

Developing Kafka deployment and publisher/consumer for ReST based services

TANU BORDIA(IMT2016002) | VISHESH RUPARELIA(IMT2016006)

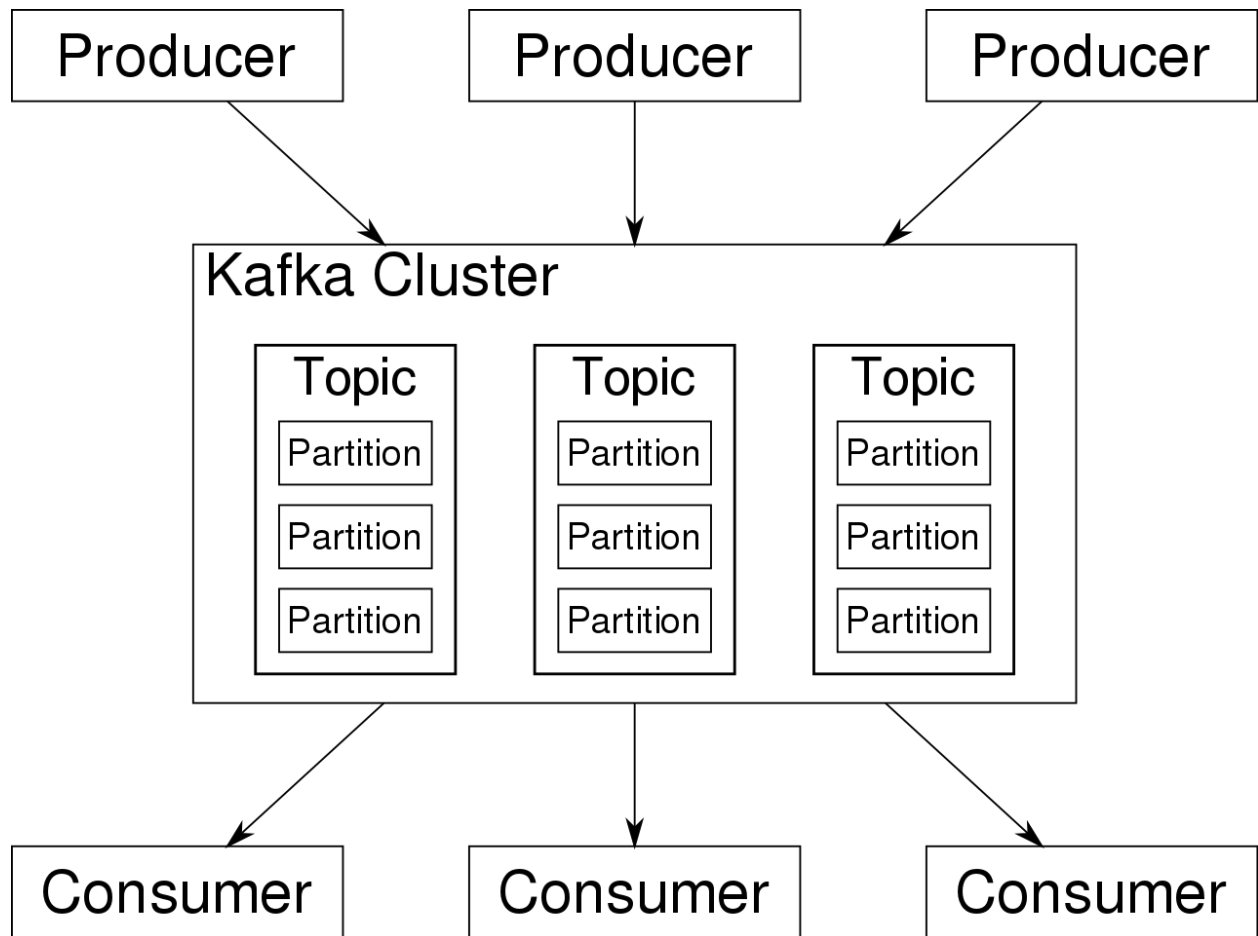
GITHUB: <https://github.com/visheshruparelia/hackernewsKafka>

Introduction

Apache Kafka is a community-based streaming event system that can handle trillions of events a day. Kafka is based on a distributed commit log abstraction, though it was thought of as a messaging queue in the beginning. Kafka has grown rapidly from messaging queue to a full-fledged event streaming platform since it was developed and launched by LinkedIn in 2011.

There are many arbitrary processes called producers that send key-valued messages to Kafka, which it stores. Within different "topics," data can be partitioned into different "partitions." In a partition, messages are ordered based on their offsets and are indexed and stored with timestamp. Other processes called consumers read the messages from these partitions.

Brokers are clusters of one or more servers, Kafka runs on these. Partition of all topics is distributed across the cluster nodes. Also, partitions are copied to multiple brokers. This architecture enables Kafka to deliver massive messaging streams in a tolerant manner.



Three major APIs used in our project were:

- **Producer API** – Permits an application to publish streams of records.
- **Consumer API** – Permits an application to subscribe to topics and processes streams of records.
- **Connector API** – Executes the reusable producer and consumer APIs that can link the topics to the existing applications.

