

Confluent Stream Processing using Apache Kafka® Streams & ksqlDB

Exercise Book

Version 6.0.0-v1.0.0



CONFLUENT

tejaswin.enugunta@walgreens.com

Table of Contents

Lab 01 Introduction	1
Copyright & Trademarks	1
a. The Lab Environment & Sample Solution	2
b. Continued Learning After Class	6
Lab 02 Motivation and Concepts for Streams	7
a. Streams Basics	7
Lab 03 Kafka Streams Architecture	16
Lab 04 Writing Kafka Streams Applications	17
a. Anatomy of a Kafka Streams App	17
b. Working with JSON	27
c. Joining Two Streams	32
d. Using the Processor API	38
Lab 05 Testing Kafka Streams Applications	43
a. Building Unit Tests	43
b. Integration Tests using Embedded Kafka	49
Lab 06 Introduction to ksqlDB	53
a. Introduction to ksqlDB	53
Lab 07 Using ksqlDB	64
a. Using ksqlDB	64
b. Using the ksqlDB REST API	76
Lab 08 Scaling a Kafka Streams Application	82
a. Scaling a Kafka Streams Application	82
Lab 09 Securing a Kafka Streams Application	91
a. Securing a Kafka Streams Application	91
Lab 10 Monitoring Kafka Streams Applications	100
a. Getting Metrics from a Kafka Streams Application	100
b. Using JConsole to monitor a Streams App	104
c. Monitoring a Kafka Streams App in Confluent Control Center	109
Appendix A: Running All Labs with Docker	114
Running Labs in Docker for Desktop	114
Running the Exercise Applications	115

Lab 01 Introduction

Copyright & Trademarks

Copyright © Confluent, Inc. 2014-2020. [Privacy Policy](#) | [Terms & Conditions](#).

Apache, Apache Kafka, Kafka and the Kafka logo are trademarks of the
[Apache Software Foundation](#)

tejaswin.renugunta@walgreens.com

Alternative Lab Environments

As an alternative you can also download the VM to your laptop and run it in VirtualBox. Make sure you have the newest version of VirtualBox installed. Download the VM from this link:

- <https://confluent-training-images-us-east-1.s3.amazonaws.com/training-ubuntu-18-04-may2020.ova>

If you have installed Docker for Desktop on your Mac or Windows 10 Pro machine then you can run the labs there. But please note that your trainer might not be able to troubleshoot any potential problems if you are running the labs locally. If you choose to do this, follow the instructions at → [Running Labs in Docker for Desktop](#).

a. The Lab Environment & Sample Solution

Welcome to your lab environment! You are connected as user **training**, password **training**.

tejaswin.renugunta@walgreens.com

If you haven't already done so, you should open the **Exercise Guide** that is located on the lab virtual machine. To do so, open the **Confluent Training Exercises** folder that is located on the lab virtual machine desktop. Then double-click the shortcut that is in the folder to open the **Exercise Guide**.



Copy and paste works best if you copy from the Exercise Guide on your lab virtual machine.

- Standard Ubuntu keyboard shortcuts will work: **Ctrl+C** → Copy, **Ctrl+V** → Paste
- In a Terminal window: **Ctrl+Shift+C** → Copy, **Ctrl+Shift+V** → Paste.

If you find these keyboard shortcuts are not working you can use the right-click context menu for copy and paste.

To make things easier for you and your trainer, we encourage using the preferred folder structure for our exercises. We also have working sample solutions for all exercises of this book, which you can use as a guideline in case your sample is not working or you have a different approach than the one suggested in the book and want to compare it with the default solution.

1. Open a terminal window.
2. Clone the **GitHub** repository with the code into the folder **~/confluent-streams**:

Bash:

```
$ git clone --depth 1 --branch 6.0.0-v1.0.0 \
  https://github.com/confluentinc/training-ksql-and-streams-src.git \
  ~/confluent-streams
```

Powershell:

```
PS> git clone --depth 1 --branch 6.0.0-v1.0.0 `
  https://github.com/confluentinc/training-ksql-and-streams-src.git `
  ~/confluent-streams
```

3. **Windows users only:** Create a folder **.gradle** in your home folder:

```
PS> mkdir ~/.gradle
```



If you chose to select another folder for the labs then note that many of our samples assume that the lab folder is **~/confluent-streams**. You will have to adjust all commands that include this default lab folder to fit your specific environment.

4. Navigate to the **confluent-streams** folder:

```
$ cd ~/confluent-streams
```

5. Start the Kafka cluster:

```
$ docker-compose up -d zookeeper kafka
```

You should see something similar to this:

```
Creating network "confluent-streams_kafka-net" with the default
driver
Creating kafka          ... done
Creating zookeeper      ... done
```



In the first steps of each exercise, you may launch the containers needed for the exercise with a **docker-compose up** command. Simply typing **docker-compose up -d** will start all of the containers defined in the docker-compose.yml file. You can start fewer containers by specifying only those you want to run, for example: **docker-compose up -d zookeeper kafka**.

The majority of the exercises use the docker-compose.yml file in the ~/confluent-streams directory. The **docker-compose up** command will search up the directory hierarchy until it finds a docker-compose.yml file, so the one in the confluent-streams directory will usually be used. The exception is the docker-compose.yml file used in the security exercise as this has additional security settings. See the comments at the beginning of Lab 09 Securing a Kafka Streams Application.

If at any time you want to get your environment back to a clean state use **docker-compose down** to end all of your containers. Then return to your last **docker-compose up** to get back to the beginning of an exercise.

Exercises do not need to be completed in order. You can start from the beginning of any exercise at any time.

If you want to completely clear out your docker environment use the script on the VM at **~/docker-nuke.sh**. The nuke script will forcefully end all of your running docker containers.

6. Monitor the cluster with:

```
$ docker-compose ps
```

Name	Command	State	Ports
-----	-----	-----	-----
kafka	/etc/confluent/docker/run	Up	0.0.0.0:9092-
>9092/tcp			
zookeeper	/etc/confluent/docker/run	Up	2181/tcp, 2888/tcp,
3888/tcp			

All services should have **State** equal to **Up**.

7. You can also observe the stats of Docker on your VM:

```
$ docker stats
```

CONTAINER ID	NAME	CPU %	MEM USAGE / LIMIT	MEM %	NET
I/O	BLOCK I/O	PIDS			
ab9c97077e94	zookeeper	0.14%	88.14MiB / 9.737GiB	0.88%	
106kB / 130kB	0B / 0B	48			
ff47bece9e4f	kafka	1.17%	421.8MiB / 9.737GiB	4.23%	
646kB / 522kB	0B / 0B	77			

Cleanup

1. Press **Ctrl+C** to exit the Docker statistics.
2. Shut down your Kafka cluster with the **docker-compose down -v** command.

b. Continued Learning After Class

Once the course ends, the VM in Content Raven will terminate and you will no longer have access to it. However, you can still download the VM onto your own machine or use Docker locally to revisit these materials. We encourage you to bring up your own test environment, explore configuration files, inspect scripts, and perform tests. Here are some activities we encourage to reinforce your learning:

- Revisit the exercises in this manual
- Summarize and discuss the student handbook with your peers
- Consult the README in this public repository for more resources and your own development playground: <https://github.com/confluentinc/training-ksql-and-streams-src>

Lab 02 Motivation and Concepts for Streams

a. Streams Basics

During the instruction module, you learned basic concepts about stream processing. In this lab exercise, you will work with streams as objects in Java (or optionally, Scala).

This lab contains the same exercise once using a Scala REPL (Read–Eval–Print Loop) and once using a Java 10 REPL. Please select the version you prefer. We're starting with **Java** but if you prefer **Scala** then you can find it here: [Optional: Playing with Scala](#)

Because this exercise does not use the kafka cluster, you do not need to run a docker-compose up command.

Playing with Java

1. Start the Java 11 REPL with the following command:

```
$ jshell
```

You should be greeted by a message similar to this:

```
May 30, 2018 8:19:38 AM java.util.prefs.FileSystemPreferences$1 run
INFO: Created user preferences directory.
| Welcome to JShell -- Version 11.0.4
| For an introduction type: /help intro

jshell>
```

2. Define a **stream** of numbers in Java.:

```
jshell> var numbers = Stream.of(1,2,3,4,5)
```

The output produced by the Java REPL should look like this:

```
numbers ==> java.util.stream.ReferencePipeline$Head@6acdbdf5
```

Now this is not very helpful if we want to see what's generated. But a stream is the base of what we want to do in the following.

3. If ever we want to see the output of an operation on a Java stream we can collect the elements into a list, e.g.:

```
jshell> var output = Stream.of(1,2,3,4,5).toArray()
```

Note how we're using the `.toArray()` function at the end of the above expression. This results in the following output:

```
output ==> Object[5] { 1, 2, 3, 4, 5 }
```



Java streams **cannot be reused**. As soon as you call any terminal operation the stream is closed. This is contrary to what we get from Kafka and ksqldb. There, streams can be reused over and over again.

4. Filter this list for odd numbers only:

```
jshell> var odds = Stream.of(1,2,3,4,5).  
    filter(x -> x%2==1).  
    toArray()
```

which yields this output:

```
odds ==> Object[3] { 1, 3, 5 }
```

As we can see, the output is the odd numbers only. But the important thing to note, that the output of the `filter` function is still a **stream** and we need to convert it to an array with `toArray()` to see a user friendly output.

5. Let's use the `map` function to change each value of the stream:

```
jshell> var squares = Stream.of(1,2,3,4,5).  
    map(x -> x*x).  
    toArray()
```

resulting in this output:

```
squares ==> Object[5] { 1, 4, 9, 16, 25 }
```

Once again the output of the **map** function is a stream too.

6. We can of course combine multiple stream functions:

```
jshell> var oddSquares = Stream.of(1,2,3,4,5).  
    filter(x -> x%2==1).  
    map(x -> x*x).  
    toArray()
```

```
oddSquares ==> Object[3] { 1, 9, 25 }
```

7. Now we can do even more sophisticated mappings like this:

```
jshell> var nestedStream = Stream.of(1,2,3,4,5).  
    map(x -> Stream.of(x-1,x,x+1))
```

8. The previous step gives us a **nested stream**. We can use the function **flatMap** to flatten that nested stream and make it a simple stream:

```
jshell> var flattened = nestedStream.  
    flatMap(x -> x).  
    toArray()
```

which results in (after collection):

```
flattened ==> Object[15] { 0, 1, 2, 1, 2, 3, 2, 3, 4, 3, 4, 5, 4, 5,  
6 }
```

Why would that be interesting? We will see in a moment where it makes sense...

9. Define a **list** of words:

```
jshell> var words = List.of("Kafka", "is", "cool", "Kafka",  
"Streams", "and", "Kafka", "Topics", "Streams")
```

10. Now let's do some aggregation. First let's count the elements in the stream:

```
jshell> var count = words.  
    stream().  
    count()
```

which gives:

```
count ==> 9
```



Note how the last command didn't return a stream but a single value. So, the input format (stream) is different than the output format (single value).

11. Now let's count the occurrence of each distinct word by grouping by word:

```
jshell> var wordCount = words.  
    stream().  
    collect(Collectors.groupingBy(  
        Function.identity(), Collectors.counting()  
    ))
```

this results in a **map** as follows:

```
wordCount ==> {Topics=1, and=1, cool=1, Kafka=3, is=1, Streams=2}
```



The groupingBy function produced a **Map** out of a stream. A Map is a set of (key,value) pairs.

12. Let's count distinct words in a list of lines with sentences. First we define a list of lines:

```
jshell> var lines = List.of("Kafka is a really cool technology",  
    "Many enterprises use Kafka and Kafka Streams",  
    "I want to hear more about Kafka Streams and ksqldb",  
    "ksqldb is for Kafka what SQL is for databases")
```

Then we split all lines into words and flatten the resulting nested list and count the words:

```
jshell> var wordCount = lines.  
    stream().  
    flatMap(x -> Stream.of(x.split(" "))).  
    collect(Collectors.groupingBy(  
        Function.identity(), Collectors.counting()  
    ))
```

the output of the above looks like this:

```
wordCount ==> {ksqlDB=2, databases=1, a=1, more=1, use=1, want=1, ...  
Many=1, hear=1, Streams=2}
```

The formula is a bit convoluted due to the fact that the result of the split function is an array and needs to be converted to a stream prior to be used by the **flatMap** function via **Stream.of(...)**.

Cleaning up

1. Exit the Java REPL by typing **Ctrl+D**

Optional: Playing with Scala

1. Open a terminal window and install Scala, then run Scala interactively. If you are prompted for a password for the training account enter training as the password.

```
$ sudo apt install -y scala ; scala
```

You will see an output that ends with this:

```
Welcome to Scala 2.11.12 (OpenJDK 64-Bit Server VM, Java 11.0.4).  
Type in expressions for evaluation. Or try :help.  
  
scala>
```

2. Define a list of numbers in Scala. This list simulates a **stream**:

```
scala> val numbers = List(1, 2, 3, 4, 5)
```

The output produced by the Scala REPL should look like this:

```
list: List[Int] = List(1, 2, 3, 4, 5)
```

3. Filter this list for odd numbers only:

```
scala> val odds = numbers.filter(x => x%2==1)
```

which yields this output:

```
odds: List[Int] = List(1, 3, 5)
```

As we can see, the output is the odd numbers only. But the important thing to note, that it is still a list.

4. Let's use the **map** function to change each value of the list:

```
scala> val squares = numbers.map(x => x*x)
```

resulting in this output:

```
squares: List[Int] = List(1, 4, 9, 16, 25)
```

Once again the output is a list too.

5. Now we can do even more sophisticated mappings like this:

```
scala> val nestedList = numbers.map(x => List(x-1,x,x+1))
```

with output:

```
nestedList: List[List[Int]] = List(List(0, 1, 2), List(1, 2, 3),  
List(2, 3, 4), List(3, 4, 5), List(4, 5, 6))
```

6. The previous step delivered us a nested list. We can use the function **flatMap** to flatten that nested list and make it a simple list:

```
scala> val flattened = nestedList.flatMap(x => x)
flattened: List[Int] = List(0, 1, 2, 1, 2, 3, 2, 3, 4, 3, 4, 5, 4, 5, 6)
```

Why would that be interesting? We will see in a moment where it makes sense...

7. Define a list of words:

```
scala> val words = List("Kafka", "is", "cool", "Kafka", "Streams",
"and", "Kafka", "Topics", "Streams")
words: List[String] = List(Kafka, is,...
```

8. Now let's do some aggregation. First let's count the elements in the list:

```
scala> val count = words.count(x => true)
count: Int = 9
```



Note how the last command didn't return a list but a single value. So, the input format (list) is different than the output format (single value).

9. Now let's count the occurrence of each distinct word by grouping by word:

```
scala> val wordCount = words.groupBy(x => x).mapValues(x => x.size)
```

resulting in:

```
wordCount: scala.collection.immutable.Map[String,Int] = Map(is -> 1,
Streams -> 2, cool -> 1, Kafka -> 3, Topics -> 1, and -> 1)
```



The `groupBy` function produced a `Map` out of a `List`. A `Map` is a set of (key,value) pairs. With the **mapValues** function we then access the **value** part of each pair.

10. Let's count distinct words in a list of lines with sentences. First we define a list of lines:

```
scala> val lines = List("Kafka is a really cool technology",  
    "Many enterprises use Kafka and Kafka Streams",  
    "I want to hear more about Kafka Streams and ksqlDB",  
    "ksqlDB is for Kafka what SQL is for databases")
```

Then we split all lines into words and flatten the resulting nested list and count the words:

```
scala> val wordCount = lines.flatMap(x => x.split(" ")).groupBy(word  
=> word).mapValues(x => x.size)
```

The last result is once again a Map where the key is the word and the value is the number of occurrences of this word in all sentences:

```
wordCount: scala.collection.immutable.Map[String,Int] = Map(for -> 2,  
is -> 3, Streams -> 2, want -> 1, what -> 1, a -> 1, I -> 1, hear ->  
1, to -> 1, technology -> 1, enterprises -> 1, cool -> 1, ksqlDB ->  
2, Kafka -> 5, use -> 1, Many -> 1, databases -> 1, more -> 1, really  
-> 1, about -> 1, and -> 2, SQL -> 1)
```

Joining two streams



There is no such thing as a **join** function in Scala but we can use the Scala **zip** function to "simulate" joining two lists/streams.

