



EUROPYTHON
2021 Jul 26–Aug 1 online

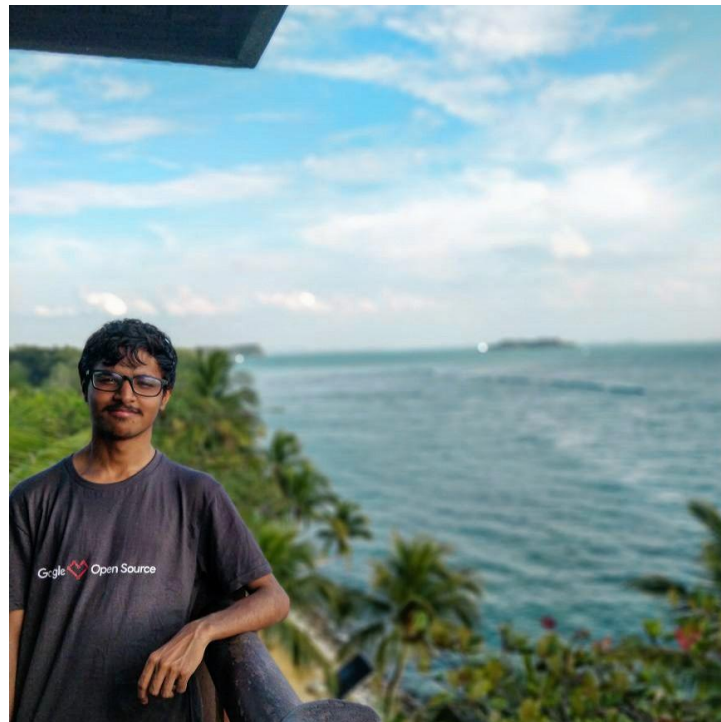
HIGH PERFORMANCE DATA PROCESSING WITH PYTHON, KAFKA AND ELASTICSEARCH

EUROPYTHON ONLINE 2021

Harshit Prasad / harshitprasad.com

INTRODUCTION

- Software Engineer at **Grofers - India**.
- Bachelors in Electronics and Communication Engineering from the **LNM Institute of Information Technology**
- Past GSoC'16 - **FOSSASIA**
- Past GSoC'17 - **CERN**
- Open Source Contributor.



