

DataEng: Data Transport Activity

[this lab activity references tutorials at confluence.com]

Make a copy of this document and use it to record your results. Store a PDF copy of the document in your git repository along with your code before submitting for this week. For your code, you create several producer/consumer programs or you might make various features within one program. There is no one single correct way to do it. Regardless, store your code in your repository.

The goal for this week is to gain experience and knowledge of using a streaming data transport system (Kafka). Complete as many of the following exercises as you can. Proceed at a pace that allows you to learn and understand the use of Kafka with python.

Submit: [In-class Activity Submission Form](#)

A. Initialization

1. Get your cloud.google.com account up and running
 - a. Redeem your GCP coupon
 - b. Login to your GCP console
 - c. Create a new, separate VM instance
2. Follow the Kafka tutorial from project assignment #1
 - a. Create a separate topic for this in-class activity
 - b. Make it “small” as you will not want to use many resources for this activity. By “small” I mean that you should choose medium or minimal options when asked for any configuration decisions about the topic, cluster, partitions, storage, anything. GCP/Confluent will ask you to choose the configs, and because you are using a free account you should opt for limited resources where possible.
 - c. Get a basic producer and consumer working with a Kafka topic as described in the tutorials.
3. Create a sample breadcrumb data file (named bcsample.json) consisting of a sample of 1000 breadcrumb records. These can be any records because we will not be concerned with the actual contents of the breadcrumb records during this assignment.
4. Update your producer to parse your sample.json file and send its contents, one record at a time, to the kafka topic.
5. Use your consumer.py program (from the tutorial) to consume your records.

B. Kafka Monitoring

1. Find the Kafka monitoring console for your topic. Briefly describe its contents. Do the measured values seem reasonable to you?
2. Use this monitoring feature as you do each of the following exercises.

Answer: Yes the values seem reasonable, although they don't match each other's actual values. It displays bytes that were produced in seconds by the producer and the bytes that were consumed by the consumer in bytes per second. For example, The producer might create 3000 bytes but the consumer updates slightly lesser bytes.

C. Kafka Storage

1. Run the linux command “wc bcsample.json”. Record the output here so that we can verify that your sample data file is of reasonable size.

```
(confluent-exercise) ctawde@in-class-assign-2:~/in-class-assign$ wc bcsample.json
35001  85069 1071472 bcsample.json
```

2. What happens if you run your consumer multiple times while only running the producer once?

Answer: The consumer keeps waiting for new messages from the producer

3. Before the consumer runs, where might the data go, where might it be stored?

Answer: Data in Kafka is stored in topics. Topics are partitioned and Each partition is further divided into segments. Each segment has a log file to store the actual message and an index file to store the position of the messages in the log file.

4. Is there a way to determine how much data Kafka/Confluent is storing for your topic? Do the Confluent monitoring tools help with this?

Answer: The Confluent monitoring tool only tells how much data (in bytes) were produced by producers per second and how much data (in bytes) was consumed by consumers per second

5. Create a “topic_clean.py” consumer that reads and discards all records for a given topic. This type of program can be very useful during debugging.

D. Multiple Producers

1. Clear all data from the topic
2. Run two versions of your producer concurrently, have each of them send all 1000 of your sample records. When finished, run your consumer once. Describe the results.

Answer: The producers produces 2000 records and the consumer consumes all 2000 records

E. Multiple Concurrent Producers and Consumers

1. Clear all data from the topic
2. Update your Producer code to include a 250 msec sleep after each send of a message to the topic.
3. Run two or three concurrent producers and two concurrent consumers all at the same time.
4. Describe the results.

Answer: Both the consumers keep consuming the produced messages simultaneously, the consumers also seem a little slow as if waiting 250 milliseconds to fetch a record.

F. Varying Keys

1. Clear all data from the topic

So far you have kept the “key” value constant for each record sent on a topic. But keys can be very useful to choose specific records from a stream.

2. Update your producer code to choose a random number between 1 and 5 for each record's key.
3. Modify your consumer to consume only records with a specific key (or subset of keys).
4. Attempt to consume records with a key that does not exist. E.g., consume records with key value of “100”. Describe the results
5. Can you create a consumer that only consumes specific keys? If you run this consumer multiple times with varying keys then does it allow you to consume messages out of order while maintaining order within each key?

