

As soon as a company has **real-time data streaming needs**, a streaming platform must be put in place. Apache Kafka is one of the most popular **data streaming processing platforms** in the industry today.

## What are the use cases of Apache Kafka?

The use cases of Apache Kafka are many. These include stream processing for different business applications. Apache Kafka makes up the storage mechanism for some of the prominent stream processing frameworks, e.g., Apache Flink, Samza.

- \* Messaging systems
- \* Activity Tracking
- \* Gather metrics from many different locations, i.e., IoT devices
- \* Application logs analysis
- \* De-coupling of system dependencies
- \* Integration with Big Data technologies like Spark, Flink, Storm, Hadoop.
- \* Event-sourcing store

## How is Kafka concretely being used within the industry?

**UBER** : The speed and flexibility of Kafka allows Uber to adjust their pricing models to the constantly evolving events in the real world (number of available drivers and their position, users and their position, weather event, other events), and bill users the right amount to manage offer and demand.

**Netflix** : It has integrated Kafka as the core component of its data platform. They refer to it internally as their Keystone data pipeline. As part of Netflix's Keystone, Kafka handles billions of events a day. Just to give an idea about the huge amount of data that Kafka can handle, Netflix sends about 5 hundred billion events and 1.3 petabytes of data per day into Kafka.

Kafka is at the core of lots of the services we enjoy on a daily basis from some of the world's largest tech companies such as Uber, Netflix, Airbnb, LinkedIn, Apple & Walmart.

## Kafka Setup

- **Download**

<https://kafka.apache.org/downloads>

Download scala latest version in binary downloads.

- **Setup Kafka in Path**

We are doing this to access Kafka Binaries with ease.

- nano .zshrc
- PATH="\$PATH:/Users/shivamjaiswal/kafka\_2.13-3.0.0/bin"
- (We will not be required to provide full path for accessing Kafka binaries)

See complete setup guidance here :

<https://www.conduktor.io/kafka/how-to-install-apache-kafka-on-mac>

## Start Zookeeper

```
zookeeper-server-start.sh ~/kafka_2.13-3.4.0/config/zookeeper.properties
```

## Start Kafka

```
kafka-server-start.sh ~/kafka_2.13-3.4.0/config/server.properties
```

## Changing the Kafka and Zookeeper data storage directory

- edit the zookeeper.properties file at ~/kafka\_2.13-3.0.0/config/zookeeper.properties and set the following to your heart's desire dataDir=/your/path/to/data/zookeeper.
- edit the server.properties file at ~/kafka\_2.13-3.0.0/config/server.properties and set the following to your heart's desire log.dirs=/your/path/to/data/kafka.

**Note:- Zookeeper is going to be completely removed from Kafka 4.0.**

## Install Kafka with KRaft (without Zookeeper) on Mac OS X

1. Install Java JDK version 11
  2. Download Apache Kafka from <https://kafka.apache.org/downloads> under 'Binary Downloads'
  3. Extract the contents on your Mac
  4. Generate a cluster ID and format the storage using kafka-storage.sh
  5. Start Kafka using the binaries
  6. Setup the \$PATH environment variables for easy access to the Kafka binaries
1. generate a new ID for your cluster.

```
~/kafka_2.13-3.4.0/bin/kafka-storage.sh random-uuid
```

this returns a UUID, for example bs8tpjusSZuLirvAj700Wg

2. Format your storage directory (replace <uuid> by your UUID obtained above)

```
~/kafka_2.13-3.4.0/bin/kafka-storage.sh format -t <uuid> -c ~/kafka_2.13-3.0.0/config/kraft/server.properties
```

This will format the directory that is in the log.dirs in the config/kraft/server.properties file (by default /tmp/kraft-combined-logs)

3. Now you can launch the broker itself in daemon mode by running this command

```
~/kafka_2.13-3.0.0/bin/kafka-server-start.sh ~/kafka_2.13-3.4.0/config/Kraft/server.properties
```

## Kafka Message Anatomy :

Kafka messages are created by the producer. A Kafka message consists of the following elements:

**Key** : Key is optional in the Kafka message and it can be null. A key may be a string, number, or any object and then the key is serialized into binary format.

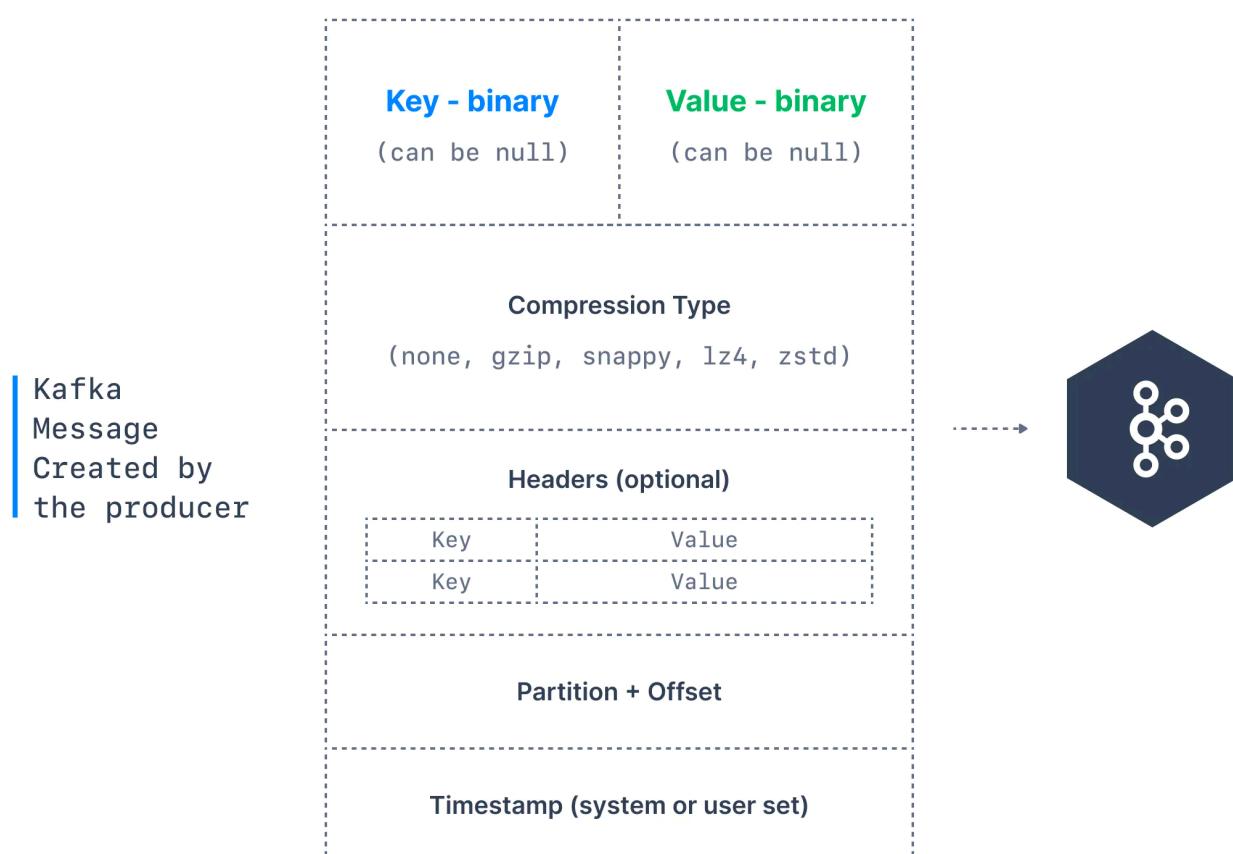
**Value** : The value represents the content of the message and can also be null. The value format is arbitrary and is then also serialized into binary format.

