

Chicago Java Users Group
&
Chicago Advanced Analytics Meetup

June 8th 2017



Kafka Streams For Java Enthusiasts



Slim Baltagi



Agenda

1. Apache Kafka: a Streaming Data Platform
2. Overview of Kafka Streams
3. Writing, deploying and running your first Kafka Streams application
4. Code and Demo of an end-to-end Kafka-based Streaming Data Application
5. Where to go from here for further learning?

Batch



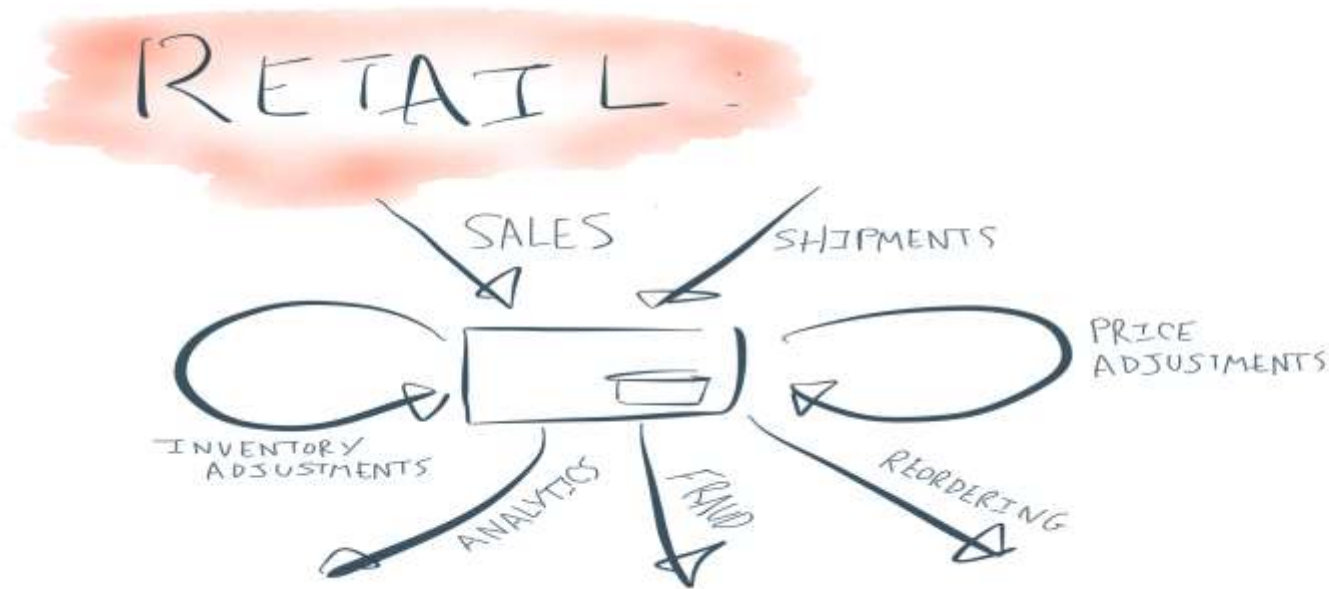
Streaming



1. Apache Kafka: a Streaming Data Platform

➤ Most of **what a business does** can be thought as **event streams**. They are in a

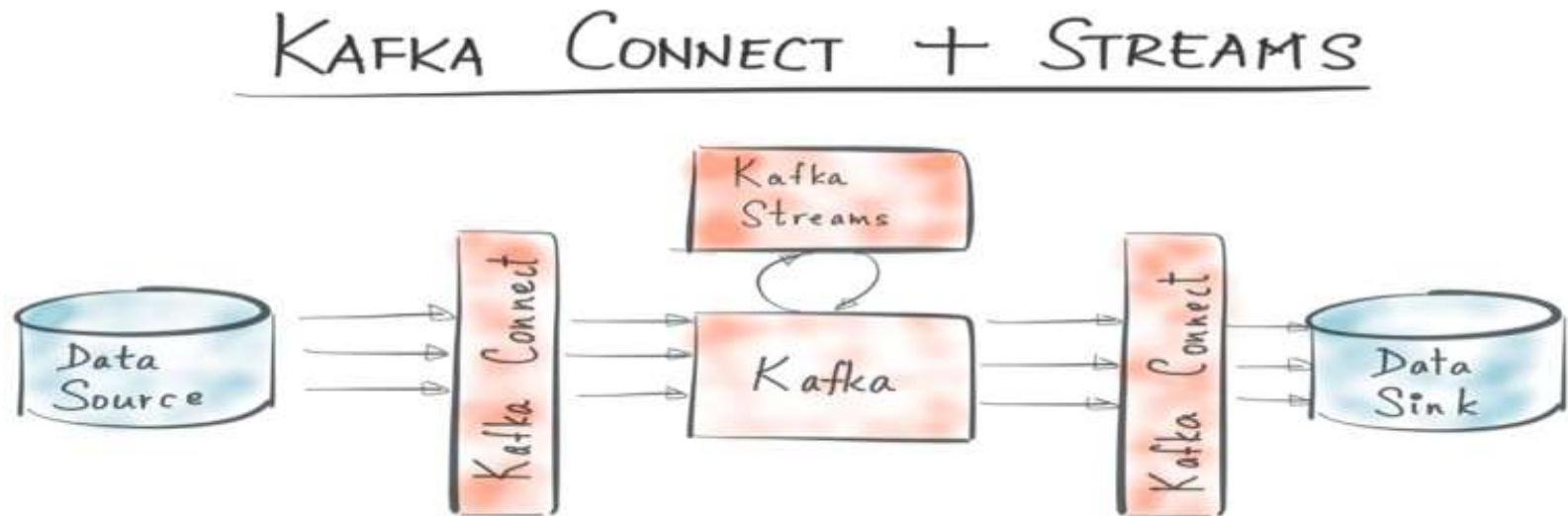
- **Retail system**: orders, shipments, returns, ...
 - **Financial system**: stock ticks, orders, ...
 - **Web site**: page views, clicks, searches, ...
 - **IoT**: sensor readings, ...
- and so on.



1. Apache Kafka: a Streaming Data Platform

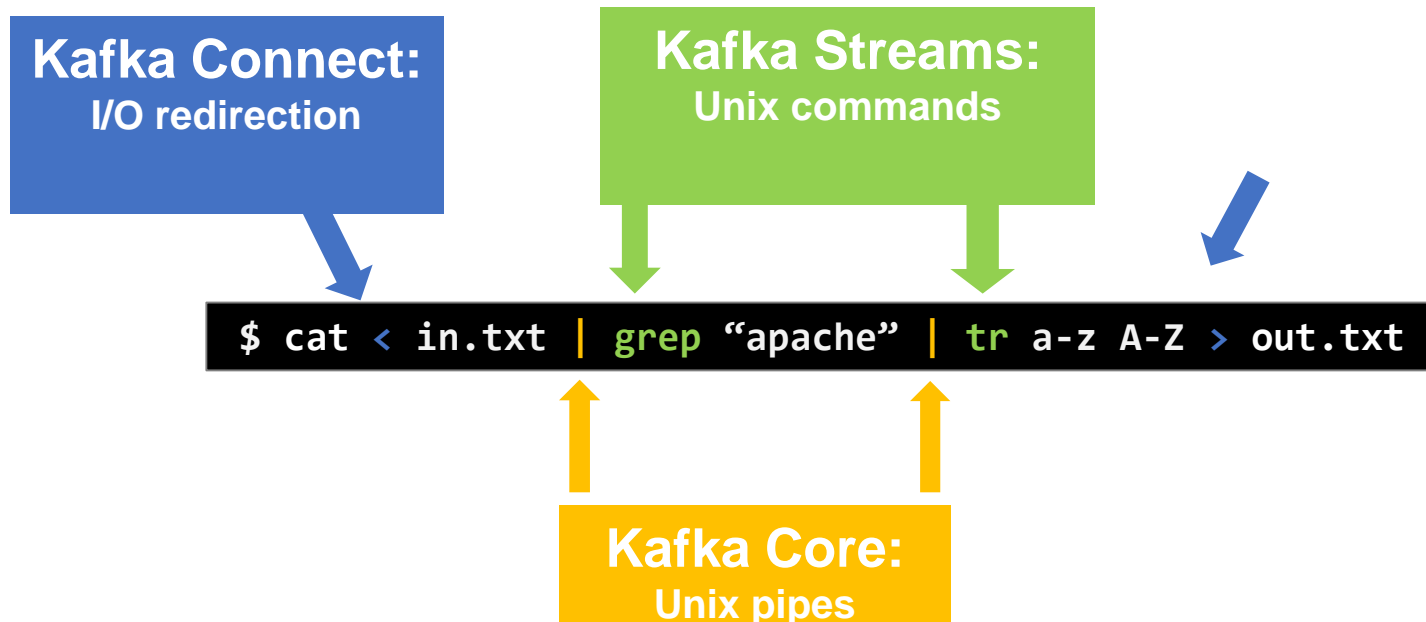
➤ **Apache Kafka** is an open source streaming data platform (a new category of software!) with **3 major components**:

1. **Kafka Core**: A **central hub** to **transport** and **store** event streams in real-time.
2. **Kafka Connect**: A **framework** to **import** event streams from other source data systems into Kafka and **export** event streams from **Kafka** to destination data systems.
3. **Kafka Streams**: A **Java library** to **process** event streams live as they occur.



1. Apache Kafka: a Streaming Data Platform

Unix Pipelines Analogy



- **Kafka Core:** is the **distributed, durable equivalent** of Unix pipes. Use it to connect and compose your large-scale data applications.
- **Kafka Streams** are the commands of your Unix pipelines. Use it to transform data stored in Kafka.
- **Kafka Connect** is the I/O redirection in your Unix pipelines. Use it to get your data into and out of Kafka.