1. Create a list of names:

```
scala> val names = List("Ann", "Sue", "Rolf", "Jack", "Tammy")
```

2. Create a list of colors:

```
scala> val colors = List("yellow", "blue", "green", "black", "white",  
    "pink")
```

3. Join or merge the lists:


```
scala> val tuples = names zip colors
```

which yields this:

```
tuples: List[(String, String)] = List((Ann,yellow), (Sue,blue),  
  (Rolf,green), (Jack,black), (Tammy,white))
```

Note how the extra color "pink" is ignored since there is nothing to match in the names list.

Cleaning up

1. Exit Scala by pressing **Ctrl+D**.

Conclusion

In this lab, we used Scala or Java to learn and apply the basics of streams and stream operations. We have played with operations that convert a stream into another stream, as well as with operations that convert a stream into a table. The latter operations are called aggregation functions. Examples of aggregation functions include count, group, sum, etc.

Note that all these stream operations in JVM based languages do not scale beyond a single instance. We will find that streaming applications powered by Kafka have nearly unlimited scalability.



STOP HERE. THIS IS THE END OF THE EXERCISE.

Lab 03 Kafka Streams Architecture



Module 03 Kafka Streams Architecture does not have an associated exercise.

Please continue to Lab 04 Writing Kafka Streams Applications.

tejaswin.renugunta@walgreens.com

Lab 04 Writing Kafka Streams Applications

This lab contains 4 exercises:

- Anatomy of a Kafka Streams App
- Working with JSON
- Joining Two Streams
- Using the Processor API

a. Anatomy of a Kafka Streams App

In this exercise, you will create a Kafka Streams application and deploy it using Gradle (or, optionally, Maven). The purpose of this exercise is to illustrate the structure of Kafka Streams application code and the routine of deploying code with build tools.

The application itself reads data from a topic whose keys are integers and whose values are sentence strings. The application transforms the input so that the strings are lower-case and output to a new topic.

Prerequisites

Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)

Preparing the Kafka Cluster

To be able to run Kafka Streams applications we need a working Kafka cluster. We will run one consisting of a single broker and zookeeper.

1. Locate the file **docker-compose.yml** in the ~/confluent-streams directory.



This **docker-compose.yml** file will be used to run a simple Kafka cluster as a backend for our Kafka Streams applications. If you're curious, please open the file and analyze its contents.

2. Navigate to the folder **~/confluent-streams/labs/streams-writing**:

```
$ cd ~/confluent-streams/labs/streams-writing
```

3. From within the **streams-writing** folder run the cluster with the following command:

```
$ docker-compose up -d zookeeper kafka
Creating network "confluent-streams_kafka-net" with the default
driver
Creating kafka          ... done
Creating zookeeper      ... done
```

4. Double check that the cluster is up and running:

```
$ docker-compose ps
```

you should see something similar to this:

Ports	Name	Command	State

kafka		/etc/confluent/docker/run	Up
0.0.0.0:9092->9092/tcp			
zookeeper		/etc/confluent/docker/run	Up
2181/tcp, 2888/tcp, 3888/tcp			

Make sure all services have **State=Up**.

5. Create an input topic called **lines-topic** in Kafka:

```
$ kafka-topics \  
  --create \  
  --bootstrap-server kafka:9092 \  
  --replication-factor 1 \  
  --partitions 1 \  
  --topic lines-topic
```

6. Create an output topic called **lines-lower-topic** in Kafka:

```
$ kafka-topics \  
  --create \  
  --bootstrap-server kafka:9092 \  
  --replication-factor 1 \  
  --partitions 1 \  
  --topic lines-lower-topic  
Created topic "lines-lower-topic".
```

7. Let's now check what topics are in Kafka using this command:

```
$ kafka-topics \  
  --bootstrap-server kafka:9092 \  
  --list
```

We should see something like this:

```
__confluent.support.metrics  
lines-lower-topic  
lines-topic
```



If you need to delete a topic, say the one with name **<topic name>**, (e.g. to start over) you can use this command:

```
$ kafka-topics \  
  --bootstrap-server kafka:9092 \  
  --delete \  
  --topic <topic name>
```

Now we are ready to create, build and run our first **Kafka Streams** application. We will first build and run the application using Java in conjunction with **Gradle**. Optionally, we will also see how to build and run the application with **Maven**.

Authoring the Kafka Streams Application using Java & Gradle

1. Navigate to the folder **gradle-sample** within the **streams-writing** folder, and launch VS Code:



Be certain to include the period in the **code .** command below. That indicates starting VS Code in the current directory - otherwise some references may not correctly resolve.

```
$ cd ~/confluent-streams/labs/streams-writing/gradle-sample
$ code .
```

2. In this folder locate the **build.gradle** file and open it to analyze its content.



The external dependencies for a simple **Kafka Streams** app are the **kafka-clients** and **kafka-streams** libraries. We also add the **slf4j-log4j12** dependency for logging purposes.

3. In VS Code, open the file **MapSample.java** in the subfolder **src/main/java/streams** and inspect its content. From here, you can challenge yourself to complete the TODOs, or you can move forward to see a step-by-step walkthrough of the code. If you decide to challenge yourself, you can always peek at the corresponding subfolder in **~/confluent-streams/solutions/** if you get stuck.
4. Now we start to add actual **Kafka Streams** application logic. We will start with the configuration part. Please add this code snippet to the **main** method of the class (right after the initial **println**):

```
Properties settings = new Properties();
settings.put(StreamsConfig.APPLICATION_ID_CONFIG, "map-sample-
v0.1.0");
settings.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "kafka:9092");
```

We're providing an ID to our application and tell it where to find the Kafka cluster. This is the minimal configuration needed!

5. Next we will define the topology of our **Kafka Streams** application. Add this snippet right after the configuration part:

```
final Serde<String> stringSerde = Serdes.String();
StreamsBuilder builder = new StreamsBuilder();
KStream<String, String> lines = builder
    .stream("lines-topic", Consumed.with(stringSerde, stringSerde));
KStream<String, String> transformed = lines
    .mapValues(value -> value.toLowerCase());
transformed.to("lines-lower-topic", Produced.with(stringSerde,
stringSerde));
Topology topology = builder.build();
```

We're defining a builder for the topology, use it to create a **KStream** from the Kafka topic **line-topic** using **String** Serdes for both key and value of the messages. Then we're using a **mapValues** function on the **KStream** to convert the **value** into all lower case. Finally we're writing the result into the Kafka topic **line-lower-topic**. Then we build the topology.

6. With the settings and the topology at hand we can now create the **Streams** app:

```
KafkaStreams streams = new KafkaStreams(topology, settings);
```

Our application will now start consuming data from the input topic, transform it and write it to the output topic.

7. To have our application terminate in an orderly way when requested without leaving any resource leaks behind we add a **shutdown hook** at the end of the main method:

```
final CountdownLatch latch = new CountdownLatch(1);
Runtime.getRuntime().addShutdownHook(new Thread(() -> {
    System.out.println("### Stopping Map Sample Application ###");
    streams.close();
    latch.countDown();
}));

try{
    streams.start();
    latch.await();
} catch (final Throwable e) {
    System.exit(1);
}
System.exit(0);
```

This **Shutdown Hook** will be executed when the application receives a **SIG_TERM** signal. The **CountDownLatch** is used as a best practice to avoid rare cases of deadlock. Notice the use of the **await()** method after the application starts.

- Within the project folder **gradle-sample** locate the folder for the Java resources **src/main/resources**. In this folder we have a file **log4j.properties**. This file is used to configure the logger for our Kafka Streams app.



Use **INFO** instead of **WARN** for the **rootLogger** if you want to be more verbose in the logs.

And that's all we need. Our first **Kafka Streams** app is ready to go! Use **Run → Start Debugging** to run your code.



The first time you run the debugger it may take extra time while resources are downloaded.

```
PROBLEMS 2 OUTPUT DEBUG CONSOLE TERMINAL C#
Installing C# dependencies...
Platform: linux, x86_64, name=ubuntu, version=18.04

Downloading package 'OmniSharp for Linux (x64)' (30900 KB)..... Done!
Installing package 'OmniSharp for Linux (x64)'

Downloading package '.NET Core Debugger (linux / x64)' (58664 KB)..... Done!
Installing package '.NET Core Debugger (linux / x64)'

Finished
|
```

In the VS Code DEBUG CONSOLE tab you should see something like this:

```
*** Starting Map Sample Application ***
2018-07-04 13:47:36 WARN ConsumerConfig:287 - The configuration
'admin.retries' was supplied but isn't a known config.
```

At this time you can safely ignore the WARN log item.



At this time nothing will happen since we did not yet produce any data in the **lines-topic** in Kafka. Let your Kafka Streams app remain running in the VS Code debugger. This will be our next task **after** we have shown how to build the same app using Maven instead of Gradle. If you want to skip the **Maven** part then go to [Producing some Input Data](#).

Optional: Build and Run the application with Java & Maven

Here we're basically showing the same steps as in the previous section, except we will now use **Maven** instead of **Gradle**.

1. Navigate to the **maven-sample** folder, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-writing/maven-sample
$ code .
```

2. Locate the **pom.xml** file, open it and analyze its content.



this is the minimum Maven file required to compile, package and run a Java application that has external dependencies on the 2 required libraries **kafka-clients** and **kafka-streams**. It also uses the **slf4j-log4j12** library for logging.

3. Locate the file **MapSample.java** in the subfolder **src/main/java/streams** and open it. Double check that it looks similar to the one you created in the Gradle example.
4. Note the same file **log4j.properties** in the folder **src/main/resources** as in the Gradle example.



Use **INFO** instead of **WARN** for the **rootLogger** if you want to be more verbose in the logs.

5. Solution code has been provided in this exercise, so our first **Kafka Streams** app is ready to go! Use **Run** → **Start Debugging** to run your code.

You should see something like this:

```
*** Starting Map Sample Application ***
2018-07-04 13:53:12 WARN ConsumerConfig:287 - The configuration
'admin.retries' was supplied but isn't a known config.
```

At this time you can safely ignore the WARN log item.



At this time nothing will happen since we did not yet produce any data in the **lines-topic** in Kafka. Let's do this next. Let your Kafka Streams app remain running in the VS Code debugger.

Producing some Input Data

We're going to use the **kafka-console-producer** tool to create some input data in the topic **lines-topic**.

1. Open a new terminal window and navigate to the **streams-writing** folder:

```
$ cd ~/confluent-streams/labs/streams-writing
```

2. From a terminal window execute this command:

```
$ cat << EOF | kafka-console-producer \
  --bootstrap-server kafka:9092 \
  --property "parse.key=true" \
  --property "key.separator=:" \
  --topic lines-topic
1:"Kafka powers the Confluent Streaming Platform"
2:"Events are stored in Kafka"
3:"Confluent contributes to Kafka"
EOF
```

This writes 3 entries into the topic called **lines-topic** using **String** serializers for both key and value.

Reading the Transformed Messages

Here we're using the **kafka-console-consumer** tool to read from the output topic. We're again using String deserializers for key and value.

1. In the same terminal window, run:

```
$ kafka-console-consumer \  
  --bootstrap-server kafka:9092 \  
  --from-beginning \  
  --topic lines-lower-topic
```

You should see this:

```
"kafka powers the confluent streaming platform"  
"events are stored in kafka"  
"confluent contributes to kafka"
```



You might need to be a bit patient until the messages appear. It can take a few seconds to a minute or so depending on the performance of your computer...

Cleanup

1. Terminate the Kafka console consumer by pressing **Ctrl+C**.
2. Terminate your **Kafka Streams** application with **Run** → **Stop Debugging**. If you used the **./gradlew run** command instead, you can terminate with **Ctrl+C**.
3. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

We have created our first complete **Kafka Streams** application. We built the application using Gradle, and then again with Maven. We have used the command line tools provided by Kafka to produce input data and display the transformed output data.



STOP HERE. THIS IS THE END OF THE EXERCISE.

tejaswin.renugunta@walgreens.com

b. Working with JSON

The purpose of this exercise is to learn how to create serializers and deserializers for custom Java objects.

In this case, we will create a Serde for an object that records temperature data using the **KafkaJsonSerializer** and **KafkaJsonDeserializer** helper classes. The application itself reads temperature data from an input topic, filters for temperatures higher than 25 degrees, and outputs that data to a new output topic.

Prerequisites

1. Navigate to this lab's folder **~/confluent-streams/labs/streams-writing**:

```
$ cd ~/confluent-streams/labs/streams-writing
```

2. If it is not already running, start the Kafka cluster:

```
$ docker-compose up -d zookeeper kafka
```

Do not proceed until all services are up and running; test with:

```
$ docker-compose ps
```

and assert that all services are in state **Up**.

3. Create an input topic called **temperatures-topic** and an output topic called **high-temperatures-topic** in Kafka:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic temperatures-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic high-temperatures-topic
```

Writing the Kafka Streams App

1. Open a new terminal window and navigate to the **streams-writing/json-sample** folder, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-writing/json-sample
$ code .
```

2. Locate the file **build.gradle** in this folder and open it to analyze its content.



Compared to the **build.gradle** file in the previous exercise we have added the **kafka-json-serializer** library for the JSON serializer/deserializer

3. In the subfolder **src/main/java/streams** locate the file **JsonSample.java** and familiarize yourself with the code. It basically does the configuration of the **Kafka Streams** app but the interesting code is missing. As before, you can challenge yourself to implement the missing code or follow the step-by-step instructions in this book.
4. To create a **Serde** (Serializer/Deserializer) for JSON formatted data, add this code to the function **getJsonSerde()** after the **TODO** comment:

```

Map<String, Object> serdeProps = new HashMap<>();
serdeProps.put("json.value.type", TempReading.class);

final Serializer<TempReading> temperatureSerializer = new
KafkaJsonSerializer<>();
temperatureSerializer.configure(serdeProps, false);

final Deserializer<TempReading> temperatureDeserializer = new
KafkaJsonDeserializer<>();
temperatureDeserializer.configure(serdeProps, false);

return Serdes.serdeFrom(temperatureSerializer,
temperatureDeserializer);

```

We're basically using the two helper classes **KafkaJsonSerializer** and **KafkaJsonDeserializer** to create a serializer and a deserializer which in turn we then use to create a **Serde**.

We will use this Serde to serialize and deserialize our **TempReading** POJO.

5. The final thing to do is to define the **Topology** for our application. We want to keep it simple and just filter the input topic **temperatures-topic** for high temperatures (>25 degrees) and output the result to the output topic **high-temperatures-topic**. Add this code to the **getTopology()** function after the **TODO** comment:

```

builder.stream("temperatures-topic", Consumed.with(stringSerde,
temperatureSerde))
    .filter((key,value) -> value.temperature > 25)
    .to("high-temperatures-topic", Produced.with(stringSerde,
temperatureSerde));
return builder.build();

```



Note how the filter function uses the **value** which is an object of type **TempReading**.

6. Note the file **log4j.properties** in the folder **src/main/resources** which is used to configure logging for our application.
7. Use **Run → Start Debugging** to run your code. Let your Kafka Streams app remain running in the VS Code debugger.

You should get this output:

```
*** Starting JSON Sample Application ***  
...
```

Creating input Data

1. Switch back to the terminal.
2. Use the following command to generate some temperature readings in JSON format:

```
$ cat << EOF | kafka-console-producer \  
  --bootstrap-server kafka:9092 \  
  --property "parse.key=true" \  
  --property "key.separator=:" \  
  --topic temperatures-topic  
"S1":{"station":"S1", "temperature": 10.2, "timestamp": 1}  
"S1":{"station":"S1", "temperature": 11.2, "timestamp": 2}  
"S1":{"station":"S1", "temperature": 11.1, "timestamp": 3}  
"S1":{"station":"S1", "temperature": 12.5, "timestamp": 4}  
"S2":{"station":"S2", "temperature": 15.2, "timestamp": 1}  
"S2":{"station":"S2", "temperature": 21.7, "timestamp": 2}  
"S2":{"station":"S2", "temperature": 25.1, "timestamp": 3}  
"S2":{"station":"S2", "temperature": 27.8, "timestamp": 4}  
EOF
```



Run this command repeatedly to generate more messages...

Reading the Output

1. Use the following command to read the output generated:

```
$ kafka-console-consumer \  
  --bootstrap-server kafka:9092 \  
  --from-beginning \  
  --topic high-temperatures-topic
```

You should get this output showing only readings with temperature higher than 25 degrees:

```
{"station":"S2","temperature":25.1,"timestamp":3}  
{"station":"S2","temperature":27.8,"timestamp":4}
```


Cleanup

1. Terminate the Kafka console consumer by pressing **Ctrl+C**.
2. Terminate your **Kafka Streams** application with **Run** → **Stop Debugging**.
3. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this sample we have built a **Kafka Streams** application that uses custom serializer and deserializer to work with data that is JSON formatted.

Please proceed to the next exercise: Joining Two Streams.

tejaswin.renugunta@walgreens.com

c. Joining Two Streams

In event-driven architecture, it is important to think about what event will trigger an output. Different kinds of joins will trigger outputs under different conditions. The purpose of this exercise is to experiment with the behavior of various stream-stream joins.

The streaming application itself performs a stream-stream join. It takes a string value from a "left stream" and a string value from a "right stream" and concatenates them together inside of brackets. The output is produced to a new stream. Remember that all stream-stream joins must be windowed since streams are unbounded. This application will use a tumbling window of 5 minutes.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. If your Kafka cluster is not already running, start it with:

```
$ cd ~/confluent-streams/labs/streams-writing
$ docker-compose up -d zookeeper kafka
```

3. Create two input topics called **left-topic** and **right-topic** and an output topic called **joined-topic** in Kafka:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic left-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic right-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic joined-topic
```



Remember that joins require the input topics to have the same number of partitions so that all input records with the same key, from both sides of the join, are delivered to the same stream task during processing. (called co-partitioning).

