



Avro producer with python

Create Avro Producers for Confluent Kafka with Python



Mrugank Ray · [Follow](#)

5 min read · Feb 21



3



Overview:

This article will teach you how to create an Avro producer using the Confluent kafka library in python.

Setup:

Clone [Big Data Cluster](#) repo. This will set up an environment for producing messages in Avro format for Kafka.

Run the following command to initialize a docker cluster

```
sudo docker-compose -f kafka-zookeper.yaml up
```

NOTE: You need docker on your computer to run the following commands. If you already have a kafka cluster running, you can skip this step.

Dependency:

Running avro producer will require some dependencies to be installed on your system. In our project directory, let's create a **requirements.txt** file and paste the following text in there.

```
1 confluent-kafka==1.9.2
2 fastavro==1.4.7
3 requests==2.27.1
```

[avro_producer_requirements.txt](#) hosted with ❤ by GitHub

[view raw](#)

requirements.txt

Let's install the dependencies using the following command in the project directory

```
pip install -r requirements.txt
```

Source:

We're going to use Wikimedia's event streams as the source. To learn more visit [https://stream.wikimedia.org/?](https://stream.wikimedia.org/)

[doc#/streams/get_v2_stream_recentchange](#). We use this source because it's publicly available and easy to use.

For reading event streams, we will use **requests** library available for python. Requesting <https://stream.wikimedia.org/v2/stream/recentchange> gives you events that look like this.

```
:ok
event: message
id: [{"topic":"eqiad.mediawiki.recentchange","partition":0,"timestamp":1676754367001}, {"topic":"codfw.mediawiki.recentchange","partition":0,"offset":-1}]
data: {"$schema":"/mediawiki/recentchange/1.0.0","meta":{"uri":"https://tr.wikipedia.org/wiki/Kullan%C4%B1c%C4%81_mesaj:Ahmetlii","request_id":"58827397-ee60-40fa-84e6-7654c1954653","id":"1b3f419b-a7d7-4b31-bfda-4310f9605243","dt":"2023-02-18T21:06:07Z","domain":"tr.wikipedia.org","stream":"mediawiki.recentchange","topic":"eqiad.mediawiki.recentchange","partition":0,"offset":4544755407}, "id":57511524,"type":"edit","namespace":3,"title":"Kullanici_mesaj:Ahmetlii","comment": "[[Kullanici_mesaj:Ahmetlii|mesaj]] ([[Kullanici_mesaj:Ahmetlii|mesaj]]) tarafından yapılan 29274421 sayılı değişiklik geri alınıyor. ","timestamp":1676754367,"user":"Ahmetlii","bot":false,"minor":true,"length": {"old":62798,"new":65535}, "revision": {"old":29274421,"new":29274427}, "server_url": "https://tr.wikipedia.org", "server_name": "tr.wikipedia.org", "server_script_path": "/w", "wiki": "trwiki", "parsedcomment": "<a href=\"/wiki/Kullan%C4%B1c%C4%81:Ahmetlii\" title=\"Kullanici:Ahmetlii\">Ahmetlii</a> <a href=\"/wiki/Kullan%C4%B1c%C4%81_mesaj:Ahmetlii\" title=\"Kullanici_mesaj:Ahmetlii\">mesaj</a>} tarafından yapılan 29274421 sayılı değişiklik geri alınıyor."}

event: message
id: [{"topic":"eqiad.mediawiki.recentchange","partition":0,"timestamp":1676754369001}, {"topic":"codfw.mediawiki.recentchange","partition":0,"offset":-1}]
data: {"$schema":"/mediawiki/recentchange/1.0.0","meta":{"uri":"https://es.wikipedia.org/wiki/Cultura_vir%C3%BAA","request_id":"aeblac22-3966-4f99-9064-810da1b235e5","id":"6e81dc8f-1b13-448b-94e5-7c827da69e84","dt":"2023-02-18T21:06:09Z","domain":"es.wikipedia.org","stream":"mediawiki.recentchange","topic":"eqiad.mediawiki.recentchange","partition":0,"offset":4544755408}, "id":282521856,"type":"edit","namespace":0,"title":"Cultura virú","comment": "/ Capital", "timestamp":1676754369,"user": "190.108.83.130","bot":false,"minor":false,"length": {"old":4894,"new":6308}, "revision": {"old":149364159,"new":149364200}, "server_url": "https://es.wikipedia.org", "server_name": "es.wikipedia.org", "server_script_path": "/w", "wiki": "eswiki", "parsedcomment": "<span dir=\"auto\"><span class=\"autocomment\"><a href=\"/wiki/Cultura_vir%C3%BA#Capital\" title=\"Cultura virú\">Capital</a></span></span>"}
```