Answer: No we cannot consume only specific keys, Also if we use the same key the order of messages is maintained whereas using different keys the order of messages is not maintained.

G. Producer Flush

The provided tutorial producer program calls “`producer.flush()`” at the very end, and presumably your new producer also calls `producer.flush()`.

1. What does `Producer.flush()` do?

Answer: `produce()` method is asynchronous, what it does is it adds the messages on an internal queue which is then later sent to the broker. Invoking `flush()` makes the client wait before all the messages have been actually sent to the broker.

2. What happens if you do not call `producer.flush()`?

Answer: if we do not call `producer.flush()`, there is no guarantee that the messages produced by the `produce()` method will be delivered to the broker.

3. What happens if you call `producer.flush()` after sending each record?

Answer: If you add `flush()` after each `produce()` call you are effectively implementing a sync producer, which will make sure every message once produced (invoked `produce()`) will be delivered to the broker synchronously.

4. What happens if you wait for 2 seconds after every 5th record send, and you call `flush` only after every 15 record sends, and you have a consumer running concurrently? Specifically, does the consumer receive each message immediately? only after a flush? Something else?

Answer: The consumer seems to be receiving 5 records after every flush
Producer:

```
ctawde@in-class-assign-2: ~/in-class-assign - Google Chrome
ssh.cloud.google.com/projects/cloud-f21-chinmay-tawde-ctawde/zones/us-central1-a/instances/in-class-assign-2?authuser=0&hl=en_US&projectNumber=698099379340&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot2...
a.placeholder.com/150/6b45ea4"
Producing record: 4631dd1-4742-47a7-a30b-1cd6cf61f8f39 ("albumId": 8, "id": 368, "title": "quaerat labore aut ducimus incidunt ex", "url": "https://via.placeholder.com/600/94182d", "thumbnailUrl": "https://via.placeholder.com/150/94182d")
Producing record: e3542e89-e016-4a3e-9e2d-d930b48f721a ("albumId": 8, "id": 369, "title": "neque perspiciatis sint vero non qui", "url": "https://via.placeholder.com/600/77e4a2", "thumbnailUrl": "https://via.placeholder.com/150/77e4a2")
Producing record: 4609d17-e353-457e-82b6-81ed1bf8f691a ("albumId": 8, "id": 370, "title": "rerum non quia dolore", "url": "https://via.placeholder.com/600/6d53ce", "thumbnailUrl": "https://via.placeholder.com/150/6d53ce")
Produced record to topic simple_photo partition [5] @ offset 9508
Producing record: 9397a0b1-477b-495b-a35b-8fda6b34a696 ("albumId": 8, "id": 372, "title": "ratione omnis fugiat sit fuga", "url": "https://via.placeholder.com/600/9c0b1e", "thumbnailUrl": "https://via.placeholder.com/150/9c0b1e")
Producing record: 2a354e83-6d37-41bb-a779-4a2fb0845a1e ("albumId": 8, "id": 373, "title": "eum dicta deleniti porro", "url": "https://via.placeholder.com/600/6a6136", "thumbnailUrl": "https://via.placeholder.com/150/6a6136")
Producing record: 6931de3-cf69-4d10-8bf8-5f8ae3246181 ("albumId": 8, "id": 374, "title": "ullam aut consectetur libero provident et porro", "url": "https://via.placeholder.com/600/d4420e", "thumbnailUrl": "https://via.placeholder.com/150/d4420e")
Producing record: 0e03b0e2-976b-4b71-a321-08044020e5f ("albumId": 8, "id": 375, "title": "voluptas repudiandae totam dolores voluptatem tempora et assumenda ducimus", "url": "https://via.placeholder.com/600/8eb5c2", "thumbnailUrl": "https://via.placeholder.com/150/8eb5c2")
Produced record to topic simple_photo partition [4] @ offset 9501
Produced record to topic simple_photo partition [4] @ offset 9502
Produced record to topic simple_photo partition [5] @ offset 9509
Produced record to topic simple_photo partition [5] @ offset 9510
Produced record to topic simple_photo partition [3] @ offset 9564
Producing record: b7c2c2f6-65f5-4164-bd67-0df3e6639b7b ("albumId": 8, "id": 376, "title": "est exercitationem aliquam omnis quia quas qui qui dolor", "url": "https://via.placeholder.com/600/24d0d1", "thumbnailUrl": "https://via.placeholder.com/150/24d0d1")