Create the Streaming App

1. Open another terminal and navigate to the **streams-writing/join-sample** folder, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-writing/join-sample
$ code .
```

2. Open the file **build.gradle** and analyze its content.
3. Locate the file **JoinSample.java** in the folder **src/main/java/streams** and open it. Familiarize yourself with the code. It basically does the configuration of the **Kafka Streams** app but the interesting code is missing. If you would like to challenge yourself, take this time to create the streaming application logic yourself.
4. Now we define the **Topology** for our application. Add the following code snippet to the

`getTopology()` function after the **TODO** comment:

```
KStream<String, String> leftStream = builder.stream("left-topic",
    Consumed.with(stringSerde, stringSerde));
KStream<String, String> rightStream = builder.stream("right-topic",
    Consumed.with(stringSerde, stringSerde));
leftStream
    .join(rightStream,
        (leftValue, rightValue) -> "[" + leftValue + ", " +
        rightValue + "]",
        JoinWindows.of(Duration.ofMinutes(5)),
        StreamJoined.with(stringSerde, stringSerde, stringSerde)
    )
    .to("joined-topic", Produced.with(stringSerde, stringSerde));

Topology topology = builder.build();
return topology;
```

What does the above code do? Discuss with your peers.

5. Note the file `log4j.properties` in the folder `src/main/resources`, which is used to configure logging for our streams application.
6. Use **Run** → **Start Debugging** to run your code. Let your Kafka Streams app remain running in the VS Code debugger.

You should get this output:

```
*** Starting Join Sample Application ***
...
```

Creating input Data

1. Open 3 terminal windows and arrange them side by side so you can see all three of them at the same time.

We will be using the tool `kafkacat` to generate data and monitor the output:

2. In the first terminal window start the tool `kafkacat` as a producer for the `left-topic` topic:

```
kafkacat \  
-b kafka:9092 \  
-t left-topic \  
-P -K: -Z
```

3. In the second terminal window start **kafkacat** as a producer for the **right-topic** topic:

```
kafkacat \  
-b kafka:9092 \  
-t right-topic \  
-P -K: -Z
```

4. In the third terminal window run an instance of **kafkacat** as a consumer of the **joined-topic** topic:

```
kafkacat \  
-b kafka:9092 \  
-t joined-topic \  
-C -K\\t
```

5. In window 1 enter **FL:** (a record key of **FL** for Florida a **<NULL>** value for the record value) and observe the output in window 3. Hint: Nothing should happen. ...why?
6. In window 2 also enter the value **FL:** and observe the output in window 3. Hint: Nothing should happen. ...why?
7. Now in window 1 enter the value **FL:Orlando** and observe the output in window 3. Hint: Nothing should happen. ...why?
8. In window 2 enter the value **FL:Tampa** and observe the output in window 3. You should see:

```
FL  [Orlando, Tampa]
```

9. Back in window 1 enter **FL:** and observe the output in window 3. What do you see? Why?
10. Still in window 1 enter **FL:Miami** and observe the output in window 3. You should see:

```
FL  [Miami, Tampa]
```

11. Continue with window 2 and value **FL:Naples**. What do you see this time?

12. Continue to enter more values with the same key **FL**. Here are other cities in Florida to experiment with:
 - Jacksonville
 - Alachua
 - Pensacola
 - Destin
 - Fort Meyers
13. What happens if you use a different key, say **NY**? Why? Discuss the results with your peers.
14. What happens when an event falls outside of the tumbling window?

Left Join

1. Modify the **getTopology()** function and replace the **join** function with a **leftJoin** function instead.
2. Recompile and run the application.
3. This time, use a car brand for the key and different car models for values to play with input data as you have done for the (inner) join example.

What is different? Discuss with your peers if needed.

Outer Join

1. Modify the **getTopology()** function and replace the current **join** function with an **outerJoin** function instead.
2. Recompile and run the application.
3. This time, use your home country for key and different cities for values to play with input data as you have done for the (inner) join example and observe the output.

What is different? Discuss with your peers if needed.

Optional: Stream - Table Joins

1. Perform a similar experiment with a stream - table **left join** (the most common join in most streaming applications). Make sure to experiment with sending null keys and values. How are the results different from the stream - stream left join?

Cleanup

1. Terminate the first 2 instances of **kafkacat** (the producer instances) by pressing **Ctrl+D**.
2. Terminate the third instance of **kafkacat** (the consumer instance) by pressing **Ctrl+C**.
3. Terminate your **Kafka Streams** application with **Run** → **Stop Debugging**.
4. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this exercise you created a **Kafka Streams** application that joins two streams with inner, left, and outer joins. You then created data for the left and the right input stream and observed the generated output. You used the command line tool **kafkacat** to generate input data and observe output data. Consider summarizing your observations and comparing them to the information found here:

<https://docs.confluent.io/current/streams/developer-guide/dsl-api.html#kstream-kstream-join>. Especially focus on the subsection called "Semantics of stream-stream joins" with an illustrative table of the output records that are produced from a join as events flow into the left and right streams.



STOP HERE. THIS IS THE END OF THE EXERCISE.

d. Using the Processor API

The purpose of this exercise is to create an application with the lower-level Processor API. This may be required in applications that require a greater level of control over state store management and more sophisticated application logic than the Kafka Streams DSL can provide.

This exercise will give you experience creating a Kafka Streams application using the DSL and creating a new node in the topology using the Processor API through the `transform()` method.

This application uses a simple source → word count processor → sink processing topology. The source node takes in records from an input topic whose values are sentence strings. The word count processor uses each record's value to update its internal state store for word counts (key=word, value=count) and sends that state to the sink processor every second. The sink processor produces the resulting records to an output topic.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. If your Kafka cluster is not running already, start it with:

```
$ cd ~/confluent-streams/labs/streams-writing
$ docker-compose up -d zookeeper kafka
```

3. Create three topics called **lines-topic**, **lines-topic-repartition** and **word-count-topic** in Kafka:


```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-topic-repartition

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic word-count-topic
```

Create the Streaming App

1. Open another terminal and navigate to the **streams-writing/processor-sample** folder, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-writing/processor-sample
$ code .
```

2. Open the file **build.gradle** in the **streams-writing/processor-sample** folder and analyze its content. It should be quite familiar by now.
3. Open the file **WordCountTransformer.java** in subfolder **src/main/java/streams**.

Familiarize yourself with the code. We are creating an instance of type **Transformer** and overriding its methods. Thoroughly document what the **init** method does. Pay particular note to the call to the `context.schedule()` method, with its punctuation and the call to `context.forward()`. As always, you can check the corresponding code in the **solutions** folder for feedback.

4. You can challenge yourself to write the **transform** method, or continue with the next step.

5. Make it so the **transform** method of the class has these contents:

```
Long oldValue = this.kvStore.get(word);
if (oldValue == null) {
    this.kvStore.put(word, 1L);
} else {
    this.kvStore.put(word, oldValue + 1L);
}
return null;
```

This code gets the correct entry from the keystore and updates it or it creates a new entry in the keystore with a key of the incoming word and a count of 1. It leaves it to the context scheduler to forward the entries in the key value store, producing them to Kafka for durability.

6. Open the file **CustomTransformerApp.java** in the same folder, and familiarize yourself with the code. This code is entirely DSL - It defines the topology, creates the configuration, sets up the shutdown hook and starts the **Kafka Streams** app.

In the topology, it creates a stream from the input topic, uses the `flatMapValues()` method to break up the lines of input text into individual words, writes a rekeyed stream out to a repartition topic and reads it back in. Then it calls the `transform()` method we created in the `WordCountTransformer` code. And finally, it directs those results to the output topic.

1. Notice the file **log4j.properties** in the folder **src/main/resources** used to configure logging for the application.
2. Use **Run → Start Debugging** to run your code. Let your Kafka Streams app remain running in the VS Code debugger.

You should get this output:

```
*** Starting Custom Transformer App Application ***
```

Creating input Data

1. Open 2 terminal windows and arrange them side by side so you can see the two at the same time.

2. In the first terminal window start the tool **kafkacat** as a producer for the **lines-topic** topic:

```
kafkacat \  
-b kafka:9092 \  
-t lines-topic \  
-P -K :
```

3. In the second terminal window run an instance of **kafka-console-consumer** as a consumer of the **word-count-topic** topic, printing the String key and the Long value:

```
kafka-console-consumer --bootstrap-server kafka:9092 \  
--topic word-count-topic --from-beginning \  
--property print.key=true \  
--value-deserializer  
org.apache.kafka.common.serialization.LongDeserializer
```

4. In window 1 enter each of these strings one at a time:

```
kafka:Kafka powers the Confluent streaming platform  
kafka:A streaming application uses Kafka  
kafka:Everyone loves Kafka  
kafka:Many contributors to Kafka work for Confluent
```

After each line observe the output in window 2.

5. Discuss the result with your peers.

Cleanup

1. Terminate the **kafkacat** instance by pressing **Ctrl+D**.
2. Terminate the **kafka-console-consumer** instance by pressing **Ctrl+C**.
3. Terminate your **Kafka Streams** application with **Run** → **Stop Debugging**.
4. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this sample we have demonstrated the use of a custom transform written using the Processor API to count words. This is included into a Stream processing application written

entirely in the Streams DSL that takes in sentences from an input topic, process them and writes the resulting word and count pairs to an output topic.



STOP HERE. THIS IS THE END OF THE EXERCISE.

tejaswin.renugunta@walgreens.com

Lab 05 Testing Kafka Streams Applications

This lab contains 2 exercises:

- Building Unit Tests
- Integration Tests using Embedded Kafka

a. Building Unit Tests

It is essential for every application to have full test coverage of each of its components. The purpose of this exercise is to build unit tests for an existing Kafka Streams application and test the application with Gradle.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. Notice that we do not need to run our Kafka cluster - the unit and integration testing do not require it.
3. Navigate to the folder `~/confluent-streams/labs/streams-testing/simple-test`, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-testing/simple-test
$ code .
```

4. Open the file `build.gradle` inside this folder and analyze its content. Notice the `junit`, as well as the `kafka-streams-test-utils` dependencies that we use to enable testing.

Authoring the Processor

1. Have a look at the content of subfolder `src/main/java/streams`. You should find 3 Java files in it:

- `CustomMaxAggregatorSupplier.java`
 - `ConfigProvider.java`
 - `TopologyProvider.java`
2. Open the class `CustomMaxAggregatorSupplier` and have a look into the code and try to understand what's happening. This is the code we're going to test ultimately. Maybe discuss the code with your peers.
 3. Now a **Kafka Streams** application also needs some configuration. For this purpose we have the `ConfigProvider` class. Once again make sure you understand the code before you proceed.
 4. Finally we have a `TopologyProvider` class which defines the topology of the **Kafka Streams** app that we want to test. And again we invite you to analyze the code and discuss it with your peers if needed.
 5. Due to the fact that we use this topology in a test scenario it has some settings that would not be recommended in production. Can you spot them? If yes, discuss how we could improve this class to work well for both scenarios, testing and production run.
 6. Now we're ready for the actual test class. Open the class `ProcessorTest.java` located in subfolder `src/test/java/streams`. This file contains the skeleton of the test class.

Let's discuss the code:

- We are using the `TopologyTestDriver` class to test the topology. This is a helper class from the `kafka-streams-test-utils` library.
- To send records to the test driver we are using the `TestInputTopic` class.
- To verify the results of our application we are using the `TestOutputTopic` class.
- In the `setup()` method, we're using our two classes `ConfigProvider` and `TopologyProvider` to get the configuration and topology of our **Kafka Streams** application.
- With the latter two we create a test driver instance that we will use in our tests
- The processor is stateful and we pre-populate the state store with a value of `21` for the key `a`.
- In the tear down method we simply clean up by closing the test driver instance.

Writing a Test

Now we are ready to write our first test.

You may want to examine the Javadocs for the following classes we will be using:

TopologyTestDriver:

<https://kafka.apache.org/25/javadoc/org/apache/kafka/streams/TopologyTestDriver.html>

TestInputTopic:

<https://kafka.apache.org/25/javadoc/org/apache/kafka/streams/TestInputTopic.html>

TestOutputTopic:

<https://kafka.apache.org/25/javadoc/org/apache/kafka/streams/TestOutputTopic.html>

1. The first test method to the **ProcessorTest** class will test whether reading from the **result-topic** will give the current max value of **21** for the key **a** after inputting a smaller value for the same key. If so, this assures us that the first input flushed to the local state store and the result was produced to the output topic. You can choose to implement this test yourself by researching the **TopologyTestDriver**, **TestInputTopic**, and **TestOutputTopic** classes as well as the `org.hamcrest.MatcherAssert` and `org.hamcrest.CoreMatchers` classes before moving on.