2. Overview of Kafka Streams

2.1 Before Kafka Streams?

2.2 What is Kafka Streams?

2.3 Why Kafka Streams?

2.4 What are Kafka Streams key concepts?

2.5 Kafka Streams APIs and code examples?

2.1 Before Kafka Streams?

- **Before Kafka Streams**, to process the data in Kafka you have 4 options:
- **Option 1: Dot It Yourself (DIY)** – Write your own ‘stream processor’ using Kafka client libs, typically with a narrower focus.
 - **Option 2: Use a library** such as AkkaStreams-Kafka, also known as Reactive Kafka, RxJava, or Vert.x
 - **Option 3: Use an existing open source stream processing framework** such as Apache Storm, Spark Streaming, Apache Flink or Apache Samza for transforming and combining data streams which live in Kafka...
 - **Option 4: Use an existing commercial tool for** stream processing with adapter to Kafka such as IBM InfoSphere Streams, TIBCO StreamBase, ...
- Each one of the 4 options above of processing data in Kafka has advantages and disadvantages.

2.2 What is Kafka Streams?

- Available since **Apache Kafka 0.10 release in May 2016**, **Kafka Streams** is a **lightweight open source Java library** for building stream processing applications on top of Kafka.
- Kafka Streams is designed to **consume** from & **produce data** to **Kafka topics**.
- It provides a **Low-level API** for building topologies of processors, streams and tables.
- It provides a **High-Level API** for common patterns like filter, map, aggregations, joins, stateful and stateless processing.
- Kafka Streams **inherits operational characteristics** (low latency, elasticity, fault-tolerance, ...) from Kafka.
- **A library is simpler than a framework** and is easy to integrate with your existing applications and services!
- Kafka Streams **runs in your application code** and imposes no change in the Kafka cluster infrastructure, or within Kafka.

What is Kafka Streams? Java analogy

1996	1 core	java.lang
2004	multi-core	java.util.concurrent
2016	multi-machine	java.distributed org.apache.kafka.streams

2.3 Why Kafka Streams?

- Processing data in Kafka with Kafka Streams has the following advantages:
- **No need to run another framework or tool** for stream processing as Kafka Streams is already a library included in Kafka
 - **No need of external infrastructure** beyond Kafka. Kafka is already your cluster!
 - **Operational simplicity** obtained by getting rid of an additional stream processing cluster
 - As a normal library, it is **easier to integrate** with your existing applications and services
 - **Inherits Kafka features such as** fault-tolerance, scalability, elasticity, authentication, authorization
 - **Low barrier to entry**: You can quickly write and run a small-scale proof-of-concept on a single machine

2.4 Wat are Kafka Streams key concepts?

- **KStream** and **KTable** as the two basic abstractions. The distinction between them comes from how the key-value pairs are interpreted:
 - In a **stream**, each key-value is an **independent piece of information**. For example, in a stream of user addresses: Alice -> New York, Bob -> San Francisco, Alice -> Chicago, we know that Alice lived in both cities: New York and Chicago.
 - If the **table** contains a **key-value pair** for the **same key twice, the latter overwrites the mapping**. For example, a table of user addresses with Alice -> New York, Bob -> San Francisco, Alice -> Chicago means that Alice moved from New York to Chicago, not that she lives at both places at the same time.
- There's a **duality between the two concepts**: a stream can be viewed as a table, and a table as a stream. See more on this in the documentation:
<http://docs.confluent.io/current/streams/concepts.html#duality-of-streams-and-tables>

KStream vs KTable

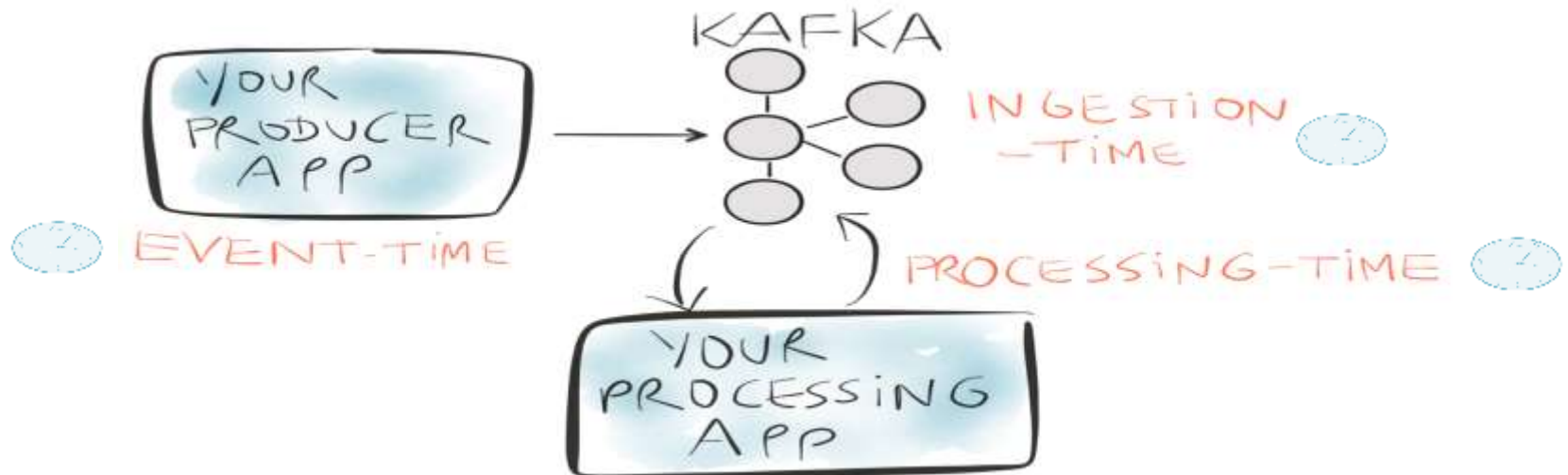
Example	When you need...	then you'd read the Kafka topic into a	so that the topic is interpreted as a	with messages interpreted as
All the cities Alice has ever lived in	All the values of a key	KStream	record stream	INSERT (append)
In what city Alice lives right now?	Latest value of a key	KTable	changelog stream	UPDATE (overwrite existing)

KStream = immutable log

KTable = mutable materialized view

2.4 What are Kafka Streams key concepts?

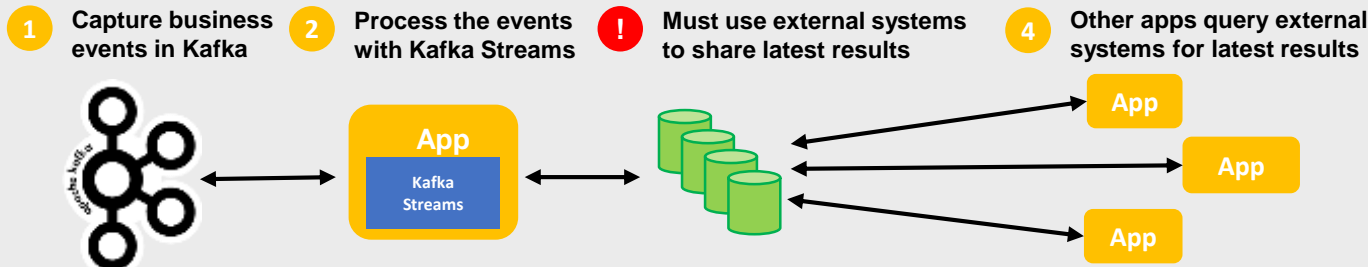
- **Event Time**: A critical aspect in stream processing is the notion of **time**, and how it is modeled and integrated.
 - **Event time**: The point in time when an **event** or **data record** occurred, i.e. was originally **created** “by the source”.
 - **Ingestion time**: The point in time when an **event** or **data record** is **stored** in a topic partition by a Kafka broker.
 - **Processing time**: The point in time when the **event** or **data record** happens to be **processed** by the stream processing application, i.e. when the record is being consumed.



2.4 What are Kafka Streams key concepts?

➤ Interactive Queries: Local queryable state

Before (0.10.0)



After (0.10.1): simplified, more app-centric architecture

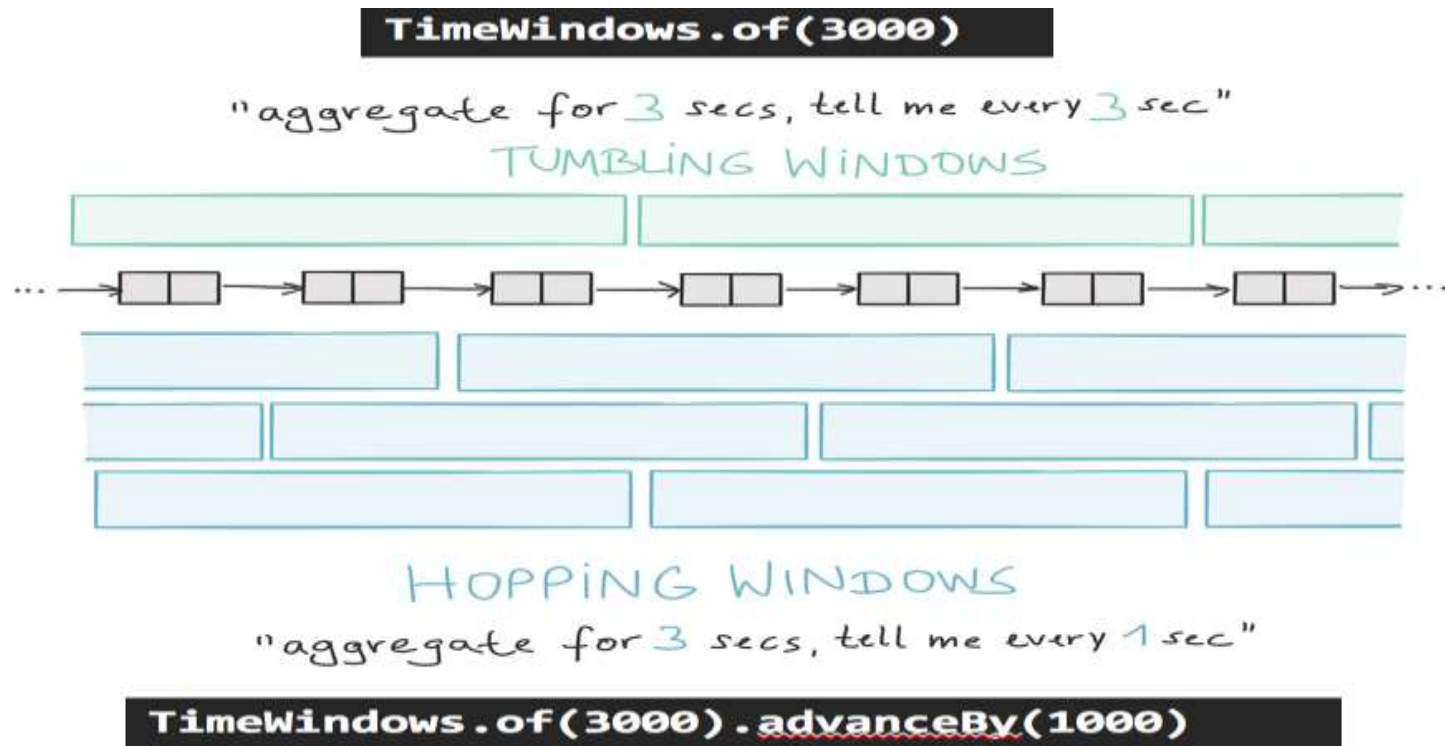


See blogs:

- **Why local state is a fundamental primitive in stream processing?** Jay Kreps, July 31st 2014
<https://www.oreilly.com/ideas/why-local-state-is-a-fundamental-primitive-in-stream-processing>
- **Unifying Stream Processing and Interactive Queries in Apache Kafka**, Eno Thereska, October 26th 2016
<https://www.confluent.io/blog/unifying-stream-processing-and-interactive-queries-in-apache-kafka/>

2.4 What are Kafka Streams key concepts?

- **Windowing**: Windowing lets you control how to group records that have the same key for stateful operations such as **aggregations** or **joins** into so-called windows.



