



- ### 1. Streams DSO

2. Processor API

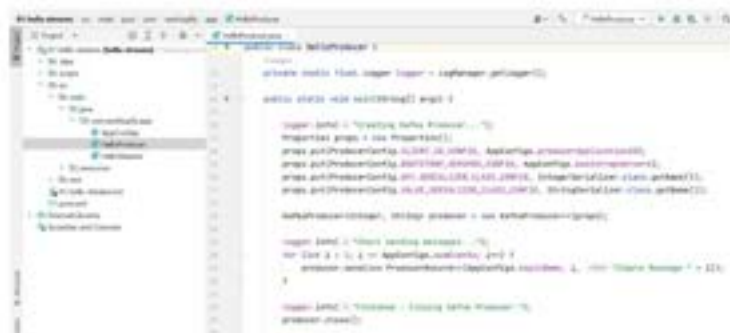
- Let's start with a simple Kafka Stream application. As our first example, we want to create a simple Kafka streams application to do the following things.

1. Connect to the Kafka cluster and start reading data from a given topic.
2. We do not want to get into the processing that data yet, but we just want to print the messages on the console.



Implementation

Step-1: Implement a AppConfigs.java class as shown below:





Apache Kafka

Step-3: Implement a Streaming class.

```

import java.util.Properties;
import org.apache.kafka.clients.consumer.ConsumerConfig;
import org.apache.kafka.clients.consumer.ConsumerRecord;
import org.apache.kafka.clients.consumer.ConsumerRecords;
import org.apache.kafka.clients.consumer.KafkaConsumer;
import org.apache.kafka.common.serialization.IntegerDeserializer;
import org.apache.kafka.common.serialization.StringDeserializer;

import java.util.concurrent.TimeUnit;

public class KafkaStreaming {

    public static void main(String[] args) {
        Properties props = new Properties();

        props.put(ConsumerConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092");
        props.put(ConsumerConfig.GROUP_ID_CONFIG, "kafka-streaming");
        props.put(ConsumerConfig.KEY_SERDE_CLASS_CONFIG, org.apache.kafka.common.serialization.IntegerDeserializer.class);
        props.put(ConsumerConfig.VALUE_SERDE_CLASS_CONFIG, org.apache.kafka.common.serialization.StringDeserializer.class);

        KafkaConsumer<Integer, String> consumer = new KafkaConsumer<>(props);
        consumer.subscribe(new String[] {"topic"});

        while (true) {
            ConsumerRecords<Integer, String> records = consumer.poll(100);

            for (ConsumerRecord<Integer, String> record : records) {
                System.out.printf("key: %d, value: %s\n", record.key(), record.value());
            }
        }
    }
}

```

At a high level, creating a Kafka Streams application is a four-step process.

Step-1: Line-17 to-21:

```

17 Properties props = new Properties();
18 props.put(StreamsConfig.APPLICATION_ID_CONFIG, "app1");
19 props.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092");
20 props.put(StreamsConfig.DEFAULT_KEY_SERDE_CLASS_CONFIG, Serdes.Integer().getClass());
21 props.put(StreamsConfig.DEFAULT_VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass());

```

- The first step is to create a Java Properties object and put the necessary configuration in the Properties object. Kafka streams API is highly configurable, and we use Java Properties object to specify those configurations. In the above example, we are setting four most basic settings.
 - ✓ APPLICATION_ID_CONFIG
 - ✓ BOOTSTRAP_SERVERS_CONFIG
 - ✓ DEFAULT_KEY_SERDE_CLASS_CONFIG
 - ✓ DEFAULT_VALUE_SERDE_CLASS_CONFIG

Every Kafka Streams application must be identified by a unique application id. If we are executing multiple instances of our Kafka Streams application, all instances of the same application must have the same application id.



Apache Kafka

The application id is used by Kafka Streams application in many ways to isolate resources for the same application and share workload among the instances of the same application.