```
@Test
public void shouldFlushStoreForFirstInput() {
    // TODO: add test code here
}
```

So far this is just standard **jUnit** test method.

2. Add code inside the above method to create an input record using the **pipeInput** method of the **TestInputTopic** class:

```
inputTopic.pipeInput("a", 1L);
```

3. We can use the **readKeyValue** method of the **TestOutputTopic** class to read the output record generated by the processor and the **assertThat** method to compare the key value with the expected result

```
assertThat(outputTopic.readKeyValue(), equalTo(new KeyValue<>("a", 21L)));
```

4. Finally we make sure that this was the only result that we got as output for the given input:

```
assertThat(outputTopic.isEmpty(), is(true));
```

5. The final test method should look like this:

```
@Test
public void shouldFlushStoreForFirstInput() {
    // TODO: add test code here
    inputTopic.pipeInput("a", 1L);
    assertThat(outputTopic.readKeyValue(), equalTo(new KeyValue<>("a", 21L)));
    assertThat(outputTopic.isEmpty(), is(true));
}
```

And that's it! We have just authored a complete test.

Running the Test(s) & Displaying the Test Report

1. In the terminal window:

```
$ ./gradlew test
```

The output of this command should look similar to this:

```
BUILD SUCCESSFUL in 4s
3 actionable tasks: 2 executed, 1 up-to-date
```

As we can see, the code compiled and the tests ran without a problem.

2. Navigate to the `simple-test/build/reports/tests/test/` folder and open the `index.html` test report in a browser.

it should look similar to this:

Test Summary

1 tests	0 failures	0 ignored	0.266s duration	100% successful
------------	---------------	--------------	--------------------	--------------------

Packages

Classes

Package	Tests	Failures	Ignored	Duration	Success rate
io.confluent.training.streams	1	0	0	0.266s	100%

Generated by Gradle 4.8 at Jun 14, 2018 7:19:08 AM

At this point, we should see 1 test that ran successfully.

Adding more Tests

1. Provided is a second test that makes sure that the store value is not updated for smaller input values

```
@Test
public void shouldNotUpdateStoreForSmallerValue() {
    inputTopic.pipeInput("a", 1L);
    assertThat(store.get("a"), equalTo(21L));
    assertThat(outputTopic.readKeyValue(), equalTo(new KeyValue<>("a", 21L)));
    assertThat(outputTopic.isEmpty(), is(true));
}
```

Discuss with your peers what this test does.

2. Add more tests and repeat the `./gradlew test` command to see new results.
 - a. Add a test that asserts that the store value is updated for a larger input value.



If you get stuck, please have a look in the solutions folder where you will find the complete sample solution.

- b. Add a test that asserts a new store value is generated if adding a value with a new key "b". Also make sure the existing store value for "a" is unchanged.
- c. Write a test that verifies that the processor **punctuates** if the **event time** advances.
- d. Write a test that verifies that the processor **punctuates** if the **wall clock time** advances.



Use the function `testDriver.advanceWallClockTime(...)` for this.

Conclusion

In this sample we have shown how to test a simple **Kafka Streams** application that uses the Processor API.

tejaswin.renugunta@walgreens.com

b. Integration Tests using Embedded Kafka

In this exercise, we are going a step further to use a full blown **embedded** Kafka single node cluster to test our **Kafka Streams** application. The application is a simple **word count** streaming app. As in the previous test, we do not need to run our Docker container Kafka cluster.

Creating the Artifacts to Test

1. Navigate to folder `~/confluent-streams/labs/streams-testing/wordcount-test`, and launch VS Code:

```
$ cd ~/confluent-streams/labs/streams-testing/wordcount-test
$ code .
```

2. Open the file `build.gradle` and analyze its content. Notice the many dependencies that we have to add to enable this kind of integration testing.



we have a few dependencies that use the **test** version of the respective **JAR** files. This can be achieved by adding `classifier: 'test'` or `classifier: 'tests'` to the dependency.

3. Open the file `TopologyProvider.java` class, located in subfolder `src/main/java/streams`, that defines the **Kafka Streams** topology that we will test. Analyze the code. Make sure you understand it. If not discuss it with your peers.

Writing the Integration Test

Now it is time to write the actual test that is leveraging the testbed that we have just prepared.

1. Locate the file `TopologyProviderTest.java` in subfolder `src/test/java/streams` and open it. It will host our test code.

Let's analyze the code a bit:

- In the `@BeforeClass` we're initializing and running an embedded single node Kafka cluster. We are preparing the cluster by creating the **input** and the **output** topic.
 - The method `shouldCountWords` contains our first test. We have clearly documented the steps:
 1. define the configuration of the streaming app to use during the test
 2. get the topology that we want to test
 3. initialize and start the streaming application
 4. produce some data for the **input** topic
 5. finally verify that the produced data in the **output** topic corresponds to the expected values
2. Inspect the method `getStreamsConfiguration`. Notice that we take the information about the bootstrap server(s) from our `CLUSTER` variable which represents our single node embedded Kafka cluster. We also provide a path to where the state store should store its information (`STATE_DIR_CONFIG`).
 3. Inspect the method `produceInputData`. Notice that we're configuring a producer and are using it to feed the data into the **input** topic.
 4. The final step is to write code that validates the results. Inspect the `verifyOutputData` method. You can challenge yourself to complete the method, or you can proceed to the next step.
 5. Make it so the `verifyOutputData` method has these contents:

```

List<KeyValue<String, Long>> expectedWordCounts = Arrays.asList(
    new KeyValue<>("hello", 1L),
    new KeyValue<>("all", 1L),
    new KeyValue<>("streams", 2L),
    new KeyValue<>("lead", 1L),
    new KeyValue<>("to", 1L),
    new KeyValue<>("join", 1L),
    new KeyValue<>("kafka", 3L),
    new KeyValue<>("summit", 1L)
);

Properties consumerConfig = new Properties();
consumerConfig.put(ConsumerConfig.BOOTSTRAP_SERVERS_CONFIG, CLUSTER
    .bootstrapServers());
consumerConfig.put(ConsumerConfig.GROUP_ID_CONFIG,
    "wordcount-lambda-integration-test-standard-consumer");
consumerConfig.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG,
    "earliest");
consumerConfig.put(ConsumerConfig.KEY_DESERIALIZER_CLASS_CONFIG,
    StringDeserializer.class);
consumerConfig.put(ConsumerConfig.VALUE_DESERIALIZER_CLASS_CONFIG,
    LongDeserializer.class);

List<KeyValue<String, Long>> actualWordCounts =
    IntegrationTestUtils.waitForMinKeyValueRecordsReceived(
        consumerConfig, outputTopic, expectedWordCounts.size());
assertThat(actualWordCounts, containsInAnyOrder(expectedWordCounts
    .toArray()));

```

We define an array with the expected data, define a consumer configuration, and then compare the produced data to the expected data.

Running the Test(s)

1. In the terminal window:

```
$ ./gradlew test
```

You should see something like this (shortened):

```
...  
> Task :compileJava  
> Task :processResources NO-SOURCE  
> Task :classes  
> Task :compileTestJava  
> Task :processTestResources NO-SOURCE  
> Task :testClasses  
> Task :test  
  
BUILD SUCCESSFUL in 46s  
3 actionable tasks: 3 executed  
...
```

2. Navigate to the **wordcount-test/build/reports/tests/test/** folder and open the **index.html** test report in a browser.

Cleanup

There is no special cleanup needed.

Conclusion

In this exercise we have written an integration test for a Kafka Streams app topology. The test bed uses an embedded version of Kafka and Zookeeper.



STOP HERE. THIS IS THE END OF THE EXERCISE.

Lab 06 Introduction to ksqlDB

a. Introduction to ksqlDB

In this lab exercise, you will use the ksqlDB CLI to slice and dice data that is being generated in real time. The purpose is to get a sense for the streaming applications that are possible using only a SQL-like syntax. We will query running streams, apply filters and maps, perform a stream-table join for data enrichment, and create new streams and tables derived from existing streams.

For more details on running ksqlDB in Docker containers, please see [Install ksqlDB with Docker](#).

Preparing the Platform

1. Open a terminal window and navigate to the folder `~/confluent-streams`

```
$ cd ~/confluent-streams
```

2. This exercise will use the `docker-compose.yml` file in `~/confluent-streams`, but this time we need to also run the `ksqldb-server` container.

Optional: Open the file in your editor and analyze its content. You do not necessarily need to understand all of it at this point.

3. Run the application with:

```
$ docker-compose up -d zookeeper kafka ksqldb-server
```

4. Wait until the app is up and running, that is all services are marked as **up**:

```
$ docker-compose ps
```

you should see something like this:

Name	Command	State	Ports
kafka	/etc/confluent/docker/run	Up	0.0.0.0:9092-
ksqldb-server	/etc/confluent/docker/run	Up	0.0.0.0:8088-
zookeeper	/etc/confluent/docker/run	Up	2181/tcp, 2888/tcp, 3888/tcp

- You can see that we have 3 containers running on our system:

- the first in the list is running a **Kafka broker**
- the second container runs an instance of **ksqlDB Server**
- and the last one runs an instance of **Zookeeper**

All these containers run on a software defined network (SDN) called **confluent-streams_kafka-net**.

Running the ksqlDB CLI

We want now to use the ksqlDB CLI and connect it with our ksqlDB server.

- Start the ksqlDB CLI

```
$ ksql http://ksqldb-server:8088
```

You will be greeted by the following screen:


```
OpenJDK 64-Bit Server VM warning: Option UseConcMarkSweepGC was
deprecated in version 9.0 and will likely be removed in a future
release.
```

```
=====
=                                     =
=                                     =
=      [ | / \ _ _ _ _ | / \ _ _ _ _ ] [ | / \ _ _ _ _ ] [ | / \ _ _ _ _ ]
=      < \ _ _ _ _ \ ( | / \ _ _ _ _ ) [ | / \ _ _ _ _ ] [ | / \ _ _ _ _ ]
=      - \| \ _ _ _ _ \| \ , - [ | / \ _ _ _ _ ] [ | / \ _ _ _ _ ]
=                                     =
= Event Streaming Database purpose-built
=       for stream processing apps
=====
```

Copyright 2017-2020 Confluent Inc.

CLI v6.0.0, Server v6.0.0 located at <http://localhost:8088>

Having trouble? Type 'help' (case-insensitive) for a rundown of how things work!

ksql>

Testing ksqlDB

Here we're going to experiment with various features of ksqlDB. We use the two topics **pageviews** and **users**.

First we will use the `ksql-datagen` tool to create data for us. This tool has a number of predefined data types - you see us requesting them here with the parameter **quickstart=users** and **quickstart=pageviews** in the commands below.

For more info on this tool, see: <https://docs.ksqldb.io/en/latest/developer-guide/test-and-debug/generate-custom-test-data/>.

1. First, open two more terminals and run each of the `ksql-datagen` commands in one of them.

```
$ ksql-datagen quickstart=users format=json topic=users \
  bootstrap-server=kafka:9092
```

```
$ ksql-datagen quickstart=pageviews format=delimited \
  topic=pageviews bootstrap-server=kafka:9092
```

1. Now, return to the ksqlDB CLI. To be able to work with data from Kafka we need to create either a **Stream** or a **Table** in ksqlDB. Let's create a stream from the topic **pageviews** using KSQL:

```
ksql> CREATE STREAM pageviews_original (  
    viewtime bigint,  
    userid varchar,  
    pageid varchar  
    ) WITH (kafka_topic='pageviews', value_format='DELIMITED');
```

This command creates a **Stream** called **pageviews_original** from the Kafka topic **pageviews**, whose record values are encoded in CSV (here called **DELIMITED**). The three values in each record value are interpreted as fields **viewtime**, **userid** and **pageid**.

The ksqlDB editor should answer with

```
Message  
-----  
Stream created  
-----
```

2. We can now describe the stream:

```
ksql> DESCRIBE pageviews_original;
```

giving us this output:

```
Name                               : PAGEVIEWS_ORIGINAL  
Field | Type  
-----  
VIEWTIME | BIGINT  
USERID   | VARCHAR(String)  
PAGEID   | VARCHAR(String)  
-----  
For runtime statistics and query details run: DESCRIBE EXTENDED  
<Stream,Table>;
```

3. Let's create a **table** from the **users** topic:

```
ksql> CREATE TABLE users_original (
    registertime BIGINT,
    gender VARCHAR,
    regionid VARCHAR,
    userid VARCHAR PRIMARY KEY
) WITH
(kafka_topic='users', value_format='JSON');
```

This command creates a **Table** from the Kafka topic **users**, whose records have values encoded in **JSON**. Since this is a table, we need a key by which the records from the source topic are grouped. In our case this is the field **userid**. The four fields **registertime**, **gender**, **regionid**, and **userid** should be part of the **JSON** value of the records in the topic.

4. Use **DESCRIBE users_original;** to get a description of the table.

```
ksql> DESCRIBE users_original;
```

giving us this output:

```
Name                               : USERS_ORIGINAL
Field                               | Type
-----
REGISTERTIME | BIGINT
GENDER       | VARCHAR(STRING)
REGIONID     | VARCHAR(STRING)
USERID       | VARCHAR(STRING) (primary key)
-----
For runtime statistics and query details run: DESCRIBE EXTENDED
<Stream,Table>;
```

5. Now use **SHOW STREAMS;** and **SHOW TABLES;** to view the list of streams and tables defined in the system.

```
ksql> SHOW STREAMS;
ksql> SHOW TABLES;
```



All ksqlDB commands need to be terminated with a semi-colon (;). ksqlDB SQL keywords such as CREATE or DESCRIBE are not case sensitive.

Querying Streams and Tables

1. Let's get some data from the **pageviews_original** stream:



The default for where to start when selecting data is latest. If data is not being continuously loaded into your stream, this may result in your **SELECT** not displaying any data for some time. Set the default to be earliest in your ksqlDB CLI session with the set command: **set 'auto.offset.reset'='earliest';**

```
ksql> SELECT * FROM pageviews_original EMIT CHANGES LIMIT 10;
```

Please note the **LIMIT** clause which limits the output to 10 records. The output should look similar to this:

VIEWTIME	USERID	PAGEID
1603202197108	User_5	Page_22
1603202197145	User_6	Page_55
1603202197145	User_1	Page_82
1603202197145	User_4	Page_86
1603202197145	User_2	Page_58
1603202197145	User_8	Page_97
1603202197145	User_6	Page_88
1603202197145	User_4	Page_21
1603202197145	User_5	Page_65
1603202197145	User_5	Page_34

Limit Reached
Query terminated



If your data is displayed in columns that are too wide, you can change the column width in your ksqlDB CLI session using the set command. For example: **SET CLI COLUMN-WIDTH 15;**

2. Let's run the same query but this time without the **LIMIT** clause:

```
ksql> SELECT * FROM pageviews_original EMIT CHANGES;
```

you will notice that the query does not stop and continues indefinitely. This is of course

expected since a stream never ends.

Hit **Ctrl+C** to end the query.

3. Now try the same with the table:

```
ksql> SELECT * FROM users_original EMIT CHANGES LIMIT 5;
```

giving us this output:

```
+-----+-----+-----+
+-----+
|USERID      |REGISTERTIME      |GENDER|
|REGIONID    |                  |      |
+-----+-----+-----+
+-----+
|User_2      |1512186975750    |FEMALE|
|Region_4    |                  |      |
|User_2      |1507817615345    |MALE  |
|Region_3    |                  |      |
|User_2      |1493582927082    |FEMALE|
|Region_3    |                  |      |
|User_5      |1515794322305    |FEMALE|
|Region_7    |                  |      |
|User_2      |1514239675179    |FEMALE|
|Region_9    |                  |      |
Limit Reached
Query terminated
```



Although tables are compacted topics and can be compared to classical database tables a query on them never ends analogous to the stream.

Filtering and Mapping Operations

Similar to what we're used to from SQL we can filter data from a stream (or from a table).

1. Let's only display records from the stream whose user ID is equal to **User_1**:

```
ksql> SELECT * FROM pageviews_original WHERE userid='User_1' EMIT CHANGES;
```

Hit **Ctrl+C** to end the query.

2. Now only records whose page is in the range 60 to 69:

```
ksql> SELECT * FROM pageviews_original
      WHERE pageid LIKE 'Page_6%'
      EMIT CHANGES;
```



The correct filter would be **Page_6_** where we want to only match one character after the 6, but ksqldb currently only supports the wildcard **%** in the filter, matching zero or more characters. See <https://docs.ksqldb.io/en/latest/developer-guide/ksqldb-reference/select-push-query/#like>

3. We can selectively output information from the stream. Only display the field **userid** and **pageid**:

```
ksql> SELECT pageid, userid
      FROM pageviews_original
      EMIT CHANGES
      LIMIT 5;
```

Joining a Stream with a Table

1. Let's enrich the **pageviews_original** stream with some data from the **users_original** table:

```
ksql> SELECT u.userid AS userid, pageid, regionid, gender
      FROM pageviews_original pv
      LEFT JOIN users_original u
      ON pv.userid = u.userid
      EMIT CHANGES
      LIMIT 10;
```

which should display something similar to this:

USERID GENDER	PAGEID	REGIONID	
User_2	Page_94	Region_6	OTHER
User_4	Page_23	Region_9	OTHER
User_6	Page_53	Region_4	OTHER
User_4	Page_90	Region_9	OTHER
User_8	Page_24	Region_1	OTHER
User_7	Page_82	Region_6	MALE
User_4	Page_51	Region_9	OTHER
User_1 FEMALE	Page_95	Region_4	
User_1 FEMALE	Page_67	Region_4	
User_3	Page_33	Region_6	MALE

Limit Reached
Query terminated

Creating derived Streams

Sometimes we want to create streams or tables that derive from existing streams and tables.

1. Create a stream that only shows pageviews by females and region 1:

```
ksql> CREATE STREAM pageviews_female_1 AS
      SELECT u.userid AS userid, pageid, regionid, gender
      FROM pageviews_original pv
      LEFT JOIN users_original u
      ON pv.userid = u.userid
      WHERE gender='FEMALE' AND regionid='Region_1'
      EMIT CHANGES;
```

2. And query this new stream:

```
ksql> SELECT * FROM pageviews_female_1 EMIT CHANGES LIMIT 5;
```

```
+-----+-----+-----+
|USERID|PAGEID|REGIONID|
|GENDER|      |         |
+-----+-----+-----+
|User_8|Page_66|Region_1|
|FEMALE|      |         |
|User_8|Page_37|Region_1|
|FEMALE|      |         |
|User_8|Page_25|Region_1|
|FEMALE|      |         |
|User_8|Page_68|Region_1|
|FEMALE|      |         |
|User_8|Page_71|Region_1|
|FEMALE|      |         |
Limit Reached
Query terminated
```

Cleaning Up

1. Exit ksqlDB with **Ctrl+D**.
2. Return to the terminals running ksql-datagen and stop them using **Ctrl+C**.
3. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this lab we have authored our very first ksqlDB queries and run them against two topics **pageviews** and **users** in Kafka. Kafka and all the other components of the ksqlDB platform ran in Docker containers to make the setup very easy and portable.



STOP HERE. THIS IS THE END OF THE EXERCISE.

tejaswin.renugunta@walgreens.com

Lab 07 Using ksqlDB

This lab has 2 exercises:

- Using ksqlDB
- Using the ksqlDB REST API

a. Using ksqlDB

In this exercise, we will explore how to apply ksqlDB's mapping, filtering, joining, and aggregating capabilities to an application that processes real-time temperature data. MQTT and Internet of Things data are perfect for real-time processing with Kafka Streams and ksqlDB. ksqlDB is an especially good choice for many of these applications because of its simplicity.

Prerequisites

1. Open a terminal window and navigate to the folder `~/confluent-streams`

```
$ cd ~/confluent-streams
```

2. Run a Kafka cluster and a ksqlDB server using this command:

```
$ docker-compose up -d zookeeper kafka ksqldb-server
Creating network "confluent-streams_kafka-net" with the default
driver
Creating ksqldb-server      ... done
Creating zookeeper          ... done
Creating kafka              ... done
```



Wait a couple of minutes until the cluster is ready.

3. Check the status with:

```
$ docker-compose ps
```

Creating Data

1. Create the **stations** topic:

```
$ kafka-topics \  
  --create \  
  --bootstrap-server kafka:9092 \  
  --replication-factor 1 \  
  --partitions 1 \  
  --topic stations \  
  --config cleanup.policy=compact
```



We do this compaction only for illustration. Normally the temperature reading stations wouldn't change too frequently to warrant compaction.

2. Create a list of stations. Here we use the kafka-console-producer command line tool to send records in to the Kafka cluster.

```
$ cat << EOF | kafka-console-producer \  
  --bootstrap-server kafka:9092 \  
  --property "parse.key=true" \  
  --property "key.separator=:" \  
  --topic stations  
1:Mombasa,Kenya  
2:Nairobi,Kenya  
3:Mogadishu,Somalia  
4:Dar es Salaam,Tanzania  
5:Pretoria,South Africa  
6:Cape Town,South Africa  
7:Bloemfontein,South Africa  
8:Diani,Kenya  
9:Embu,Kenya  
10:Johannesburg,South Africa  
EOF  
  
>>>>>>>>>>
```

3. Double check that the list of stations has been created:

```
$ kafka-console-consumer \  
  --bootstrap-server kafka:9092 \  
  --from-beginning \  
  --max-messages 7 \  
  --topic stations \  
  --property print.key=true \  
  --property key.separator=":"
```

```
1:Mombasa,Kenya  
2:Nairobi,Kenya  
3:Mogadishu,Somalia  
4:Dar es Salaam,Tanzania  
5:Pretoria,South Africa  
6:Cape Town,South Africa  
7:Bloemfontein,South Africa  
8:Diani,Kenya  
9:Embu,Kenya  
10:Johannesburg,South Africa  
Processed a total of 10 messages
```

4. Create the **temperatures** topic:

```
$ kafka-topics \  
  --create \  
  --bootstrap-server kafka:9092 \  
  --replication-factor 1 \  
  --partitions 1 \  
  --topic temperatures
```

5. Create a list of temperature readings, again using the kafka-console-producer command line tool.


```
$ ksql http://ksqldb-server:8088
```

2. Define that streams should be read from beginning:

```
ksql> SET 'auto.offset.reset' = 'earliest';  
Successfully changed local property 'auto.offset.reset' from 'null'  
to 'earliest'
```

Mapping and Filtering

1. Create a table from the topic **stations**:

```
ksql> CREATE TABLE weather_stations(  
    id VARCHAR PRIMARY KEY,  
    name VARCHAR,  
    country VARCHAR  
)  
WITH(kafka_topic='stations', value_format='DELIMITED');
```

Message

Table created

2. Run a simple query against this new table:

```
ksql> SELECT * FROM weather_stations EMIT CHANGES LIMIT 5;
```

you should see something like this:

```
+-----+-----+-----+  
| ID    | NAME          | COUNTRY    |  
+-----+-----+-----+  
| 1     | Mombasa       | Kenya    |  
| 2     | Nairobi       | Kenya    |  
| 3     | Mogadishu     | Somalia    |  
| 4     | Dar es Salaam | Tanzania   |  
| 5     | Pretoria      | South Africa |  
Limit Reached  
Query terminated
```



You can also peek into the **stations** topic using:

```
ksql> PRINT 'stations' FROM BEGINNING;
```

Press **Ctrl+C** to stop the above query

3. Now let's only output stations in Kenya:

```
ksql> SELECT * FROM weather_stations
      WHERE country='Kenya'
      EMIT CHANGES;
```

ID	NAME	COUNTRY
1	Mombasa	Kenya
2	Nairobi	Kenya
8	Diani	Kenya
9	Embu	Kenya



Press **Ctrl+C** to end the query.

4. To show only stations whose name starts with **M** and output the country in all caps use:

```
ksql> SELECT id, name, UCASE(country) AS country
      FROM weather_stations
      WHERE name LIKE 'M%'
      EMIT CHANGES;
```

ID	NAME	COUNTRY
1	Mombasa	Kenya
2	Nairobi	Kenya
3	Mogadishu	Somalia

Press **Ctrl+C** to end the query.

5. Create a stream from the topic **temperatures**:

```
ksql> CREATE STREAM mytemperatures(
    id INTEGER,
    station_id VARCHAR,
    temp DOUBLE
)
WITH(kafka_topic='temperatures', value_format='DELIMITED');
```

6. The temperatures are in degree Celsius. To output them in degree Fahrenheit use this:

```
ksql> SELECT id, station_id,
    temp AS temp_in_C,
    temp*9/5+32 as temp_in_F
    FROM mytemperatures
    EMIT CHANGES
    LIMIT 10;
```

ID	STATION_ID	TEMP_IN_C	TEMP_IN_F
6	2	18.5	65.3
5	1	23.0	73.4
20	5	21.3	70.34
15	4	35.5	95.9
17	4	34.5	94.1
9	2	18.0	64.4
24	7	17.0	62.6
13	3	33.0	91.4
19	5	21.0	69.8
25	7	18.0	64.4

Limit Reached
Query terminated

Joining

1. Having worked with the temperatures stream we would love to enrich it with information available in the stations table. We can use **LEFT JOIN** for this:


```
ksql> SELECT s.country, s.name, t.temp
      FROM mytemperatures t
      LEFT JOIN weather_stations s ON t.station_id = s.id
      EMIT CHANGES;
```

Press **Ctrl+C** to end the query.

Aggregating

1. Let's find the maximum temperature for each station:

```
ksql> SELECT station_id, MAX(temp) AS tmax
      FROM mytemperatures
      GROUP BY station_id
      EMIT CHANGES;
```

which should give:

STATION_ID	TMAX
2	18.5
1	23.0
5	21.3
4	35.5
4	35.5
2	18.5
7	17.0
3	33.0
5	21.3
7	18.0
3	33.0
1	24.5
2	18.5
7	18.0
1	26.0
3	33.0
1	26.0
7	18.0
1	26.0
4	37.5
2	18.5
2	18.5
4	37.5
2	18.5
6	23.0

Press **Ctrl+C** to end the query.

Discuss the result, and the **why**, with your peers.

HINT: Look at the environment variable **KSQL_CACHE_MAX_BYTES_BUFFERING** defined in the **docker-compose.yml**.

2. To make the above a permanent query use this command:

```
ksql> CREATE TABLE max_temperatures AS
      SELECT station_id, MAX(temp) AS tmax
      FROM mytemperatures
      GROUP BY station_id
      EMIT CHANGES;
```

Message

Created query with ID CTAS_MAX_TEMPERATURES_5

and then we can select from this new table:

```
ksql> SELECT * FROM max_temperatures EMIT CHANGES;
```

Press **Ctrl+C** to end the query.

3. Now count the measurements per station and only return those who have at least 2 measurements:

```
ksql> SELECT station_id, COUNT( * ) AS nbr
      FROM mytemperatures
      GROUP BY station_id
      HAVING COUNT( * )>=2
      EMIT CHANGES;
```

STATION_ID	NBR
4	2
2	2
5	2
7	2
3	2
1	2
2	3
7	3
1	3
3	3
1	4
7	4
1	5
4	3
2	4
2	5
4	4
2	6

Press **Ctrl+C** to end the query.



Once again the same applies as in step 1 above. We get more records output than we would "logically" expect due to the setting of **KSQL_CACHE_MAX_BYTES_BUFFERING** to **0**.

- Finally count the measurements per country and display them with the max temperature also per country:

```
ksql> SELECT s.country, COUNT( * ) AS nbr, MAX(temp) AS tmax
      FROM mytemperatures t
      LEFT JOIN weather_stations s ON t.station_id = s.id
      GROUP BY s.country
      EMIT CHANGES;
```

COUNTRY	NBR	TMAX
Kenya	1	18.5
Kenya	2	23.0
South Africa	1	21.3
Tanzania	1	35.5
Tanzania	2	35.5
Kenya	3	23.0
South Africa	2	21.3
Somalia	1	33.0
South Africa	3	21.3
South Africa	4	21.3
Somalia	2	33.0
Kenya	4	24.5
Kenya	5	24.5
South Africa	5	21.3
Kenya	6	26.0
Somalia	3	33.0
Kenya	7	26.0
South Africa	6	21.3
Kenya	8	26.0
Tanzania	3	37.5

5. We can also perform a Pull Query from this new table. A Pull query is much like a standard SQL query, in that it reports the results at a snapshot in time and then does not persist and continue to emit changes. No need to press **Ctrl+C** to end the query - it ends as soon as it returns the result.

```
ksql> SELECT * FROM max_temperatures where station_id = '2';
+-----+-----+
| STATION_ID | TMAX |
+-----+-----+
| 2          | 18.5 |
Query terminated
```

Cleanup

1. Exit the **ksqlDB CLI** by pressing **Ctrl+D**.
2. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this exercise, we have explored the various capabilities that ksqlDB offers us in a easy and convenient way. We have learned that the syntax of ksqlDB SQL strongly resembles that of ANSI SQL. We have explored stateless functions such as mapping and filtering as well as stateful operations such as joining and aggregating.

tejaswin.renugunta@walgreens.com

b. Using the ksqlDB REST API

The idea of this exercise is to show how ksqlDB server can be accessed via its RESTful API by any language that can do **HTTP POST** requests. It is possible to create stream processing applications in any such language by submitting queries to a ksqlDB server cluster. The example in this exercise is a Python application. The Python application is also equipped with Kafka client libraries to produce input data and read output data, but in the real world, there are often dedicated upstream producer and downstream consumer applications, and so the stream processing application would not even need to import Kafka client libraries. ksqlDB enables essentially any language to create real-time stream processing applications via its REST API.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. Navigate to the folder **labs/using-ksql**:

```
$ cd ~/confluent-streams/labs/using-ksql
```

3. Run a Kafka cluster and a ksqlDB server using this command:

```
$ docker-compose up -d zookeeper kafka ksqldb-server
Creating network "confluent-streams_kafka-net" with the default
driver
Creating ksqldb-server      ... done
Creating zookeeper          ... done
Creating kafka              ... done
```



Wait a couple of minutes until the cluster is ready.

Authoring a Python Client

1. Locate the folder **labs/using-ksql/ksql-rest-api** and in it the file **requirements.txt**. It has this content:

```
confluent_kafka
requests
```

These are our external dependencies. The library **confluent_kafka** contains the native Python client for Kafka. The **requests** library we use to **HTTP POST** requests to our ksqldb server

2. Again in the same folder locate and open the file **main.py**. Analyze its content. Discuss it with your peers to make sure you understand what's going on.

In essence the application does the following:

- Produce some quotes
- Call ksqldb via its REST API to generate a streaming query
- Use a Kafka consumer to consume the lowercase quotes produced by the streaming query

3. Create the **quotes** topic:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic quotes
```

4. Use **pip3** to install the requirements.

```
$ cd ~/confluent-streams/labs/using-ksql/ksql-rest-api
$ pip3 install -r requirements.txt
```

5. Run the Python client with the following command:

```
$ python3 main.py
```

You should see an output similar to this:

```

>>> Starting Python Kafka Client...
----- Writing quotes to topic 'quotes' -----
*** writing: Kafka enables the Confluent Streaming Platform
*** writing: Confluent offers a Streaming Platform powered by Kafka
*** writing: Kafka Streams are cool
*** writing: Streaming allows for real-time processing of information
*** writing: I love Kafka
----- done writing quotes -----

----- Posting to KSQL Server -----

200, [{"@type":"currentStatus","statementText":"CREATE STREAM
quotes_orig (line STRING) WITH(KAFKA_TOPIC='quotes',
VALUE_FORMAT='DELIMITED');","commandId":"stream/`QUOTES_ORIG`/create"
,"commandStatus":{"status":"SUCCESS","message":"Stream created"}
,"commandSequenceNumber":2,"warnings":[]}],

200, [{"@type":"currentStatus","statementText":"CREATE STREAM
QUOTES_LOWER WITH (KAFKA_TOPIC='QUOTES_LOWER', PARTITIONS=1,
REPLICAS=1) AS SELECT LCASE(QUOTES_ORIG.LINE) KSQL_COL_0\nFROM
QUOTES_ORIG QUOTES_ORIG\nEMIT CHANGES;","commandId":"stream/
`QUOTES_LOWER`/create","commandStatus":{"status":"SUCCESS","message":
"Created query with ID CSAS_QUOTES_LOWER_0"},
,"commandSequenceNumber":4,"warnings":[]}],

----- done posting to KSQL Server -----

>>> Starting Python Kafka Client...
----- Reading from topic 'QUOTES_LOWER' -----
Received message: kafka enables the confluent streaming platform
Received message: confluent offers a streaming platform powered by
kafka
Received message: kafka streams are cool
Received message: streaming allows for real-time processing of
information
Received message: i love kafka
<<< Ending Python Kafka Client...

```

Using the ksqlDB CLI

1. Enter the ksqlDB CLI with this command:

```
$ ksql http://ksqldb-server:8088
```

2. Show all streams:

```
$ ksql> SHOW STREAMS;
```


You should see this:

Stream Name	Kafka Topic	Format
KSQL_PROCESSING_LOG_STREAM	ksql_processing_log_topic	JSON
QUOTES_LOWER	QUOTES_LOWER	DELIMITED
QUOTES_ORIG	quotes	DELIMITED

3. Use the ksqlDB CLI Print command to list the content of the topic:

```
$ ksql> PRINT 'quotes' FROM BEGINNING;
Key format: `_(Ø)`/` - no data processed
Value format: KAFKA_STRING
rowtime: 2020/10/20 14:15:39.093 Z, key: <null>, value: Kafka enables
the Confluent Streaming Platform
rowtime: 2020/10/20 14:15:39.093 Z, key: <null>, value: Confluent
offers a Streaming Platform powered by Kafka
rowtime: 2020/10/20 14:15:39.093 Z, key: <null>, value: Kafka Streams
are cool
rowtime: 2020/10/20 14:15:39.093 Z, key: <null>, value: Streaming
allows for real-time processing of information
rowtime: 2020/10/20 14:15:39.093 Z, key: <null>, value: I love Kafka
```

Press **Ctrl+C** to end the query.