- **More concepts** in Kafka Streams documentation:
<http://docs.confluent.io/current/streams/concepts.htm>

2.5 Kafka Streams APIs and code examples?

API option 1: DSL (high level, declarative)

```
KStream<Integer, Integer> input =  
    builder.stream("numbers-topic");  
  
// Stateless computation  
KStream<Integer, Integer> doubled =  
    input.mapValues(v -> v * 2);  
  
// Stateful computation  
KTable<Integer, Integer> sumOfOdds = input  
    .filter((k,v) -> v % 2 != 0)  
    .selectKey((k, v) -> 1)  
    .groupByKey()  
    .reduce((v1, v2) -> v1 + v2, "sum-of-odds");
```

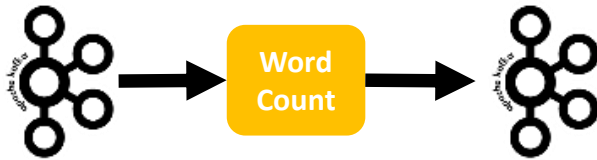
The **preferred API** for most use cases.

The DSL particularly appeals to users:

- familiar with Spark, Flink, Beam
- fans of Scala or functional programming

- If you're used to the functions that real-time processing systems like Apache Spark, Apache Flink, or Apache Beam expose, you'll be right at home in the DSL.
- If you're not, you'll need to spend some time understanding what methods like **map**, **flatMap**, or **mapValues** mean.

Code Example 1: complete app using DSL



```
1 public static void main(final String[] args) throws Exception {
2     Properties config = new Properties();
3     config.put(StreamsConfig.APPLICATION_ID_CONFIG, "wordcount-example");
4     config.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "kafka-broker1:9092");
5     config.put(StreamsConfig.KEY_SERDE_CLASS_CONFIG, Serdes.String().getClass().getName());
6     config.put(StreamsConfig.VALUE_SERDE_CLASS_CONFIG, Serdes.String().getClass().getName());
7
8     KStreamBuilder builder = new KStreamBuilder();
9     KStream<String, String> textLines = builder.stream("TextLinesTopic");
10    KStream<String, Long> wordCounts = textLines
11        .flatMapValues(value -> Arrays.asList(value.toLowerCase().split("\\W+")))
12        .groupBy((key, word) -> word)
13        .count("Counts")
14        .toStream();
15    wordCounts.to(Serdes.String(), Serdes.Long(), "WordsWithCountsTopic");
16
17    KafkaStreams streams = new KafkaStreams(builder, config);
18    streams.start();
19 }
```

**App
configuration**

**Define
processing
(here:
WordCount)**

Start processing

API option 2: Processor API (low level, imperative)

```
class PrintToConsoleProcessor
    implements Processor<K, V> {

    @Override
    public void init(ProcessorContext context) {}

    @Override
    void process(K key, V value) {
        System.out.println("Got value " + value);
    }

    @Override
    void punctuate(long timestamp) {}

    @Override
    void close() {}
}
```

Full flexibility but **more manual work**:

- The **Processor API** appeals to users:
 - **familiar with Storm, Samza**
 - Still, check out the DSL!
 - requiring **functionality that is not yet available in the DSL**
- Some people have begun using the low-level Processor API to **port their Apache Storm code to Kafka Streams**.

Code Example 2: Complete app using Processor API

```
public PrintToConsoleProcessor implements Processor<K, V> {
```

```
    @Override  
    public void init(ProcessorContext context) {  
        // No initialization needed in this case.  
    }  
}
```

Startup

```
    @Override  
    public void process(K key, V value) {  
        System.out.println("Received data record with " +  
            "key=" + key + ", value=" + value);  
    }  
}
```

Process a record

```
    @Override  
    public void punctuate(long timestamp) {  
        // No periodic actions needed in this case.  
    }  
}
```

Periodic action

```
    @Override  
    public void close() {  
        // No shutdown logic needed in this case.  
    }  
}
```

Shutdown

```
}
```

3. Writing, deploying and running your first Kafka Streams application

- **Step 1:** Ensure Kafka cluster is accessible and has data to process
- **Step 2:** Write the application code in Java or Scala
- **Step 3:** Packaging and deploying the application
- **Step 4:** Run the application

Step 1: Ensure Kafka cluster is accessible and has data to process

- Get the **input data into Kafka** via:
 - Kafka Connect (part of Apache Kafka)
 - or your own application that write data into Kafka
 - or tools such as StreamSets, Apache Nifi, ...
- **Kafka Streams** will then be used to **process the data** and write the results back to Kafka.



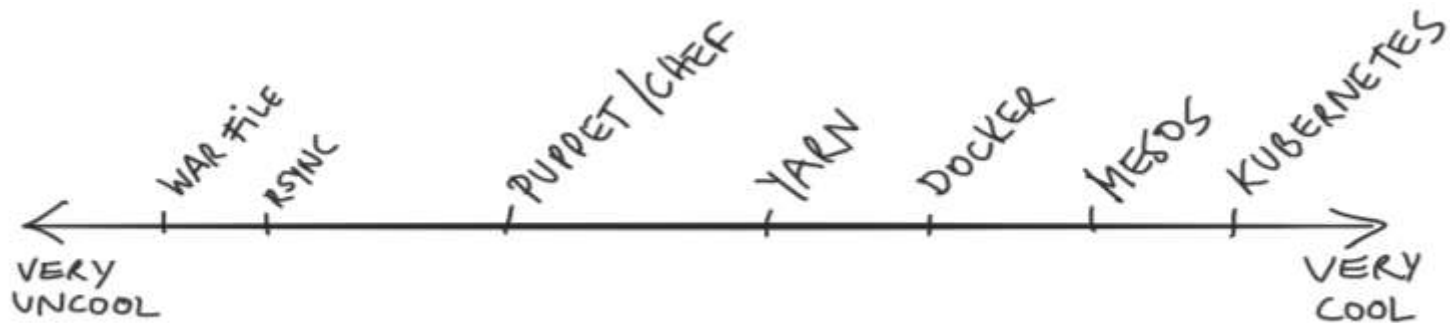
Step 2: Write the application code in Java or Scala

- How to start?
 - Learn from existing code examples:
<https://github.com/confluentinc/examples>
 - Documentation: <http://docs.confluent.io/current/streams/>
- How do I install Kafka Streams?
 - **There is no “installation”!** It's a Java library. Add it to your client applications like any other Java library.
 - Example adding 'kafka-streams' library using Maven:

```
<dependency>  
  <groupId>org.apache.kafka</groupId>  
  <artifactId>kafka-streams</artifactId>  
  <version>0.10.2.0</version>  
</dependency>
```

Step 3: Packaging and deploying the application

- How do you package and deploy your Kafka Streams apps?
 - **Whatever works for you!** Stick to what you/your company think is the best way for deploying and packaging a java application.
 - Kafka Streams integrates well with what you already use because an application that uses **Kafka Streams is a normal Java application.**



Step 4: Run the application

- **You don't need to install a cluster** as in other stream processors (Storm, Spark Streaming, Flink, ...) and **submit jobs to it!**
- **Kafka Streams runs as part of your client applications**, it does not run in the Kafka brokers.
- In production, **bundle as fat jar**, then ``java -cp my-fatjar.jar com.example.MyStreamsApp``
<http://docs.confluent.io/current/streams/developer-guide.html#running-a-kafka-streams-application>
- TIP: During development from your IDE or from CLI, the '**Kafka Streams Application Reset Tool**', available since Apache Kafka 0.10.0.1, is great for playing around.
<https://cwiki.apache.org/confluence/display/KAFKA/Kafka+Streams+Application+Reset+Tool>

Example: complete app, ready for production at large-scale!

```
public static void main(String[] args) throws Exception {
    Properties config = new Properties();
    config.put(StreamsConfig.JOB_ID_CONFIG, "wordcount-lambda-example");
    config.put(StreamsConfig.BootstrapServersConfig, "localhost:9092");
    config.put(StreamsConfig.ZooKeeperConnectConfig, "localhost:2181");
    config.put(StreamsConfig.KeySerializerClassConfig, StringSerializer.class);
    config.put(StreamsConfig.KeyDeserializerClassConfig, StringDeserializer.class);
    config.put(StreamsConfig.ValueSerializerClassConfig, StringSerializer.class);
    config.put(StreamsConfig.ValueDeserializerClassConfig, StringDeserializer.class);

    final Serializer<String> stringSerializer = new StringSerializer();
    final Deserializer<String> stringDeserializer = new StringDeserializer();
    final Serializer<Long> longSerializer = new LongSerializer();
    final Deserializer<Long> longDeserializer = new LongDeserializer();

    KStreamBuilder builder = new KStreamBuilder();
    KStream<String, String> textLines = builder.stream(stringDeserializer, stringDeserializer, "TextLinesTopic");

    KStream<String, Long> wordCounts = textLines
        .flatMapValues(value -> Arrays.asList(value.toLowerCase().split("\\W+")))
        .map((key, value) -> new KeyValue<>(value, value))
        .countByKey(stringSerializer, longSerializer, stringDeserializer, longDeserializer, "Counts")
        .toStream();
    wordCounts.to("WordsWithCountsTopic", stringSerializer, longSerializer);

    KafkaStreams streams = new KafkaStreams(builder, config);
    streams.start();
}
```

4. Code and Demo of an end-to-end Kafka-based Streaming Data Application

- 4.1 Scenario of this demo
- 4.2 Architecture of this demo
- 4.3 Setup of this demo
- 4.4 Results of this demo
- 4.5 Stopping the demo!