wikimedia event stream

Get data from stream

Let's take a look at how to get data from streams.

```

1 import requests
2 import json
3
4 init_string = 'data: '
5 source_url = 'https://stream.wikimedia.org/v2/stream/recentchange'
6
7 def avro_producer(source_url):
8     s = requests.Session()
9
10    with s.get(source_url, headers=None, stream=True) as resp:
11        for line in resp.iter_lines():
12            if line:
13                decoded_line = line.decode()
14                if decoded_line.find(init_string) >= 0:
15                    # remove data: to create a valid json
16                    decoded_line = decoded_line.replace(init_string, "")
17                    # convert to json
18                    decoded_json = json.loads(decoded_line)
19                    print(decoded_json)
20
21 avro_producer(source_url)

```

get_wikimedia_data.py hosted with ❤ by GitHub

[view raw](#)

avro_producer.py

Let's understand the code:

- Line 8 to 10: We create a session object and establish a TCP connection.
- Line 11 to 13: We loop through the events and then we need to decode the byte string to UTF-8.
- Line 14: Each event data is stored as a string. We search for events that have a “data:” string since we’re interested in these.
- Line 16 to 19: We remove the string “data:” and convert it into a valid JSON object.

Create kafka topic:

Let's create a Kafka topic. To do this, we need to log into our “kafka-broker” container. Once logged in run the following command:

```
kafka-topics --bootstrap-server kafka-broker:29092 --create --topic wikimedia --
```

This creates a topic named “wikimedia” with one partition and one replication factor.

NOTE: The kafka cluster only has one broker, so set the replication factor to 1.

Avro schema:

In this section, we will create the avro schema and register it in schema registry.

Define the Avro schema

Schemas can be created for both key and value, but we'll only create a schema for value here. The value is the data that's streaming from Wikimedia.

This [site](#) helps you create AVRO schema from JSON.

Here is [schema](#) that I've created.

Now let's make a `constants.py` file where we define constants like `SCHEMA_STR` and assign stringified schema.

```
1 SCHEMA_STR = "{\"type\":\"record\",\"name\":\"value_wikimedia\",\"namespace\":\"wikimedia\",\"fields\":[{
```

wikimedia_stringified_schema.py hosted with ❤ by GitHub

[view raw](#)

constants.py

Check out this [site](#) for more info on AVRO schema.

Add schema to the registry

There are multiple ways to add schema to the schema registry, like [REST APIs](#) and [Confluent Kafka library](#).

In this article, we're going to use the [Schema Registry Client](#).

Let's take a look at how to register a schema

```
1 # 3rd party library imported
2 from confluent_kafka.schema_registry import SchemaRegistryClient
3 from confluent_kafka.schema_registry import Schema
4
5 # import from constants
6 from constants import SCHEMA_STR
7
8 schema_registry_url = 'http://localhost:8083'
9 kafka_topic = 'test-blog'
10 schema_registry_subject = f"{kafka_topic}-value"
11
12 def get_schema_from_schema_registry(schema_registry_url, schema_registry_subject):
13     sr = SchemaRegistryClient({'url': schema_registry_url})
14     latest_version = sr.get_latest_version(schema_registry_subject)
15
16     return sr, latest_version
17
18 def register_schema(schema_registry_url, schema_registry_subject, schema_str):
19     sr = SchemaRegistryClient({'url': schema_registry_url})
20     schema = Schema(schema_str, schema_type="AVRO")
21     schema_id = sr.register_schema(subject_name=schema_registry_subject, schema=schema)
22
23     return schema_id
24
25 def update_schema(schema_registry_url, schema_registry_subject, schema_str):
26     sr = SchemaRegistryClient({'url': schema_registry_url})
27     versions_deleted_list = sr.delete_subject(schema_registry_subject)
28     print(f"versions of schema deleted list: {versions_deleted_list}")
29
30     schema_id = register_schema(schema_registry_url, schema_registry_subject, schema_str)
31     return schema_id
32
33 schema_id = register_schema(schema_registry_url, schema_registry_subject, SCHEMA_STR)
34 print(schema_id)
35
36 sr, latest_version = get_schema_from_schema_registry(schema_registry_url, schema_registry_subjec
37 print(latest_version.schema.schema_str)
```

register_schema.py hosted with ❤ by GitHub

[view raw](#)

avro_producer.py

Let's understand the code:

- Line 9: We set the kafka topic name.
- Line 10: The topic name is suffixed with “-value” for a value schema. If your schema is a key schema, it’s suffixed with “-key”.
- Line 12 to 16: We create a schema registry client instance and get the latest version of the schema using the subject. You can learn more about it [here](#)
- Line 18 to 23: We pass the schema_type of AVRO along with the stringified schema data to create the Schema object. Then we register the schema in the registry, which returns the schema ID.
- Line 25 to 31: There's no update function in confluent kafka library, so we create one function that deletes the schema and registers the updated schema.
- Line 33: We register the schema.
- Line 36: We get the latest version schema's details.

Avro producer:

Let's get started. We will break down the steps as we go along.

Send avro data to kafka

We'll update the avro_producer function to send data.

Open in app ↗

Sign up

Sign In



Search

Write



```
1 import requests
2 import json
3
4 # 3rd party library imported
5 from confluent_kafka import SerializingProducer
6 from confluent_kafka.schema_registry import SchemaRegistryClient
7 from confluent_kafka.schema_registry.avro import AvroSerializer
8 from confluent_kafka.schema_registry import Schema
9
10 # import from constants
11 from constants import SCHEMA_STR
12
13 init_string = 'data: '
14 source_url = 'https://stream.wikimedia.org/v2/stream/recentchange'
15 kafka_url = 'localhost:9092'
16 schema_registry_url = 'http://localhost:8083'
17 kafka_topic = 'wikimedia'
18 schema_registry_subject = f'{kafka_topic}-value'
19
20 def delivery_report(errmsg, msg):
21     if errmsg is not None:
22         print("Delivery failed for Message: {} : {}".format(msg.key(), errmsg))
23         return
24     print('Message: {} successfully produced to Topic: {} Partition: [{}]
25 at offset {}'.format(
26     msg.value(),
27     msg.topic(),
28     msg.partition(),
29     msg.offset()))
30
31 value_avro_serializer = AvroSerializer(schema_registry_client = sr,
32                                         schema_str = latest_version.schema.schema_str,
33                                         conf = {
34                                             'auto.register.schemas': False
35                                         })
36
37
38 # Kafka Producer
39 producer = SerializingProducer({
40     'bootstrap.servers': kafka_url,
41     'security.protocol': 'plaintext',
42     'value.serializer': value_avro_serializer,
43     'delivery.timeout.ms': 120000, # set it to 2 mins
44     'enable.idempotence': 'true'
45 })
```

```
    ..
46
47     s = requests.Session()
48
49     with s.get(source_url, headers=None, stream=True) as resp:
50         for line in resp.iter_lines():
51             if line:
52                 decoded_line = line.decode()
53                 if decoded_line.find(init_string) >= 0:
54                     # remove data: to create a valid json
55                     decoded_line = decoded_line.replace(init_string, "")
56                     # convert to json
57                     decoded_json = json.loads(decoded_line)
58
59             try:
60                 # print(decoded_line + '\n')
61                 producer.produce(topic=kafka_topic, value=decoded_json, on_delivery=delivery_report)
62
63                 # Trigger any available delivery report callbacks from previous produce()
64                 events_processed = producer.poll(1)
65                 print(f"events_processed: {events_processed}")
66
67                 messages_in_queue = producer.flush(1)
68                 print(f"messages_in_queue: {messages_in_queue}")
69             except Exception as e:
70                 print(e)
71
72     def get_schema_from_schema_registry(schema_registry_url, schema_registry_subject):
73         sr = SchemaRegistryClient({'url': schema_registry_url})
74         latest_version = sr.get_latest_version(schema_registry_subject)
75
76         return sr, latest_version
77
78     def register_schema(schema_registry_url, schema_registry_subject, schema_str):
79         sr = SchemaRegistryClient({'url': schema_registry_url})
80         schema = Schema(schema_str, schema_type="AVRO")
81         schema_id = sr.register_schema(subject_name=schema_registry_subject, schema=schema)
82
83         return schema_id
84
85     def update_schema(schema_registry_url, schema_registry_subject, schema_str):
86         sr = SchemaRegistryClient({'url': schema_registry_url})
87         versions_deleted_list = sr.delete_subject(schema_registry_subject)
88         print(f"versions of schema deleted list: {versions_deleted_list}")
89
```

```
90     schema_id = register_schema(schema_registry_url, schema_registry_subject, schema_str)
91     return schema_id
92
93 avro_producer(source_url, kafka_url, schema_registry_url, schema_registry_subject)
```