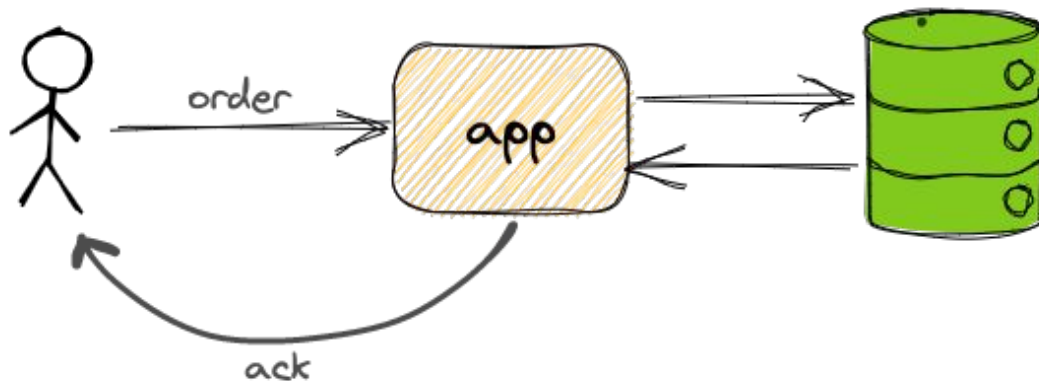
TALK OUTLINE

- Problem Statement
- Solution
- Workflow
- Code Walkthrough

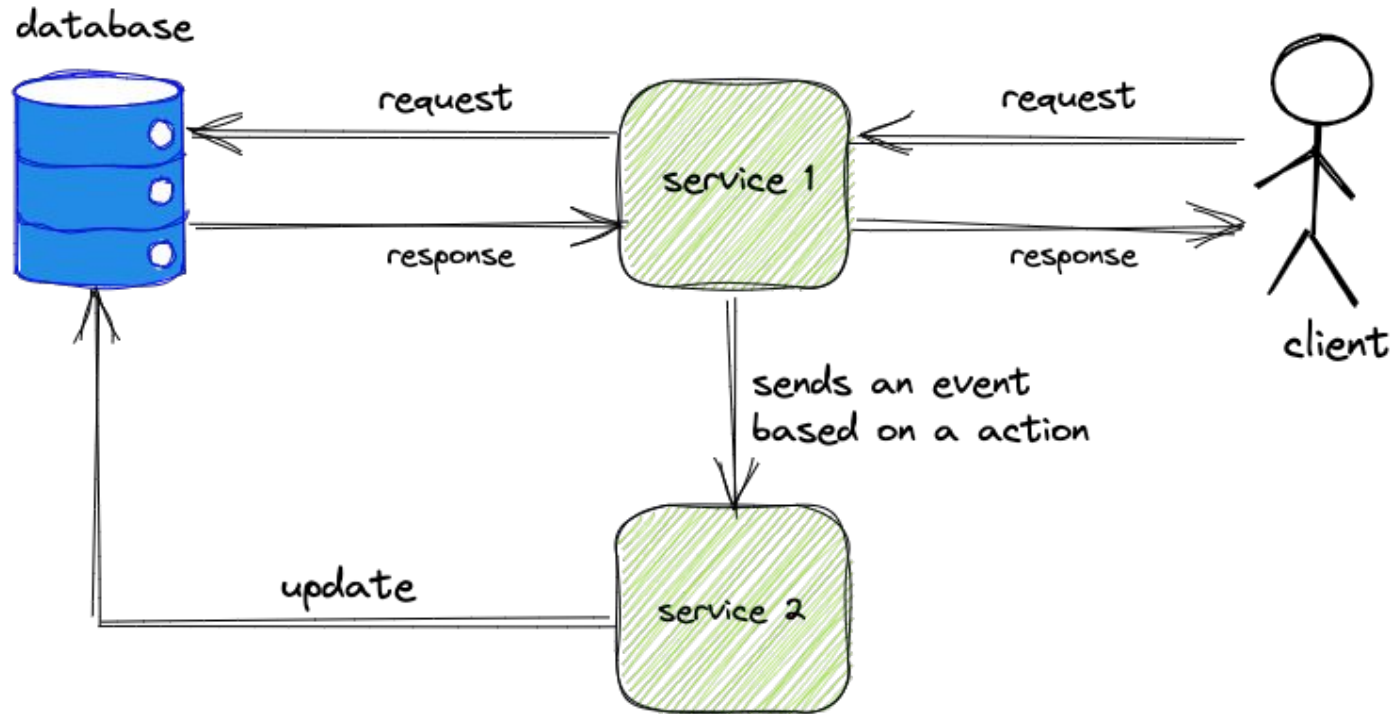