The next configuration is the bootstrap server which we already know, A bootstrap server is a comma-separated list of host/port pairs that the streams application would use to establishing the initial connection to the Kafka cluster.

The next two configurations are the Key and the Value Serdes. Serdes is a Factory class that we will be using for creating a Serde (a combination of serializer and a matching deserializer).

In this example, we are setting an Integer Serde for the key and a String Serde for the value. we might be wondering, why do we need a Serde?

While creating producers, we required a serializer and consumers needed a deserializer. But why do we need Serde (a combination of both) for the Streams API?

This example is a data consuming stream application. However, **a typical Kafka Streams application internally creates a combination of consumer and producer.** Hence, we need a serializer as well as a deserializer for the internal producer/consumer to work correctly. Therefore, we set a default Serdes for the key and a value.

Step-2: Line-23-to-28:

```
23 StreamBuilder streamBuilder = new StreamBuilder();
24 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
25 kStream.foreach((k, v) -> System.out.println("key: " + k + " value: " + v));
26 //kStream.peek((k, v) -> System.out.println("key: " + k + " value: " + v));
27
28 Topology topology = streamBuilder.build();
```

Next comes the central part of our Kafka Streams application. In this part, we define our streams computational logic that we want to execute.

For our example, the computation logic is as straightforward as following steps.

1. Open a stream to a source topic: define a Kafka stream for a Kafka topic that can be used to read all the messages.
2. Process the stream: for each message, print the Key and the Value.
3. Create a Topology: bundle your business logic into a stream topology.

1. Opening a Stream

We want to use Kafka Streams DSL for defining the above computational logic. Most of the DSL APIs are available through StreamBuilder() class.



Apache Kafka

So, the first step is to create a StreamBuilder object. After creating a builder, we can open a Kafka Stream using the stream() method on the StreamBuilder. The stream() method takes a Kafka topic name and returns a KStream object.

Once we get a KStream object, we have successfully completed the first part of our computational logic i.e. Opening a stream to a source topic.

```
23 StreamBuilder streamBuilder = new StreamBuilder();
24 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
```

2. Process the stream

KStream class provides a bunch of methods for us to build our computational logic. We just want to navigate each message (key/value pair) and apply a println. So, the example is making use of the foreach() method on KStream to implement the second part of our computational logic i.e. Process the stream.

```
StreamBuilder streamBuilder = new StreamBuilder();
KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
kStream.foreach((k, v) -> System.out.println("key: " + k + " value: " + v));
```

Instead of using foreach() method, we can also use the peek() method to work with each record in the stream.

```
23 StreamBuilder streamBuilder = new StreamBuilder();
24 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
25 //kStream.foreach((k, v) -> System.out.println("key: " + k + " value: " + v));
26 kStream.peek((k, v) -> System.out.println("key: " + k + " value: " + v));
```

The foreach() and the peek() method takes a lambda expression on key/value pair, and the logic to work with each pair is implemented in the lambda body. That's all. That's what we wanted to do as a computational logic.

3. Create a Topology

The Kafka Streams computational logic is known as a Topology and is represented by the Topology class. So, whatever we defined as computational logic, we can get all that bundled into a Topology object by making a call to the build() method on the StreamBuilder object.



Apache Kafka

```
23 StreamBuilder streamBuilder = new StreamBuilder();
24 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
25 kStream.foreach((k, v) -> System.out.println("key: " + k + " value: " + v));
```

```

17 Topology topology = streamBuilder.build();
18 KafkaStreams streams = new KafkaStreams(topology, props);
19 logger.info("Starting stream.");
20 streams.start();

```

That's all. That's what we wanted to achieve in the Step-2.

- However, we implemented the computational logic by creating some intermediate objects such as KStream, and Topology.