4.1. Scenario of this demo

➤ This demo consists of:

- reading live stream of data (tweets) from Twitter using **Kafka Connect** connector for Twitter
- storing them in Kafka broker leveraging **Kafka Core** as publish-subscribe message system.
- performing some basic stream processing on tweets in Avro format from a Kafka topic using **Kafka Streams** library to do the following:
 - **Raw word count** - every occurrence of individual words is counted and written to the **topic wordcount** (a predefined list of stopwords will be ignored)
 - **5-Minute word count** - words are counted per 5 minute window and every word that has more than 3 occurrences is written to the **topic wordcount5m**
 - **Buzzwords** - a list of special interest words can be defined and those will be tracked in the **topic buzzwords**

4.1. Scenario of this demo

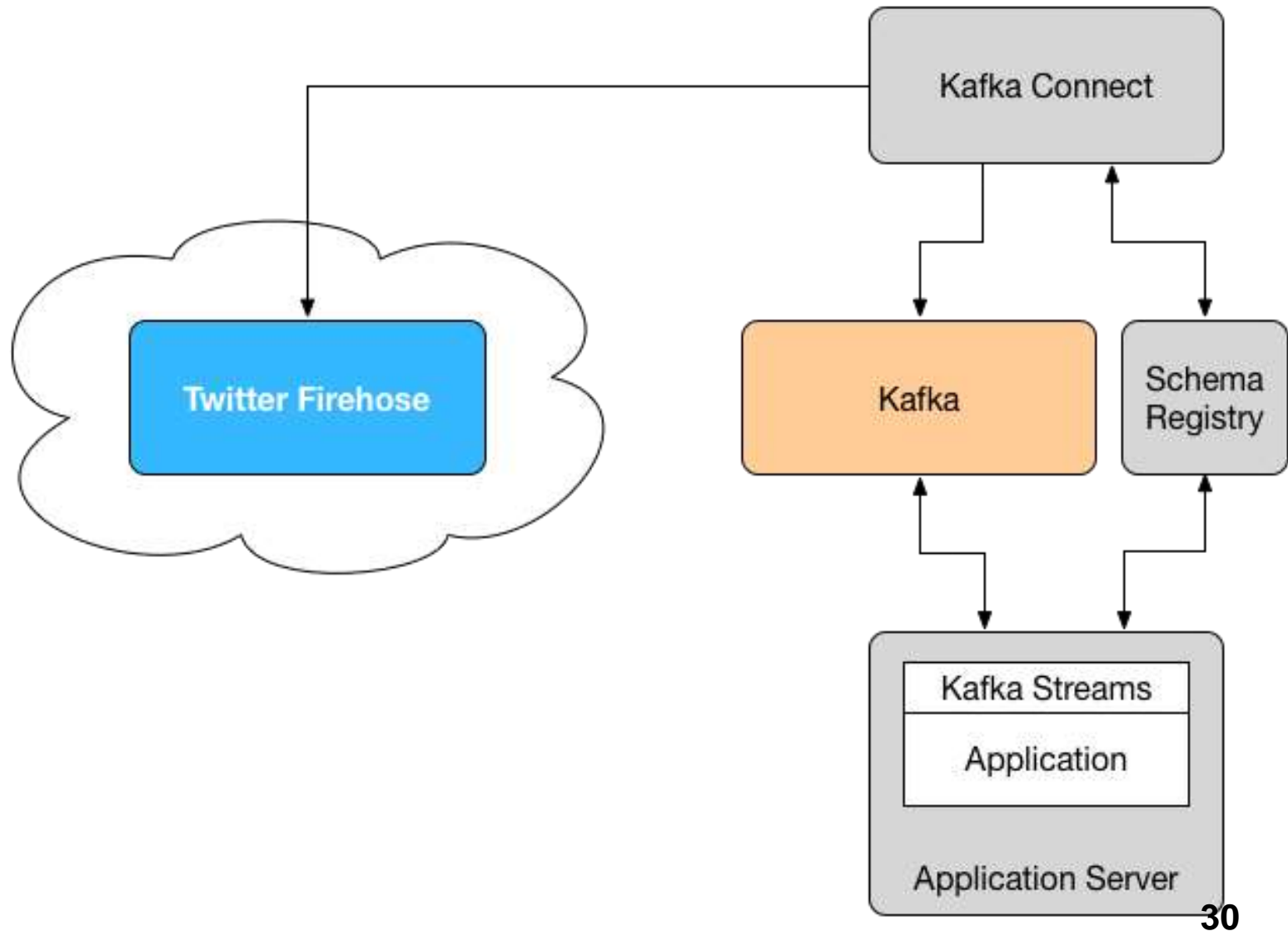
- This demo is adapted from one that was given by Sönke Liebau on July 27th 2016 from OpenCore, Germany. See blog entry titled: ‘Processing Twitter Data with Kafka Streams’

<http://www.opencore.com/blog/2016/7/kafka-streams-demo/> and related code at GitHub

<https://github.com/opencore/kafkastreamsdemo>

- In addition:
 - I’m using a **Docker container** instead of the confluent platform they are providing with their Virtual Machine defined in Vagrant.
 - I’m also using **Kafka Connect UI** from Landoop for easy and fast configuration of Twitter connector and also other Landoop’s Fast Data Web UIs.

4.2. Architecture of this demo



4.3. Setup of this demo

Step 1: Setup your Kafka Development Environment

Step 2: Get twitter credentials to connect to live data

Step 3: Get twitter live data into Kafka broker

Step 4: Write and test the application code in Java

Step 5: Run the application

Step 1: Setup your Kafka Development Environment

- The easiest way to get up and running quickly is to use a **Docker container** with all components needed.
- **First**, install Docker on your desktop or on the cloud
<https://www.docker.com/products/overview> and start it




Secure <https://www.docker.com/products/overview>

DOCKER PLATFORM DOCKER HUB DOCKER CLOUD DOCKER DATACENTER

INSTALL THE PLATFORM

Install Docker with easy to use installers for the major desktop and cloud platforms:




MAC

A native Mac application with a user interface and auto-update capabilities, that is deeply integrated with OS X native virtualization.

[Download](#)

[Learn More](#)




WINDOWS

A native Windows application with a user interface and auto-update capabilities, that is deeply integrated with Windows native virtualization.

[Download](#)

[Learn More](#)




LINUX

Install Docker on nodes which have a Linux distribution already installed.

[Install](#)

[Learn More](#)




AWS

Quickly deploy, scale, and manage Docker on AWS. Docker for AWS takes optimal advantage of the underlying infrastructure, while providing a modern Docker platform that can be used to deploy portable apps.

[Launch Stack](#)

[Learn More](#)




AZURE

Quickly deploy, scale and manage Docker on Azure. Docker for Azure takes optimal advantage of the underlying infrastructure, while providing a modern Docker platform that can be used to deploy portable apps.

[Deploy to Azure](#)

[Learn More](#)



WINDOWS SERVER

Install Docker on nodes which have Windows Server 2016 already installed.

[Learn More](#)

Step 1: Setup your Kafka Development Environment

➤ **Second**, install **Fast-data-dev**, a **Docker image for Kafka** developers which is packaging:

- **Kafka broker**
- **Zookeeper**
- **Open source version of the Confluent Platform** with its Schema registry, REST Proxy and bundled connectors
- Certified DataMountaineer **Connectors** (ElasticSearch, Cassandra, Redis, ..)
- Landoop's Fast Data **Web UIs** : schema-registry, kafka-topics, kafka-connect.
- Please note that **Fast Data Web UIs are licensed under BSL**. You should contact Landoop if you plan to use them on production clusters with more than 4 nodes.

by executing the command below, while Docker is running and you are connected to the internet:

docker run --rm -it --net=host landoop/fast-data-dev

- If you are on Mac OS X, you have to expose the ports instead:

**docker run --rm -it \
-p 2181:2181 -p 3030:3030 -p 8081:8081 \
-p 8082:8082 -p 8083:8083 -p 9092:9092 \
-e ADV_HOST=127.0.0.1 \
landoop/fast-data-dev**

- This will download the fast-data-dev Docker image from the Dock Hub.
<https://hub.docker.com/r/landoop/fast-data-dev/>
- Future runs will use your local copy.
- More details about Fast-data-dev docker image <https://github.com/Landoop/fast-data-dev>

Step 1: Setup your Kafka Development Environment

➤ Points of interest:

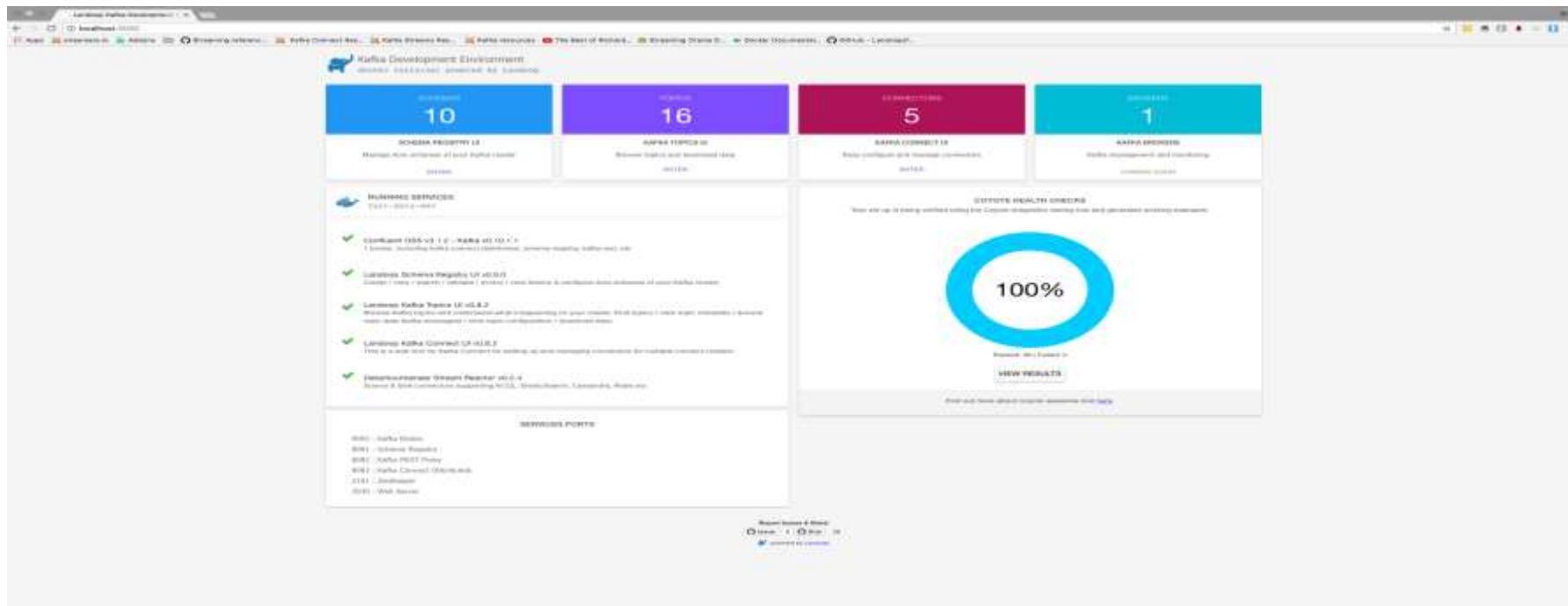
- the **-p** flag is used to publish a network port. Inside the container, ZooKeeper listens at 2181 and Kafka at 9092. If we don't publish them with -p, they are not available outside the container, so we can't really use them.
- the **-e** flag sets up environment variables.
- the last part specifies the image we want to run: **landoop/fast-data-dev**
- Docker will realize it doesn't have the landoop/fast-data-dev image locally, so it will first download it.

➤ That's it.

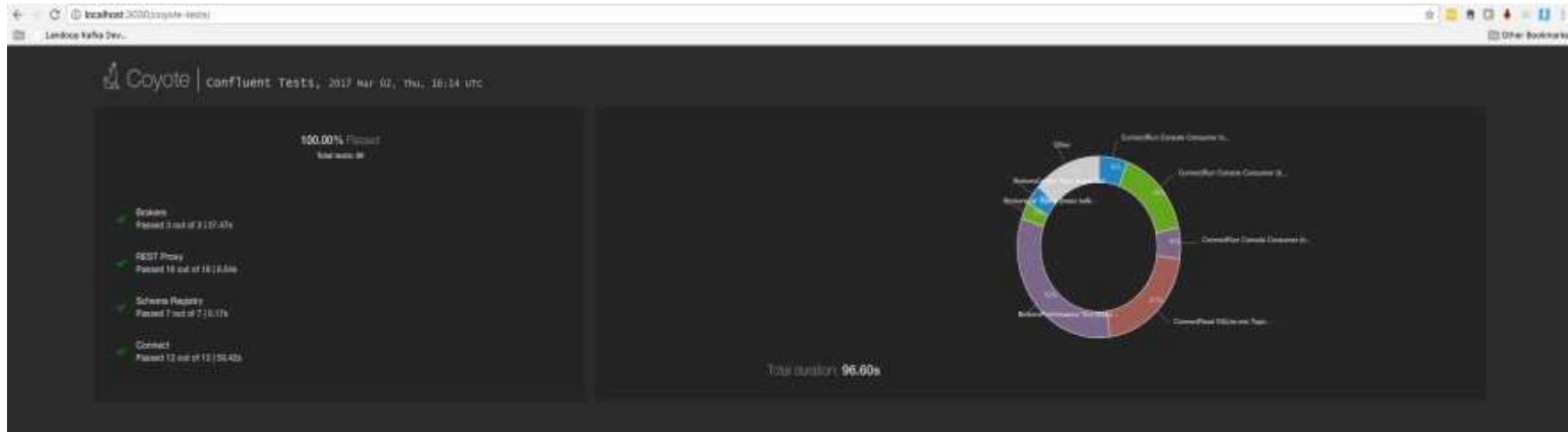
- Your **Kafka Broker** is at localhost:9092,
- your **Kafka REST Proxy** at localhost:8082,
- your **Schema Registry** at localhost:8081,
- your **Connect Distributed** at localhost:8083,
- your **ZooKeeper** at localhost:2181

Step 1: Setup your Kafka Development Environment

- At <http://localhost:3030>, you will find **Landoop's Web UIs** for:
 - Kafka Topics
 - Schema Registry
 - as well as a integration test report for connectors & infrastructure using Coyote. <https://github.com/Landoop/coyote>
- If you want to stop all services and remove everything, simply hit **Control+C**.



➤ Explore Integration test results at <http://localhost:3030/coyote-tests/>

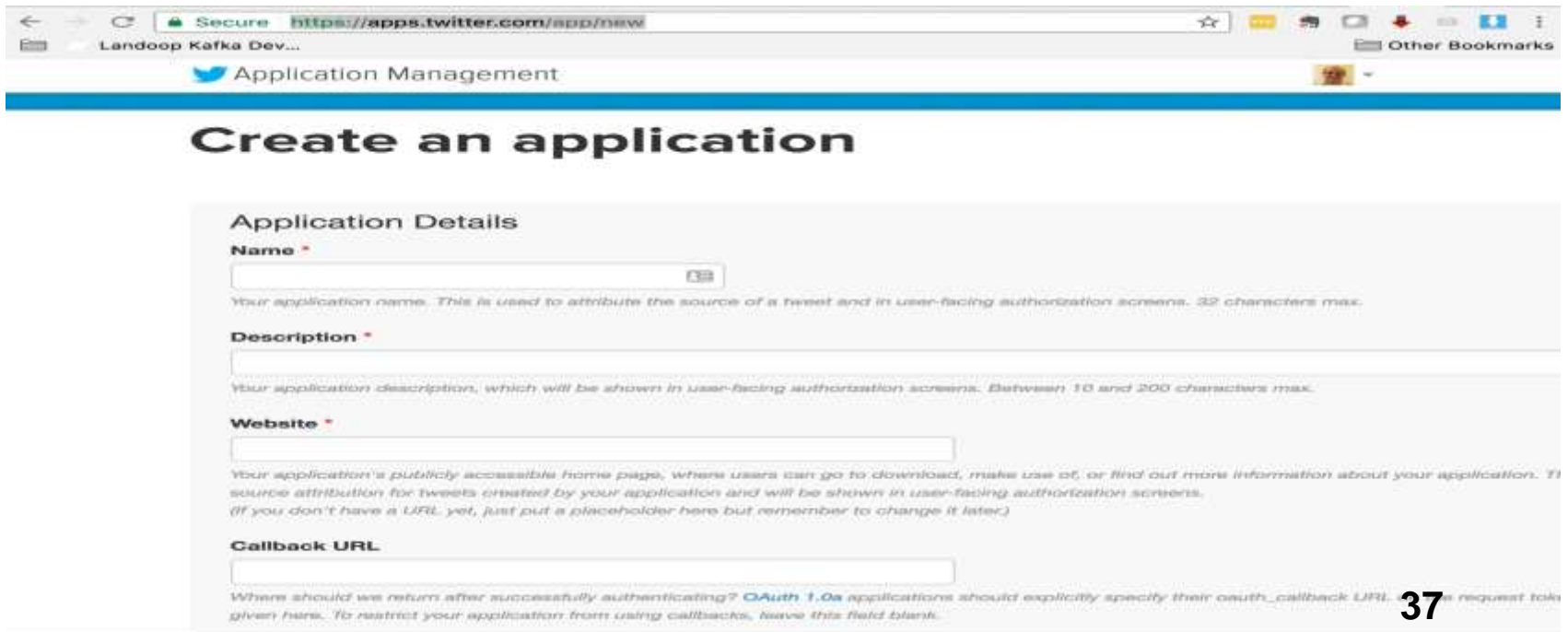


Results

Brokers					⌵
	Action	⌚ Time (sec)	📄 Command	Exit code	
+ ✓	Create Topic (basic kafka)	0.97	kafka-topics --zookeeper 127.0.0.1:2181 --topic mytest_basic_1488471144874 --partitions 1 --replication 1 --create	0	
+ ✓	List Topics (basic kafka)	2.96	kafka-topics --zookeeper 127.0.0.1:2181 --list	0	
+ ✓	Performance Test (basic kafka)	31.14	kafka-producer-perf-test --topic mytest_basic_1488471144874 --throughput 100000 --record-size 1000 --num-records 500000 --producer-group testgroup.servers="127.0.0.1:9092"	0	
Passed 3 out of 3 37.47 seconds					
REST Proxy					⌵
	Action	⌚ Time (sec)	📄 Command	Exit code	
+ ✓	List Topics (rest proxy)	0.09	curl -vs --stderr - "http://127.0.0.1:8082/topics"	0	
+ ✓	Topic Information (rest proxy)	1.29	curl -vs --stderr - "http://127.0.0.1:8082/topics/mytest_basic_1488471144874"	0	
+ ✓	Topic Partitions (rest proxy)	0.11	curl -vs --stderr - "http://127.0.0.1:8082/topics/mytest_basic_1488471144874/partitions"	0	
+ ✓	Delete Topic (basic kafka)	2.23	kafka-topics --zookeeper 127.0.0.1:2181 --topic mytest_basic_1488471144874 --delete	0	
+ ✓	Brokers from Metadata (rest proxy, schema isolated)	0.03	curl -vs --stderr - -HMYT -H "Content-Type: application/vnd.kafka.admin.v1+json" --data '{"name schema": {"type": "record", "fields": [{"name": "test", "type": "string"}]}}' "http://127.0.0.1:8082/brokers"	0	

Step 2: Get twitter credentials to connect to live data

- Now that our **single-node Kafka cluster is fully up and running**, we can proceed to preparing the input data:
- First you need to **register an application with Twitter**.
 - Second, once the application is created **copy the Consumer key and Consumer Secret**.
 - Third, **generate the Access Token Access and Secret Token** required to give your twitter account access to the new application
- Full instructions are here: <https://apps.twitter.com/app/new>



The screenshot shows a web browser window with the address bar displaying <https://apps.twitter.com/app/new>. The page title is "Application Management" and the main heading is "Create an application". The form is titled "Application Details" and contains the following fields:

- Name ***: A text input field with a character count of 32. Below the field, it says: "Your application name. This is used to attribute the source of a tweet and in user-facing authorization screens. 32 characters max."
- Description ***: A text input field with a character count of 200. Below the field, it says: "Your application description, which will be shown in user-facing authorization screens. Between 10 and 200 characters max."
- Website ***: A text input field. Below the field, it says: "Your application's publicly accessible home page, where users can go to download, make use of, or find out more information about your application. If source attribution for tweets created by your application and will be shown in user-facing authorization screens. (If you don't have a URL yet, just put a placeholder here but remember to change it later.)"
- Callback URL**: A text input field. Below the field, it says: "Where should we return after successfully authenticating? OAuth 1.0a applications should explicitly specify their oauth_callback URL. If you don't specify a callback URL, we will use the request token given here. To restrict your application from using callbacks, leave this field blank."

