.u json_string data is added to the JSON file.

Successfully consumed record with key 11 and value {'ID': 11, 'name': 'Product 9771118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 15, 19, tzinfo=da json_string data is added to the JSON file.

avro_data_consumer_4 Last Checkpoint: 16 minutes ago (autosaved)

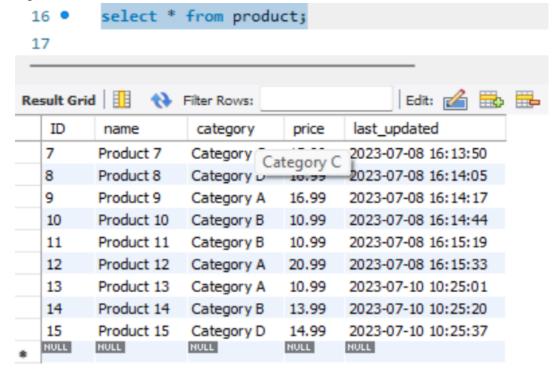
```
View
        Insert
                 Cell
                       Kernel
                                Widgets
                                           Help
               ▶ Run ■ C → Code
à 🖪
             msg.value()['price'] = round(msg.value()['price'],2)
         print('Successfully consumed record with key {} and value {}'.format(ms;
         json_string = json.dumps(msg.value(), default=datetime_encoder)
         def write_to_json_file(json_string, file_path):
             with open(file_path, 'a') as file:
                 file.write(json_string + '\n')
         # Check if the file exists
         if not os.path.isfile(file_path):
             # Create the file and write the initial data
             with open(file path, 'w') as file:
                 file.write(json_string + '\n')
             # Append the data to the existing file
             write_to_json_file(json_string, file_path)
             print("json_string data is added to the JSON file.")
         file.close()
 except KeyboardInterrupt:
     pass
 finally:
     consumer.close()
```

Successfully consumed record with key 1 and value {'ID': 1, 'name': 'Product 1' _updated': datetime.datetime(2023, 7, 8, 16, 10, 19, tzinfo=datetime.timezone.u

avro_data_consumer_5 Last Checkpoint: 16 minutes ago (autosaved)

```
Kernel
 View
          Insert
                    Cell
                                      Widgets
         4 | 4 |
2 🗗
                  ► Run ■ C → Code
                                                          [<del>]</del>
           # Check if the file exists
           if not os.path.isfile(file path):
               # Create the file and write the initial data
               with open(file path, 'w') as file:
                    file.write(json_string + '\n')
               # Append the data to the existing file
               write to json file(json string, file path)
               print("json_string data is added to the JSON file.")
           file.close()
 except KeyboardInterrupt:
      pass
 finally:
      consumer.close()
 Successfully consumed record with key 3 and value {'ID': 3, 'name': 'Product 3
 1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 11, 1, tzinfo=dateti
Successfully consumed record with key 6 and value {'ID': 6, 'name': 'Product 6'
 1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 13, 36, tzinfo=date1
 json_string data is added to the JSON file.
 Successfully consumed record with key 8 and value {'ID': 8, 'name': 'Product 8'
 1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 14, 5, tzinfo=datet
 json_string data is added to the JSON file.
 Successfully consumed record with key 10 and value {'ID': 10, 'name': 'Product 9771118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 14, 44, tzinfo=da
 json string data is added to the JSON file.
```

8. Again loaded 3 records in the table.



9. 3 new record processed to producer successfully.

er avro_data_producer Last Checkpoint: 11 minutes ago (autosaved)

```
View
          Insert
                  Cell
                         Kernel
                                 Widgets
                                            Help
42 ■
                 ► Run ■ C → Code
                                                   =====
  query = "SELECT MAX(last_updated) FROM product"
  cursor.execute(query)
  # Fetch the result
  result = cursor.fetchone()
  max_date = result[0] # Assuming the result is a single value
  # Convert datetime object to string representation
  max date str = max date.strftime("%Y-%m-%d %H:%M:%S")
  # Update the value in the config.json file
  config_data['last_read_timestamp'] = max_date_str
  with open('config.json', 'w') as file:
      json.dump(config_data, file)
  # Close the cursor and database connection
  cursor.close()
  connection.close()
  print("Data successfully published to Kafka")
  Record b'13' successfully produced to product_updates [7] at offset 154
  Record b'14' successfully produced to product_updates [2] at offset 146
  Record b'15' successfully produced to product_updates [8] at offset 159
  Data successfully published to Kafka
```