A typical application may not be creating and holding these intermediate objects. Instead, we prefer to define all these operations as a chain of methods to avoid unnecessary code clutter as shown in the below code implemented as a chain of methods.

```

StreamsBuilder builder = new StreamsBuilder();
builder.stream(topicName)
    .foreach((k, v) -> System.out.println("key = " + k + " value "
    + " + v));
KafkaStreams streams = new KafkaStreams(builder.build(),
props);
streams.start();

```

Step 3 (Line 12 - 14):

Once we have the Properties and the Topology, we are ready to instantiate the KafkaStreams and start() it. That's all we do in a typical Kafka Streams application.

```

27 KafkaStreams streams = new KafkaStreams(topology, props);
28 logger.info("Starting stream.");
29 streams.start();

```

We can summarize a Kafka Streams application in just three steps.

Configure using Properties Create Topology using StreamBuilder Create and start KafkaStreams using Properties and Topology.

- We have seen a simple Kafka Streams application. We must have realized that the core of our Kafka streams application is the Topology of our application. A typical Kafka Streams application defines its computational logic through one or more processor topologies. The next important point to notice is that our Kafka Streams application does not Kafka streams is a client library for processing



Apache Kafka

and analyzing data stored in Kafka. Being a lightweight and straightforward library, it is the most convincing feature of Kafka Streams.

- We can easily embed Kafka streams library in our Java application and execute it on a single machine or take advantage of our preferred packaging and deployment frameworks like docker.
- When we want to scale our system, we can run multiple instances of the same application on multiple machines. Yet, we do not need to get into the complexities of load balancing.
- Kafka Streams transparently handles the load balancing of multiple instances of the same application by leveraging Kafka's parallelism model.
- The core functionality of Kafka Streams library is available in two flavors.
 - Streams DSL
 - Processor API

3. Streams DSL

- Streams DSL is a high-level API that offers common data transformation functions. Kafka streams DSL is a very well thought API set that allows us to handle most of the stream processing needs.
- DSL is the recommended way to create Kafka Streams application, and it should cover most of our needs and most of the use cases.

4. Processor API

- Processor APIs are the low-level API, and they provides us with more flexibility than the DSL. Using Processor API requires little extra manual work and code on the application developer side.
- However, they allow us to handle those one-off problems that may not be possible to solve with higher level abstraction.

Let's start with a simple Kafka Stream application. As our first example, we want to create a simple Kafka streams application to do the following things.

- Connect to the Kafka cluster and start reading data from a given topic.
- We do not want to get into the processing that data yet, but we just want to print the messages on the console.



Apache Kafka

Implementation

Step-1: Implement a AppConfigs.java class as shown below:

```

1 // AppConfigs.java
2
3 import org.apache.kafka.clients.producer.ProducerConfig;
4 import org.apache.kafka.clients.producer.Producer;
5 import org.apache.kafka.clients.producer.ProducerRecord;
6
7 public class AppConfigs {
8
9     // Kafka bootstrap servers
10    static String bootstrapServers = "localhost:9092,localhost:9093";
11
12    // Kafka producer group ID
13    static String groupId = "kafka-producer-group";
14
15    // Kafka topic
16    static String topicName = "test-topic";
17
18    // Kafka topic partition count
19    static int numPartitions = 1;
20
21 }