4. Set the starting point of your queries to **earliest**:

```
ksql> SET 'auto.offset.reset' = 'earliest';
```

5. Set the column width to **50**:

```
ksql> SET CLI COLUMN-WIDTH 50;
```

6. List the content of both streams:

```
ksql> SELECT * FROM quotes_orig EMIT CHANGES LIMIT 5;
+-----+
| LINE
+-----+
| Kafka enables the Confluent Streaming Platform
| Confluent offers a Streaming Platform powered by Kafka
| Kafka Streams are cool
| Streaming allows for real-time processing of information
| I love Kafka
Limit Reached
Query terminated
```

and

```
ksql> SELECT * FROM quotes_lower EMIT CHANGES LIMIT 3;
+-----+
| KSQL_COL_0
+-----+
| kafka enables the confluent streaming platform
| confluent offers a streaming platform powered by kafka
| kafka streams are cool
Limit Reached
Query terminated
ksql>
```

7. Quit the ksqlDB CLI with **Ctrl+D**

Cleanup

1. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this exercise we have created a Kafka client application in Python that uses the ksqlDB REST API to access the ksqlDB functionality. The Python client executed the following tasks:

- write entries to an existing topic **quotes**
- post a query to ksqlDB server to create a stream from the topic **quotes**
- post another query to ksqlDB server to create a stream **quotes_lower** containing the quotes from the topic **quotes** all in lower case
- read from the topic **quotes_lower** and output the messages to the screen



STOP HERE. THIS IS THE END OF THE EXERCISE.

Lab 08 Scaling a Kafka Streams Application

a. Scaling a Kafka Streams Application

In this exercise, we're going to write a Kafka producer in either Python or Java that generates a stream of temperature readings for a set of weather stations. We are also writing a simple Kafka Streams application that will consume this topic and calculate the maximum temperature per station per time window. We will then run this application in a single instance and later scale it up to several instances. We will monitor the throughput with the **Confluent Control Center**.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)

Running the Kafka Cluster and Confluent Control Center

1. Navigate to the folder **labs/deployment**:

```
$ cd ~/confluent-streams/labs/deployment
```

2. Run the Kafka cluster:

```
$ docker-compose up -d zookeeper kafka ksqldb-server control-center
```

Wait a couple of minutes until the cluster is initialized.

3. Create the two topics **temperature-readings** and **max-temperatures**, each with 3 partitions using these commands:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 3 \
  --topic temperature-readings

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 3 \
  --topic max-temperatures
```



In any case you need to delete a topic (e.g. the input topic) use this command:

```
$ kafka-topics \
  --bootstrap-server kafka:9092 \
  --delete \
  --topic <topic-name>
```

Creating the Producer

Now it is time to create a temperature readings producer. You can do it either in Python or Java. We start with the Python producer. If you prefer Java then move ahead to → [here](#).

Create the Producer in Python

1. Navigate to the folder **labs/deployment**, and launch VS Code:

```
$ cd ~/confluent-streams/labs/deployment/temp-producer
$ code .
```

2. Locate the file **main.py**.
3. Inspect the code.
 - We are defining a few (temperature measurement) stations, their respective average temperatures and last measured temperatures
 - we're using the Confluent Python client for Kafka to create a producer

- every ~100 ms we're generating a temperature reading for one of the randomly selected station.
4. Return to the terminal window, and install the python prerequisite. Note this may already be installed from a previous exercise:

```
$ pip3 install confluent-kafka
```

5. Run the producer:

```
$ python3 main.py
```

6. In another terminal window run the **kafka-console-consumer** to display the **temperature-readings** topic:

```
$ kafka-console-consumer \  
  --bootstrap-server kafka:9092 \  
  --from-beginning \  
  --max-messages 25 \  
  --topic temperature-readings \  
  --property print.key=true \  
  --property key.separator=", "
```

You should see an output similar to this:

```
S-06, {"station": "S-06", "temperature": -1}  
S-03, {"station": "S-03", "temperature": 8}  
S-03, {"station": "S-03", "temperature": 9}  
S-06, {"station": "S-06", "temperature": 0}  
S-08, {"station": "S-08", "temperature": 31}  
S-09, {"station": "S-09", "temperature": -7}  
...
```

The **key** is the station and the **value** is a JSON object with the station and the temperature in degree Celsius.

Create the Producer in Java

1. Navigate to the folder **labs/deployment**, and launch VS Code:

```
$ cd ~/confluent-streams/labs/deployment/temp-producer
$ code .
```

2. Locate the file **build.gradle** and analyze its content. It is the build file for a simple Kafka client.
3. Locate the file **TempProducer.java** in the subfolder **src/main/java/streams** and open it. Analyze the file and make sure you understand the code. If necessary discuss with your peers.
4. Notice the file **log4j.properties** in the folder **src/main/resources** that configures logging for the producer.
5. Use **Run → Start Debugging** in VS Code or **./gradlew run** in the terminal to run the Java producer.

You should see an output similar to this:

```
The record is: S-06, {"station": "S-06", "temperature": -1}
The record is: S-03, {"station": "S-03", "temperature": 8}
The record is: S-03, {"station": "S-03", "temperature": 9}
The record is: S-06, {"station": "S-06", "temperature": 0}
The record is: S-08, {"station": "S-08", "temperature": 31}
The record is: S-09, {"station": "S-09", "temperature": -7}
...
```

The **key** is the station and the **value** is a JSON object with the station and temperature in degree Celsius.

6. From another terminal window, run the **kafka-console-consumer** to display the **temperature-readings** topic:

```
$ kafka-console-consumer \
  --bootstrap-server kafka:9092 \
  --from-beginning \
  --max-messages 25 \
  --topic temperature-readings \
  --property print.key=true \
  --property key.separator=", "
```

You should see an output similar to this:

```
S-06, {"station": "S-06", "temperature": -1}
S-03, {"station": "S-03", "temperature": 8}
S-03, {"station": "S-03", "temperature": 9}
S-06, {"station": "S-06", "temperature": 0}
S-08, {"station": "S-08", "temperature": 31}
S-09, {"station": "S-09", "temperature": -7}
...
```

Writing the Kafka Streams Application

1. In the folder **deployment/streams-app** locate the file **build.gradle** and analyze its content.



In addition to the usual libraries we also load the **monitoring-interceptors** library to be able to integrate with the Confluent Control Center and the **kafka-json-serializer** library for the JSON serde.

2. Locate the file **StreamsApp.java** in the subfolder **src/main/java/streams** and open it. Make sure you understand the code in it. If not discuss it with your peers.
3. Open a terminal window and navigate to the **streams-app** folder:

```
$ cd ~/confluent-streams/labs/deployment/streams-app
```

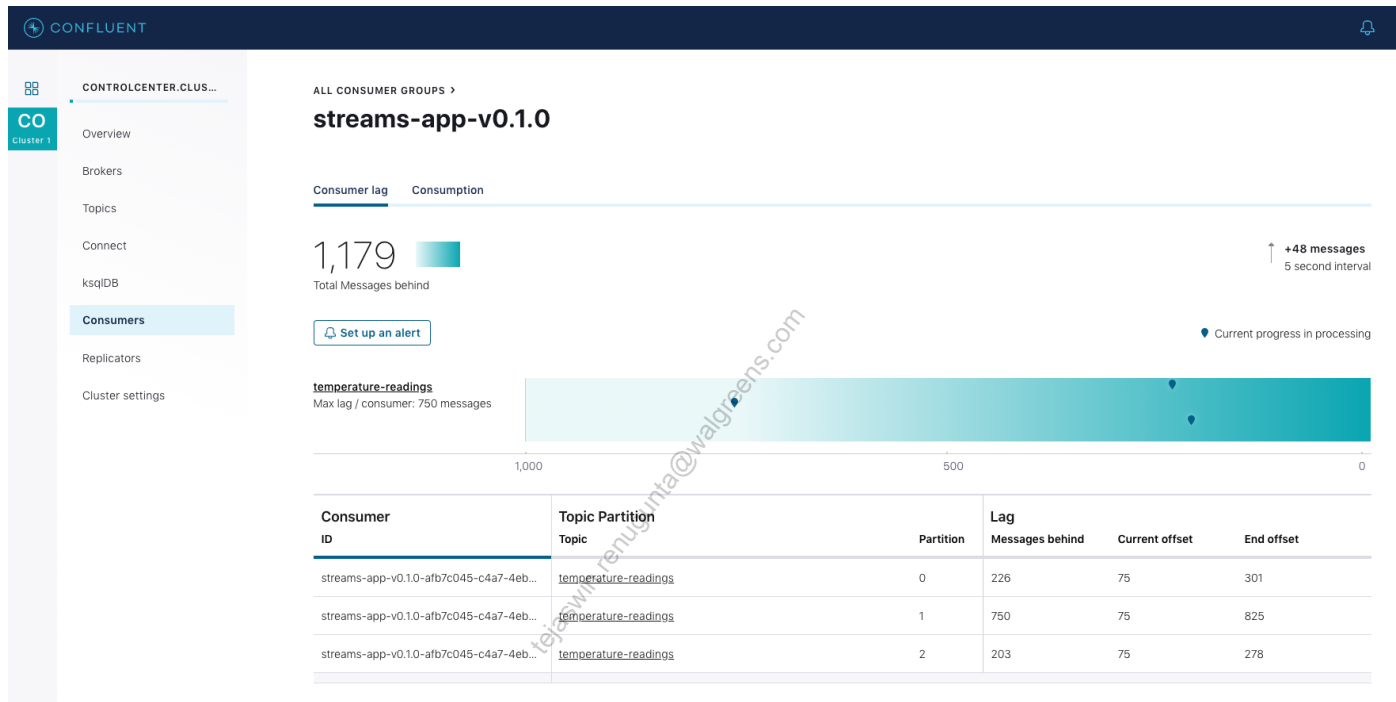
4. You may choose to launch VS Code with **code .** to build and run the application. Or simply use gradle with **./gradlew run**.
5. Run an instance of **kafka-console-consumer** to display the **max-temperatures** topic. Note it may take some time for max temperatures to appear:

```
$ kafka-console-consumer \
  --bootstrap-server kafka:9092 \
  --from-beginning \
  --topic max-temperatures
```




Is the output surprising to you? Refresh yourself on the nature of the **commit.interval.ms** property and how it relates to the output of KTables to Kafka topics.

6. Open Confluent Control Center at <http://localhost:9021>
7. In Control Center, under **Consumers**, monitor the consumer lag for the **streams-app-v.0.1.0**. Note that the consumer group falls more and more behind:

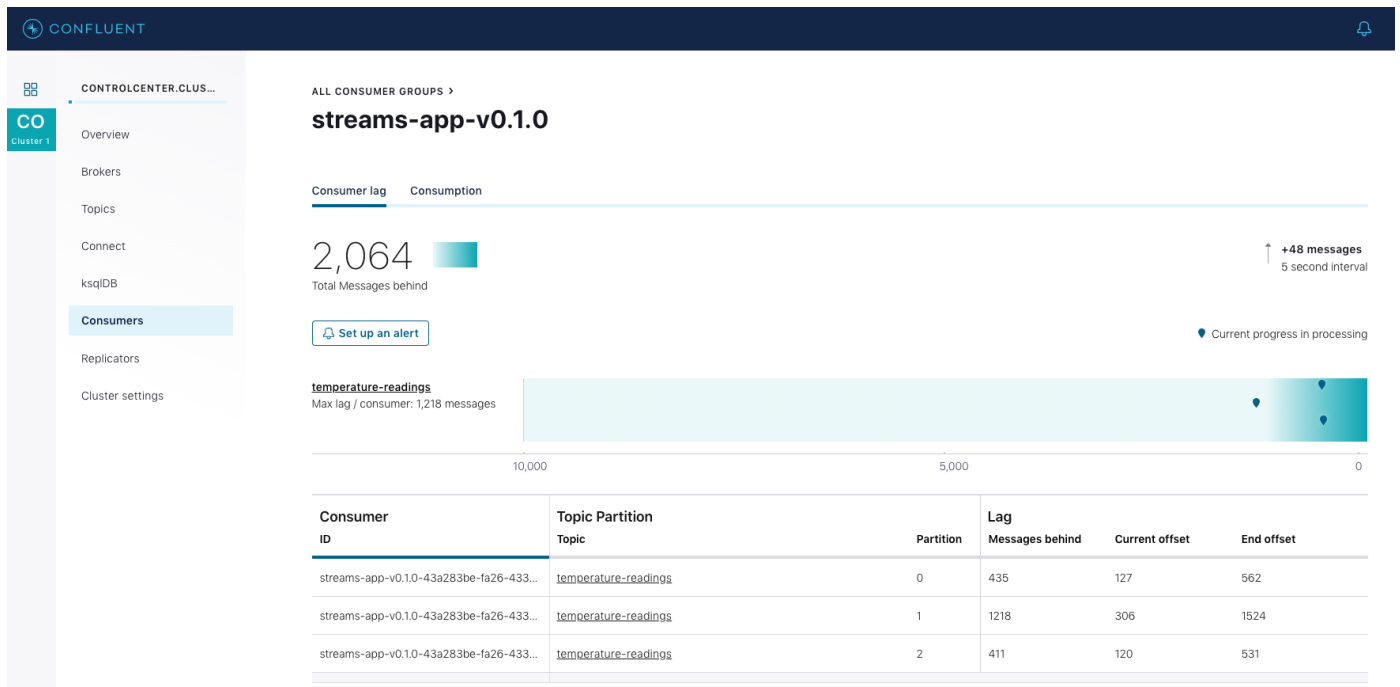


Scaling the Kafka Streams Application

1. Open another terminal window to run another instance of your Kafka Streams application:

```
$ cd ~/confluent-streams/labs/deployment/streams-app
$ ./gradlew run
```

2. In **Confluent Control Center** observe how the throughput of the streams app nearly doubles.
3. Note that we have now two consumer instances listed (recognizable by their ID):



4. Also observe that the consumer lag increases more slowly...
5. Now scale the streams app to 3 instances and again monitor an increase in throughput and reduction in consumer lag.
6. Finally scale the app again, this time to 4 instances. Monitor the throughput after scaling the app. What are you observing? Explain your observation. See the conclusion for an explanation of what happens here.

Optional: Using the ksqldb CLI

We can achieve the same results using ksqldb!

1. Use the ksqldb CLI to play with the data:

```
$ ksql http://ksqldb-server:8088
```

2. Set the starting point of your queries to **earliest**:

```
ksql> SET 'auto.offset.reset' = 'earliest';
```

3. Create a stream for the source topic:

```
ksql> CREATE STREAM temperatures(station STRING, temperature INTEGER)
      WITH(KAFKA_TOPIC='temperature-readings',
      VALUE_FORMAT='JSON');
```

4. Create a table that shows the maximum, minimum, and average temperatures per station per minute:

```
ksql> CREATE TABLE temp_agg_per_min AS
      SELECT station,
             max(temperature) AS max,
             min(temperature) AS min,
             sum(temperature) / count( * ) AS avg
      FROM temperatures
      WINDOW TUMBLING (SIZE 1 MINUTE)
      GROUP BY station;
```

5. Inspect the aggregated temperature data as new records flow in from the producer.