**Compression Type** : Kafka messages may be compressed. The compression type can be specified as part of the message. Options are none, gzip, lz4, snappy, and zstd

**Headers** : There can be a list of optional Kafka message headers in the form of key-value pairs. It is common to add headers to specify metadata about the message, especially for tracing.

**Partition + Offset** : Once a message is sent into a Kafka topic, it receives a partition number and an offset id. The combination of [topic+partition+offset](#) uniquely identifies the message.

**Timestamp** : A timestamp is added either by the user or the system in the message.



## Kafka Terms :

### **Kafka Broker**

- A single Kafka server is called a Kafka Broker.
- Kafka brokers store data in a directory on the server disk they run on.

### **Kafka Cluster**

An ensemble of Kafka brokers working together is called a Kafka cluster.

### **Bootstrap Server**

- A client that wants to send or receive messages from the Kafka cluster **may connect to any broker in the cluster**. Every broker in the cluster has metadata about all the other brokers and will help the client connect to them as well, and **therefore any broker in the cluster is also called a bootstrap server**.
- The bootstrap server will return metadata to the client that consists of a list of all the brokers in the cluster. Then, when required, the client will know which exact broker to connect to to send or receive data, and accurately find which brokers contain the relevant topic-partition.
- In practice, it is common for the Kafka client to reference at least two bootstrap servers in its connection URL

### **Topic**

Similar to how databases have tables to organize and segment datasets, Kafka uses the concept of topics to organize related messages.

- A topic is identified by its name.
- Kafka topics can contain any kind of message in any format, and the sequence of all these messages is called a **data stream**.
- Data in Kafka topics is deleted after one week by default (also called the default message retention period), and this value is configurable. This mechanism of deleting old data ensures a Kafka cluster does not run out of disk space by recycling topics over time.
- Kafka topics are **immutable**: once data is written to a partition, it cannot be changed

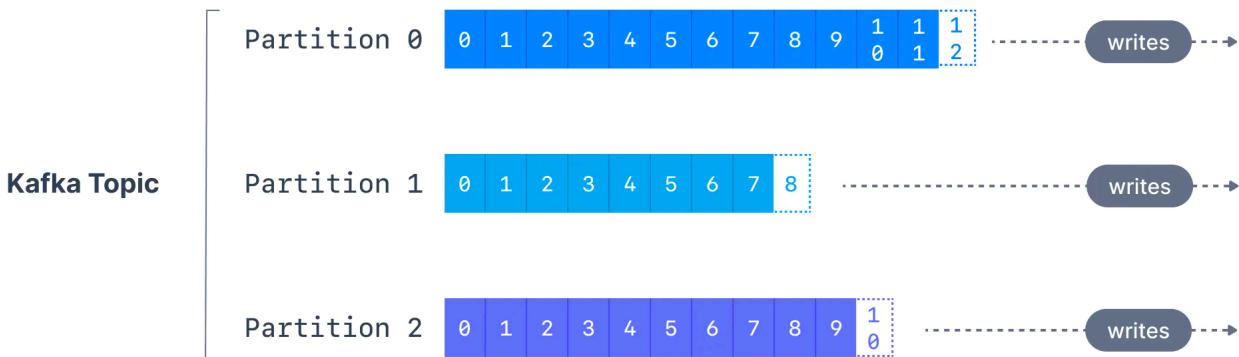
Unlike database tables, Kafka topics are not query-able. Instead, we have to create Kafka producers to send data to the topic and Kafka consumers to read the data from the topic in order.

## Partitions

- Topics are broken down into a number of partitions.
- Partitions are 0 indexed.
- partitions for a given topic will be distributed among the brokers evenly, to achieve load balancing and scalability.

## Offset

- The offset is an integer value that Kafka adds to each message as it is written into a partition.
- Offsets are 0 indexed.
- Each message in a given partition has a unique offset.
- Even though we know that messages in Kafka topics are deleted over time (as seen above), the offsets are not re-used. They are continually incremented in a never-ending sequence



## Replication Factor

- Replication means that data is written down not just to one broker, but many.
- The replication factor is a topic setting and is specified at topic creation time.
- replication factor of 1 means no replication, replication factor of 3 is a commonly used .

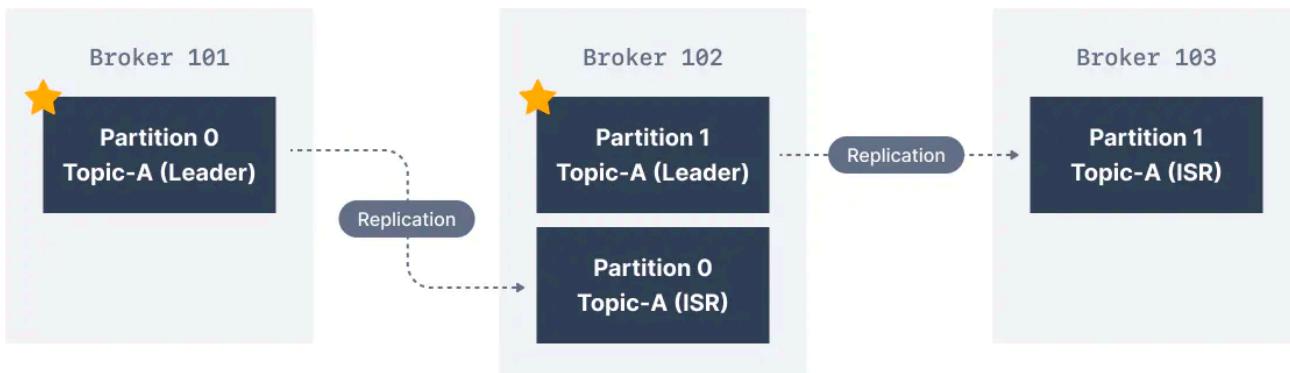
## Kafka Partitioner

A Kafka partitioner is a code logic that takes a record and determines to which partition to send it into. In the default Kafka partitioner, the keys are hashed using the **murmur2 algorithm**.

Messages sent to a Kafka topic that do not respect the agreed-upon serialization format are called poison pills.

## Kafka Partitions Leader and Replicas

- For a given topic-partition, one Kafka broker is designated by the cluster to be responsible for sending and receiving data to clients. That broker is known as the leader broker of that topic partition. Any other broker that is storing replicated data for that partition is referred to as a replica.
- Therefore, each partition has one leader and multiple replicas.
- An ISR is a replica that is up to date with the leader broker for a partition. Any replica that is not up to date with the leader is out of sync.