```

Step-2: Implement a Producer as shown below:

```

1 // Producer.java
2
3 import org.apache.kafka.clients.producer.ProducerConfig;
4 import org.apache.kafka.clients.producer.Producer;
5 import org.apache.kafka.clients.producer.ProducerRecord;
6
7 public class Producer {
8
9     // Kafka bootstrap servers
10    static String bootstrapServers = "localhost:9092,localhost:9093";
11
12    // Kafka producer group ID
13    static String groupId = "kafka-producer-group";
14
15    // Kafka topic
16    static String topicName = "test-topic";
17
18    // Kafka topic partition count
19    static int numPartitions = 1;
20
21    // Kafka producer
22    static Producer producer;
23
24    // Kafka producer
25    static Producer producer;
26
27    // Kafka producer
28    static Producer producer;
29
30    // Kafka producer
31    static Producer producer;
32
33    // Kafka producer
34    static Producer producer;
35
36    // Kafka producer
37    static Producer producer;
38
39    // Kafka producer
40    static Producer producer;
41
42    // Kafka producer
43    static Producer producer;
44
45    // Kafka producer
46    static Producer producer;
47
48    // Kafka producer
49    static Producer producer;
50
51    // Kafka producer
52    static Producer producer;
53
54    // Kafka producer
55    static Producer producer;
56
57    // Kafka producer
58    static Producer producer;
59
60    // Kafka producer
61    static Producer producer;
62
63    // Kafka producer
64    static Producer producer;
65
66    // Kafka producer
67    static Producer producer;
68
69    // Kafka producer
70    static Producer producer;
71
72    // Kafka producer
73    static Producer producer;
74
75    // Kafka producer
76    static Producer producer;
77
78    // Kafka producer
79    static Producer producer;
80
81    // Kafka producer
82    static Producer producer;
83
84    // Kafka producer
85    static Producer producer;
86
87    // Kafka producer
88    static Producer producer;
89
90    // Kafka producer
91    static Producer producer;
92
93    // Kafka producer
94    static Producer producer;
95
96    // Kafka producer
97    static Producer producer;
98
99    // Kafka producer
100   static Producer producer;
101
102 }

```



Apache Kafka

Step-3: Implement a Streaming class.

```

1 // Streaming.java
2
3 import org.apache.kafka.clients.producer.ProducerConfig;
4 import org.apache.kafka.clients.producer.Producer;
5 import org.apache.kafka.clients.producer.ProducerRecord;
6
7 public class Streaming {
8
9     // Kafka bootstrap servers
10    static String bootstrapServers = "localhost:9092,localhost:9093";
11
12    // Kafka producer group ID
13    static String groupId = "kafka-producer-group";
14
15    // Kafka topic
16    static String topicName = "test-topic";
17
18    // Kafka topic partition count
19    static int numPartitions = 1;
20
21    // Kafka producer
22    static Producer producer;
23
24    // Kafka producer
25    static Producer producer;
26
27    // Kafka producer
28    static Producer producer;
29
30    // Kafka producer
31    static Producer producer;
32
33    // Kafka producer
34    static Producer producer;
35
36    // Kafka producer
37    static Producer producer;
38
39    // Kafka producer
40    static Producer producer;
41
42    // Kafka producer
43    static Producer producer;
44
45    // Kafka producer
46    static Producer producer;
47
48    // Kafka producer
49    static Producer producer;
50
51    // Kafka producer
52    static Producer producer;
53
54    // Kafka producer
55    static Producer producer;
56
57    // Kafka producer
58    static Producer producer;
59
60    // Kafka producer
61    static Producer producer;
62
63    // Kafka producer
64    static Producer producer;
65
66    // Kafka producer
67    static Producer producer;
68
69    // Kafka producer
70    static Producer producer;
71
72    // Kafka producer
73    static Producer producer;
74
75    // Kafka producer
76    static Producer producer;
77
78    // Kafka producer
79    static Producer producer;
80
81    // Kafka producer
82    static Producer producer;
83
84    // Kafka producer
85    static Producer producer;
86
87    // Kafka producer
88    static Producer producer;
89
90    // Kafka producer
91    static Producer producer;
92
93    // Kafka producer
94    static Producer producer;
95
96    // Kafka producer
97    static Producer producer;
98
99    // Kafka producer
100   static Producer producer;
101
102 }

```

At a high level, creating a Kafka Streams application is a four-step process.

Step-1: Line 17 to 21

```

17 Properties props = new Properties();
18 props.put(StreamsConfig.APPLICATION_ID_CONFIG, "AppConfig.streaming");
19 props.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092,localhost:9093");
20 props.put(StreamsConfig.DEFAULT_KEY_SERIALIZER_CONFIG, Serdes.String().getClass());
21 props.put(StreamsConfig.DEFAULT_VALUE_SERIALIZER_CONFIG, Serdes.String().getClass());

```

➤ The first step is to create a Java Properties object and put the necessary configuration in the Properties object. Kafka streams API is highly configurable, and we use Java Properties object to

specify these configurations. In the above example, we are setting four most basic settings.

