
Kafka

Fernando Anselmo

<http://fernandoanselmo.orgfree.com/wordpress/>

Versão 1.01 em 7 de outubro de 2023

Resumo

Kafka[1] ou "Apache Kafka" (aqui chamarei apenas de Kafka) é parte do Ecosistema Hadoop e uma plataforma de streaming de eventos distribuída pela comunidade capaz de lidar com trilhões por dia. Inicialmente concebido como uma simples fila de mensagens, Kafka é baseado em uma abstração de um log de confirmação distribuído. Desde que foi criado e aberto pelo LinkedIn em 2011, evoluiu rapidamente da fila de mensagens para uma plataforma de streaming de eventos completa.

1 Parte inicial

Kafka é uma plataforma de streaming distribuída. Kafka pode ser usado para construir aplicativos de streaming em tempo real que podem transformar os fluxos de dados ou deduzir alguma inteligência deles.



Figura 1: Logo do Apache Kafka

Kafka fornece um banco de dados de streaming de vários sistemas de origem. Kafka propõe ter os seguintes recursos mencionados para seus bancos de dados de streaming:

- Análise em tempo real de fluxos de *Big Data*
- Carrega recursos tolerantes a falhas
- Sistema durável, escalável e rápido
- Permite o rastreamento de dados de IoT

Kafka é integrado a pipelines de dados de streaming que compartilham dados entre sistemas e/ou aplicativos e também é integrado aos sistemas e aplicativos que consomem esses dados. Oferece suporte a uma variedade de casos de uso em que alta taxa de transferência e escalabilidade são vitais. Minimiza a necessidade de integrações ponto a ponto para compartilhamento de dados em determinados aplicativos, ele pode reduzir a latência para milissegundos. Isso significa que os dados estão disponíveis para os usuários mais rapidamente, o que pode ser vantajoso em casos de uso que exigem disponibilidade de dados em tempo real, como operações de TI e comércio eletrônico.

1.1 Como funcionam os bancos de dados de streaming?

Quando um fluxo atinge um banco de dados de streaming em tempo real, é processado imediatamente. Esses dados podem ser usados por um aplicativo após sua análise.

As entradas de dados para um banco de dados de streaming são chamadas de "fluxos de dados". Esses fluxos de dados são eventos de sequências e são imutáveis. Os dados de entrada para um banco de dados são categorizados em duas camadas:

- A primeira é a streaming
- A segunda é criada pelo usuário com base no comportamento desses fluxos, que podem ser chamados de "Estatísticas de eventos".

Essa análise de fluxo de entrada é armazenada em colunas e tabelas da mesma forma que é armazenada em um banco de dados relacional tradicional. A imagem a seguir apresenta o fluxograma para o funcionamento de um banco de dados de streaming:



Figura 2: Representação do Banco de Dados de Streaming

Os dados reservados em um banco de dados de streaming podem ser informações do usuário de plataformas de mídia social, um arquivo de log gerado pela Web, tendências do usuário de comércio eletrônico, relatórios de atividade em um jogo ou telemetria de vários aplicativos em um *data centers*. Esses dados são processados sequencialmente e incrementalmente e, em seguida, usados para fazer análises como regressão, filtragem, amostragem e correlação.

Essa análise de dados em tempo real abre vários casos de uso para diferentes setores. As empresas podem fazer uso dessa análise e tomar decisões relevantes com base nos resultados da análise. Considere um exemplo para uma organização, onde a análise de mídia social é feita por meio dos recursos de um banco de dados de streaming. A organização pode analisar facilmente o comportamento e a atividade do usuário. Essa observação pode ajudar a empresa a dar novos passos a partir dos dados analisados para melhorar a eficiência.

1.2 Criar o contêiner Docker

A forma mais simples para obtermos o Kafka é através de um contêiner no Docker, deste modo podemos ter várias versões do banco instalada e controlar mais facilmente qual banco está ativo ou não. E ainda colhemos o benefício adicional de não termos absolutamente nada deixando sujeira em nosso sistema operacional ou áreas de memória.

Precisamos do Zookeeper para termos o Kafka, sendo assim, vamos começar com o contêiner deste:

Baixar a imagem oficial do Zookeeper:

```
$ docker pull zookeeper
```

Criar o contêiner do Zookeeper:

```
$ docker run -d -p 2181:2181 --name meu-zookeeper zookeeper
```

Devemos esperar o Zookeeper está levantado e rondando sem problemas para executar o Kafka, para verificarmos isso:

```
$ docker logs meu-zookeeper
```

Qual o papel do Zookeeper? É usado para selecionar automaticamente um líder para uma partição. Em caso de desligamento de algum broker, é realizada uma eleição, pelo ZooKeeper, para a posição de líder das partições (que caíram com o broker). Também metadados como, em qual corretor uma partição líder está vivendo. Esses detalhes são mantidos pelo ZooKeeper. Produtores que transmitem dados para tópicos ou Consumidores que leem os dados do fluxo de tópicos entrem em contato com o ZooKeeper para obter o broker mais próximo ou menos ocupado.