## Ordering

If a topic has more than one partition, Kafka guarantees the order of messages within a partition, but there is no ordering of messages across partitions while reading.

## Producers

- Applications that send data into topics are known as Kafka producers. Applications typically integrate a Kafka client library to write to Apache Kafka.
- A Kafka producer sends messages to a topic, and messages are distributed to partitions according to a mechanism such as key hashing. In case the key

(`key=null`) is not specified by the producer, messages are distributed evenly across partitions in a topic(Round Robin Fashion).

- For a message to be successfully written into a Kafka topic, a producer must specify a level of acknowledgment (acks).
- We need to use appropriate serialiser to send message to Kafka.

## Consumers

- Applications that read data from Kafka topics are known as consumers.
- Reads data from lower to higher offset.
- Consumers can read from one or more partitions at a time in Apache Kafka, and data is read in order **within each partition**.
- By default, Kafka consumers will only consume data that was produced after it first connected to Kafka.
- Data being consumed must be deserialized in the same format it was serialized in.

## Consumer group

- Kafka Consumers that are part of the same application and therefore performing the same "logical job" can be grouped together as a Kafka consumer group.
- Consumers automatically use a `GroupCoordinator` and a `ConsumerCoordinator` to assign consumers to a partition and ensure the load balancing is achieved across all consumers in the same group.
- Each topic partition is only assigned to one consumer within a consumer group, but a consumer from a consumer group can be assigned multiple partitions.

## Consumer Offsets

- Kafka brokers use an internal topic named `_consumer_offsets` that keeps track of what messages a given **consumer group** last successfully processed.
- the consumer will regularly **commit** the latest processed message, also known as **consumer offset**.
- The process of committing offsets is not done for every message consumed (because this would be inefficient), and instead is a periodic process.
- If rebalancing happens ,assigned consumer starts reading from consumer offset.

- A consumer will not be able to consume messages from beginning if consumer offset is already set for respective consumer group.
- Reset offset cannot happen when the consumer is running.

## Delivery semantics for consumers

- By default, Java consumers automatically commit offsets (controlled by the `enable.auto.commit=true` property) every `auto.commit.interval.ms` (5 seconds by default) when `.poll()` is called.
- A consumer may opt to commit offsets by itself (`enable.auto.commit=false`).
- Depending on when it chooses to commit offsets, there are delivery semantics available to the consumer. The three delivery semantics are explained below.

### → At most once

Offsets are committed as soon as the message is received.

If the processing goes wrong, the message will be lost (it won't be read again).

### → At Least once

Offsets are committed after the message is processed.

If the processing goes wrong, the message will be read again.

This can result in duplicate processing of messages. Therefore, it is best practice to make sure data processing is idempotent (i.e. processing the same message twice won't produce any undesirable effects)

### → Exactly once

This can only be achieved for Kafka topic to Kafka topic workflows using the transactions API. The Kafka Streams API simplifies the usage of that API and enables exactly once using the

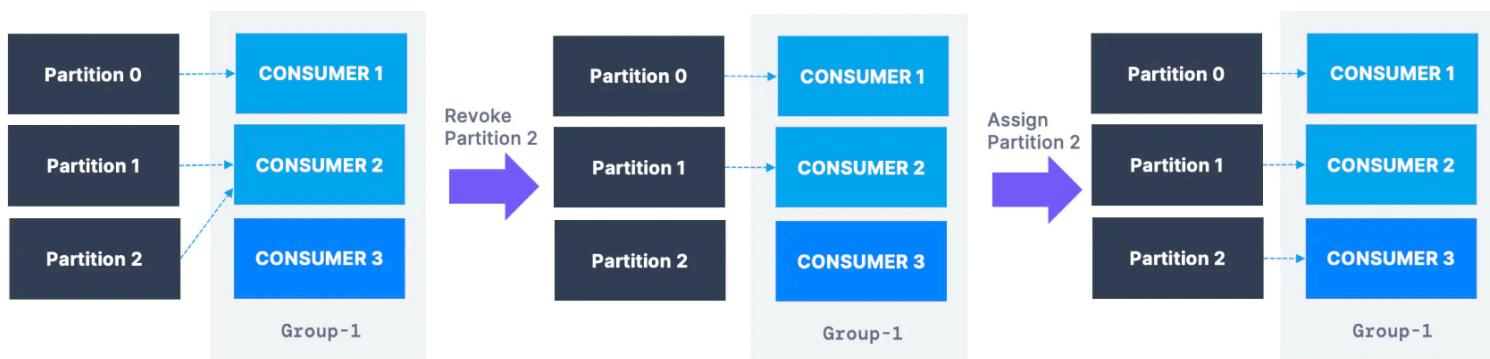
setting `processing.guarantee=exactly_once_v2` (`exactly_once` on Kafka < 2.5).

For Kafka topic to External System workflows, to effectively achieve exactly once, you must use an idempotent consumer. In practice, at least once with idempotent processing is the most desirable and widely implemented mechanism for Kafka consumers.

## Consumer Rebalance

Consumer rebalances happen for the following events:

- Number of partitions change for any of the subscribed topics
- A subscribed topic is created or deleted
- An existing member of the consumer group is shutdown or fails
- A new member is added to the consumer group
- Rebalancing techniques
  1. **Eager rebalance** : all partitions get revoked from consumer and get reassigned.
  2. **Co-operative rebalancing** : only subset of partition(which are required to rebalance) get revoked hence other consumers keeps on reading the data from unrequired partitions.



## Consumer Configurations

### 1. `partition.assignment.strategy`

**RangeAssignor**: assign partitions on a per-topic basis (can lead to imbalance)

**RoundRobin**: assign partitions across all topics in a round-robin fashion, optimal balance.

**StickyAssignor**: balanced like RoundRobin, and then minimises partition movements when consumers join/leave the group in order to minimise movements

**CooperativeStickyAssignor**: rebalance strategy is identical to StickyAssignor but supports cooperative rebalances and therefore consumers can keep on consuming from the topic.

### 2. `group.instance.id`

it makes the consumer a static member. I.e., upon leaving consumer has `session.timeout.ms` to join back and get back its partition, without triggering a rebalance.

## Producer Configurations

### 1. min.insync.replicas

it is used together with acks=all , and it assures acks from minimum x brokers.i.e., `min.insync.replicas=x{1,2,3,4....}`.

### 2. acks

Kafka producers only write data to the current leader broker for a partition.It must also specify a level of acknowledgment `acks` to specify if the message must be written to a minimum number of replicas before being considered a successful write.producers consider messages as "written successfully" for

- `acks = 0` : the moment the message was sent without waiting for the broker to accept it at all.
- `acks =1` : when the message was acknowledged by only the leader.
- `acks=all(-1)` : when the message is accepted by all in-sync replicas (ISR). Leader partition ask for acks to other ISR partitions and responds to producer together with its acks.

**default values is acks = all (-1)**

### 3. retries

The retries setting determines how many times the producer will attempt to send a message before marking it as failed. The default values is `MAX_INT`, i.e., `2147483647` for Kafka  $\geq 2.1$ .