Producing record: 599ab0ee-e2f2-485b-ba9b-6saacc09293d ("albumId": 8, "id": 377, "title": "illum architecto rerum rerum", "url": "https://via.placeholder.com/600/bf47cb", "thumbnailUrl": "https://via.placeholder.com/150/bf47cb")
Producing record: e613cb73-c647-46ba-9b50-dde37c67814c ("albumId": 8, "id": 378, "title": "veritatis quos vel omnis error", "url": "https://via.placeholder.com/600/c74808", "thumbnailUrl": "https://via.placeholder.com/150/c74808")
Producing record: 279eef2c-910b-40f3-919a-d922af45640d ("albumId": 8, "id": 379, "title": "quaerat rerum non", "url": "https://via.placeholder.com/600/ea74e", "thumbnailUrl": "https://via.placeholder.com/150/ea74e")
Producing record: 9f4aa3cf-d93b-41ad-9d0c-f41dd286b06 ("albumId": 8, "id": 380, "title": "voluptates earum dolor perferendis et", "url": "https://via.placeholder.com/600/6be8c1", "thumbnailUrl": "https://via.placeholder.com/150/6be8c1")
Produced record to topic simple_photo partition [3] @ offset 9565
Produced record to topic simple_photo partition [3] @ offset 9566
Produced record to topic simple_photo partition [1] @ offset 9491
Produced record to topic simple_photo partition [0] @ offset 9477
Produced record to topic simple_photo partition [4] @ offset 9503
Producing record: 16c7243b-d061-4c6b-9b52-71c578c72124 ("albumId": 8, "id": 381, "title": "sed quo et et nemo earum omnis quia", "url": "https://via.placeholder.com/600/627b42", "thumbnailUrl": "https://via.placeholder.com/150/627b42")
Producing record: d348ae6-06b8-447a-9409-482f015ef3d ("albumId": 8, "id": 382, "title": "iusto nam atque facilis est eos", "url": "https://via.placeholder.com/600/36f93e", "thumbnailUrl": "https://via.placeholder.com/150/36f93e")
Producing record: 03e4f03d-f95b-4bb8-b582-a54230ad1aeb ("albumId": 8, "id": 383, "title": "doloribus est assumenda eligendi cum asperiores earum vel", "url": "https://via.placeholder.com/600/6f3ee", "thumbnailUrl": "https://via.placeholder.com/150/6f3ee")
Producing record: 59946414-8119-4967-9f93-0c402e86dd0c ("albumId": 8, "id": 384, "title": "aut quia ad earum consequatur", "url": "https://via.placeholder.com/600/d94fb7", "thumbnailUrl": "https://via.placeholder.com/150/d94fb7")
Producing record: 60869d92-6486-446a-9b50-dde37c67814c ("albumId": 8, "id": 385, "title": "blanditiis labore fugiat eum esse dolores inventore", "url": "https://via.placeholder.com/600/696ef", "thumbnailUrl": "https://via.placeholder.com/150/696ef")
Produced record to topic simple_photo partition [3] @ offset 9567
Produced record to topic simple_photo partition [3] @ offset 9568
Produced record to topic simple_photo partition [0] @ offset 9478
Produced record to topic simple_photo partition [0] @ offset 9479
Produced record to topic simple_photo partition [5] @ offset 9511
Producing record: c7f4a3cf-d93b-41ad-9d0c-f41dd286b06 ("albumId": 8, "id": 386, "title": "sequi autem fugiat ab incidunt mollitia", "url": "https://via.placeholder.com/600/6b51f3", "thumbnailUrl": "https://via.placeholder.com/150/6b51f3")
Producing record: 8f8d4f50-9b52-4604-94d9-eb15ef1ef6b6 ("albumId": 8, "id": 387, "title": "et quam explicabo molestiae fugiat ipsa eum nesciunt quae", "url": "https://via.placeholder.com/600/747986", "thumbnailUrl": "https://via.placeholder.com/150/747986")
Producing record: 7710b387-68bd-476c-9a73-f13570e5c83 ("albumId": 8, "id": 388, "title": "quos tempore nihil rerum rerum aut libero", "url": "https://via.placeholder.com/600/8661f8", "thumbnailUrl": "https://via.placeholder.com/150/8661f8")
Producing record: c7f4a3cf-d93b-41ad-9d0c-f41dd286b06 ("albumId": 8, "id": 389, "title": "sapiente illum vel adipisci aliquid quia", "url": "https://via.placeholder.com/600/122741", "thumbnailUrl": "https://via.placeholder.com/150/122741")
Producing record: 73e395df-3db9-480f-a7b4-bf6d33c635c8 ("albumId": 8, "id": 390, "title": "reprehenderit nesciunt delectus", "url": "https://via.placeholder.com/600/7df63c", "thumbnailUrl": "https://via.placeholder.com/150/7df63c")
Ctrl+back (most recent call last):
  File ~/sample_producer.py, line 74, in <module>
    time.sleep(2)
KeyboardInterrupt
(confluent-exercise) ctawde@in-class-assign-2: ~/in-class-assign$
```