- ✓ APPLICATION_ID_CONFIG
- ✓ BOOTSTRAP_SERVERS_CONFIG
- ✓ DEFAULT_KEY_SERDE_CLASS_CONFIG
- ✓ DEFAULT_VALUE_SERDE_CLASS_CONFIG

Every Kafka Streams application must be identified by a unique application id. If we are executing multiple instances of our Kafka Streams application, all instances of the same application must have the same application id.



Apache Kafka

The application id is used by Kafka Streams application in many ways to isolate resources for the same application and share workload among the instances of the same application.

The next configuration is the bootstrap server which we already know. A bootstrap server is a comma-separated list of host/port pairs that the streams application would use to establish the initial connection to the Kafka cluster.

The next two configurations are the Key and the Value Serdes. Serdes is a Factory class that we will be using for creating a Serde (a combination of serializer and a matching deserializer).

In this example, we are setting an Integer Serde for the key and a String Serde for the value. We might be wondering, why do we need a Serde?

While creating producers, we required a serializer and consumers needed a deserializer. But why do we need Serde (a combination of both) for the Streams API?

This example is a data consuming stream application. However, **a typical Kafka Streams application internally creates a combination of consumer and producer**. Hence, we need a serializer as well as a deserializer for the internal producer/consumer to work correctly. Therefore, we set a default Serdes for the key and a value.

Step 2: Line 23 to 28

```
23 StreamBuilder streamBuilder = new StreamBuilder();
24 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
25 kStream.foreach((k, v) -> System.out.println("Key: " + k + " Value: " + v));
26 //stream.close(k, v) -> System.out.println("Key: " + k + " Value: " + v);
27
28 Topology topology = streamBuilder.build();
```

Next comes the central part of our Kafka Streams application. In this part, we define our streams computational logic that we want to execute.

For our example, the computation logic is as straightforward as following steps.

4. Open a stream to a source topic: define a Kafka stream for a Kafka topic that can be used to read all the messages.
5. Process the stream: for each message, print the Key and the Value.
6. Create a Topology: bundle your business logic into a stream topology.

4. Opening a Stream

We want to use Kafka Streams DSL for defining the above computational logic. Most of the DSL APIs are available through StreamBuilder() class.



Apache Kafka

So, the first step is to create a StreamBuilder object. After creating a builder, we can open a Kafka Stream using the stream() method on the StreamBuilder. The stream() method takes a Kafka topic name and returns a KStream object.

Once we get a KStream object, we have successfully completed the first part of our computational logic i.e. Opening a stream to a source topic.

```
14 StreamBuilder streamBuilder = new StreamBuilder();
15 KStream<Integer, String> kStream = streamBuilder.stream(AppConfigs.topicName);
```

5. Process the stream

KStream class provides a bunch of methods for us to build our computational logic. We just want to navigate each message (key/value pair) and apply a println. So, the example is making use of the

`foreach()` method on `KStream` to implement the second part of our computational logic i.e. Process the stream.

```
StreamsBuilder streamsBuilder = new StreamsBuilder();
KStream<Integer, String> kStream = streamsBuilder.stream(AppConfigs.topicName);
kStream.foreach((k, v) -> System.out.println("Key= " + k + " Value= " + v));
```

Instead of using `foreach()` method, we can also use the `peek()` method to work with each record in the stream.

```
25 StreamsBuilder streamsBuilder = new StreamsBuilder();
26 KStream<Integer, String> kStream = streamsBuilder.stream(AppConfigs.topicName);
27 //kStream.foreach((k, v) -> System.out.println("Key= " + k + " Value= " + v));
30 kStream.peek((k, v) -> System.out.println("Key= " + k + " Value= " + v));
```

The `foreach()` and the `peek()` method takes a lambda expression on key/value pair, and the logic to work with each pair is implemented in the lambda body. That's all. That's what we wanted to do as a computational logic.