The number 37 is visible in the bottom right corner of the page.

Step 3: Get twitter live data into Kafka broker

➤ First, **create a new Kafka Connect for Twitter**

The screenshot displays the Kafka Connect web console. On the left, a sidebar shows '3 Connectors' and a 'NEW' button. Below this is a search bar and a list of installed connectors: 'logs-broker', 'logs-rest-proxy', and 'logs-schema-registry', each with a status of '1 x' and a 'New' button. At the bottom of the sidebar, it shows 'Kafka Connect : /api/kafka-connect' and 'Kafka Connect Version : 0.10.1.0'. The main area is titled 'New Connector' and contains a search bar. It is divided into two columns: 'Sources' and 'Sinks'. Under 'Sources', there are four options: 'twitter' (with a Twitter logo and description), 'ftp' (with a folder icon and description), 'file' (with a folder icon and description), and 'bloomberg' (with a Bloomberg logo and description). Under 'Sinks', there are three options: 'cassandra' (with a Cassandra logo and description), 'influxDB' (with an InfluxDB logo and description), and 'elastic' (with an Elastic logo and description). Below these, there is a section for 'hbase' (with an HBase logo and description).

Step 3: Get twitter live data into Kafka broker

- Second, **configure this Kafka Connect for Twitter** to write to the **topic *twitter*** by entering your own track.terms and also the values of twitter.token, twitter.secret, twitter.consumerkey and twitter.consumer.secret

The screenshot displays the Kafka Connect web interface. On the left, a sidebar shows '3 Connectors' and a 'NEW' button. Below this is a search bar and a list of connectors: 'logs-broker', 'logs-rest-proxy', and 'logs-schema-registry', each with a status icon and a '1 x' count. At the bottom of the sidebar, it shows 'Kafka Connect : /api/kafka-connect' and 'Kafka Connect Version : 0.10.1.0'. The main panel is titled 'New Connector (Source): twitter'. It features a Twitter logo and a description: 'Subscribe to feeds using the Twitter API and stream data into a kafka topic.' Below this, there are checkboxes for 'Show Hints' and 'Editor mode' (which is checked). A 'FINISH' button is visible. At the bottom, there is a checkbox for 'Show cURL command' and a 'SHOW OPTIONAL FIELDS' button. The configuration details are listed in a code block:

```
1 name=twitter-source
2 topic=twitter
3 tasks.max=1
4 connector.class=com.eneco.trading.kafka.connect.twitter.TwitterSourceConnector
5 track.terms=#ApacheSpark,#ApacheFlink,#ApacheKafka
6 track.locations=
7 track.follow=
8 twitter.token=
9 twitter.secret=
10 twitter.consumerkey=
11 twitter.consumersecret=
12 language=
13 batch.size=100
```

Step 3: Get twitter live data into Kafka broker

➤ **Kafka Connect for Twitter** is now configured to write data to the **topic *twitter***.

The screenshot displays the Kafka Connect dashboard. On the left, a list of 6 connectors is shown, each with a status icon, a checkmark, and a '1 x' label. The connectors are: logs-broker, logs-connect-distributed, logs-rest-proxy, logs-schema-registry, logs-zookeeper, and twitter-source. The twitter-source connector is highlighted with a blue background. Below the list, the Kafka Connect API endpoint and version are displayed: /api/kafka-connect and 0.10.1.0.

The main dashboard area shows three summary cards: SINK CONNECTORS (0), SOURCE CONNECTORS (6), and TOPICS USED BY CONNECTORS (6). Below these, the 'Connect topology' section shows a diagram of the connectors and their connections. The topology includes: logs-schema-registry, logs-rest-proxy, logs-broker, logs-zookeeper, twitter-source, and logs-connect-distributed. The twitter-source connector is connected to the twitter topic.

KAFKA CONNECT

6 Connectors [NEW](#)

Search connectors

- logs-broker 1 x [icon] → [icon]
- logs-connect-distributed 1 x [icon] → [icon]
- logs-rest-proxy 1 x [icon] → [icon]
- logs-schema-registry 1 x [icon] → [icon]
- logs-zookeeper 1 x [icon] → [icon]
- twitter-source 1 x [icon] → [icon]

Kafka Connect : /api/kafka-connect
Kafka Connect Version : 0.10.1.0

Dashboard [EXPORT CONFIG](#)

SINK CONNECTORS: 0

SOURCE CONNECTORS: 6

TOPICS USED BY CONNECTORS: 6

Connect topology

- logs-schema-registry → logs-schema-registry
- logs-rest-proxy → logs-rest-proxy
- logs-broker → logs-broker
- logs-zookeeper → logs-zookeeper
- twitter-source → twitter
- logs-connect-distributed → logs-connect-distributed

Step 3: Get twitter live data into Kafka broker

➤ Data is now being written to the **topic twitter**.

localhost:3030/kafka-topics-ui/#/cluster/fast-data-dev/topic/n/twitter/table

Landoop Kafka Dev...

Other Bookmarks

KAFKA TOPICS

12 Topics

☐ System Topics

Search topics

logs-connect-distributed
1 Replication • 1 Partition

logs-rest-proxy
1 Replication • 1 Partition

logs-schema-registry
1 Replication • 1 Partition

logs-zookeeper
1 Replication • 1 Partition

twitter
1 Replication • 1 Partition

< 1 2 >

Kafka Rest : /api/kafka-rest-proxy
Kafka Brokers : 1
kafka-topics-ui: 0.8.2

twitter

Filter

TOPIC

TABLE

RAW DATA

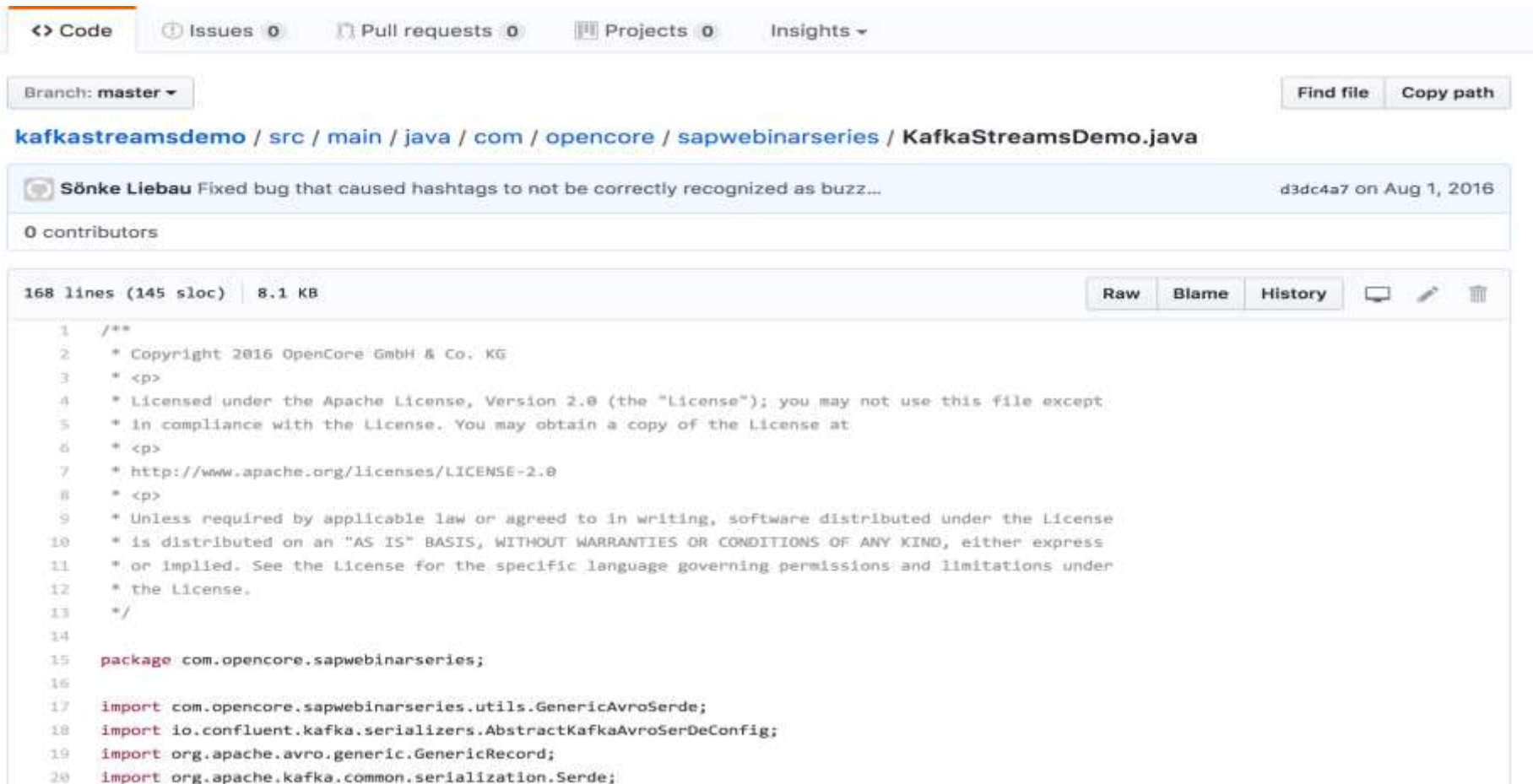
offset	partition	key	Value				
			id	created_at.string	user.id	user.name.string	user.scr
0	0		837388933633728500	2017-03-02T19:47:34.000+0000	113254439	Phil Tibitoski ☐ GDC	PTIbz
1	0		837388934124556300	2017-03-02T19:47:34.000+0000	22311026	TnT's Daddy	Scotttam
2	0		837388941938593800	2017-03-02T19:47:36.000+0000	23662055	David Bitner	bitnerd
3	0		837388944312582100	2017-03-02T19:47:36.000+0000	257823267	Nostradamass	blood_on
4	0		837388944987811800	2017-03-02T19:47:37.000+0000	2830655980	Fred	mrchas42
5	0		837388947529531400	2017-03-02T19:47:37.000+0000	761222977623056400	Kevin W	kevincaux
6	0		837388948586430500	2017-03-02T19:47:37.000+0000	31245235	Some Nigga	FukYoTwi
7	0		837388953904894000	2017-03-02T19:47:39.000+0000	787703409969336300	F	neko759
8	0		837388956618666000	2017-03-02T19:47:39.000+0000	743188217265000400	iSnouKeR7w7_	yiSnouKe
9	0		837388956882722800	2017-03-02T19:47:39.000+0000	269501813	rose 🇺🇸	lexiechief