Consumer :

```
ctawde@in-class-assign-2: ~/in-class-assign - Google Chrome
ssh.cloud.google.com/projects/cloud-f21-chinmay-tawde-ctawde/zones/us-central1-a/instances/in-class-assign-2?authuser=0&hl=en_US&projectNumber=698099379340&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true
a.placeholder.com/150/bf47cb"
Consumed record with key b'16c7243b-d061-4c6b-9b50-dde37c67814c' and value b'{"albumId": 8, "id": 378, "title": "veritatis quos vel omnis error", "url": "https://via.placeholder.com/600/c74808", "thumbnailUrl": "https://via.placeholder.com/150/c74808"}', and data is {'albumId': 8, 'id': 378, 'title': 'veritatis quos vel omnis error', 'url': 'https://via.placeholder.com/600/c74808', 'thumbnailUrl': 'https://via.placeholder.com/150/c74808'}
Consumed record with key b'279eef2c-910b-40f3-919a-d922af45640d' and value b'{"albumId": 8, "id": 379, "title": "quaerat rerum non", "url": "https://via.placeholder.com/600/ea74e", "thumbnailUrl": "https://via.placeholder.com/150/ea74e"}', and data is {'albumId': 8, 'id': 379, 'title': 'quaerat rerum non', 'url': 'https://via.placeholder.com/600/ea74e', 'thumbnailUrl': 'https://via.placeholder.com/150/ea74e'}
Consumed record with key b'b7c2c2f6-65f5-4164-bd67-0df3e6639b7b' and value b'{"albumId": 8, "id": 376, "title": "est exercitationem aliquam omnis quia quas qui qui dolor", "url": "https://via.placeholder.com/600/24d0d1", "thumbnailUrl": "https://via.placeholder.com/150/24d0d1"}', and data is {'albumId': 8, 'id': 376, 'title': 'est exercitationem aliquam omnis quia quas qui qui dolor', 'url': 'https://via.placeholder.com/600/24d0d1', 'thumbnailUrl': 'https://via.placeholder.com/150/24d0d1'}
Consumed record with key b'9f4aa3cf-d93b-41ad-9d0c-f41dd286b06' and value b'{"albumId": 8, "id": 380, "title": "voluptates earum dolor perferendis et", "url": "https://via.placeholder.com/600/6be8c1", "thumbnailUrl": "https://via.placeholder.com/150/6be8c1"}', and data is {'albumId': 8, 'id': 380, 'title': 'voluptates earum dolor perferendis et', 'url': 'https://via.placeholder.com/600/6be8c1', 'thumbnailUrl': 'https://via.placeholder.com/150/6be8c1'}
Waiting for message or event/error in poll()
Consumed record with key b'4631dd1-4742-47a7-a30b-1cd6cf61f8f39' and value b'{"albumId": 8, "id": 368, "title": "quaerat labore aut ducimus incidunt ex", "url": "https://via.placeholder.com/600/94182d", "thumbnailUrl": "https://via.placeholder.com/150/94182d"}', and data is {'albumId': 8, 'id': 368, 'title': 'quaerat labore aut ducimus incidunt ex', 'url': 'https://via.placeholder.com/600/94182d', 'thumbnailUrl': 'https://via.placeholder.com/150/94182d'}
Consumed record with key b'0e03b0e2-976b-4b71-a321-08044020e5f' and value b'{"albumId": 8, "id": 375, "title": "voluptas repudiandae totam dolores voluptatem tempora et assumenda ducimus", "url": "https://via.placeholder.com/600/8eb5c2", "thumbnailUrl": "https://via.placeholder.com/150/8eb5c2"}', and data is {'albumId': 8, 'id': 375, 'title': 'voluptas repudiandae totam dolores voluptatem tempora et assumenda ducimus', 'url': 'https://via.placeholder.com/600/8eb5c2', 'thumbnailUrl': 'https://via.placeholder.com/150/8eb5c2'}