avro_producer.py hosted with ❤️ by GitHub

[view raw](#)

Let's understand the code:

- Line 20 to 24: It's a callback function that's called when the message is delivered successfully or failed. In case the message doesn't get delivered, we log the error to console. In case the message gets delivered, we log the key, topic, partition and offset.
- Line 28 to 36: Once we get the latest schema, we create an avro serializer and make "auto.register.schemas" False in the config so that the library won't create a new schema.
- Line 39 to 45: We set up a producer and some configurations. You can read about SerializingProducer [here](#). "delivery.timeout.ms" sets the maximum time a producer can take to deliver a message (including retries). We have set it to 120000 milliseconds or 2 mins. When we set "enable.idempotence" to true, messages are produced exactly once and in the original order. You can read more about config [here](#).
- Line 61: This function is async in nature and calls the delivery_report function once messages have been delivered or failed.
- Line 64: This polls the producer for events and gets back the number of events processed. We can optionally pass max waiting time.
- Line 67: This function waits for all messages in the Producer queue to be delivered. It returns the number of messages still in the queue. We can optionally pass a waiting time. If messages are delivered within the waiting time then it would return zero.

Output

```
Message: None successfully produced to Topic: wikimedia Partition: [0] at offset 579
events_processed: 1
messages_in_queue: 0
Message: None successfully produced to Topic: wikimedia Partition: [0] at offset 580
events_processed: 1
messages_in_queue: 0
```

Avro producer

Avro Consumer:

We will discuss two ways to consume the messages.

CLI

We will use CLI to consume the messages from Kafka topic. To do this, we need to log into our “schema-registry” container. Once logged in run the following command:

```
kafka-console-consumer \
  --bootstrap-server kafka-broker:29092 \
  --topic wikipedia \
  --from-beginning \
  --property schema.registry.url=http://schema-registry:8083
```

This consumer will read data from the beginning of the topic and write all the data to the console.

```
File Edit View Search Terminal Help
a.org"},"server_script_path":{"string":"/w"},"server_url":{"string":"https://www.wikidata.org"},"timestamp":{"int":1676723263},"title":{"string":"Q116817758"},"type":{"string":"edit"},"user":{"string":"Piastu"},"wiki":{"string":"wikidatawiki"}}
{"bot":{"boolean":false},"comment":{"string":"added [[Category:Looms in art]]"}, "id":{"int":2120600418}, "length":{"wikimedia.Length":{"new":{"int":1096}, "old":{"int":1070}}}, "meta":{"wikimedia.Meta":{"domain":{"string":"commons.wikimedia.org"}}, "dt":{"string":"2023-02-18T12:27:43Z"}, "id":{"string":"b70d1a6f-ea73-4d3e-9daf-03dafdd9bb6b"}, "offset":{"long":4543820769}, "partition":{"int":0}, "request_id":{"string":"fa3d1463-5b29-4128-8602-6c2228e67d98"}, "stream":{"string":"mediawiki.recentchange"}, "topic":{"string":"eqiad.mediawiki.recentchange"}, "uri":{"string":"https://commons.wikimedia.org/wiki/File:Heva_Coomans_-_Penelope_awaiting_Odysseus.jpg"}}, "minor":{"boolean":true}, "namespace":{"int":6}, "parsedcomment":{"string":"added <a href=\"/wiki/Category:Looms_in_art\" title=\"Category:Looms in art\">Category:Looms in art</a>"}, "patrolled":{"boolean":true}, "revision":{"wikimedia.Revision":{"new":{"int":733572250}, "old":{"int":733572150}}}, "schema":null, "server_name":{"string":"commons.wikimedia.org"}, "server_script_path":{"string":"/w"}, "server_url":{"string":"https://commons.wikimedia.org"}, "timestamp":{"int":1676723263}, "title":{"string":"File:Heva Coomans - Penelope awaiting Odysseus.jpg"}, "type":{"string":"edit"}, "user":{"string":"Kilom691"}, "wiki":{"string":"commonswiki"}}
{"bot":{"boolean":false}, "comment":{"string":"[[Spezial:LintErrors/missing-end-tag|End-Tag fehlt]] kursiv nicht beidseitig geschlossen"}, "id":{"int":333123150}, "length":{"wikimedia.Length":{"new":{"int":26039}, "old":{"int":26031}}}, "meta":{
```