```
ksql> SELECT station, max, min, avg FROM temp_agg_per_min EMIT
      CHANGES;
```

Press **CTL-C** to terminate the query.

6. Exit ksqlDB with **Ctrl+D**.

Cleanup

1. Stop the producer, consumers, and stream application with **Ctrl+C** in the terminal or **Run → Stop Debugging** in VScode.
2. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

In this exercise we have created a Kafka Streams application that processed an input topic and produced an output topic. First we ran only one application instance and then we scaled the application up to several instances. We noticed a significant boost in throughput until the number of instances was greater than the number of partitions of the input topic. At this point the additional application instances were sitting there idle.



In the solutions folder, the Java producer and Kafka Streams app include Dockerfiles so that they can be deployed as containers. A separate **docker-compose.services.yml** has also been provided to start the microservices.



STOP HERE. THIS IS THE END OF THE EXERCISE.

tejaswin.renugunta@walgreens.com

Lab 09 Securing a Kafka Streams Application

a. Securing a Kafka Streams Application

In this exercise, we will interact with a secure Kafka cluster in several ways. The cluster is configured to use SASL-PLAIN for authentication and SSL for transport encryption. A certificate creation script is used to create keystores and truststores for clients and brokers so that they can authenticate with each other. There is a client `.properties` that will be used to create topics, to produce to an input topic, and to consume from an output topic. There is another `.properties` file to configure the security for a Kafka Streams application. In practice, SSL certificates must be carefully managed, credentials should be created separately for different clients, and permissions should be applied to secure credentials. These areas are beyond the scope of the course. Rather, the purpose of the exercise is to provide hands-on experience with the configurations necessary to connect a Kafka Streams application to a secured cluster.

Preparing the Project

Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)



If you have a cluster running from an earlier lab, shut it down now. This lab depends on the settings in the local `docker-compose.yml` file. You can check if you have a cluster running with `docker-compose ps` and you can shutdown a running cluster with `docker-compose down -v`.

Preparing the Certificates

First we need to generate all the necessary certificates and credentials that are used to create a secure (single node) Kafka cluster.

1. In a terminal window, navigate to the folder `security/secure-`

app/scripts/security:

```
$ cd ~/confluent-streams/labs/security/secure-app/scripts/security
```

2. Take a few minutes to inspect the contents of the files in this directory and discuss them with your peers.
3. Run the **certs-create.sh** script to generate the necessary certs and credential files:

```
$ ./certs-create.sh
```

Running the Cluster

1. In a terminal window navigate to folder **security/secure-app**, and launch VS Code:

```
$ cd ~/confluent-streams/labs/security/secure-app  
$ code .
```

2. From this folder open the file **docker-compose.yml** and analyze it. Specifically focus on the settings of the security relevant environment variables. Discuss the content with your peers if you don't fully understand it.
3. Start the cluster:

```
$ docker-compose up -d zookeeper kafka control-center
```

and wait a couple of minutes until the cluster is initialized. You may observe the progress by using:

```
$ docker-compose ps
```

or following the logs, e.g.:

```
$ docker-compose logs -f kafka
```

Press **Ctrl+C** to stop following the log.

Creating the Application

1. In the **secure-app** folder, locate the file **build.gradle** and analyze its content. There should be no surprises at this point.
2. Open the file **SecureAppSample.java** located in the subfolder **src/main/java/streams** and inspect its contents. As you can see, we're using a class **ConfigProvider** to get us the configuration for the **Kafka Streams** app and a class **TopologyProvider** to get us the topology for the application. Other than that all is well known and familiar code by now.
3. Now let's have a look at the file **ConfigProvider.java**. Notice that this code gets the configuration settings from a **.properties** file, which is preferable to hard-coding properties into our application.
4. In the project sub-folder **scripts/security**, open the file called **secureapp-sample.properties** analyze its content.
 - Note that we're using **SASL** to authenticate our application with the secured Kafka cluster.
 - The selection of the **protocol** (**SASL_SSL**) also has an implication on which port we're using for the communication with the brokers; **9091** in this case. The brokers are also listening to **SSL** on port **11091**, if we wanted to do mutual SSL instead of SASL + SSL. In that case, we would also have to configure **ssl.keystore** location and password on the clients so that Brokers could authenticate them.
 - We need to define the **trust store** and its password as well as the login module and the credentials.
 - We're also securing the monitoring interceptors so that the app can be monitored from **Confluent Control Center**.
5. Copy the security files to the location in **/etc/kafka/secrets/** - remember our user ID is training and the password which will be requested is also training:

```
$ sudo mkdir -p /etc/kafka/secrets/  
$ sudo cp scripts/security/* /etc/kafka/secrets/
```

6. Finally have a look at the file **TopologyProvider.java** which should have this content:

```

package streams;

import org.apache.kafka.streams.Topology;
import org.apache.kafka.streams.StreamsBuilder;

public class TopologyProvider {
    public Topology getTopology() {
        final StreamsBuilder builder = new StreamsBuilder();
        builder.stream("secure-input").to("secure-output");
        return builder.build();
    }
}

```



Since this sample is all about securing a **Kafka Streams** application we have chosen a very simple **Topology**, it's basically a **NO-OP** topology. We read input from topic **secure-input** and directly and unmodified output all the messages to the topic **secure-output**.

- Before we can run our **Kafka Streams** sample application we need to first manually create the topics **secure-input** and **secure-output** since we have configured the Kafka cluster to **disable** auto-create of topics (see setting in the **docker-compose.yml** file). The **AdminClient** that creates topics must also authenticate with the Kafka cluster, so we must use the **--command-config** option to point to the **.properties** file that contains the security configurations.

```

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9091 \
  --replication-factor 1 \
  --partitions 3 \
  --topic secure-input \
  --command-config /etc/kafka/secrets/client_security.properties

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9091 \
  --replication-factor 1 \
  --partitions 3 \
  --topic secure-output \
  --command-config /etc/kafka/secrets/client_security.properties

```


You will see some warnings such as



```
2018-10-15 18:12:39 WARN  ProducerConfig:287 - The
configuration
'confluent.monitoring.interceptor.security.protocol' was
supplied but isn't a known config.
```

You can safely ignore these warnings. The properties are necessary monitoring producers and consumers in Confluent Control Center, but are not recognized by the core open source Apache Kafka **ProducerConfig** and the **ConsumerConfig** classes.

8. Now we're ready to run the app. Use **Run** → **Start Debugging** to run your **Kafka Streams** application.

Creating Input Data

To see the sample streams application working we need to generate some data.

1. In the `~/confluent-streams/labs/security/secure-app/scripts/security` folder, there is a file called `client_security.properties` with the following content:

```

bootstrap.servers=kafka:9091
security.protocol=SASL_SSL
ssl.truststore.location=/etc/kafka/secrets/kafka.client.truststore.jks
ssl.truststore.password=confluent
sasl.mechanism=PLAIN
sasl.jaas.config=org.apache.kafka.common.security.plain.PlainLoginModule required \
    username=\"client\" \
    password=\"client-secret\";

# authenticate the monitor interceptor with Kafka.
confluent.monitoring.interceptor.bootstrap.servers=kafka:9091
confluent.monitoring.interceptor.security.protocol=SASL_SSL
confluent.monitoring.interceptor.ssl.truststore.location=/etc/kafka/secrets/kafka.client.truststore.jks
confluent.monitoring.interceptor.ssl.truststore.password=confluent
confluent.monitoring.interceptor.sasl.mechanism=PLAIN
confluent.monitoring.interceptor.sasl.jaas.config=org.apache.kafka.common.security.plain.PlainLoginModule required \
    username=\"client\" \
    password=\"client-secret\";

```

In an earlier step, we copied this file, along with others, to **/etc/kafka/secrets/**. This file will be used by the client tools that we're going to use in a moment to authenticate themselves, as it was earlier when we created the secure-input and secure-output topics.

2. Open a new terminal window, navigate to the **secure-app** folder:

```
$ cd ~/confluent-streams/labs/security/secure-app
```

3. To run the **kafka-console-producer** that we will use to feed some data to the topic **secure-input** use this command:

```

$ NS=io.confluent.monitoring.clients.interceptor && \
kafka-console-producer \
    --bootstrap-server kafka:9091 \
    --topic secure-input \
    --producer.config /etc/kafka/secrets/client_security.properties \
    --producer-property \
    interceptor.classes="${NS}.MonitoringProducerInterceptor"

```



- We're configuring the producer for monitoring via Control Center through the parameter **--producer-property**
- We're passing the **properties** file with the security settings through the parameter **--producer.config** to the producer.

4. Open another terminal window, navigate to the **secure-app** folder:

```
$ cd ~/confluent-streams/labs/security/secure-app
```

5. To list the data produced by our sample application to the topic **secure-output** use this command:

```
$ NS=io.confluent.monitoring.clients.interceptor && \
kafka-console-consumer \
  --group secure-console-consumer \
  --bootstrap-server kafka:9091 \
  --topic secure-output \
  --from-beginning \
  --consumer.config /etc/kafka/secrets/client_security.properties \
  --consumer-property \
  interceptor.classes="${NS}.MonitoringConsumerInterceptor"
```



- We're configuring the consumer for monitoring via Control Center through the parameter **--consumer-property**
- We're passing the **properties** file with the security settings through the parameter **--consumer.config** to the consumer

6. In the terminal where the producer is running enter a few lines such as:

```
Kafka is powering the Confluent streaming platform
Real time streaming is exceedingly important
Confluent offers support for Kafka Streaming and ksqlDB
In my company we build Kafka Streams applications
```

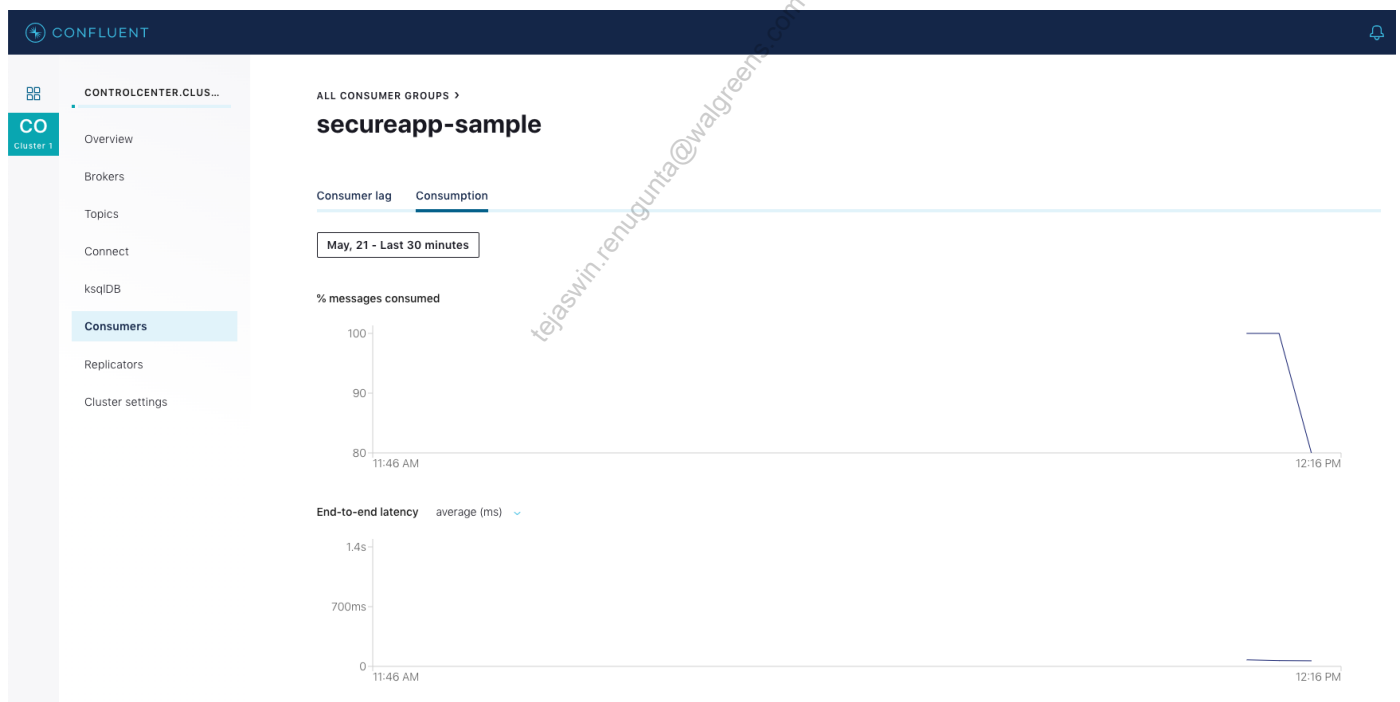
7. Observe how output is generated in the terminal window where the **kafka-console-consumer** is running. Specifically notice that the order of the output relative to the input might be changed (if you enter values fast enough) due to the fact that we have 3

partitions per topic and ordering is only guaranteed within a partition but not globally per topic!

8. **Optional:** Add a lot more data using the **kafka-console-producer**. For example:

```
$ NS=io.confluent.monitoring.clients.interceptor
$ for i in {1 .. 100}; do
    sleep 1
    seq 1000 | kafka-console-producer \
    --bootstrap-server kafka:9091 \
    --topic secure-input \
    --producer.config /etc/kafka/secrets/client_security.properties \
    --producer-property \
    interceptor.classes="${NS}.MonitoringProducerInterceptor"
done
```

9. In **Confluent Control Center** go to **Consumers** → **secureapp-sample** and you should see something like this:



Cleanup

1. Quit both the producer and consumer by pressing **Ctrl+C**.
2. Quit the sample **Kafka Streams** application with **Run** → **Stop Debugging**.
3. Shutdown the Kafka cluster:

```
$ docker-compose down -v
```

Conclusion

In this example we have shown how to run a secure Kafka cluster and then build a **Kafka Streams** application that integrates with this cluster using **SASL** for authentication and SSL for encryption. The application logic was trivial yet that is not the point of this example. The important fact is the integration with the Kafka cluster security settings.



STOP HERE. THIS IS THE END OF THE EXERCISE.

Lab 10 Monitoring Kafka Streams Applications

This lab contains 3 exercises:

- Getting Metrics from a Kafka Streams Application
- Using JConsole to monitor a Streams App
- Monitoring a Kafka Streams App in Confluent Control Center

a. Getting Metrics from a Kafka Streams Application

The purpose of this exercise is to learn how to expose metrics in a Kafka Streams application. This particular application sends metrics to standard output every 10 seconds. In practice, these metrics can be exposed to external sources for aggregation and analysis, as we will see in the exercises that follow.

Preparing the application

In this exercise we're going to use the **word count** exercise from a previous lab.

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. Navigate to the folder for this lab:

```
$ cd ~/confluent-streams/labs/monitoring
```

3. Run the Kafka cluster:

```
$ docker-compose up -d zookeeper kafka control-center
```

4. Create the two topics called **lines-topic** and **word-count-topic** in Kafka:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic word-count-topic
```

Building & Running the Application

1. In a terminal window navigate to the **word-count** folder, and launch VS Code:

```
$ cd ~/confluent-streams/labs/monitoring/word-count
$ code .
```

2. Open the file **build.gradle** in folder **word-count** and analyze its content. It should be quite familiar by now.
3. Notice the four Java files in subfolder **src/main/java/streams**:

WordCountSample.java	Main class of the application. This class makes use of the 3 following classes
ConfigProvider.java	Defines the configuration for the streams application
TopologyProvider.java	Defines the topology of the stream application
MetricsReporter.java	Defines how the list of metrics is output every 10 seconds by the app

Analyze their code. Specifically note the use of the **MetricsReporter** class. Make sure you understand what's going on in the code.

4. Use **Run** → **Start Debugging** in VS Code or **./gradlew run** in the terminal to run your code.

The application will print out the list of metrics to the terminal every 10 seconds. It should look similar to this (shortened for readability):

```
--- Application Metrics ---
MetricName [name=count, group=kafka-metrics-count, description=total
number of registered metrics, tags={client-id=wordCount-3e02d8d9-
490b-492a-9bf0-cf85629e7fcf-StreamThread-1-producer}], 81.0
MetricName [name=io-time-ns-avg, group=producer-metrics,
description=The average length of time for I/O per select call in
nanoseconds., tags={client-id=wordCount-3e02d8d9-490b-492a-9bf0-
cf85629e7fcf-StreamThread-1-producer}], 612250.0
```

Producing Input Data

To see how the metrics change, we can produce some input data that the application will process. We use the **kafka-console-producer** tool for this job.