Consumed record with key b'9397a0b1-477b-495b-a35b-8fda6b34a696' and value b'{"albumId": 8, "id": 372, "title": "ratione omnis fugiat sit fuga", "url": "https://via.placeholder.com/600/9c0b1e", "thumbnailUrl": "https://via.placeholder.com/150/9c0b1e"}', and data is {'albumId': 8, 'id': 372, 'title': 'ratione omnis fugiat sit fuga', 'url': 'https://via.placeholder.com/600/9c0b1e', 'thumbnailUrl': 'https://via.placeholder.com/150/9c0b1e'}
Consumed record with key b'2a354e83-6d37-41bb-a779-4a2fb0845a1e' and value b'{"albumId": 8, "id": 373, "title": "eum dicta deleniti porro", "url": "https://via.placeholder.com/600/6a6136", "thumbnailUrl": "https://via.placeholder.com/150/6a6136"}', and data is {'albumId': 8, 'id': 373, 'title': 'eum dicta deleniti porro', 'url': 'https://via.placeholder.com/600/6a6136', 'thumbnailUrl': 'https://via.placeholder.com/150/6a6136'}
Consumed record with key b'60869d92-6486-446a-9b50-dde37c67814c' and value b'{"albumId": 8, "id": 385, "title": "blanditiis labore fugiat eum esse dolores inventore", "url": "https://via.placeholder.com/600/696ef", "thumbnailUrl": "https://via.placeholder.com/150/696ef"}', and data is {'albumId': 8, 'id': 385, 'title': 'blanditiis labore fugiat eum esse dolores inventore', 'url': 'https://via.placeholder.com/600/696ef', 'thumbnailUrl': 'https://via.placeholder.com/150/696ef'}
Consumed record with key b'16c7243b-d061-4c6b-9b52-71c578c72124' and value b'{"albumId": 8, "id": 381, "title": "sed quo et et nemo earum omnis quia", "url": "https://via.placeholder.com/600/627b42", "thumbnailUrl": "https://via.placeholder.com/150/627b42"}', and data is {'albumId': 8, 'id': 381, 'title': 'sed quo et et nemo earum omnis quia', 'url': 'https://via.placeholder.com/600/627b42', 'thumbnailUrl': 'https://via.placeholder.com/150/627b42'}
Consumed record with key b'd348ae6-06b8-447a-9409-482f015ef3d' and value b'{"albumId": 8, "id": 382, "title": "iusto nam atque facilis est eos", "url": "https://via.placeholder.com/600/36f93e", "thumbnailUrl": "https://via.placeholder.com/150/36f93e"}', and data is {'albumId': 8, 'id': 382, 'title': 'iusto nam atque facilis est eos', 'url': 'https://via.placeholder.com/600/36f93e', 'thumbnailUrl': 'https://via.placeholder.com/150/36f93e'}
Consumed record with key b'a0346148-8119-4967-9f93-0c402e86dd0c' and value b'{"albumId": 8, "id": 384, "title": "aut quia ad earum consequatur", "url": "https://via.placeholder.com/600/d94fb7", "thumbnailUrl": "https://via.placeholder.com/150/d94fb7"}', and data is {'albumId': 8, 'id': 384, 'title': 'aut quia ad earum consequatur', 'url': 'https://via.placeholder.com/600/d94fb7', 'thumbnailUrl': 'https://via.placeholder.com/150/d94fb7'}
Waiting for message or event/error in poll()
Consumed record with key b'9f4aa3cf-d93b-41ad-9d0c-f41dd286b06' and value b'{"albumId": 8, "id": 380, "title": "sequi autem fugiat ab incidunt mollitia", "url": "https://via.placeholder.com/600/6b51f3", "thumbnailUrl": "https://via.placeholder.com/150/6b51f3"}', and data is {'albumId': 8, 'id': 380, 'title': 'sequi autem fugiat ab incidunt mollitia', 'url': 'https://via.placeholder.com/600/6b51f3', 'thumbnailUrl': 'https://via.placeholder.com/150/6b51f3'}