Problem Statement

Developing Kafka deployment and publisher/consumer for ReST based services

1. Deploy kafka
2. Assume as application is exposing data using ReST services - for that any standard API available on open network can be used.
3. There is a publisher that calls this ReST api to read the data (ensure that the ReST APIs are returning data in JSON Format) and put it on kafka topic using Java APIs.
4. There are three consumers who are reading this data from this topic. i) one of them is printing it locally ii) Sending it using ReST apis to a remote applications. lii) Storing the data in ElasticSerach DB.

Solution

Overview:

A producer script will pull new stories posted on <https://news.ycombinator.com/> and publish it to the kafka bus under the topic name "story". Three consumers pull data from this bus and process and store it accordingly.

Language: Python

Library Dependencies: haxor, kafka-python, flask, elasticsearch

Installing dependencies:

1. cd into the root directory of the project.
2. Run "pip3 install -r requirements.txt"
3. **Install Docker:**
 - a. <https://docs.docker.com/v17.12/install/linux/docker-ce/ubuntu/#install-using-the-repository>
4. **Install Docker-compose:**
 - a. <https://docs.docker.com/compose/install/>

Steps to run:

1. Kafka, Zookeeper, ElasticSearch, Kibana:

- a. cd into root directory of the project
- b. Run "sudo docker-compose up --no-recreate"

Outcomes:

- Kafka will run on port 9092
- Zookeeper will run on port 2181
- ElasticSearch will run on port 9200
- Kibana will run on port 5601

2. Producer:

- a. cd into the root directory of the project
- b. Run "python3 producer.py"

- **Outcome:** The script will keep publishing data to a topic named "story" in the kafka bus.

3. Local Consumer:

- a. cd into local_consumer folder
- b. Run "python3 consumer.py"

- **Outcome:** The script will pull the data from the bus and append it to consumerdata.csv

4. Remote Consumer:

- a. cd into flask_app folder
- b. Run "python3 consumerflask.py"

- **Outcome:** Webpage will be hosted on <http://localhost:5000>

5. Elastic Search Consumer:

- a. cd into elastic_consumer folder
 - b. Run "python3 consumer.py"
- **Outcome:** Data will be indexed to Elasticsearch under index "hn_story" and can be viewed from Kibana console at <http://localhost:5601>

The following section will also cover results and observations:

Documentation:

1. **HackerNews API:** Exposes data continuously

Documentation: <https://github.com/HackerNews/API>

```
while(True):
    story=hn.new_stories(limit=1)
    if(len(story)==0 or story[0].item_id<=latest_id):
        continue
    else:
        start=latest_id+1
        for i in range(start,story[0].item_id+1):
            curr_story=hn.get_item(i)
            if(curr_story.item_type!='story'):
                continue
            title=curr_story.title
            id=curr_story.item_id
            by=curr_story.by
            time=curr_story.submission_time
            a=time.strftime("%m/%d/%Y, %H:%M:%S")
            data={'id':id,'title':title,'by':by,'time':a}
            latest_id=id
            print('publishing {}'.format(data))
            producer.send('story',value=data)
            print('\n')
        sleep(5)
```

Will call "new_stories" method every 5 seconds.

-
2. **Kafka Producer:** A Kafka client that publishes records to the Kafka cluster.

The producer consists of a pool of buffer space that holds records that haven't yet been transmitted to the server as well as a background I/O thread that is responsible for turning these records into requests and transmitting them to the cluster.

File location: producer.py

Command to run: python3 producer.py

It keeps pulling new stories posted on the HackerNews forum

```
from kafka import KafkaProducer

producer = KafkaProducer(bootstrap_servers=['localhost:9092'],
                          value_serializer=lambda x:
                              dumps(x).encode('utf-8'))
```

When run, Kafka producer gets data from the HackerNews and sends it to the consumer.

```
producer.send('story', value=data)
```

3. **Kafka Consumers:** Consume records from a Kafka cluster.

There are 3 consumer groups:

-
- a. **CSVConsumer:** Consumes the data and stores it in a CSV file on the local system which gets updated continuously. (Local Consumer)

```
consumer = KafkaConsumer(  
    'story',  
    bootstrap_servers=['localhost:9092'],  
    auto_offset_reset='earliest',  
    enable_auto_commit=True,  
    group_id='CSVConsumer',  
    value_deserializer=lambda x: loads(x.decode('utf-8')))
```