### 4. retry.backoff.ms

time to wait before next retry

**default values is 100ms**

### 5. delivery.timeout.ms

retries are bounded by a timeout

**default values is 120000ms == 2 min**

### 6. max.in.flight.requests.per.connection

how many produce request(message batches) can be made in parallel

**default values is 5, set it to 1 to ensure ordering of messages.**

### 7. enable.idempotence =true

**Note : For idempotent producer** ordering will be same even for higher values and duplicates are not introduced due to network retries.This producers are default from Kafka 3.0.

if we enable idempotence `enable=idempotence=true`, then it is required for `max.in.flight.requests.per.connection` to be less than or equal to 5 with message ordering preserved for any allowable value!!

## 8. `compression.type`

can be at broker or producer level , producer level is recommended.

Note : at consumer end we don't need to decompress the batches , consumers do it by itself.

**default values is none.**

## 9. `linger.ms`

how long to wait until we send a batch. Adding a small number for example 5ms helps add more messages in a batch at the expense of latency.

**default values is 0.**

## 10. `batch.size`

if a batch is filled before linger.ms, increase the batch size.

smaller batches leads to more request and higher latency so prefer to have larger batch size.

**default values is 16kb**

## 11. `buffer.memory`

if the producer produces faster than the broker can take, the records will be buffered in memory

**default size is 33554432 (32 MB):the size of send buffer**

this buffer will fill up over time and empty back down when the throughput to the broker increases

## 12. `max.block.ms`

the time the `.send()` will block until throwing the exception.

Exception occurs when:

if the buffer is full , then the `.send()` method will start to block(won't return right away).

if broker is down or overloaded

if 60 sec has elapsed

`flush()` and `close()` are required to ensure the producer is shut down after the message is sent to Kafka.

## Seek and Assign

to read specific messages from specific partitions, the `.seek()` and `.assign()` API may help you. These APIs are also helpful to replay data from a specific offset.

### **Sticky Partitioner(used for performance improvement)**

which means the producer that receives messages sent in time close to each other will try to fill a batch into ONE partition before switching to creating a batch for another partition.

Keys become useful when a user wants to introduce ordering and ensure the messages that share the same key end up in the same partition.

# Kafka Streams

**Stream** : a sequence of immutable data records, that fully ordered , can be replayed and is fault tolerant(think of Kafka topic as parallel).

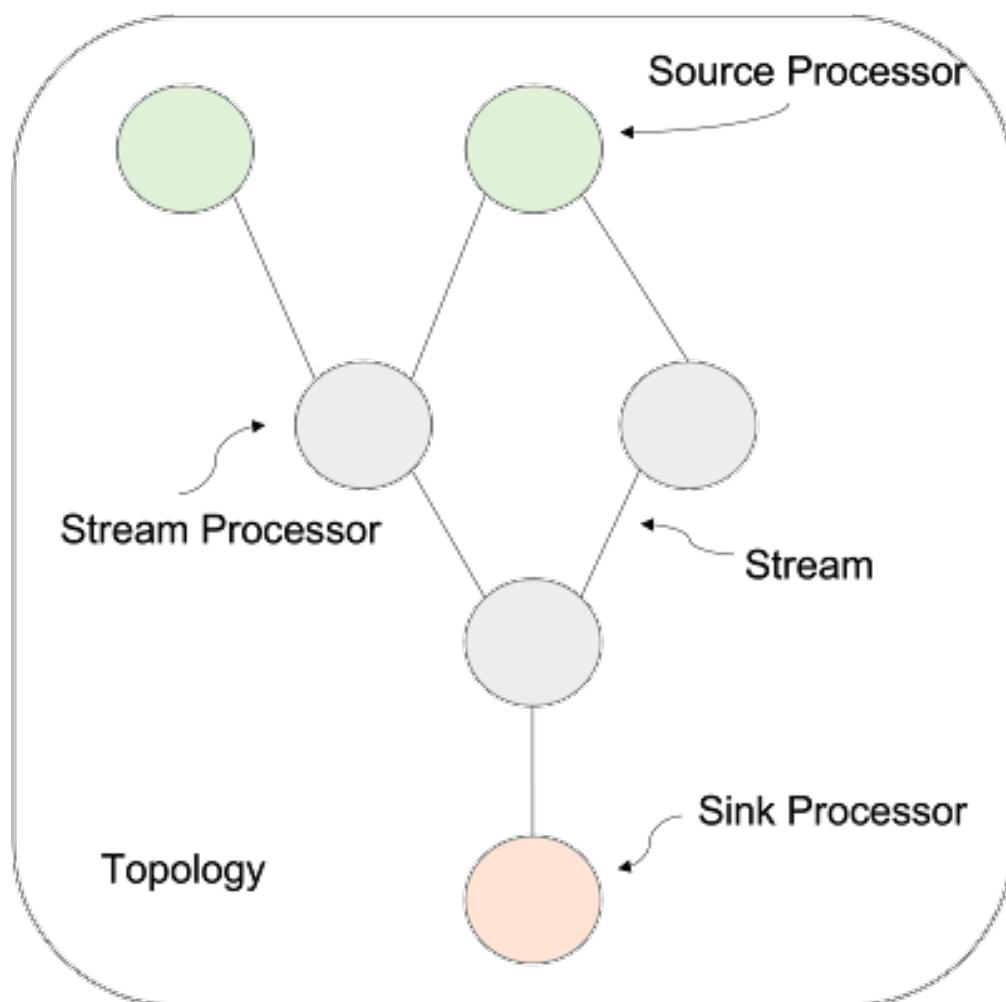
**Stream Processor** : It is node in the processor topology(graph).It transform incoming stream, record by record, and may create a new stream from it.

**Topology** : It is a graph of processors chain together by streams.

**Source Processor** : It is a special processor that takes data directly from Kafka topics. It has no predecessors in a topology and does not transform the data.

**Sink Processor** : It does not have children , it sends stream data directly to Kafka topic.

In between every other Processor, streams not necessary persisted into Kafka.It lives only within Kafka stream application.



**Internal Topics** : running a Kafka Streams may eventually create internal intermediary topics.Two types are :

**1. Repartitioning topics** : Incase we start transforming the key of our stream , a repartitioning will happen at some processor.

**2. Changelog topics** : Incase we perform aggregations, Kafka stream will save compacted data in these topic.

Note : we should not mess with internal topics, they are for internal use of Kafka streams, i.e., save, restore and repartition data.