Consumed record with key b'8f8d4f50-9b52-4604-94d9-eb15ef1ef6b6' and value b'{"albumId": 8, "id": 387, "title": "et quam explicabo molestiae fugiat ipsa eum nesciunt quae", "url": "https://via.placeholder.com/600/747986", "thumbnailUrl": "https://via.placeholder.com/150/747986"}', and data is {'albumId': 8, 'id': 387, 'title': 'et quam explicabo molestiae fugiat ipsa eum nesciunt quae', 'url': 'https://via.placeholder.com/600/747986', 'thumbnailUrl': 'https://via.placeholder.com/150/747986'}
Consumed record with key b'7710b387-68bd-476c-9a73-f13570e5c83' and value b'{"albumId": 8, "id": 388, "title": "quos tempore nihil rerum rerum aut libero", "url": "https://via.placeholder.com/600/8661f8", "thumbnailUrl": "https://via.placeholder.com/150/8661f8"}', and data is {'albumId': 8, 'id': 388, 'title': 'quos tempore nihil rerum rerum aut libero', 'url': 'https://via.placeholder.com/600/8661f8', 'thumbnailUrl': 'https://via.placeholder.com/150/8661f8'}
Consumed record with key b'c7f4a3cf-d93b-41ad-9d0c-f41dd286b06' and value b'{"albumId": 8, "id": 389, "title": "sapiente illum vel adipisci aliquid quia", "url": "https://via.placeholder.com/600/122741", "thumbnailUrl": "https://via.placeholder.com/150/122741"}', and data is {'albumId': 8, 'id': 389, 'title': 'sapiente illum vel adipisci aliquid quia', 'url': 'https://via.placeholder.com/600/122741', 'thumbnailUrl': 'https://via.placeholder.com/150/122741'}
Consumed record with key b'73e395df-3db9-480f-a7b4-bf6d33c635c8' and value b'{"albumId": 8, "id": 390, "title": "reprehenderit nesciunt delectus", "url": "https://via.placeholder.com/600/7df63c", "thumbnailUrl": "https://via.placeholder.com/150/7df63c"}', and data is {'albumId': 8, 'id': 390, 'title': 'reprehenderit nesciunt delectus', 'url': 'https://via.placeholder.com/600/7df63c', 'thumbnailUrl': 'https://via.placeholder.com/150/7df63c'}
Waiting for message or event/error in poll()
```

Reference: <https://lifesaver.codes/answer/is-producer-flush-a-must-137>

H. Consumer Groups

1. Create two consumer groups with one consumer program instance in each group.
2. Run the producer and have it produce all 1000 messages from your sample file.
3. Run each of the consumers and verify that each consumer consumes all of the 50 messages.
4. Create a second consumer within one of the groups so that you now have three consumers total.
5. Rerun the producer and consumers. Verify that each consumer group consumes the full set of messages but that each consumer within a consumer group only consumes a portion of the messages sent to the topic.

I. Kafka Transactions

6. Create a new producer, similar to the previous producer, that uses transactions.
7. The producer should begin a transaction, send 4 records in the transactions, then wait for 2 seconds, then choose True/False randomly with equal probability. If True then finish the transaction successfully with a commit. If False is picked then cancel the transaction.
8. Create a new transaction-aware consumer. The consumer should consume the data. It should also use the Confluent/Kafka transaction API with a "read_committed" isolation level. (I can't find evidence of other isolation levels).
9. Transaction across multiple topics. Create a second topic and modify your producer to send two records to the first topic and two records to the second topic before randomly committing or canceling the transaction. Modify the consumer to consume from the two queues. Verify that it only consumes committed data and not uncommitted or canceled data.