Console

Confluent control center

It's a browser-based UI. Open your browser and go to "<http://localhost:9021>".

To start a consumer go to “Cluster 1 -> Topics” in the left panel and then select topic name which in our case is wikimedia. Go to the Messages tab and enter 0 for Offset. This will read data from the beginning of the topic.

	value.bot.boolean	value.comment.str...	value.id.int	value.length.wikim...	value.length.wikim...	value.meta.wikim...	valu...
Bytes in/sec 0	false	[[File:Brighton 202...	2120600622	null	null	null	8 min ago null
Bytes out/sec 0	false	[[File:Uutelan kana...	2120600621	null	null	null	null
Message fields	true	[[Catégorie:Entrepr...	486551504	null	null	null	null
topic	false	[[File:Uutelan kana...	2120600620	null	null	null	null
partition	false	[[Benutzer:Aka/Fe...	333123160	null	null	null	null
offset							

Confluent control center

Kafka

Kafka Producer

Python

Data Engineering

Data Streaming



Written by Mrugank Ray

7 Followers

Follow

More from Mrugank Ray



 Mrugank Ray

Real-time AVRO Data Analysis with Spark Streaming and Confluent...

Overview:

12 min read · Mar 25

 25

 1

 +



 Mrugank Ray

Unlocking the Vault of Tech Interview Questions

As a Senior Engineer who has gone through several job interviews in the tech industry an...

10 min read · Jun 18

 16



 +

 Mrugank Ray

PostgreSQL vs MySQL | Choose What's Best For Your Project

Introduction :There are many options available when it comes to database system...

8 min read · Sep 25, 2022

 1 1

See all from Mrugank Ray

Recommended from Medium





Ismael Sánchez Chaves in C# Programming

Schema Registry in Kafka: Avro, JSON and Protobuf

The importance of having a structured data schema for messaging-based systems

7 min read · May 15



24



1



Lists



Coding & Development

11 stories · 258 saves



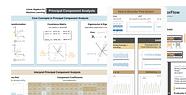
Predictive Modeling w/ Python

20 stories · 561 saves



New_Reading_List

174 stories · 176 saves



Practical Guides to Machine Learning

10 stories · 637 saves



Jay Kim

[Docker] Docker Compose example for Kafka, Zookeeper, and Schema...

1 Kafka Broker Setup

2 min read · Jun 1



Prateek

Avro Example 5 -use Date and Timestamp

How to use Date and Timestamp in Avro, lets check that out —

3 min read · Jun 22

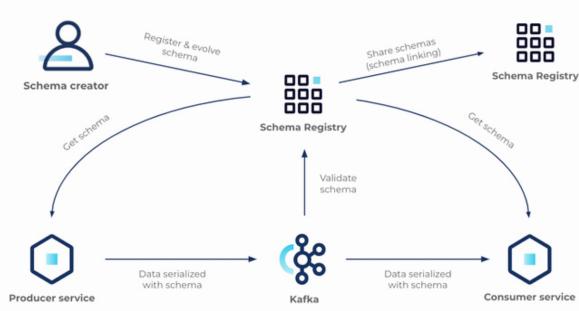


Ahmed Uz Zaman in Dev Genius

PySpark and Kafka Streaming: Reading and Writing Data in...

A Comprehensive Guide to Reading and Writing Data in Different Formats from Kafka...

· 4 min read · May 7



JSON
SERIALIZATION



Kamini Kamal in Level Up Coding

Schema Registry in Kafka ecosystem

The Schema Registry is an integral component in the Kafka ecosystem, providing...

4 min read · Jul 2



• 11 min read · Jun 17



[See more recommendations](#)