Estando tudo OK, vamos baixar a imagem do Kafka da **Confluent**:

```
$ docker pull confluentinc/cp-kafka
```

A **Confluent**[2] é uma empresa fundada pelos desenvolvedores do Kafka. Fornecem conectores adicionais ao Kafka por meio de versões Aberta e Comercial do Confluent Kafka. Também oferecem suporte por meio de qualquer Kafka Cluster ou configuração de aplicativo, caso necessite em sua organização.

E criar o contêiner do Kafka:

```
$ docker run -d --name meu-kafka -p 9092:9092
-e KAFKA_ZOOKEEPER_CONNECT=[meu-ip]:2181
-e KAFKA_ADVERTISED_LISTENERS=PLAINTEXT://[meu-ip]:9092
-e KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR=1 confluentinc/cp-kafka
```

E para verificar se está rondando sem problemas:

```
$ docker logs meu-kafka
```

2 Produzir e Consumir mensagens

Vamos usar o **NodeJS** para exemplificar como podemos verificar o funcionamento do nosso serviço do Kafka. O primeiro passo é criar uma pasta e inicializá-la com o NodeJS: `$ npm init -y`

Precisamos também da biblioteca de comunicação com o Kafka: `$ npm install kafkajs`

O **tópico** é uma categoria na qual os fluxos de eventos/registros são armazenados no cluster Kafka.

Um cluster Kafka pode ter vários tópicos. Criamos um programa chamado "topic.js" com o seguinte conteúdo:

```
1 const { Kafka } = require("kafkajs");
2
3 run();
4
5 async function run() {
6   try {
7     const kafka = new Kafka({
8       "clientId": "myapp",
9       "brokers": ["192.168.15.5:9092"]
10    });
11    const admin = kafka.admin();
12    console.log("Conectando...");
13    await admin.connect();
14    console.log("Conectado!");
15    await admin.createTopics({
16      "topics": [{
17        "topic": "Usuarios",
18        "numPartitions": 2
19      }]
20    });
21    console.log("Sucesso");
22    await admin.disconnect();
23    console.log("Desconectado!");
24  } catch (ex) {
25    console.error('Alguma coisa aconteceu ${ex}')
26  } finally {
27    process.exit(0);
28  }
29 }
```

E rodamos este com o seguinte comando: `$ node topic.js`

Se olharmos no "Docker Logs" veremos que foi criado um tópico chamado "Usuarios" com duas partições. Agora podemos criar um produtor de conteúdo.

Produtores são aplicativos que enviam fluxos de dados para tópicos no Kafka Cluster. Um produtor pode enviar o fluxo de registros para vários tópicos. A API do Apache Kafka Producer permite que um aplicativo se torne um produtor. Criamos um programa chamado "producer.js" com o seguinte conteúdo:

```
1 const { Kafka } = require("kafkajs");
2 const usuario = process.argv[2];
3
4 run();
5
6 async function run() {
7   try {
8     const kafka = new Kafka({
9       "clientId": "myapp",
10      "brokers": ["192.168.15.5:9092"]
11    });
12    const producer = kafka.producer();
13    console.log("Conectando...");
14    await producer.connect();
```

```

15 console.log("Produtor Conectado!");
16
17 // A-M = Part0 ou N-Z = Part1
18 const particao = usuario[0] < "N" ? 0 : 1;
19
20 const result = await producer.send({
21   "topic": "Usuarios",
22   "messages": [
23     {
24       "value": usuario,
25       "partition": particao
26     }
27   ]
28 });
29 console.log('Mensagem Enviada: ${JSON.stringify(result)}');
30 await producer.disconnect();
31 console.log("Produtor Desconectado!");
32 } catch (ex) {
33   console.error('Alguma coisa aconteceu ${ex}')
34 } finally {
35   process.exit(0);
36 }
37 }

```

E rodamos este com o seguinte comando: `$ node producer.js nomequalquer`

A primeira letra do nome vai definir em qual partição esse usuário vai entrar entre as letras "A" e "M" será a partição 0 e de "N" a "Z" na 1. Rodamos esse comando algumas vezes adicionando nomes de usuários para a nossa lista e vamos finalmente tratar do consumidor.

