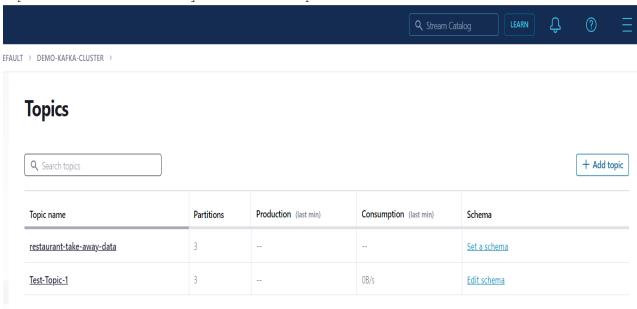
KAFKA ASSIGNMENT

-Download restaurant data

Link: https://github.com/shashank-mishra219/Confluent-Kafka-Setup/blob/main/restaurant_orders.csv

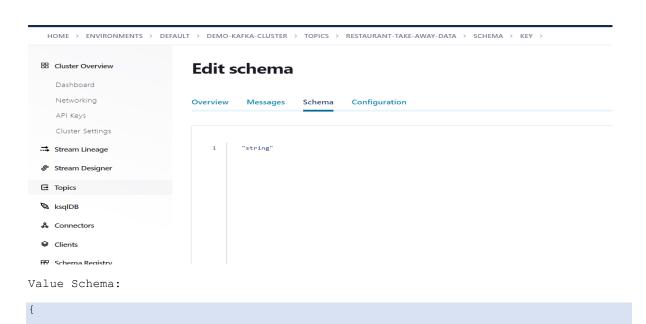
- 1. Setup Confluent Kafka Account
- 2. Create one kafka topic named as "restaurant-take-away-data" with 3 partitions

Topic 'restaurant-take-away-data ' with 3 partitions is created

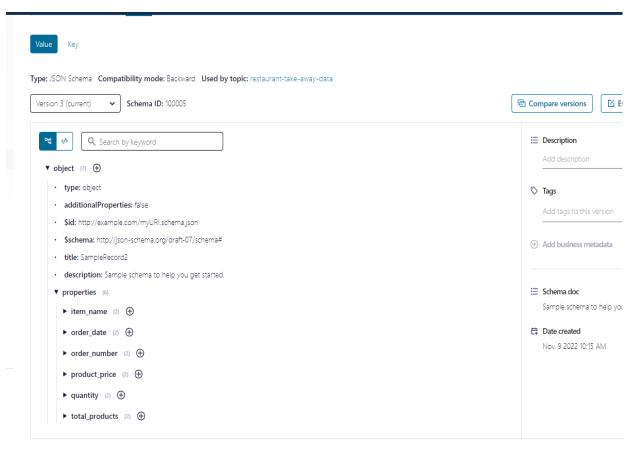


-Setup key (string) & value (json) schema in the confluent schema registry Key Schema:

"string"



```
"$id": "http://example.com/myURI.schema.json",
"$schema": "http://json-schema.org/draft-07/schema#",
"additionalProperties": false,
"description": "Sample schema to help you get started.",
"properties": { "item_name": {
   "description": "The type(v) type is used.",
   "type": "string"
 },
 "order_date": {
   "description": "The type(v) type is used.",
   "type": "string"
 "order_number": {
   "description": "The type(v) type is used.",
   "type": "number"
 "product price": {
   "description": "The type(v) type is used.",
   "type": "number"
 "quantity": {
   "description": "The type(v) type is used.",
   "type": "number"
 },
 "total products": {
   "description": "The type(v) type is used.",
   "type": "number"
},
"title": "SampleRecord2",
"type": "object"
```



-Write a kafka producer program (python or any other language) to read data records from restaurant data csv file, make sure schema is not hardcoded in the producer code, read the latest version of schema and schema_str from schema registry and use it for data serialization.

From producer code, publish data in Kafka Topic one by one and use dynamic key while publishing the records into the Kafka Topic