10. All 3(with id as 13,14,15) records consumed by consumers 1,4 and 3 respectively.

r avro data consumer 1 Last Checkpoint: 22 minutes ago (autosaved)

```
View
       Insert
               Cell
                      Kernel
                               Widgets
                                          Help
              ► Run ■ C → Code
                                                 -----
            with open(file path, 'w') as file:
                file.write(json string + '\n')
        else:
            # Append the data to the existing file
            write_to_json_file(json_string, file_path)
            print("json_string data is added to the JSON file.")
        file.close()
except KeyboardInterrupt:
   pass
finally:
   consumer.close()
Successfully consumed record with key 2 and value {'ID': 2, 'name': 'Product 2', 'category':
1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 10, 36, tzinfo=datetime.timezone.u
json_string data is added to the JSON file.
Successfully consumed record with key 5 and value {'ID': 5, 'name': 'Product 5', 'category':
updated': datetime.datetime(2023, 7, 8, 16, 13, 21, tzinfo=datetime.timezone.utc)}
json_string data is added to the JSON file.
Successfully consumed record with key 7 and value {'ID': 7, 'name': 'Product 7', 'category':
1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 13, 50, tzinfo=datetime.timezone.u
json_string data is added to the JSON file.
Successfully consumed record with key 13 and value {'ID': 13, 'name': 'Product 13', 'categor'
ast_updated': datetime.datetime(2023, 7, 10, 10, 25, 1, tzinfo=datetime.timezone.utc)}
json string data is added to the JSON file.
```

If avro data consumer 4 Last Checkpoint: 21 minutes ago (autosaved)

```
View
         Insert
                 Cell
                        Kernel
                                Widgets
                                           Help
4 6
        ↑ ↓
                ► Run ■ C → Code
                                              v 📟
              msg.value()['price'] = round(msg.value()['price'],2)
          print('Successfully consumed record with key {} and value {}'.format(msg.key(), msg.value()))
          json_string = json.dumps(msg.value(), default=datetime_encoder)
          def write_to_json_file(json_string, file_path):
              with open(file_path, 'a') as file:
                  file.write(json_string + '\n')
          # Check if the file exists
          if not os.path.isfile(file path):
              # Create the file and write the initial data
              with open(file_path, 'w') as file:
                  file.write(json string + '\n')
              # Append the data to the existing file
              write_to_json_file(json_string, file_path)
              print("json_string data is added to the JSON file.")
          file.close()
  except KeyboardInterrupt:
     pass
  finally:
      consumer.close()
  Successfully consumed record with key 1 and value {'ID': 1, 'name': 'Product 1', 'category': 'category
```

updated': datetime.datetime(2023, 7, 8, 16, 10, 19, tzinfo=datetime.timezone.utc)}
Successfully consumed record with key 14 and value {'ID': 14, 'name': 'Product 14', 'category': 'catego

'r avro_data_consumer_3 Last Checkpoint: 22 minutes ago (autosaved)

```
View
         Insert
                 Cell
                        Kernel
                                 Widgets
                                           Help
4
                ▶ Run | ■ | C | ▶ | | Code
  #####
          # Check if the file exists
          if not os.path.isfile(file_path):
              # Create the file and write the initial data
              with open(file_path, 'w') as file:
                  file.write(json_string + '\n')
          else:
              # Append the data to the existing file
              write_to_json_file(json_string, file_path)
              print("json_string data is added to the JSON file.")
          file.close()
  except KeyboardInterrupt:
     pass
  finally:
     consumer.close()
  Successfully consumed record with key 4 and value {'ID': 4, 'name': 'Product 4', 'category': 'category
  1118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 13, 5, tzinfo=datetime.timezone.utc)}
  Successfully consumed record with key 9 and value {'ID': 9, 'name': 'Product 9', 'category': 'category
  _updated': datetime.datetime(2023, 7, 8, 16, 14, 17, tzinfo=datetime.timezone.utc)}
  json_string data is added to the JSON file.
  Successfully consumed record with key 11 and value {'ID': 11, 'name': 'Product 11', 'category': 'cate@
  9771118164, 'last_updated': datetime.datetime(2023, 7, 8, 16, 15, 19, tzinfo=datetime.timezone.utc)}
  json_string data is added to the JSON file.
  Successfully consumed record with key 15 and value {'ID': 15, 'name': 'Product 15', 'category': 'categ
  9771118164, 'last_updated': datetime.datetime(2023, 7, 10, 10, 25, 37, tzinfo=datetime.timezone.utc)}
  json string data is added to the JSON file.
```