### We can write a topology by two ways

**High Level DSL** : it is the simplest. It has all the operations we need to perform.

**Low Level Processor Api** : It is an imperative API,Can be used to implement the most complex logic.(Rarely Used)

### Some extra Configuration for Kafka Streams

#### 1. `application.id`

=> Consumer group.id = `application.id`

=> default client.id prefix

=> prefix to internal changelog topics

#### 2. `default.[key/value].serde`

for serialisation and deserialisation of data

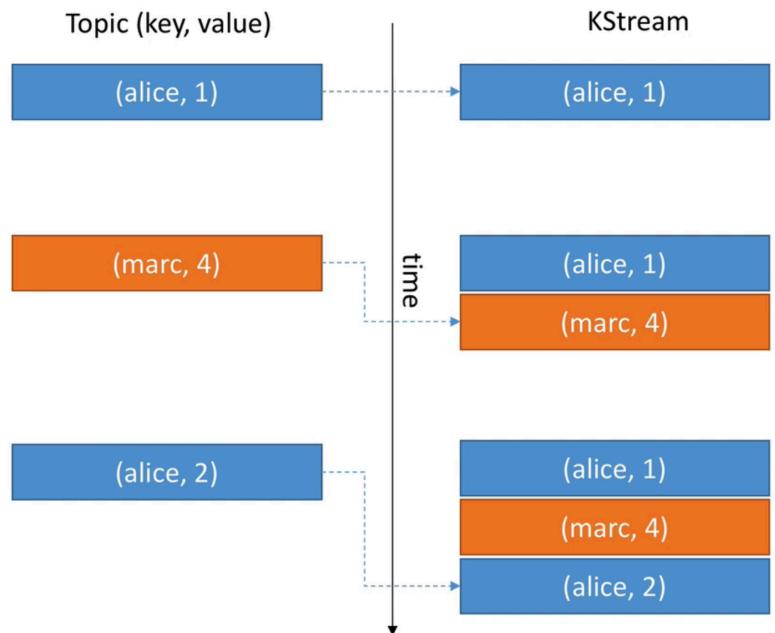
**Scale our Application** : run multiple instances for Kafka streams .

/\*-----End of Basic Doc-----\*/



# KStreams

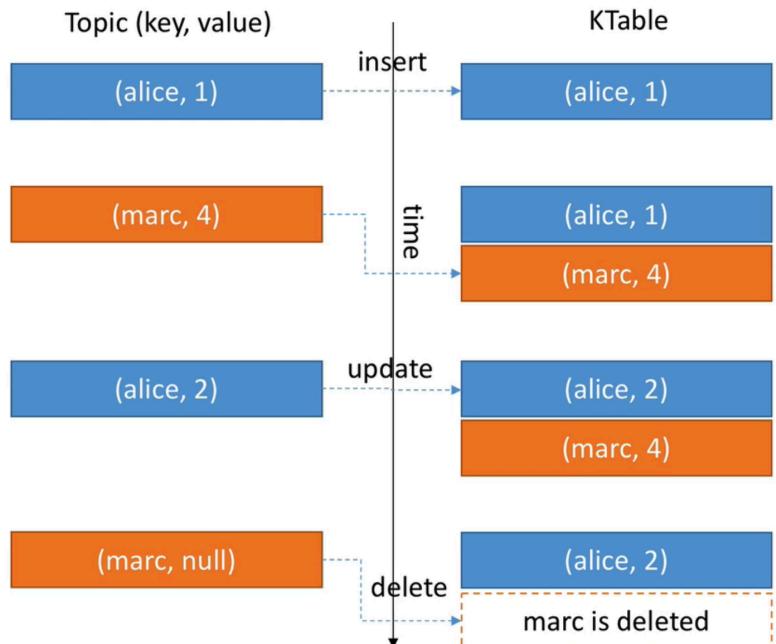
- All **inserts**
  - Similar to a log
  - Infinite
  - Unbounded data streams



# KTables



- All **upserts** on non null values
  - Deletes on null values
  - Similar to a table
  - Parallel with log compacted topics



# When to use KStream vs KTable ?



- **KStream** reading from a topic that's not compacted
- **KTable** reading from a topic that's log-compacted (aggregations)
  
- **KStream** if new data is partial information / transactional
- **KTable** more if you need a structure that's like a “database table”, where every update is self sufficient (think – total bank balance)

# Stateless vs Stateful Operations



- **Stateless** means that the result of a transformation only depends on the data-point you process
  - Example: a “**multiply value by 2**” operation is stateless because it doesn't need memory of the past to be achieved.
  - 1 => 2
  - 300 => 600
- **Stateful** means that the result of a transformation also depends on an external information – the **state**
  - Example: a **count operation** is stateful because your app needs to know what happened since it started running in order to know the computation result
  - hello => 1
  - hello => 2



# MapValues / Map

- Takes one record and produces one record
- **MapValues**
  - Is only affecting values
  - == does not change keys
  - == does not trigger a repartition
  - For KStreams and KTables
- **Map**
  - Affects both keys and values
  - Triggers a re-partitions
  - For KStreams only

```
// Java 8+ example, using lambda expressions KStream<byte[], String>
uppercased = stream.mapValues(value -> value.toUpperCase());
```



# Filter / FilterNot

- Takes one record and produces zero or one record
- **Filter**
  - does not change keys / values
  - == does not trigger a repartition
  - For KStreams and KTables
- **FilterNot**
  - Inverse of Filter

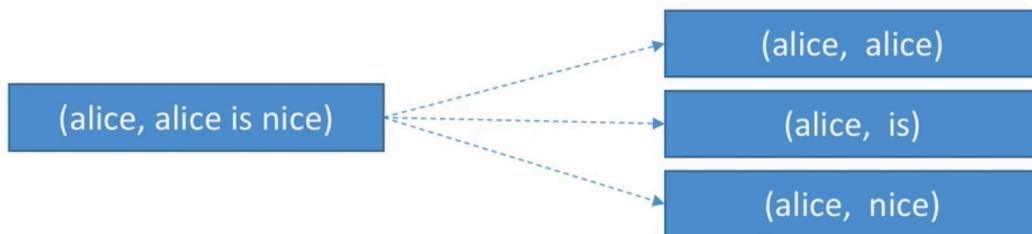
```
// A filter that selects (keeps) only positive numbers
KStream<String, Long> onlyPositives = stream.filter((key, value) -> value > 0);
```



# FlatMapValues / FlatMap

- Takes one record and produces zero, one or more record
- **FlatMapValues**
  - does not change keys
  - == does not trigger a repartition
  - For KStreams only
- **FlatMap**
  - Changes keys
  - == triggers a repartitions
  - For KStreams only