PROBLEM STATEMENT

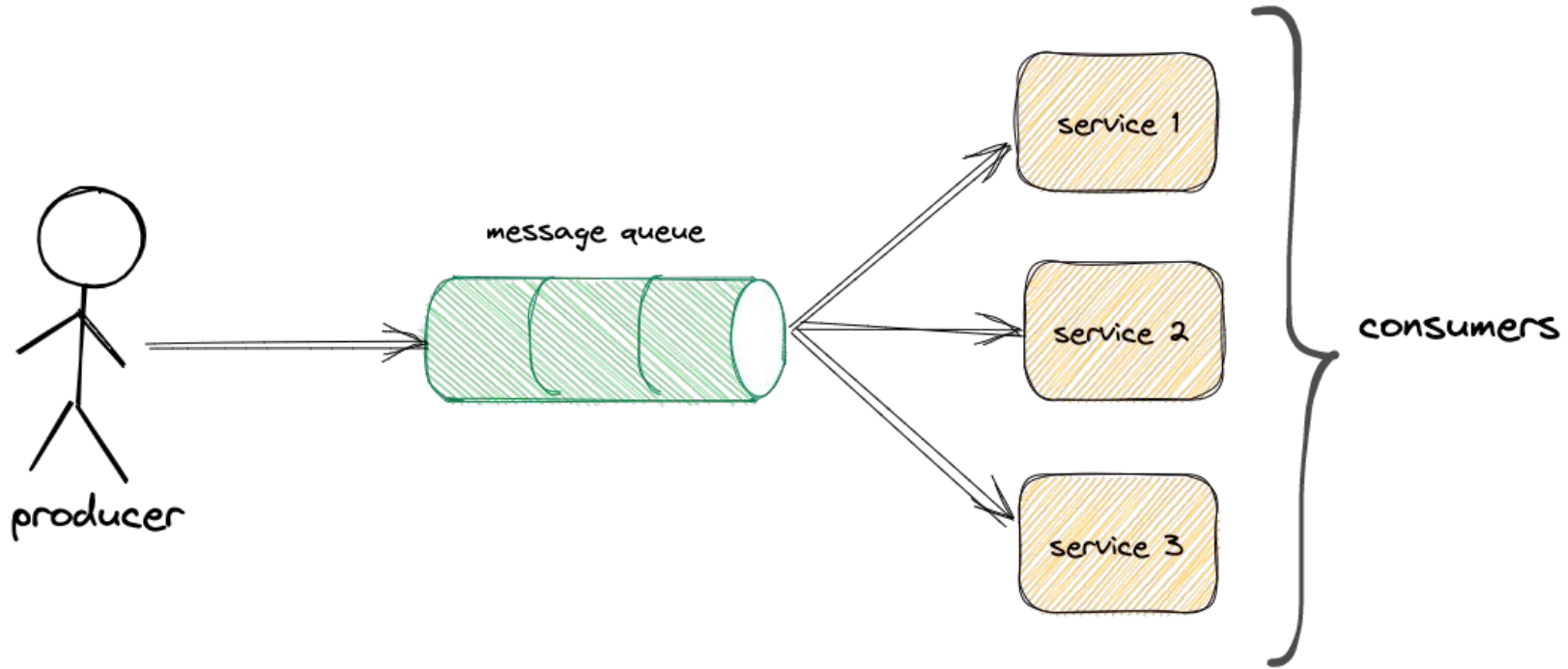
- Your own application
- Avg order rate X
- Peak hour order rate $X * 1000$