< 1 2 3 4 5 6 7 >

Step 4: Write and test the application code in Java

- Instead of writing our own code for this demo, we will be leveraging an existing code from GitHub by Sonke Liebau:

<https://github.com/opencore/kafkstreamsdemo>



```
1  /**
2   * Copyright 2016 OpenCore GmbH & Co. KG
3   * <p>
4   * Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except
5   * in compliance with the License. You may obtain a copy of the License at
6   * <p>
7   * http://www.apache.org/licenses/LICENSE-2.0
8   * <p>
9   * Unless required by applicable law or agreed to in writing, software distributed under the License
10  * is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express
11  * or implied. See the License for the specific language governing permissions and limitations under
12  * the License.
13  */
14
15  package com.opencore.sapwebinarseries;
16
17  import com.opencore.sapwebinarseries.utils.GenericAvroSerde;
18  import io.confluent.kafka.serializers.AbstractKafkaAvroSerDeConfig;
19  import org.apache.avro.generic.GenericRecord;
20  import org.apache.kafka.common.serialization.Serde;
```


Step 4: Write and test the application code in Java

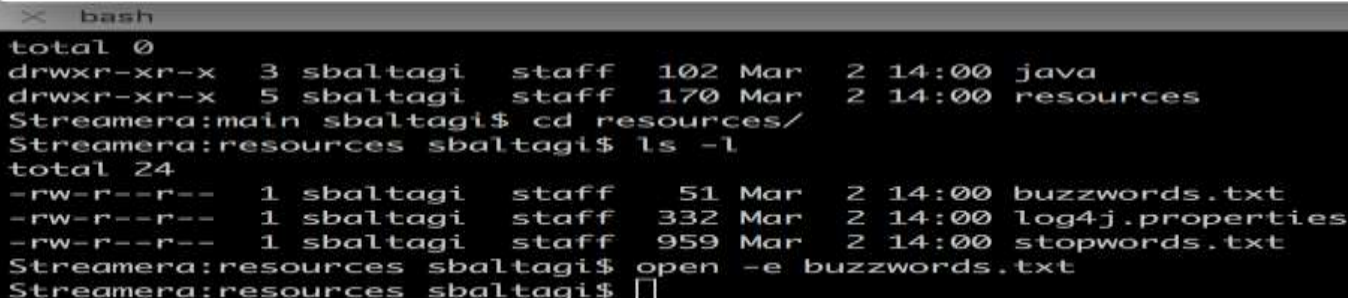
- **git clone** <https://github.com/opencore/kafkstreamsdemo>

```
Streamera:~ sbaltagi$ git clone https://github.com/opencore/kafkstreamsdemo
Cloning into 'kafkstreamsdemo'...
remote: Counting objects: 34, done.
remote: Total 34 (delta 0), reused 0 (delta 0), pack-reused 34
Unpacking objects: 100% (34/34), done.
Streamera:~ sbaltagi$
```

- Edit the buzzwords.txt file with your own works and probably one of the twitter terms that you are watching live:



```
chicago
kafka
meetup
streaming|
```



```
total 0
drwxr-xr-x  3 sbaltagi  staff  102 Mar  2 14:00 java
drwxr-xr-x  5 sbaltagi  staff  170 Mar  2 14:00 resources
Streamera:main sbaltagi$ cd resources/
Streamera:resources sbaltagi$ ls -l
total 24
-rw-r--r--  1 sbaltagi  staff    51 Mar  2 14:00 buzzwords.txt
-rw-r--r--  1 sbaltagi  staff   332 Mar  2 14:00 log4j.properties
-rw-r--r--  1 sbaltagi  staff   959 Mar  2 14:00 stopwords.txt
Streamera:resources sbaltagi$ open -e buzzwords.txt
Streamera:resources sbaltagi$
```

Step 5: Run the application

- The next step is to run the Kafka Streams application that processes twitter data.
- First, install Maven <http://maven.apache.org/install.html>
- Then, **compile the code into a fat jar with Maven.**

\$ *mvn package*

```
Streamera:kafkastreamsdemo sbaltagi$ pwd
/Users/sbaltagi/kafkastreamsdemo
Streamera:kafkastreamsdemo sbaltagi$ ls -l
total 40
-rw-r--r--  1 sbaltagi  staff   11357 Mar  2 14:00 LICENSE
-rw-r--r--  1 sbaltagi  staff    2135 Mar  2 14:00 README.md
-rw-r--r--  1 sbaltagi  staff    2910 Mar  2 14:00 pom.xml
drwxr-xr-x  3 sbaltagi  staff     102 Mar  2 14:00 src
Streamera:kafkastreamsdemo sbaltagi$ mvn package
[INFO] Scanning for projects...
[INFO]
[INFO] -----
[INFO] Building KafkaStreamsDemo 1.0-SNAPSHOT
[INFO] -----
Downloading: http://packages.confluent.io/maven/org/apache/kafka/kafka-streams/0.10.0.0/kafka-s
treams-0.10.0.0.pom
```

```
Component: org.apache.maven.artifact.handler.ArtifactHandlerejb-client is already defined. Skipping.
Component: org.apache.maven.artifact.handler.ArtifactHandlererrar is already defined. Skipping.
Component: org.apache.maven.artifact.handler.ArtifactHandlerpar is already defined. Skipping.
Component: org.apache.maven.artifact.handler.ArtifactHandlerejb3 is already defined. Skipping.
Component: org.apache.maven.artifact.factory.ArtifactFactory is already defined. Skipping.
Component: org.apache.maven.artifact.resolver.ArtifactCollector is already defined. Skipping.
Component: org.apache.maven.plugin.registry.MavenPluginRegistryBuilder is already defined. Skipping.
Component: org.apache.maven.shared.io.download.DownloadManagerdefault is already defined. Skipping.
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 11.475 s
[INFO] Finished at: 2017-03-02T14:24:58-06:00
[INFO] Final Memory: 37M/422M
[INFO] -----
Streamera:kafkastreamsdemo sbaltagi$
```


Step 5: Run the application

- Two jar files will be created in the **target** folder:
1. KafkaStreamsDemo-1.0-SNAPSHOT.jar – Only your project classes
 2. KafkaStreamsDemo-1.0-SNAPSHOT-**jar-with-dependencies**.jar – Project and dependency classes in a single jar.

```
Streamera:kafkastreamsdemo sbaltagi$ pwd
/Users/sbaltagi/kafkastreamsdemo
Streamera:kafkastreamsdemo sbaltagi$ ls -l
total 40
-rw-r--r--  1 sbaltagi  staff  11357 Mar  2 14:00 LICENSE
-rw-r--r--  1 sbaltagi  staff   2135 Mar  2 14:00 README.md
-rw-r--r--  1 sbaltagi  staff   2910 Mar  2 14:00 pom.xml
drwxr-xr-x  3 sbaltagi  staff    102 Mar  2 14:00 src
drwxr-xr-x  8 sbaltagi  staff    272 Mar  2 14:24 target
Streamera:kafkastreamsdemo sbaltagi$ cd target/
Streamera:target sbaltagi$ ls -l
total 31672
-rw-r--r--  1 sbaltagi  staff 16203673 Mar  2 14:24 KafkaStreamsDemo-1.0-SNAPSHOT-jar-with-dependencies.jar
-rw-r--r--  1 sbaltagi  staff   11537 Mar  2 14:24 KafkaStreamsDemo-1.0-SNAPSHOT.jar
drwxr-xr-x  2 sbaltagi  staff     68 Mar  2 14:24 archive-tmp
drwxr-xr-x  6 sbaltagi  staff    204 Mar  2 14:24 classes
drwxr-xr-x  3 sbaltagi  staff    102 Mar  2 14:24 generated-sources
drwxr-xr-x  3 sbaltagi  staff    102 Mar  2 14:24 maven-archiver
Streamera:target sbaltagi$
```

Step 5: Run the application

➤ Then

```
java -cp target/KafkaStreamsDemo-1.0-SNAPSHOT-  
jar-with-dependencies.jar  
com.opencore.sapwebinarseries.KafkaStreamsDemo
```

➤ **TIP:** During development: from your IDE, from CLI ...
Kafka Streams Application Reset Tool, **available**
since Apache Kafka 0.10.0.1, is great for playing
around.