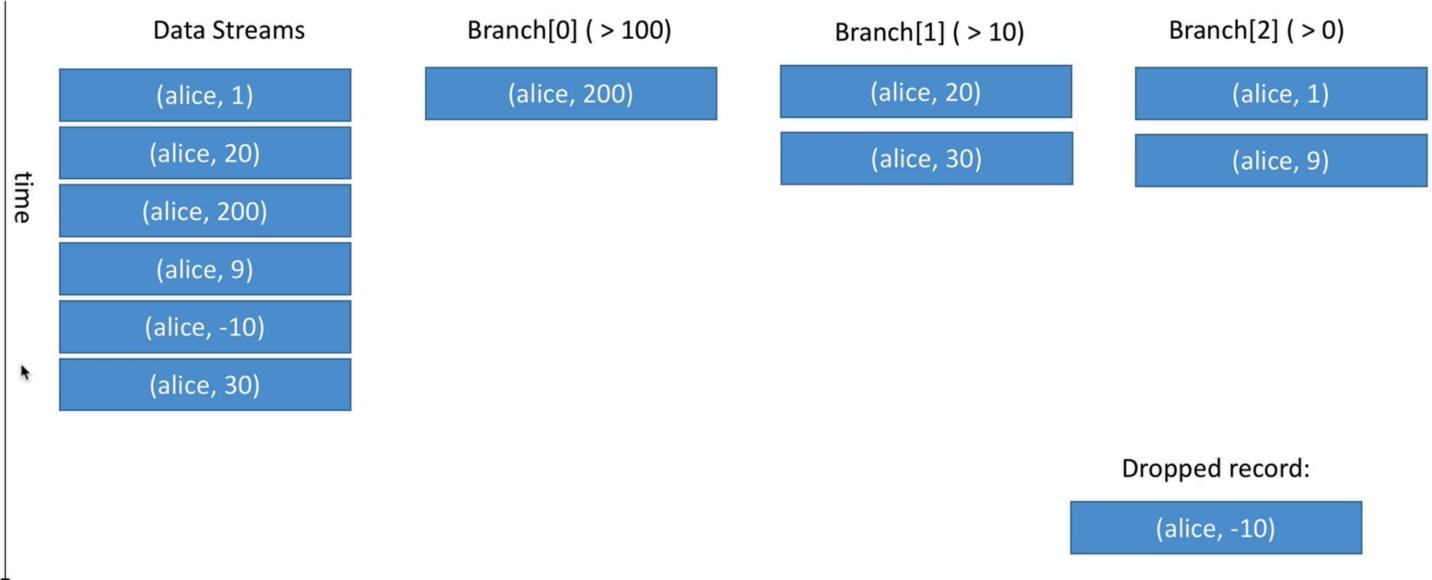
```
// Split a sentence into words.
words = sentences.flatMapValues(value -> Arrays.asList(value.split("\\s+")));
```



## KStream Branch

- Branch (split) a KStream based on one or more predicates
- Predicates are evaluated in order, if no matches, records are dropped
- You get multiple KStreams as a result

```
KStream<String, Long>[] branches = stream.branch(
    (key, value) -> value > 100, /* first predicate */
    (key, value) -> value > 10, /* second predicate */
    (key, value) -> value > 0/* third predicate */
);
```



## KStream SelectKey

- Assigns a new Key to the record (from old key and value)
- == marks the data for re-partitioning
- Best practice to isolate that transformation to know exactly where the partitioning happens.

```
// Use the first letter of the key as the new key
rekeyed = stream.selectKey((key, value) -> key.substring(0, 1))
```



# Reading from Kafka

- You can read a topic as a KStream, a KTable or a GlobalKTable

```
KStream<String, Long> wordCounts = builder.stream(
    Serdes.String(), /* key serde */
    Serdes.Long(), /* value serde */
    "word-counts-input-topic" /* input topic */);
```

```
KTable<String, Long> wordCounts = builder.table(
    Serdes.String(), /* key serde */
    Serdes.Long(), /* value serde */
    "word-counts-input-topic" /* input topic */);
```

```
GlobalKTable<String, Long> wordCounts = builder.globalTable(
    Serdes.String(), /* key serde */
    Serdes.Long(), /* value serde */
    "word-counts-input-topic" /* input topic */);
```

# Writing to Kafka

- You can write any KStream or KTable back to Kafka
- If you write a KTable back to Kafka, think about creating a log compacted topic!
- To: Terminal operation – write the records to a topic

```
stream.to("my-stream-output-topic");
table.to("my-table-output-topic");
```

- Through: write to a topic and get a stream / table from the topic

```
KStream<String, Long> newStream = stream.through("user-clicks-topic");
KTable<String, Long> newTable = table.through("my-table-output-topic");
```

# Streams marked for re-partition

- As soon as an operation can possibly change the key, the stream will be marked for repartition:
  - Map
  - FlatMap
  - SelectKey
- So only use these APIs if you need to change the key, otherwise use their counterparts:
  - MapValues
  - FlatMapValues
- Repartitioning is done seamlessly behind the scenes but will incur a performance cost (read and write to Kafka)

# Refresher on Log Compaction

- Log Compaction can be a huge improvement in performance when dealing with KTables because eventually records get discarded
- This means less reads to get to the final state (less time to recover)
- Log Compaction has to be enabled by you on the topics that get created (source or sink topics)

# Log Cleanup Policy: Compact

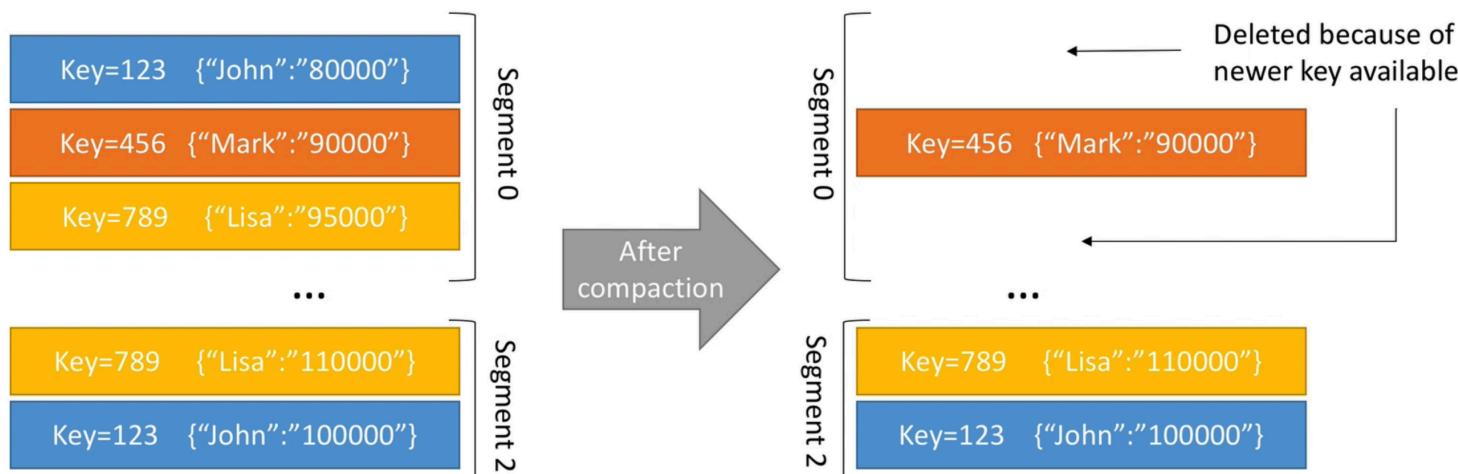


- Log compaction ensures that your log contains at least the last known value for a specific key within a partition
- Very useful if we just require a SNAPSHOT instead of full history (such as for a data table in a database)
- The idea is that we only keep the latest “update” for a key in our log

## Log Compaction: Example



- Our topic is: employee-salary
- We want to keep the most recent salary for our employees



# Log Compaction Guarantees

- **Any consumer that is reading from the head of a log will still see all the messages sent to the topic**
- Ordering of messages it kept, log compaction only removes some messages, but does not re-order them
- The offset of a message is immutable (it never changes). Offsets are just skipped if a message is missing
- Deleted records can still be seen by consumers for a period of delete.retention.ms (default is 24 hours).