4. Create a Topology

The Kafka Streams computational logic is known as a Topology and is represented by the `Topology` class. So, whatever we defined as computational logic, we can get all that bundled into a `Topology` object by making a call to the `build()` method on the `StreamsBuilder` object.



```
25 StreamsBuilder streamsBuilder = new StreamsBuilder();
26 KStream<Integer, String> kStream = streamsBuilder.stream(AppConfigs.topicName);
27 kStream.foreach((k, v) -> System.out.println("Key= " + k + " Value= " + v));
28 Topology topology = streamsBuilder.build();
29 KafkaStreams streams = new KafkaStreams(topology, props);
30 logger.info("Starting stream.");
31 streams.start();
```

That's all. That's what we wanted to achieve in the Step-2.

However, we implemented the computational logic by creating some intermediate objects such as `KStream`, and `Topology`.

A typical application may not be creating and holding these intermediate objects. Instead, we prefer to define all these operations as a chain of methods to avoid unnecessary code clutter as shown in the below code implemented as a chain of methods.

```
StreamsBuilder builder = new StreamsBuilder();
builder.stream(topicName)
    .foreach((k, v) -> System.out.println("Key= " + k + " Value= " + v));
KafkaStreams streams = new KafkaStreams(builder.build(), props);
streams.start();
```

Step 3 (Line 12 - 14):

Once we have the Properties and the Topology, we are ready to instantiate the `KafkaStreams` and start it. That's all we do in a typical Kafka Streams application.

```
29 KafkaStreams streams = new KafkaStreams(topology, props);
30 logger.info("Starting stream.");
31 streams.start();
```

We can summarize a Kafka Streams application in just three steps.

1. Configure using Properties
2. Create Topology using `StreamsBuilder`
3. Create and start `KafkaStreams` using Properties and Topology.

We have seen a simple Kafka Streams application. We must have realized that the core of our Kafka streams application is the Topology of our application.



A typical Kafka Streams application defines its computational logic through one or more processor topologies.

The next important point to notice is that our Kafka Streams application does not execute inside the Kafka cluster or the Kafka Broker. Instead, it runs in a separate JVM instance, or in a different machine entirely.

For reading and writing data to the Kafka cluster, it creates a remote connection using the bootstrap servers configuration. The simplicity to run our Kafka Streams application allows us to use our favourite packaging and deployment infrastructure such as Docker and Kubernetes container platforms.

Finally, we can run a single instance of our application or alternatively, we can start multiple instances of the same application to scale our production deployment horizontally.

Kafka Streams API is designed to help us elastically scale and parallelize our application by simply starting more and more instances.

These additional instances not only contribute to share workload but also provide automatic fault-tolerance.

Streams Topology

A processor topology or simply topology defines a step-by-step computational logic of the data processing that needs to be performed by a stream processing application. We can easily represent these steps using a DAG (Directed Acyclic Graph).



Hello Streams topology is reasonably straightforward. However, it helps us to visualize and uncover several concepts.

Source Processor



Apache Kafka

Every Kafka streams application or the topology starts by subscribing to some Kafka topics, and that's fair because we need to consume a data stream to begin our processing.

So, our topology starts by creating a source processor. Creating a source processor is straightforward. The below figure shows the code snippet from the Hello Streams example that creates a source processor.

```
KStream<String, String> kStream = builder.stream(topicName);
```

Source Processor

As shown above, a source processor internally implements a Kafka Consumer to consume data from one or multiple Kafka topics, and it produces a KStream object that will be used by its down-stream processors as an input stream.

Once we have a source processor and our first KStream object, we are ready to add more processors in the topology. Adding a new processor node is as simple as calling a transformation method on the KStream object.

For example, we take kStream that was returned by the source node and make a call to the kStream.foreach() method. This operation simply adds a foreach() processor node to the topology.



TechHubVault