Producer Code:

```
import argparse
from uuid import uuid4
from six.moves import input
from confluent_kafka import Producer
from confluent_kafka.serialization import StringSerializer, SerializationContext,
MessageField
from confluent_kafka.schema_registry import SchemaRegistryClient
from confluent_kafka.schema_registry.json_schema import JSONSerializer
#from confluent_kafka.schema_registry import *
import pandas as pd
from typing import List

FILE_PATH = "/Users/namrt/OneDrive/Desktop/Kafka/restaurant_orders.csv"
columns=['order_number', 'order_date', 'item_name', 'quantity', 'product_price',
'total_products']
API KEY = '3DKHZA77TCORTOVM'
```

```
API SECRET KEY = '0c+PCR+sp1cEL0ns6tBTZ18bA9Z+AmpYEPSgQ4DfnsiwUKPRSR2vZs2Cx6f8IfoL'
BOOTSTRAP_SERVER = 'pkc-ymrq7.us-east-2.aws.confluent.cloud:9092'
SSL MACHENISM = 'PLAIN'
SCHEMA_REGISTRY_API_KEY = 'G6WDT4VOCBSLA4CO'
SCHEMA_REGISTRY_API_SECRET =
'28e5JcKaaXek7dbw6SCOoaX0T0UNnuYUg3wopsPqRHy9LrEfhqGB9pmquUbT1Pqe'
def sasl_conf():
def schema config():
class Restaurant:
    def __init__(self,record:dict):
        for k,v in record.items():
        self.record=record
    def dict_to_restaurant(data:dict,ctx):
        return Restaurant(record=data)
    def __str__(self):
def get restaurant instance(file path):
```

```
df=pd.read_csv(file_path)
    df=df.iloc[:,:]
def restaurant_to_dict(restaurant:Restaurant, ctx):
            operation.
    Returns:
        dict: Dict populated with user attributes to be serialized.
    return restaurant.record
def delivery_report(err, msg):
    Reports the success or failure of a message delivery.
    Args:
        print("Delivery failed for User record {}: {}".format(msg.key(), err))
        msg.key(), msg.topic(), msg.partition(), msg.offset()))
def main(topic):
    schema_registry_conf = schema_config()
    producer = Producer(sasl conf())
```

```
Restaurant producer ison.py >
                                                                                                                                                                                                                                                                                                                         PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL JUPYTER
{'order_number': 15253, 'order_date': '08/06/2019 19:38',
                                                                                                                                     'item_name': 'Mango Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
                                                                                                                                   item_name: 'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products':
'item_name': 'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 4}
'item_name': 'Red Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 9}
'item_name': 'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 'item_name': 'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 5}
'item_name': 'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 5}
'item_name': 'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 5}
'item_name': 'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 5}
'item_name': 'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 5}
 'order_number': 15251, 'order_date': '08/06/2019 19:22',
'order_number': 15249, 'order_date': '08/06/2019 19:13',
   'order_number': 15244, 'order_date':
'order_number': 15240, 'order_date':
                                                                                      '08/06/2019 17:43', '08/06/2019 17:05',
  'order_number': 15244,
   'order_number': 15224,
                                                      'order_date':
'order_date':
                                                                                      '07/06/2019 19:26'
                                                                                        '07/06/2019 18:48',
   'order_number': 15217, 'order_date':
'order_number': 15216, 'order_date':
                                                                                      '07/06/2019 18:28',
'07/06/2019 18:21',
                                                                                                                                     'item_name':
'item_name':
                                                                                                                                                                 'Onion Chutney',
'Mango Chutney',
                                                                                                                                                                                                         'quantity': 2, 'product_price': 0.5, 'total_products': 7} 'quantity': 1, 'product_price': 0.5, 'total_products': 14]
                                                                                                                                                                                                        'quantity'
   'order_number': 15216,
                                     15216, 'order_date':
15213, 'order_date':
                                                                                      '07/06/2019 18:21',
'07/06/2019 18:12',
                                                                                                                                     'item name':
                                                                                                                                                                 'Onion Chutney',
                                                                                                                                                                                                        'quantity': 3,
                                                                                                                                                                                                                                           'product_price': 0.5, 'total_products': 14]
'product_price': 0.5, 'total_products': 14]
                                                                                                                                                                 'Onion Chutney', 'quantity':
'Mint Sauce', 'quantity': 2,
'Mango Chutney', 'quantity':
                                                                                                                                                                                                       'quantity': 1, 'product_price': 0.5, 'total_products': 14
uantity': 2, 'product_price': 0.5, 'total_products': 8}
'quantity': 2, 'product_price': 0.5, 'total_products': 8}
'quantity': 2, 'product_price': 0.5, 'total_products': 8}
uantity': 2, 'product_price': 0.5, 'total_products': 5}
                                                                                                                                      item_name
                                     15213, 'order_date':
15211, 'order_date':
   order_number':
                                                                                      '07/06/2019 18:12',
'07/06/2019 18:04',
                                                                                                                                     'item_name':
'item_name':
                                                                                                                                                                'Onion Chutney', 'quantity':
'Mint Sauce', 'quantity': 2,
                                     15211, 'order_date':
15210, 'order_date':
    order number':
                                                                                       '07/06/2019 18:04'