# Log Compaction Myth Busting

- It doesn't prevent you from pushing duplicate data to Kafka
  - De-duplication is done after a segment is committed
  - Your consumers will still read from head as soon as the data arrives
- It doesn't prevent you from reading duplicate data from Kafka
  - Same points as above
- Log Compaction can fail from time to time
  - It is an optimization and if the compaction thread might crash
  - Make sure you assign enough memory to it and that it gets triggered

# KStream & KTable Duality (from confluent docs)

- **Stream as Table:** A stream can be considered a changelog of a table, where each data record in the stream captures a state change of the table.
- **Table as Stream:** A table can be considered a snapshot, at a point in time, of the latest value for each key in a stream (a stream's data records are key-value pairs).

## Transforming a KTable to a KStream

- It is sometimes helpful to transform a KTable to a Kstream in order to keep a changelog of all the changes to the Ktable (see last lecture on Kstream / Ktable duality)
- This can be easily achieved in one line of code!

```
KTable<byte[], String> table = ...;

// Also, a variant of `toStream` exists that allows you
// to select a new key for the resulting stream.
KStream<byte[], String> stream = table.toStream();
```

# Transforming a KStream to a KTable

- Two ways:

- Chain a `groupByKey()` and an aggregation step (`count`, `aggregate`, `reduce`)

```
KTable<String, Long> table = usersAndColours.groupByKey()
    .count();
```

- Write back to Kafka and read as KTable

```
// write to Kafka
stream.to("intermediary-topic");

// read from Kafka as a table
KTable<String, String> table = builder.table("intermediary-topic");
```

## Stateful Operations (KStream and KTable)

### KTable GroupBy



- GroupBy allows you to perform more aggregations within a KTable
- We have used it in the previous section during our Favourite Colour Example!
- It triggers a repartition because the key changes.

```
// Group the table by a new key and key type
KGroupedTable<String, Integer> groupedTable = table.groupBy(
    (key, value) -> KeyValue.pair(value, value.length()),
    Serdes.String(), /* key (note: type was modified) */
    Serdes.Integer() /* value (note: type was modified) */ );
```

# KGroupedStream / KGroupedTable Count

---

- As a reminder, KGroupedStream are obtained after a `groupBy/groupByKey()` call on a KStream
- `Count` counts the number of record by grouped key.
- If used on KGroupedStream:
  - Null keys or values are ignored
- If used on KGroupedTable:
  - Null keys are ignored
  - Null values are treated as “delete” (= tombstones)

## KGroupedStream Aggregate

---



- You need an initializer (of any type), an adder, a Serde and a State Store name (name of your aggregation)
- Example: Count total string length by key

```
// Aggregating a KGroupedStream (note how the value type changes from String to Long)
KTable<byte[], Long> aggregatedStream = groupedStream.aggregate(
    () -> 0L, /* initializer */
    (aggKey, newValue, aggValue) -> aggValue + newValue.length(), /* adder */
    Serdes.Long(), /* serde for aggregate value */
    "aggregated-stream-store" /* state store name */);
```

# KGroupedTable Aggregate



- You need an initializer (of any type), an adder, a subtractor, a Serde and a State Store name (name of your aggregation)
- Example: Count total string length by key

```
// Aggregating a KGroupedStream (note how the value type changes from String to Long)
KTable<byte[], Long> aggregatedStream = groupedStream.aggregate(
    () -> 0L, /* initializer */
    (aggKey, newValue, aggValue) -> aggValue + newValue.length(), /* adder */
    (aggKey, oldValue, aggValue) -> aggValue - oldValue.length(), /* subtractor */
    Serdes.Long(), /* serde for aggregate value */
    "aggregated-table-store" /* state store name */);
```

# KGroupedStream / KGroupedTable Reduce

- Similar to [Aggregate](#), but the result type has to be the same as an input:
- (Int, Int) => Int (example: a \* b)
- (String, String) => String (example concat(a, b))

```
// Reducing a KGroupedStream
KTable<String, Long> aggregatedStream = groupedStream.reduce(
    (aggValue, newValue) -> aggValue + newValue, /* adder */
    "reduced-stream-store" /* state store name */);

// Reducing a KGroupedTable
KTable<String, Long> aggregatedTable = groupedTable.reduce(
    (aggValue, newValue) -> aggValue + newValue, /* adder */
    (aggValue, oldValue) -> aggValue - oldValue, /* subtractor */
    "reduced-table-store" /* state store name */);
```

# KGroupedStream / KGroupedTable Reduce

- Similar to [Aggregate](#), but the result type has to be the same as an input:
- $(\text{Int}, \text{Int}) \Rightarrow \text{Int}$  (example:  $a * b$ )
- $(\text{String}, \text{String}) \Rightarrow \text{String}$  (example  $\text{concat}(a, b)$ )

```
// Reducing a KGroupedStream
KTable<String, Long> aggregatedStream = groupedStream.reduce(
    (aggValue, newValue) -> aggValue + newValue, /* adder */
    "reduced-stream-store" /* state store name */);

// Reducing a KGroupedTable
KTable<String, Long> aggregatedTable = groupedTable.reduce(
    (aggValue, newValue) -> aggValue + newValue, /* adder */
    (aggValue, oldValue) -> aggValue - oldValue, /* subtractor */
    "reduced-table-store" /* state store name */);
```

Activ

## KStream Peek



- [Peek](#) allows you to apply a side-effect operation to a KStream and get the same KStream as a result.
- A side effect could be:
  - printing the stream to the console
  - Statistics collection
- Warning: It could be executed multiple times as it is side effect (in case of failures)

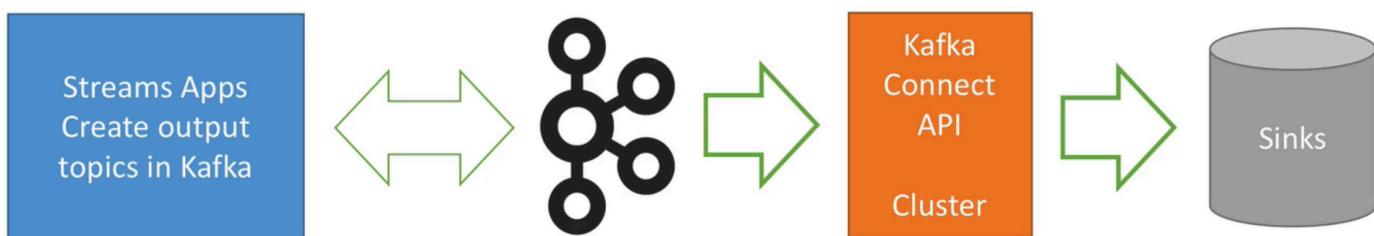
```
KStream<byte[], String> stream = ...;

// Java 8+ example, using lambda expressions
KStream<byte[], String> unmodifiedStream = stream.peek(
    (key, value) -> System.out.println("key=" + key + ", value=" + value));
```