Os consumidores são aplicativos que se alimentam de fluxos de dados de tópicos no Kafka Cluster. Um consumidor pode receber fluxo de registros de vários tópicos por meio de assinatura. A API de consumidor do Apache Kafka permite que um aplicativo se torne um consumidor. Criamos um programa chamado "consumer.js" com o seguinte conteúdo:

```

1 const { Kafka } = require("kafkajs");
2
3 run();
4
5 async function run() {
6   try {
7     const kafka = new Kafka({
8       "clientId": "myapp",
9       "brokers": ["192.168.15.5:9092"]
10    });
11    const consumer = kafka.consumer({ "groupId": "teste" });
12    console.log("Conectando...");
13    await consumer.connect();
14    console.log("Consumidor Conectado!");
15    consumer.subscribe({
16      "topic": "Usuarios",
17      "fromBeginning": true
18    })
19
20    await consumer.run({
21      "eachMessage": async result => {
22        console.log('Mensagem recebida ${result.message.value} da partição

```

```

    ${result.partition}');
23     }
24   });
25   console.log("Consumidor Rodando...");
26 } catch (ex) {
27   console.error('Alguma coisa aconteceu ${ex}');
28 }
29 }

```

Devemos observar que o comando `process.exit(0)` foi removido, assim como o comando para "desconectar", pois o consumidor deve ficar sempre rodando e recebendo as mensagens enviadas pelo produtor. Caso precise interrompê-lo usar o comando `CRTL+C`.

E rodamos este com o seguinte comando: `$ node consumer.js`

Todas as mensagens caem para esse consumidor, se abrirmos uma janela e executamos o produtor veremos que pouco importa se usamos a partição 0 ou 1 as mensagens caem para esse consumidor. Como mudar essa situação? simples abrimos mais uma janela e iniciamos um segundo consumidor. Agora o primeiro recebe mensagens somente da partição 0 enquanto que o segundo da 1.

3 Sistema em Java

Nada disso teria graça senão usássemos Java para demonstrar o poder do Kafka, para isso vamos utilizar o Spring Tool Suite no qual criamos um projeto *Maven Project*.

A grande vantagem de criarmos um projeto tipo Maven que podemos inserir todas as bibliotecas que necessitamos sem termos o menor trabalho, basta apenas buscamos estas do repositório central do Maven através de sua declaração no arquivo **pom.xml**:

```

1 <project xmlns="http://maven.apache.org/POM/4.0.0"
2   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
3   xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
4     http://maven.apache.org/xsd/maven-4.0.0.xsd">
5   <modelVersion>4.0.0</modelVersion>
6   <groupId>kafkaproj</groupId>
7   <artifactId>kafkaproj</artifactId>
8   <version>0.0.1-SNAPSHOT</version>
9   <name>Projeto em Kafka</name>
10  <description>Exemplo do projeto em Kafka</description>
11  <dependencies>
12    <dependency>
13      <groupId>org.apache.kafka</groupId>
14      <artifactId>kafka_2.13</artifactId>
15      <version>3.1.0</version>
16    </dependency>
17    <dependency>
18      <groupId>org.apache.kafka</groupId>
19      <artifactId>kafka-clients</artifactId>
20      <version>3.1.0</version>
21    </dependency>
22    <dependency>
23      <groupId>com.typesafe.scala-logging</groupId>
24      <artifactId>scala-logging-slf4j_2.11</artifactId>

```

```
24     <version>2.1.2</version>
25   </dependency>
26 </dependencies>
27 </project>
```