File Location: local_consumer/consumer.py

Output Location: local_consumer/consumerdata.csv

Output Screenshot:

| | A | B | C | D |
|----|----------|---|-------------|----------------------|
| 1 | id | title | by | time |
| 2 | 21655958 | Firefox Replay | nachtigall | 11/28/2019, 15:29:47 |
| 3 | 21655951 | Bill Gates-backed Heliogen targets 1,500°C solar thermal | nclzz | 11/28/2019, 15:28:15 |
| 4 | 21655945 | QuestDB – Fast open source relational time series database | beagle3 | 11/28/2019, 15:27:03 |
| 5 | 21655940 | Use Tensor Flow Lite and a Particle Xenon to Build an ML-Gesture Wand | adunk | 11/28/2019, 15:25:58 |
| 6 | 21655936 | Test Automation Without Coding | dsheiko | 11/28/2019, 15:25:22 |
| 7 | 21655926 | Biological enzymes as source of hydrogen fuel | conse_lad | 11/28/2019, 15:23:59 |
| 8 | 21655903 | Replace Windows 7 with Linux | dagurp | 11/28/2019, 15:19:40 |
| 9 | 21655897 | How the FCC Helped Pave the Way for Predatory Prison Telecoms | vezycash | 11/28/2019, 15:18:52 |
| 10 | 21655892 | That Uplifting Tweet You Just Shared? A Russian Troll Sent It | edoloughlin | 11/28/2019, 15:18:14 |
| 11 | 21655890 | The blockers and catalysts of digital transformation | mooreds | 11/28/2019, 15:17:11 |

- b. **FlaskConsumer:** It is a flask app which consumes data from the bus under the group name “FlaskConsumer” and publishes the data to localhost:5000. Two threads run, one is the main flask app and the other one keeps pulling data from the bus.

File Location: flask_app/consumerflask.py

Output: Can be accessed at <http://localhost:5000>, refresh page every time to see new records.

Output Screenshot:

Stories

| Id | Title | Author | Time |
|----------|--|-----------------|----------------------|
| 21663269 | Just like I imaged it would happen in C&C when I was a kid | kakadoodle | 11/29/2019, 15:11:19 |
| 21663256 | Collaborative Organigrams App | javieranton | 11/29/2019, 15:08:32 |
| 21663248 | AWS Transit Gateway Examined | lucyinkedup | 11/29/2019, 15:05:54 |
| 21663245 | Singapore government must realise human error also a security breach | Jarqwiz | 11/29/2019, 15:05:31 |
| 21663234 | OpenSTF – Remote control all your Smartphone devices | samueladam | 11/29/2019, 15:03:51 |
| 21663228 | I Interviewed David Ngo, Head of Product and Eng. Metallic at CommvaultGO 2019 | dezb | 11/29/2019, 15:02:21 |
| 21663225 | Vectornator For Mac – New Graphic Design Software | VladShvets96 | 11/29/2019, 15:01:18 |
| 21663218 | Ask HN: De-identifying data without destroying quality | monkeydust | 11/29/2019, 14:59:34 |
| 21663204 | Software maintenance categories and what they mean for your users | roxyabercrombie | 11/29/2019, 14:57:29 |
| 21663176 | Cinematic video with a consumer drone and ND filter set | malik9 | 11/29/2019, 14:52:09 |

c. ElasticSearch Consumer:

Using “elasticsearch” module available in python, the script consumes data from the bus as the group “ElasticConsumer” and indexes it to the ElasticSearch DB instance running on localhost:9200 under the index name “hn_story” and doc type “test” with the id of the data being the id of the HN story.

A kibana instance connects with this elastic search instance and displays the data at localhost:5601.

NOTE: When opening kibana for the first time, an index named “hn_story” needs to be created.

File Location: elastic_consumer/consumer.py

```
⇒ python3 consumer.py
21656119 indexed.
21656115 indexed.
21656104 indexed.
21656091 indexed.
21656090 indexed.
21656082 indexed.
21656069 indexed.
21656061 indexed.
21656057 indexed.
21656039 indexed.
21656119 indexed.
21656115 indexed.
21656104 indexed.
21656091 indexed.
21656090 indexed.
21656082 indexed.
21656069 indexed.
21656061 indexed.
21656057 indexed.
21656039 indexed.
```

Kibana Output:

| 40 hits | | | | | | | | | |
|---------|----------------|---|-------------------|---------------|-------------|------------------|-----------|--|--|
| _source | | | | | | | | | |
| > | id: 21,620,522 | title: They Know My Secrets: Tales from Atomic's Software Tester | by: philk10 | _id: 21620522 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,643 | title: FlowTrace - remote app inspecting(Ruby). Declare steps and receive params | by: kirill_shevch | _id: 21620643 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,641 | title: Rats laugh when tickled - this and other reasons not to grow up quickly | by: TheSpine | _id: 21620641 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,640 | title: Scientists may have found a fifth force of nature | by: palebt | _id: 21620640 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,626 | title: Critical perspectives on provable security: Fifteen years of another look papers | by: JetSpiegel | _id: 21620626 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,613 | title: Show HN: Pmfx-shader - Cross platform shader system | by: polymonster | _id: 21620613 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,611 | title: Fibonacci Series and Indian Mathematics | by: inception44 | _id: 21620611 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,595 | title: Finding Unicorns: A Quick Lesson from Adam Sandler | by: rebase-vc | _id: 21620595 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,594 | title: Creating a Bloom Filter with Go | by: Insanity | _id: 21620594 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,586 | title: The Hunter Memorial Library | by: mcgin | _id: 21620586 | _type: test | _index: hn_story | _score: 0 | | |
| > | id: 21,620,540 | title: Mixing Hue and Innr Smart Lights | by: edent | _id: 21620540 | _type: test | _index: hn_story | _score: 0 | | |