# What if I want to write the result to an external System?

- Although it is theoretically doable to do it using [Kafka Streams library](#), it is NOT recommended
- The recommend way of doing so is Kafka Streams to transform the data and then using **Kafka Connect API** (see my other course)



# Kafka Administration

## Factors impacting Kafka performance

### Disks: I/O



- Reads are done sequentially (as in not randomly), therefore make sure you should a disk type that corresponds well to the requirements
- Format your drives as XFS (easiest, no tuning required)
- If read / write throughput is your bottleneck
  - it is possible to mount multiple disks in parallel for Kafka
  - The config is `log.dirs=/disk1/kafka-logs,/disk2/kafka-logs,/disk3/kafka-logs...`
- Kafka performance is constant with regards to the amount of data stored in Kafka.
  - Make sure you expire data fast enough (default is one week)
  - Make sure you monitor disk performance

## Factors impacting Kafka performance

### Network



- Latency is key in Kafka
  - Ensure your Kafka instances are your Zookeeper instances are geographically close!!!
  - Do not put one broker in Europe and the other broker in the US
  - Having two brokers live on the same rack is good for performance, but a big risk if the rack goes down.
- Network bandwidth is key in Kafka
  - Network will be your bottleneck.
  - Make sure you have enough bandwidth to handle many connections, and TCP requests.
  - Make sure your network is high performance
- Monitor network usage to understand when it becomes a bottleneck

# Factors impacting Kafka performance RAM



- Kafka has amazing performance thanks to the page cache which utilizes your RAM
- Understanding RAM in Kafka means understanding two parts:
  - The Java HEAP for the Kafka process
  - The rest of the RAM used by the OS page cache
- Let's understand how both of those should be sized
- Overall, your Kafka production machines should have at least 8GB of RAM to them (the more the better – it's common to have 16GB or 32GB per broker)

Ex: 16 GB of RAM

Kafka Heap  
(4G to start)

OS page  
cache  
(free memory  
automatically  
assigned by the  
OS)

Activate Windows  
Go to Settings to activate Win

## Factors impacting Kafka performance RAM – Java Heap



- When you launch Kafka, you specify Kafka Heap Options (KAFKA\_HEAP\_OPTS environment variable)
- I recommend to assign a MAX amount (-Xmx) of 4GB to get started to the kafka heap:
- **`export KAFKA_HEAP_OPTS="-Xmx4g"`**
- Don't set -Xms (starting heap size):
  - Let heap grow over time
  - Monitor the heap over time to see if you need to increase Xmx
- Kafka should keep a low heap usage over time, and heap should increase only if you have more partitions in your broker

Ex: 16 GB of RAM

Kafka Heap  
(4G to start)

OS page  
cache  
(free memory  
automatically  
assigned by the  
OS)

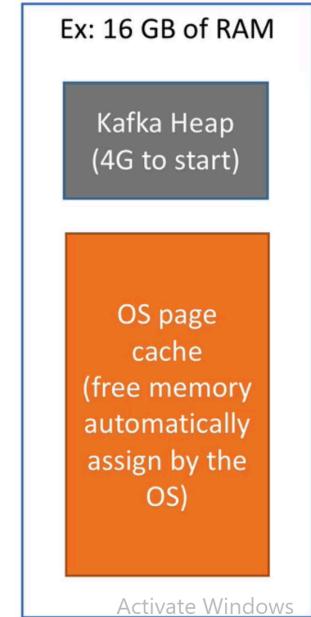
Activate Windows

# Factors impacting Kafka performance

## RAM – OS Page Cache



- The remaining RAM will be used automatically for the Linux **OS Page Cache**.
- This is used to buffer data to the disk and this is what gives Kafka an amazing performance
- You don't have to specify anything!
- Any un-used memory will automatically be leveraged by the Linux Operating System and assign memory to the page cache
- Note: Make sure swapping is disabled for Kafka entirely  
**vm.swappiness=0 or vm.swappiness=1  
(default is 60 on Linux)**



# Factors impacting Kafka performance

## CPU

- CPU is usually not a performance bottle neck in Kafka because Kafka does not parse any messages , but can become one in some situations
- If you have SSL enabled, Kafka has to encrypt and decrypt every payload, which adds load on the CPU
- Compression can be CPU bound if you force Kafka to do it. Instead, if you send compressed data, make sure your producer and consumers are the ones doing the compression work (that's the default setting anyway)
- Make sure you monitor Garbage Collection over time to ensure the pauses are not too long



# Factors impacting Kafka performance

## Operating System (OS)

- Use Linux or Solaris, running production Kafka clusters on Windows is not recommended.
- Increase the file descriptor limits (at least 100,000 as a starting point)  
<https://unix.stackexchange.com/a/8949/170887>
- Make sure only Kafka is running on your Operating System. Anything else will just slow the machine down.

## Factors impacting Kafka performance

### Other

- Make sure you have enough file handles opened on your servers, as Kafka opens 3 file descriptor for each topic-partition-segment that lives on the Broker.
- Make sure you use Java 8
- You may want to tune the GC implementation: (see in the resources)
- Set Kafka quotas in order to prevent unexpected spikes in usage

## Running Kafka on AWS in Production

- Separate your instances between different availability zones
- Use `st1` EBS volumes for the best price / performance ratio
- Mount multiple EBS volumes to the same broker if you need to scale
- Use `r4.xlarge` or `m4.2xlarge` if you're using EBS (these instances are EBS optimized). You can use something smaller but performance may degrade
- Setup DNS names for your brokers / fixed IPs so that your clients aren't affected if you recycle your EC2 instances

# Advanced Kafka Configuration

## Important parameters to be aware of



- `auto.create.topics.enable=true` => set to false in production
- `background.threads=10` => increase if you have a good CPU
- `delete.topic.enable=false` => your choice
- `log.flush.interval.messages` => don't ever change. Let your OS do it
- `log.retention.hours=168` => Set in regards to your requirements
- `message.max.bytes=1000012` => increase if you need more than 1MB
- `min.insync.replicas=1` => set to 2 if you want to be extra safe
- `num.io.threads=8` => ++if your network io is a bottleneck
- `num.network.threads=3` => ++if your network is a bottleneck
  
- `num.recovery.threads.per.data.dir=1` => set to number of disks
- `num.replica.fetchers=1` => increase if your replicas are lagging
- `offsets.retention.minutes=1440` => after 24 hours you lose offsets
- `unclean.leader.election.enable=true` => false if you don't want data loss
- `zookeeper.session.timeout.ms=6000` => increase if you timeout often
- `broker.rack=null` => set your to availability zone in AWS
- `default.replication.factor=1` => set to 2 or 3 in a kafka cluster
- `num.partitions=1` => set from 3 to 6 in your cluster
- `quota.producer.default=10485760` => set quota to 10MBs
- `quota.consumer.default=10485760` => set quota to 10MBs

Next :

Automate Deployment  
Setup monitoring