Assim de forma simples temos as bibliotecas do Kafka e do Log do Scala que é utilizada pela classe *ShutdownableThread* que usaremos para manter o consumidor ativo. Vamos começar com a classe que produz nossas mensagens, criamos uma classe chamada **ExProdutor** com a seguinte codificação:

```
1 package kafkaproj;
2
3 import java.util.Properties;
4 import java.util.concurrent.ExecutionException;
5
6 import org.apache.kafka.clients.producer.Callback;
7 import org.apache.kafka.clients.producer.KafkaProducer;
8 import org.apache.kafka.clients.producer.ProducerRecord;
9 import org.apache.kafka.clients.producer.RecordMetadata;
10
11 public class ExProdutor extends Thread {
12     private final KafkaProducer<Integer, String> producer;
13     private final String topic;
14     private final boolean isAsync;
15
16     public static final String KAFKA_SERVER_URL = "ENDERECO IP";
17     public static final int KAFKA_SERVER_PORT = 9092;
18     public static final String CLIENT_ID = "SimplesProdutor";
19
20     public ExProdutor(String topic, Boolean isAsync) {
21         Properties properties = new Properties();
22         properties.put("bootstrap.servers", KAFKA_SERVER_URL + ":" + KAFKA_SERVER_PORT);
23         properties.put("client.id", CLIENT_ID);
24         properties.put("key.serializer",
25             "org.apache.kafka.common.serialization.IntegerSerializer");
26         properties.put("value.serializer",
27             "org.apache.kafka.common.serialization.StringSerializer");
28         producer = new KafkaProducer<>(properties);
29         this.topic = topic;
30         this.isAsync = isAsync;
31     }
32
33     @Override
34     public void run() {
35         for (int messageNo = 1; messageNo < 100; messageNo++) {
36             String messageStr = "Esta é a mensagem no." + messageNo;
37             long startTime = System.currentTimeMillis();
38             if (isAsync) {
39                 producer.send(new ProducerRecord<>(topic, messageNo,
40                     messageStr), new DemoCallBack(startTime, messageNo, messageStr));
41             } else {
42                 try {
43                     producer.send(new ProducerRecord<>(topic, messageNo, messageStr)).get();
44                     System.out.println("Enviando mensagem: (" + messageNo + ", " + messageStr +
45                         ")");
46                 } catch (InterruptedException e) {
47                     e.printStackTrace();
48                 }
49             }
50         }
51     }
52 }
```

```

45         Thread.currentThread().interrupt();
46     } catch (ExecutionException e) {
47         e.printStackTrace();
48     }
49 }
50 }
51 }
52 }
53
54 class DemoCallBack implements Callback {
55
56     private final long startTime;
57     private final int key;
58     private final String message;
59
60     public DemoCallBack(long startTime, int key, String message) {
61         this.startTime = startTime;
62         this.key = key;
63         this.message = message;
64     }
65
66     public void onCompletion(RecordMetadata metadata, Exception exception) {
67         long elapsedTime = System.currentTimeMillis() - startTime;
68         if (metadata != null) {
69             System.out.println(
70                 "message(" + key + ", " + message + ") sent to partition(" +
71                 metadata.partition() + "), " +
72                 "offset(" + metadata.offset() + ") in " + elapsedTime + " ms");
73         } else {
74             exception.printStackTrace();
75         }
76     }
77 }

```

Para executá-la uma simples classe que contém o método `main()` chamada **GeraMensagens**:

```

1 package kafkaproj;
2
3 public class GeraMensagens {
4     public static final String TOPIC = "testTopic";
5
6     public static void main(String[] args) {
7         boolean isAsync = false;
8         ExProdutor producerThread = new ExProdutor(TOPIC, isAsync);
9         producerThread.start();
10    }
11 }

```

Para consumirmos as mensagens geradas, criamos uma classe chamada **ExConsumidor** com a seguinte codificação:

```

1 package kafkaproj;
2
3 import java.time.Duration;
4 import java.util.Collections;
5 import java.util.Properties;

```



```

6
7 import org.apache.kafka.clients.consumer.ConsumerConfig;
8 import org.apache.kafka.clients.consumer.ConsumerRecord;
9 import org.apache.kafka.clients.consumer.ConsumerRecords;
10 import org.apache.kafka.clients.consumer.KafkaConsumer;
11
12 import kafka.utils.ShutdownableThread;
13
14 public class ExConsumidor extends ShutdownableThread {
15     private final KafkaConsumer<Integer, String> consumer;
16     private final String topic;
17
18     public static final String KAFKA_SERVER_URL = "ENDERECO IP";
19     public static final int KAFKA_SERVER_PORT = 9092;
20     public static final String CLIENT_ID = "SimplesConsumidor";
21
22     public ExConsumidor(String topic) {
23         super("KafkaConsumidor", false);
24         Properties props = new Properties();
25         props.put(ConsumerConfig.BootstrapServersConfig, KAFKA_SERVER_URL + ":" +
26             KAFKA_SERVER_PORT);
27         props.put(ConsumerConfig.GroupIdConfig, CLIENT_ID);
28         props.put(ConsumerConfig.EnableAutoCommitConfig, "true");
29         props.put(ConsumerConfig.AutoCommitIntervalMsConfig, "1000");
30         props.put(ConsumerConfig.SessionTimeoutMsConfig, "30000");
31         props.put(ConsumerConfig.KeyDeserializerClassConfig,
32             "org.apache.kafka.common.serialization.IntegerDeserializer");
33         props.put(ConsumerConfig.ValueDeserializerClassConfig,
34             "org.apache.kafka.common.serialization.StringDeserializer");
35
36         consumer = new KafkaConsumer<>(props);
37         this.topic = topic;
38     }
39
40     @Override
41     public void doWork() {
42         consumer.subscribe(Collections.singletonList(this.topic));
43         ConsumerRecords<Integer, String> records = consumer.poll(Duration.ofMillis(100));
44         for (ConsumerRecord<Integer, String> registro : records) {
45             System.out.println("Recebendo mensagem: (" + registro.key() + ", " +
46                 registro.value() + ") at offset " + registro.offset());
47         }
48     }
49
50     @Override
51     public String name() {
52         return null;
53     }
54
55     @Override
56     public boolean isInterruptible() {
57         return false;
58     }
59 }