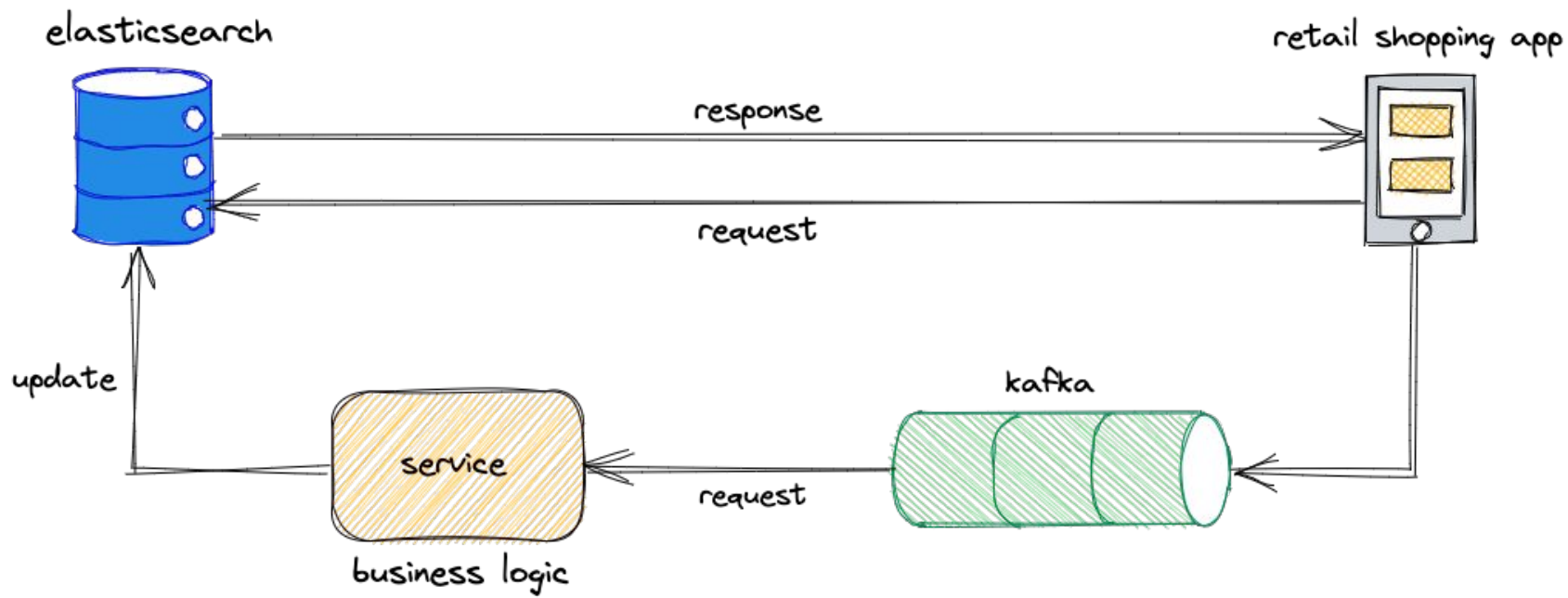
PROBLEM STATEMENT



PRODUCER CONSUMER MODEL



WORKFLOW



WHAT KIND OF DATA ARE WE UPDATING ?



```
{
  "_index" : "retail-catalogue",
  "_type" : "_doc",
  "_id" : "10",
  "_version" : 3,
  "_seq_no" : 4,
  "_primary_term" : 2,
  "found" : true,
  "_source" : {
    "product" : {
      "id" : 10,
      "name" : "Funny Farm House Ketchup",
      "category" : "Dips and Ketchups",
      "price" : 15.0,
      "stock" : 5
    },
    "update_ts" : 1625931472
  }
}
```


CODE WALKTHROUGH



- Python based implementation of Producer Consumer model.
- Walkthrough of Producer Consumer implementation - like producing data, pushing it to Kafka and then again consuming from Kafka to send an update to the Elasticsearch datastore.
- Data Streamer:
 - Producer
 - Consumer



DATA PRODUCER

```
1 app = FastAPI()
2
3 loop = asyncio.get_event_loop()
4
5 aioproducer = AIOKafkaProducer(loop=loop,
6                               client_id='retail-product',
7                               bootstrap_servers=KAFKA_INSTANCE)
8
9
10 @app.on_event("startup")
11 async def startup_event():
12     await aioproducer.start()
13
14
15 @app.on_event("shutdown")
16 async def shutdown_event():
17     await aioproducer.stop()
18
19
20 @app.post("/producer/{topicname}", response_model=ProducerResponse)
21 async def kafka_produce(msg: ProducerMessage, topicname: str):
22     """
23     Produce a message into <topicname>
24     This will produce a message into a Apache Kafka topic.
25     """
26     message_id = f"{msg.name}_{uuid.uuid4()}"
27     logger.debug(msg)
28
29     await aioproducer.send(topicname, json.dumps(msg.dict()).encode("ascii"))
30     response = ProducerResponse(name=msg.name,
31                                message_id=message_id,
32                                topic=topicname,
33                                timestamp=msg.timestamp)
34     logger.info(response)
35     return response
```

Initialization of AIOKafka Producer
with KAFKA_INSTANCE

Producer sends message to Kafka.
This message holds information
related to product.



GET http://0.0.0.0:8001/ping POST http://127.0.0.1:8001/produce... GET http://0.0.0.0:8000/ping GET http://127.0.0.1:8000/consumer... + ...

Untitled Request

POST http://127.0.0.1:8001/producer/retail-product

Params Authorization Headers (9) **Body** Pre-request Script Tests Settings

none form-data x-www-form-urlencoded raw binary GraphQL **JSON** ▾

```
1 {
2   "name": "Funny Farm House Ketchup",
3   "category": "Dips and Ketchups",
4   "price": 15,
5   "stock": 3,
6   "product_id": 10,
7   "timestamp": ""
8 }
```

Body Cookies Headers (4) Test Results

Pretty Raw Preview Visualize JSON ▾

```
1 [
2   {
3     "name": "Funny Farm House Ketchup",
4     "message_id": "Funny Farm House Ketchup_e5c0f1b0-ac1e-44c7-92c1-1f86728a36dc",
5     "topic": "retail-product",
6     "timestamp": "2021-07-12 15:17:20.314894"
7   }
8 ]
```



DATA CONSUMER

```
1 app = FastAPI()
2
3 loop = asyncio.get_event_loop()
4
5 es = AsyncElasticsearch(hosts=[ELASTICSEARCH_INSTANCE])
6
7 @app.get("/consumer/{topicname}")
8 async def kafka_consume(topicname: str):
9     """
10    Produce a message into <topicname>
11    This will produce a message into a Apache Kafka topic.
12    """
13    consumer = AIOKafkaConsumer(topicname,
14                                loop=loop,
15                                client_id='retail-catalogue',
16                                bootstrap_servers=KAFKA_INSTANCE,
17                                enable_auto_commit=True,
18                                value_deserializer=kafka_json_deserializer)
19
20    await consumer.start()
21    logger.debug('starting consumer task')
22
23    retrieved_requests = []
24    try:
25        result = await consumer.getmany(timeout_ms=20000)
26        logger.info(f"Got {len(result)} messages in {topicname}.")
27
28        for _, messages in result.items():
29            if messages:
30                for message in messages:
31                    retrieved_requests.append({"value": message.value})
32
33        for request in retrieved_requests:
34            value = request.get('value')
35            ....
```

Initialization of AIOKafka Consumer
with KAFKA_INSTANCE

Consume messages in a batch from
"retail-product" topic.