<https://cwiki.apache.org/confluence/display/KAFKA/Kafka+Streams+Application+Reset+Tool>

4.4. Results of this demo

➤ Once the above is running, the following topics will be populated with data :

- **Raw word count** - Every occurrence of individual words is counted and written to the topic *wordcount* (a predefined list of stopwords will be ignored)
- **5-Minute word count** - Words are counted per 5 minute window and every word that has more than three occurrences is written to the topic *wordcount5m*
- **Buzzwords** - a list of special interest words can be defined and those will be tracked in the topic *buzzwords* - the list of these words can be defined in the file *buzzwords.txt*

4.4. Results of this demo

- **Accessing the data generated by the code** is as simple as starting a **console consumer** which is shipped with Kafka
- You need first to enter the container to use any tool as you like:
docker run --rm -it --net=host landoop/fast-data-dev bash
- Use the following command to check the topics:
 - ***kafka-console-consumer --topic wordcount --new-consumer --bootstrap-server 127.0.0.1:9092 --property print.key=true***
 - ***kafka-console-consumer --topic wordcount5m --new-consumer --bootstrap-server 127.0.0.1:9092 --property print.key=true***
 - ***kafka-console-consumer --topic buzzwords --new-consumer --bootstrap-server 127.0.0.1:9092 --property print.key=true***

4.4. Results of this demo

5. docker

```

@SandraTXAS: Mayor Rahm Emanuel and Democrats allow Chicago citizens to be murdered in their neighborhoods for decades, yet blam... en...SATRA
SandraTXAS
null 2017-03-02T21:28:48.000+0000 Marsha GanzMarshaGanz$Ft. Lauderdale, FL @Science4UsSays: Bwhahahahaha #thursdaythoughts #ICE1
7 #Chicago #davinci https://t.co/89NqAgfHwn @MarshaGanz https://t.co/324aVsDUXot thursdaythoughts
ICE17Chicagodavinci4pic.twitter.com/324aVsDUXothttps://twitter.com/Science4UsSays/status/837413032342999040/photo/1
photo,https://t.co/324aVsDUXobuff.ly/2lyHLAH,http://buff.ly/2lyHLAH,https://t.co/89NqAgfHwn
"Jenna and SciefusScience4UsSaysMarsha GanzMarshaGanz
z
null 2017-03-02T21:29:10.000+0000 Billy BlokeeBayDownUnder8Australian Capital TerritoryjChicago girl writes 'I hate this life' d
ays before death https://t.co/kJ6Zmq7F1F https://t.co/jzvOWF0qwyen4pic.twitter.com/jzvOWF0qwyhttps://twitter.com/eBayDownUnder/status/837414503738322
944/photo/1
photo,https://t.co/jzvOWF0qwydlvr.it/NWnpJ3*http://dlvr.it/NWnpJ3.https://t.co/kJ6Zmq7F1F

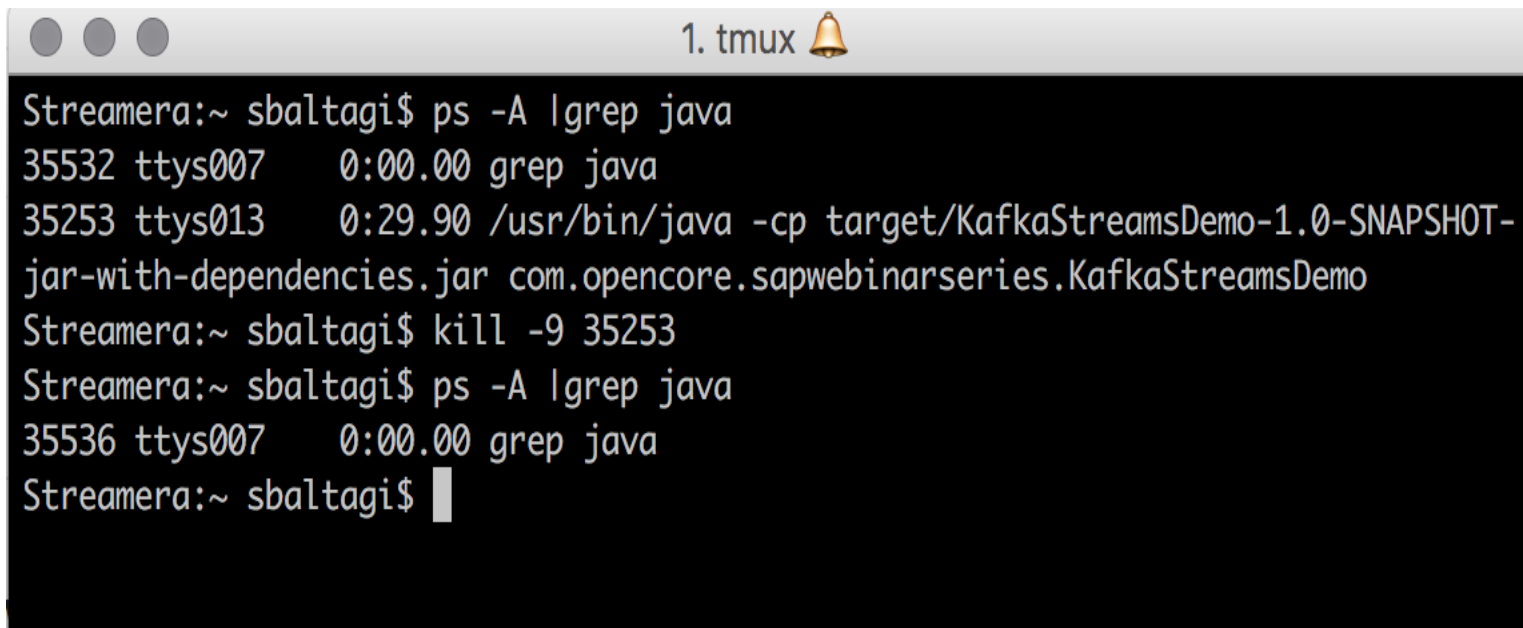
```

docker	docker	docker
https://t.co/89nqagfHwn 2	Thu Mar 02 15:30:00 CST 2017-democrats 5	Thu Mar 02 15:30:00 CST 2017-chicago 91
@marshoganz 2	Thu Mar 02 15:30:00 CST 2017-allow 5	Thu Mar 02 15:30:00 CST 2017-chicago 92
https://t.co/324avsduxo 2	Thu Mar 02 15:30:00 CST 2017-chicago 99	Thu Mar 02 15:30:00 CST 2017-chicago 93
chicago 2263	Thu Mar 02 15:30:00 CST 2017-citizens 5	Thu Mar 02 15:30:00 CST 2017-chicago 94
girl 155	Thu Mar 02 15:30:00 CST 2017-murdered 5	Thu Mar 02 15:30:00 CST 2017-chicago 95
writes 51	Thu Mar 02 15:30:00 CST 2017-neighborhoods 5	Thu Mar 02 15:30:00 CST 2017-meetup 8
'i 9	Thu Mar 02 15:30:00 CST 2017-decades, 5	Thu Mar 02 15:30:00 CST 2017-chicago 96
hate 134	Thu Mar 02 15:30:00 CST 2017-yet 5	Thu Mar 02 15:30:00 CST 2017-chicago 97
life' 1	Thu Mar 02 15:30:00 CST 2017-blam_ 5	Thu Mar 02 15:30:00 CST 2017-chicago 98
days 38	Thu Mar 02 15:30:00 CST 2017-rt 124	Thu Mar 02 15:30:00 CST 2017-meetup 9
death 18	Thu Mar 02 15:30:00 CST 2017-#chicago 17	Thu Mar 02 15:30:00 CST 2017-chicago 99
https://t.co/kj6zmq7f1f 1	Thu Mar 02 15:30:00 CST 2017-chicago 100	Thu Mar 02 15:30:00 CST 2017-#chicago 17
https://t.co/jzvowf0qwy 1	Thu Mar 02 15:30:00 CST 2017-days 3	Thu Mar 02 15:30:00 CST 2017-chicago 100

4.5. Stopping the demo!

➤ To stop the Kafka Streams Demo application:

- **\$ ps -A | grep java**
- **\$ kill -9 PID**

A terminal window titled "1. tmux" with a bell icon. The terminal shows the following commands and output:

```
Streamera:~ sbaltagi$ ps -A |grep java
35532 ttys007    0:00.00 grep java
35253 ttys013    0:29.90 /usr/bin/java -cp target/KafkaStreamsDemo-1.0-SNAPSHOT-
jar-with-dependencies.jar com.opencore.sapwebinarseries.KafkaStreamsDemo
Streamera:~ sbaltagi$ kill -9 35253
Streamera:~ sbaltagi$ ps -A |grep java
35536 ttys007    0:00.00 grep java
Streamera:~ sbaltagi$
```

➤ If you want to **stop all services in fast-data-dev Docker** image and remove everything, simply hit **Control+C**.

5. Where to go from here for further learning?

➤ **Kafka Streams code examples**

- Apache Kafka
<https://github.com/apache/kafka/tree/trunk/streams/examples/src/main/java/org/apache/kafka/streams/examples>
- Confluent <https://github.com/confluentinc/examples/tree/master/kafka-streams>

➤ **Source Code** <https://github.com/apache/kafka/tree/trunk/streams>

➤ **Kafka Streams Java docs**

<http://docs.confluent.io/current/streams/javadocs/index.html>

➤ **First book on Kafka Streams (MEAP)**

- Kafka Streams in Action <https://www.manning.com/books/kafka-streams-in-action>

➤ **Kafka Streams download**

- Apache Kafka <https://kafka.apache.org/downloads>
- Confluent Platform <http://www.confluent.io/download>

5. Where to go from here for further learning?

- **Kafka Users mailing list** <https://kafka.apache.org/contact>
- **Kafka Streams at Confluent Community on Slack**
 - <https://confluentcommunity.slack.com/messages/streams/>
- **Free ebook:**
 - Making Sense of Stream processing by Martin Klepmann <https://www.confluent.io/making-sense-of-stream-processing-ebook-download/>
- **Kafka Streams documentation**
 - Apache Kafka <http://kafka.apache.org/documentation/streams>
 - Confluent <http://docs.confluent.io/3.2.0/streams/>
- **All web resources related to Kafka Streams**

Thank you!

Let's keep in touch!



[@SlimBaltagi](https://twitter.com/SlimBaltagi)



<https://www.linkedin.com/in/slimbaltagi>



sbaltagi@gmail.com