```

E finalmente para executá-la uma simples classe que contém o método `main()` chamada **Consume-Mensagens**:

```

1 package kafkaproj;
2
3 public class ConsumeMensagens {
4     public static void main(String[] args) {
5         ExConsumidor consumerThread = new ExConsumidor("testTopic");
6         consumerThread.start();
7     }
8 }

```

Cada vez são geradas 100 mensagens para o "testTopic", o padrão de mensagem pode ser qualquer coisa que desejamos, o que mais impressiona nesse sistema é a velocidade como essas mensagens são consumidas de forma quase instantânea. Não pense que é porquê estamos na mesma máquina, Kafka foi criado para ser um sistema extremamente performático.

Imagine então que através dele sistemas totalmente heterogêneos e que antes não trocavam informações agora podem fazê-lo sem o menor problema. E em linguagens totalmente diferentes. Mais fontes de informação podem ser obtidas em diversos sites que apresenta tutoriais completos sobre o Apache Hive como a Tutorials Point[3].

4 Conclusão

Kafka é ideal para justificar a necessidade que o mundo tem gerado um volume cada vez maior de informações. Um cenário comum é que estes dados, uma vez gerados, precisam ser transportados para diversas aplicações. Kafka é uma plataforma de streaming de dados distribuído que pode publicar, assinar, armazenar e processar streams de registros em tempo real. Projetado para lidar com fluxos de dados de várias fontes e entregá-los a vários consumidores. Em resumo: mover grandes quantidades de dados – não apenas do ponto A ao B, mas dos pontos A ao Z e a qualquer outro lugar que necessitemos e tudo ao mesmo tempo.

É uma alternativa natural a um sistema de mensagens corporativo tradicional. Começou como um simples sistema interno desenvolvido pelo LinkedIn para lidar com 1,4 trilhão de mensagens por dia, mas agora é uma solução de streaming de dados de código aberto com aplicação para uma variedade de necessidades corporativas.

Kafka pode lidar com milhões de dados por segundo, o que o torna adequado para desafios de Big Data. No entanto, Kafka também faz sentido para empresas que atualmente não estão lidando com cenários de dados tão extremos. Em muitos casos de uso de processamento de dados, como Internet das Coisas (IoT) e mídia social, os dados estão aumentando exponencialmente e podem sobrecarregar rapidamente um aplicativo que estamos criando com base no volume de dados atual. Em termos de processamento de dados, podemos considerar a escalabilidade, e isso significa planejar a maior proliferação dos dados.

Sou um entusiasta do mundo **Open Source** e novas tecnologias. Qual a diferença entre Livre e Open Source? Livre significa que esta apostila é gratuita e pode ser compartilhada a vontade. Open Source além de livre todos os arquivos que permitem a geração desta (chamados de arquivos fontes) devem ser disponibilizados para que qualquer pessoa possa modificar ao seu prazer, gerar novas, complementar ou fazer o que quiser. Os fontes da apostila (que foi produzida com o LaTeX) está disponibilizado no GitHub [6]. Veja ainda outros artigos que publico sobre tecnologia através do meu Blog Oficial [4].

Referências

- [1] Página do Apache Kafka
<https://kafka.apache.org/>
- [2] Página da Confluent
<https://www.confluent.io/what-is-apache-kafka/>
- [3] Tutorials Point sobre Kafka
https://www.tutorialspoint.com/apache_kafka/index.htm
- [4] Fernando Anselmo - Blog Oficial de Tecnologia
<http://www.fernandoanselmo.blogspot.com.br/>
- [5] Encontre essa e outras publicações em
<https://cetrex.academia.edu/FernandoAnselmo>
- [6] Repositório para os fontes da apostila
<https://github.com/fernandoans/publicacoes>