                                                                                                                                     'item name':
                                                                                                                                                                 'Onion Chutney',
                                     15210, 'order_date':
15204, 'order_date':
                                                                                      '07/06/2019 17:49',
'06/06/2019 18:59',
                                                                                                                                     'item_name':
'item_name':
                                                                                                                                                                                                                                          'product_price': 0.5, 'total_products': 5]
'product_price': 0.5, 'total_products': 6]
   'order number':
                                                                                                                                                                                                        'quantity':
                                                                                                                                                                  'Mango Chutney',
                                                                                                                                                                                                         'quantity'
                                                                                       '05/06/2019 21:08',
                                                                                                                                                                                                                                   1, 'product_price': 0
'product_price': 0.5,
    order number':
                                     15199, 'order_date':
15199, 'order_date':
                                                                                                                                     'item name':
                                                                                                                                                                 'Mango Chutney', 'quantity':
'Mint Sauce', 'quantity': 1,
                                                                                                                                                                                                        'quantity':
                                                                                                                                                                                                                                           'product price': 0.5,
                                                                                                                                                                                                                                                                                            'total products': 8)
                                                                                                                                                                'Onion Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 8}
'Mango Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 8}
'Mango Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 7}
'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 7}
'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
'Lime Pickle', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
'Lime Pickle', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
                                     15199, 'order_date':
15199, 'order_date':
                                                                                      '05/06/2019 21:08',
'05/06/2019 21:08',
                                                                                                                                     'item_name':
'item_name':
   order_number':
                                    15197, 'order_date':
15197, 'order_date':
15197, 'order_date':
15191, 'order_date':
15191, 'order_date':
                                                                                       '05/06/2019 20:44',
    order number':
                                                                                                                                     'item name':
                                                                                        '05/06/2019 20:44',
                                                                                       '05/06/2019 11:29',
                                                                                                                                     'item_name':
'item_name':
   'order number':
                                                                                        '05/06/2019 11:29',
   'order_number': 15186, 'order_date':
'order_number': 15186, 'order_date':
                                                                                                                                                                                                        'quantity': 1, 'product_price': 0.5, 'total_products': 6}
                                                                                       '04/06/2019 20:08',
                                                                                                                                                                 'Onion Chutney',
                                                                                                                                     'item name':
                                                                                                                                                                 'Red Sauce', 'quantity': 1,
'Mint Sauce', 'quantity': 1,
  order_number': 15174, 'order_date':
'order_number': 15174, 'order_date':
'order_number': 15174, 'order_date':
'order_date':
'order_date':
'order_date':
                                                                                       '04/06/2019 20:08'
                                                                                                                                     'item_name'
                                                                                                                                                                                                                                     'product_price': 0.5,
                                                                                                                                                                                                                                                                                      'total products': 6}
                                                                                                                                   item_name: Munt Sauce , quantity: 1, product price: 0.5, 'total_products': 6}
'item_name': 'Mango Chutney', 'quantity': 2, 'product_price': 0.5, 'total_products': 1:
'item_name': 'Mint Sauce', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
'item_name': 'Mango Chutney', 'quantity': 1, 'product_price': 0.5, 'total_products': 4'
'item_name': 'Lime Pickle', 'quantity': 1, 'product_price': 0.5, 'total_products': 11}
                                                                                      '04/06/2019 17:58',
'04/06/2019 13:38',
'04/06/2019 12:10',
   'order number': 15171, 'order date':
```