GET http://0.0.0.0:8001/ping ● POST http://127.0.0.1:8001/produce... ● GET http://0.0.0.0:8000/ping ● GET http://127.0.0.1:8000/consumer... ● + ...

Untitled Request

GET http://127.0.0.1:8000/consumer/retail-product

Params Authorization Headers (7) Body Pre-request Script Tests Settings

Query Params

KEY	VALUE
Key	Value

Body Cookies Headers (4) Test Results

Pretty Raw Preview Visualize JSON

```
1 {
2   "topic": "retail-product",
3   "timestamp": "2021-07-12 20:47:21.067802",
4   "product_name": "Funny Farm House Ketchup",
5   "product_id": 10,
6   "success": true
7 }
```



DATA UPDATE

Get the product doc which has to be update by “_id”

Update the product price and stock

Send update request to Elasticsearch to update the following product doc.

```
1     ....
2     product_id = value.get('product_id')
3     price = value.get('price')
4     stock = value.get('stock')
5
6     logger.info("data before sending to elasticsearch for update")
7     try:
8         data = await es.get(index='retail-catalogue', id=product_id)
9         logger.info("data: ", data['_source'])
10    except TransportError as e:
11        raise e
12
13    logger.info("after update, new data in elasticsearch")
14    data['_source']['product']['_id'] = product_id
15    data['_source']['product']['price'] = price
16    data['_source']['product']['stock'] = stock
17
18    product_name = data['_source']['product']['name']
19    product_id = data['_source']['product']['_id']
20
21    update_data = data['_source']
22    data_for_update = {'doc': {'product': update_data['product']}}
23
24    logger.info("data going for update in elasticsearch")
25    try:
26        await es.update(index='retail-catalogue',
27                        body=data_for_update,
28                        id=product_id)
29
30        response = ConsumerResponse(topic=topicname,
31                                    timestamp=str(datetime.now()),
32                                    product_name=product_name,
33                                    product_id=product_id,
34                                    success=True)
35
36        logger.info(response)
37    except TransportError as e:
38        raise e
39
40    return response
41
42 except Exception as e:
43     logger.error(
44         f"Error when trying to consume request on topic {topicname}: {str(e)}"
45     )
46 finally:
47     await consumer.stop()
```

RESULT



```
{
  "_index" : "retail-catalogue",
  "_type" : "_doc",
  "_id" : "10",
  "_version" : 3,
  "_seq_no" : 4,
  "_primary_term" : 2,
  "found" : true,
  "_source" : {
    "product" : {
      "id" : 10,
      "name" : "Funny Farm House Ketchup",
      "category" : "Dips and Ketchups",
      "price" : 15.0,
      "stock" : 5
    }
  },
  "update_ts" : 1625931472
}
```



```
{
  "_index" : "retail-catalogue",
  "_type" : "_doc",
  "_id" : "10",
  "_version" : 4,
  "_seq_no" : 5,
  "_primary_term" : 2,
  "found" : true,
  "_source" : {
    "product" : {
      "id" : 10,
      "name" : "Funny Farm House Ketchup",
      "category" : "Dips and Ketchups",
      "price" : 15.0,
      "stock" : 3
    }
  },
  "update_ts" : 1625931472
}
```



CONCLUSION AND LEARNINGS

- We learnt about how Producer Consumer models work and their advantages to process millions of data.
- Deep dive in example of Producer Consumer working and how it can be used with as a background service.
- Understanding of how datastore is being updated and achieving high performance and other benefits using messaging queue.
- Python code implementation can be found in my git repository - <https://github.com/harshit98/Retail-Updates-Streamer>
- You can reach me out on Matrix / Email / LinkedIn.

THAT'S ALL FOLKS, THANK YOU !



EUROPYTHON

2021 Jul 26 - Aug 1 Online



harshit98



harshit-prasad



HarshitPrasad8



harshitprasad.com