1. Return to the terminal window and create a list of input sentences that will be randomly produced to the input topic:

```
$ INPUTS=('Kafka powers the Confluent streaming platform' \
'All Streams come from Kafka'\
'Streams will all flow to Kafka'\
'Follow the streams to Kafka Summit' \
'Check out Confluent Cloud')
```

2. Run **Kafkacat** to produce a steady flow of sentences to the input topic:

```
$ while true; do
    MESSAGE=${INPUTS[${RANDOM} % ${#INPUTS[@]}]}
    echo ${MESSAGE} | kafkacat -P \
        -b kafka:9092 \
        -t lines-topic
    sleep 0.1
done
```

3. Observe the metric values printed by the Word Count application.
4. In another terminal window run **kafka-console-consumer** to report the output of our

sample app:

```
$ kafka-console-consumer --bootstrap-server kafka:9092 \  
  --topic word-count-topic \  
  --property print.key=true \  
  --value-deserializer  
org.apache.kafka.common.serialization.LongDeserializer
```

Cleanup

1. Quit the Producer with **Ctrl+C**.
2. Quit the Consumer with **Ctrl+C**.
3. Quit our sample app with **Run → Stop Debugging**.
4. Shut down your Kafka cluster with the **docker-compose down -v** command.

tejaswin.renugunta@walgreens.com

b. Using JConsole to monitor a Streams App

In this exercise, we will use **JConsole** to monitor the various metrics a simple **Kafka Streams** application exposes.

We will be using an application we created in an earlier lab that reads data from a topic whose keys are integers and whose values are sentence strings. The input values are transformed to lower-case and output to a new topic.

1. Navigate to the module folder:

```
$ cd ~/confluent-streams/labs/monitoring
```

2. Run the Kafka cluster:

```
$ docker-compose up -d zookeeper kafka
```

3. Create an input topic called **lines-topic** in Kafka:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-topic
```

4. Create the output topic called **lines-lower-topic** in Kafka:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-lower-topic
```

5. In a terminal window navigate to the project folder **jmx-sample** and build the artifact:

```
$ cd ~/confluent-streams/labs/monitoring/jmx-sample
$ ./gradlew build
```

6. Examine **build.gradle** and observe it has been configured to expose JMX metrics on port **4444** that we can then use to attach **JConsole**:

```
$ cat build.gradle
...
applicationDefaultJvmArgs = [
    "-Dcom.sun.management.jmxremote",
    "-Dcom.sun.management.jmxremote.authenticate=false",
    "-Dcom.sun.management.jmxremote.ssl=false",
    "-Djava.rmi.server.hostname=127.0.0.1",
    "-Dcom.sun.management.jmxremote.rmi.port=4444",
    "-Dcom.sun.management.jmxremote.port=4444"]
...
```

7. Still in the project folder **jmx-sample** run your **Kafka Streams** application:

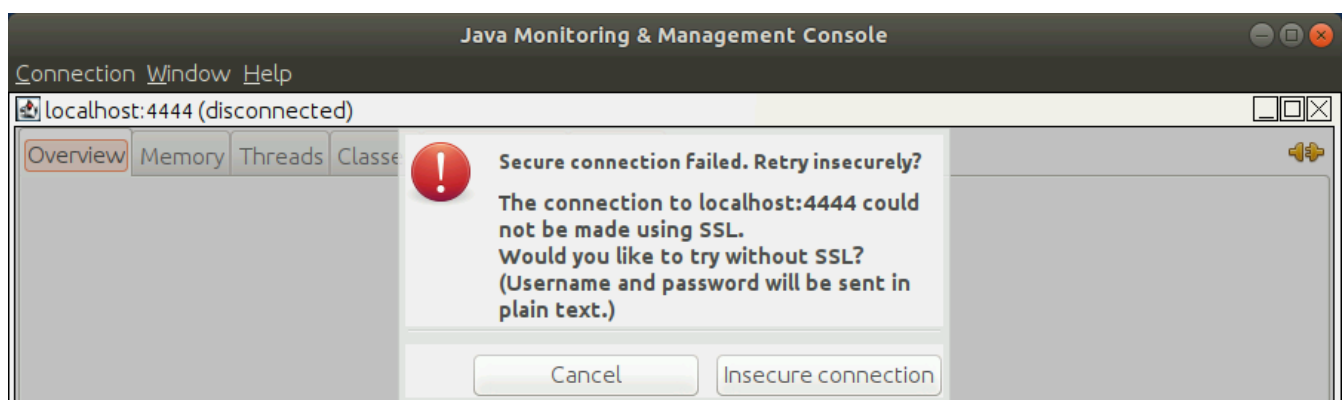
```
$ ./gradlew run
```

8. Observe the JMX metrics of the **Kafka Streams** application:

- a. Open a **jconsole** connection to port **4444** which is the JMX port for the **Kafka Streams** application.

```
$ jconsole localhost:4444 &
```

- b. Select **Insecure connection** when asked



- c. Navigate to the **MBeans** tab
- d. Explore the node under **kafka-streams** as indicated in the below image:

The screenshot shows the JMX console for 'localhost:4444'. The left pane displays a tree view of metrics. The right pane shows a table of attribute values.

Name	Value
commit-latency-avg	NaN
commit-latency-max	NaN
commit-rate	0.0
commit-ratio	0.0
commit-total	0.0
poll-latency-avg	NaN
poll-latency-max	NaN
poll-rate	0.0
poll-ratio	1.0
poll-records-avg	NaN
poll-records-max	NaN
poll-total	0.0
process-latency-avg	NaN
process-latency-max	NaN
process-rate	0.0
process-ratio	0.0
process-records-avg	0.0
process-records-max	0.0
process-total	0.0
punctuate-latency-avg	NaN
punctuate-latency-max	NaN
punctuate-rate	0.0
punctuate-ratio	0.0
punctuate-total	0.0
task-closed-rate	0.0
task-closed-total	0.0
task-created-rate	0.0
task-created-total	1.0

Initially some of the numbers, such as **process total** will be zero.

9. Open a new terminal window and navigate to the **monitoring** folder:

```
$ cd ~/confluent-streams/labs/monitoring
```

10. Now create some data that will be consumed and processed by the **Kafka Streams** application:

```
$ cat << EOF | kafka-console-producer \
  --bootstrap-server kafka:9092 \
  --property "parse.key=true" \
  --property "key.separator=:" \
  --topic lines-topic
1:"Kafka powers the Confluent Streaming Platform"
2:"Events are stored in Kafka"
3:"Confluent contributes to Kafka"
EOF
```

11. Observe how the values of the metrics in **JConsole** change (you will need to click the **Refresh** button):

The screenshot shows the JConsole interface for a Kafka consumer. The left pane displays the MBeans tree, with the following structure expanded:

- kafka-metrics-count
 - Attributes
 - count
 - stream-metrics
 - map-sample-v0.1.0-5d7803ce-d935-4946-8e7e-8f0000000002
 - Attributes
 - topology-description
 - state
 - alive-stream-threads
 - commit-id
 - version
 - application-id

The right pane displays the 'Attribute values' table:

Name	Value
commit-latency-avg	17.0
commit-latency-max	17.0
commit-rate	0.0283077619883372
commit-ratio	0.0
commit-total	1.0
poll-latency-avg	114.0
poll-latency-max	114.0
poll-rate	0.0224804981678394
poll-ratio	1.0
poll-records-avg	3.0
poll-records-max	3.0
poll-total	1.0
process-latency-avg	7.0
process-latency-max	14.0
process-rate	0.06746424395070612
process-ratio	0.0
process-records-avg	0.008310249307479225
process-records-max	3.0
process-total	3.0
punctuate-latency-avg	NaN
punctuate-latency-max	NaN
punctuate-rate	0.0
punctuate-ratio	0.0
punctuate-total	0.0
task-closed-rate	0.0
task-closed-total	0.0
task-created-rate	0.0
task-created-total	1.0

A 'Refresh' button is located at the bottom right of the interface.

12. **Optional:** add more data to the topic and monitor the attributes and how the values

change.

Cleanup

1. Quit **JConsole**.
2. Stop the sample Kafka Streams application by pressing **Ctrl+C**.
3. Shut down your Kafka cluster with the **docker-compose down -v** command.

Conclusion

We have used Java code to retrieve the metrics directly from the **KafkaStreams** object and also configured the **Kafka Streams** sample application to expose JMX data that we then explored using **JConsole**.

tejaswin.renugunta@walgreens.com

c. Monitoring a Kafka Streams App in Confluent Control Center

In this exercise we're going to reuse the **word count** processor API example application from a previous exercise and extend it so that it can be monitored in **Confluent Control Center**.

Prerequisites

1. Please make sure you have prepared your lab environment as described here: → [Lab Environment](#)
2. In a terminal window navigate to the **processor-sample** folder:

```
$ cd ~/confluent-streams/labs/monitoring/processor-sample
```

3. Open the **build.gradle** file and observe **monitoring-interceptors** in the list of dependencies:

```
compile group: "io.confluent", name: "monitoring-interceptors",  
version: "6.0.0"
```

This dependency contains the interceptors classes that we will use to configure our application for monitoring via **Confluent Control Center**.

4. Open this application's root directory in VS Code.

```
$ code .
```

5. Open the file **ProcessorSample.java** (in folder **src/main/java/streams**) which contains the **main** function of the sample application. Notice these lines that have been added to the **getConfig** function:

```
settings.put(StreamsConfig.producerPrefix(ProducerConfig.INTERCEPTOR_CLASSES_CONFIG),

"io.confluent.monitoring.clients.interceptor.MonitoringProducerInterceptor");
settings.put(StreamsConfig.consumerPrefix(ConsumerConfig.INTERCEPTOR_CLASSES_CONFIG),

"io.confluent.monitoring.clients.interceptor.MonitoringConsumerInterceptor");
```



Please note the **prefix** used for both consumer and producer interceptors.

With interceptors configured, we can monitor our app through **Confluent Control Center**.

6. Start the Kafka cluster and the **Confluent Control Center** server with:

```
$ docker-compose up -d zookeeper kafka control-center
```

As usual, you are encouraged to inspect the content of the **docker-compose** file in the `~confluent-streams` directory to make sure you understand all the settings and discuss it with your peers. Wait a couple of minutes until the cluster is initialized. Open **Confluent Control Center** at <http://localhost:9021> and wait until it displays the system health.

7. Create the input and output topics called **lines-topic** and **word-count-topic** respectively. For the moment let's give them each 1 partition and replication factor 1:

```
$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic lines-topic

$ kafka-topics \
  --create \
  --bootstrap-server kafka:9092 \
  --replication-factor 1 \
  --partitions 1 \
  --topic word-count-topic
```

8. Use **Run** → **Start Debugging** in VS Code or **./gradlew run** in the terminal to run your

streams app.

Producing Data

1. Run the **kafka-console-producer** that we will use to feed some data to the topic **lines-topic**:

```
$ NS=io.confluent.monitoring.clients.interceptor && \
kafka-console-producer \
  --bootstrap-server kafka:9092 \
  --topic lines-topic \
  --producer-property \
  "interceptor.classes=${NS}.MonitoringProducerInterceptor"
```

2. Open another terminal tab and run the **kafka-console-consumer** for the **word-count-topic** topic:

```
$ NS=io.confluent.monitoring.clients.interceptor && \
kafka-console-consumer \
  --group word-count-consumer \
  --bootstrap-server kafka:9092 \
  --topic word-count-topic \
  --property print.key=true \
  --consumer-property \
  "interceptor.classes=${NS}.MonitoringConsumerInterceptor"
```

Note the use of the group name **word-count-consumer** and the **MonitoringConsumerInterceptor** class to enable monitoring of the consumer in Control Center.

3. In the terminal window where the producer runs, enter a few lines of text such as:

```
Kafka is powering the Confluent streaming platform
Streaming in real-time is more and more important
Our company will invest in real-time streaming
For this we need to know loads about Kafka and Kafka Streams
ksqlDB is a simpler alternative to Kafka Streams
Everybody loves ksqlDB since no programming is required
```

and observe the output in the terminal window where the consumer runs.



Add more data at will so that there is some activity ongoing.

4. In **Confluent Control Center** (<http://localhost:9021>) navigate to **Consumers**. You should see something like this:

All consumer groups

Search consumer groups

Consumer group ID	Messages behind	Number of consumers	Number of topics
secure-console-consumer	0	1	1
_confluent-controlcenter-5-5-0-1-command	0	1	1
processor-sample-v0.1.0	0	1	1
_confluent-controlcenter-5-5-0-1	215	1	15

Click on our Kafka Streams (processor-sample-v0.1.0) application to investigate its metrics.

Cleanup

1. To stop the producer hit **Ctrl+C**.
2. To stop the consumer hit **Ctrl+C**.
3. To stop the **Kafka Streams** application with **Run** → **Stop Debugging**.
4. To stop the Kafka Cluster and delete the orphaned volumes execute this command:

```
$ cd ~/confluent-streams/labs/monitoring/processor-sample
$ docker-compose down -v
```



STOP HERE. THIS IS THE END OF THE EXERCISE.

tejaswin.renugunta@walgreens.com

Appendix A: Running All Labs with Docker

Running Labs in Docker for Desktop

If you have installed Docker for Desktop on your Mac or Windows 10 Pro machine you are able to complete the course by building and running your applications from the command line.

- Increase the memory available to Docker Desktop to a minimum of 6 GiB. See the advanced settings for [Docker Desktop for Mac](#), and [Docker Desktop for Windows](#).
- Follow the instructions at → [The Lab Environment & Sample Solutions](#) to **git clone** the source code, in each exercise follow the instructions to launch the cluster containers with **docker-compose** on your host machine. The exercise source code will now be on your host machine where you can edit the source code with any editor.
- Begin the exercises by first opening a bash shell on the tools container. All the command line instructions will work from the tools container. This container has been preconfigured with all of the tools you use in the exercises, e.g. **kafka-topics** and **python**.

```
$ docker-compose exec tools bash
bash-4.4#
```

- At the time of writing, the Python prerequisite **confluent_kafka** won't install with the version of Python 3 available in the tools container. You should use **pip** and **python** rather than **pip3** and **python3**.
- At the time of writing, Maven is not installed in the tools container, so the optional maven exercise would require a hefty 300 MB download with **apt-get install maven**. It might be best to skip this optional exercise and simply use it as reference if you use Maven in your day-to-day work.
- Anywhere you are instructed to open additional terminal windows you can **exec** additional bash shells on the tools container with the same command as above on your host machine.

- Any subsequent **docker** or **docker-compose** instructions should be run on your host machine.

Running the Exercise Applications

From the **tools** container you can use command line alternatives to the VS Code steps used in the instructions. Complete the exercise code with an editor on your host machine, then use the following command line instructions to build and run the applications from the exercise directory.

- For Java applications: **./gradlew run**
- For Python applications: **python main.py**

Where you are instructed to use **Run** → **Stop Debugging** in VS code, use **Ctrl+C** to end the running exercise.

Our **docker-compose.yml** file sets the working directory inside the tools container to the **~/confluent-streams directory**. Be sure to change the working directory to each exercise directory as stated in each exercise instructions.

```
$ docker-compose exec tools bash
bash-4.4# pwd
/root/confluent-streams
```

To build and run the **Anatomy of a Kafka Streams App** exercise. First make the source code updates to **~/confluent-streams/labs/streams-writing/gradle-sample/src/main/java/streams/MapSample.java** on your host machine. Then enter into the bash shell on your tools container:

```
bash-4.4# cd labs/streams-writing/gradle-sample
bash-4.4# ./gradlew run
```

In the **Monitoring Kafka Streams Applications** exercise, you must use **jconsole** on port 4444 on the host to view JMX metrics. If **jconsole** is not already installed on your host system, install using the password **training**:

```
$ sudo apt install -y openjdk-11-jdk
```

Run a new tools container with a port mapping to expose JMX metrics to the host:

```
$ docker-compose run -p 4444:4444 tools  
bash-4.4# cd jmx-sample && ./gradlew run
```

The application will now expose metrics to port 4444 in the container, which is mapped to port 4444 on the host. Running **jconsole** on the host on port 4444 will now pick up the metrics exposed by the application.

tejaswin.renugunta@walgreens.com