-Write kafka consumer code and create two copies of same consumer code and save it with different names (kafka_consumer_1.py & kafka_consumer_2.py), again make sure lates schema version and schema_str is not hardcoded in the consumer code, read it automatically from the schema registry to deserialize the data.

Now test two scenarios with your consumer code:

a.) Use "group.id" property in consumer config for both consumers and mention different group ids in kafka consumer 1.py & kafka consumer 2.py,

apply "earliest" offset property in both consumers and run these two consumers from two different terminals. Calculate how many records each consumer consumed and printed on the terminal

b.) Use "group.id" property in consumer config for both consumers and mention same group ids in kafka consumer 1.py & kafka consumer 2.py,

apply "earliest" offset property in both consumers and run these two consumers from two different terminals. Calculate how many records each consumer consumed and printed on the terminal

Consumer Code:

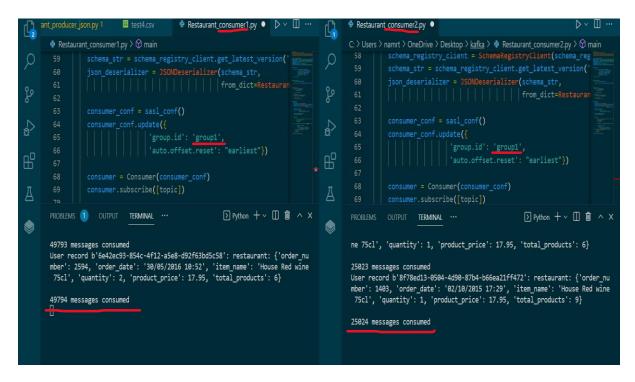
```
API KEY = '3DKHZA77TCORTQVM'
API SECRET KEY = '0c+PCR+sp1cEL0ns6tBTZ18bA9Z+AmpYEPSgQ4DfnsiwUKPRSR2vZs2Cx6f8IfoL'
SSL MACHENISM = 'PLAIN'
SCHEMA REGISTRY API SECRET =
'28e5JcKaaXek7dbw6SCOoaX0T0UNnuYUg3wopsPqRHy9LrEfhqGB9pmquUbT1Pqe'
def sasl_conf():
                'sasl.password': API SECRET KEY
    return sasl conf
```

```
def schema_config():
class Restaurant:
    def __init__(self,record:dict):
        for k,v in record.items():
        self.record=record
    def dict to restaurant(data:dict,ctx):
        return Restaurant(record=data)
    def str (self):
def main(topic):
    schema_registry_conf = schema_config()
                                         from dict=Restaurant.dict to restaurant)
    consumer = Consumer(consumer conf)
   while True:
```

For testing, we need to change group id only, rest code will remain same.

```
consumer_conf.update({
    'group.id': 'group1',
    'auto.offset.reset': "earliest"})
```

a) group id for consumer1 and consumer2 is same. Task gets distributed among the number of consumers we have. Load sharing results in fast processing.



Total messages: 74818 Consumer1: 49794 Consumer2: 25024

```
b) Group if for consumer1 and consumer2 is different.
Both will perform their own task.
Total messages: 74818
Consumer1: 74818
Consumer2: 74818
```

-Once above questions are done, write another kafka consumer to read data from kafka topic and from the consumer code create one csv file "output.csv" and append consumed records output.csv file

```
from confluent kafka.serialization import SerializationContext, MessageField
API KEY = '3DKHZA77TCORTQVM'
ENDPOINT_SCHEMA_URL = 'https://psrc-zj6ny.us-east-2.aws.confluent.cloud'
API_SECRET_KEY = '0c+PCR+sp1cEL0ns6tBTZ18bA9Z+AmpYEPSgQ4DfnsiwUKPRSR2vZs2Cx6f8IfoL'
BOOTSTRAP_SERVER = 'pkc-ymrq7.us-east-2.aws.confluent.cloud:9092'
SSL MACHENISM = 'PLAIN'
'28e5JcKaaXek7dbw6SCOoaX0T0UNnuYUg3wopsPqRHy9LrEfhqGB9pmquUbT1Pqe'
def sasl_conf():
                'bootstrap.servers':BOOTSTRAP SERVER,
def schema_config():
```

```
class Restaurant:
    def __init__(self,record:dict):
        for k,v in record.items():
        self.record=record
    def dict_to_restaurant(data:dict,ctx):
        return Restaurant(record=data)
    def __str__(self):
def main(topic):
    schema_registry_conf = schema_config()
    schema str = schema registry client.get latest version('restaurant-take-away-data-
                                         from dict=Restaurant.dict to restaurant)
                     'group.id': 'group1',
    consumer = Consumer(consumer conf)
    consumer.subscribe([topic])
    columns=['order number', 'order date', 'item name', 'quantity', 'product price',
 total products']
   while True:
        try:
                    writer.writerows(final)
```

```
■ output.csv •
output.csv
       order_number,order_date,item_name,quantity,product_price,total_products
       16118,03/08/2019 20:25,Plain Papadum,2,0.8,6
       16118,03/08/2019 20:25,Mango Chutney,1,0.5,6
       16117,03/08/2019 20:17, Tandoori Chicken (1/4),1,4.95,7
       16117,03/08/2019 20:17,Saag Paneer,1,5.95,7
       16116,03/08/2019 20:09,Aloo Chaat,1,4.95,5
       16116,03/08/2019 20:09, Lamb Biryani, 1,9.95,5
       16115,03/08/2019 20:01,Chicken Pakora,1,5.95,7
       16114,03/08/2019 19:44, Special Fried Rice, 2, 3.95, 2
       16113,03/08/2019 19:42,Pilau Rice,1,2.95,5
       16112,03/08/2019 19:41,Plain Papadum,2,0.8,4
       16111,03/08/2019 19:29,Saag Aloo,1,5.95,4
       16110,03/08/2019 19:28,Aloo Gobi,1,5.95,8
       16110,03/08/2019 19:28, Chicken Biryani, 1,9.95,8
       16109,03/08/2019 19:26,Plain Papadum,4,0.8,7
```