



Shaun Thomas

mini sLab | Kafka in fintech: An introduction to Apache Kafka

<https://www.meetup.com/Sytac-sLab/events/256434066/>

Introduction

Kafka

Kafka® is used for building real-time data pipelines and streaming apps. It is horizontally scalable, fault-tolerant, wicked fast, and runs in production in thousands of companies.

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



HOME
INTRODUCTION
QUICKSTART
USE CASES
DOCUMENTATION
PERFORMANCE
POWERED BY
PROJECT INFO
ECOSYSTEM
CLIENTS
EVENTS
CONTACT US
APACHE

[Download](#)

@apachekafka

PUBLISH & SUBSCRIBE

[Read and write streams of data like a messaging system.](#)

[Learn more »](#)

PROCESS

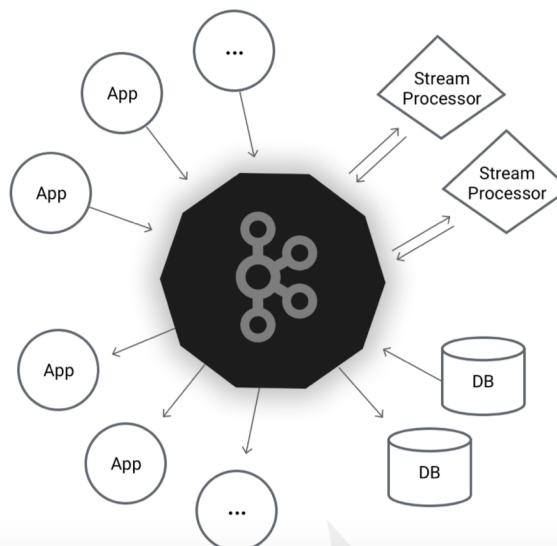
[Write scalable stream processing applications that react to events in real-time.](#)

[Learn more »](#)

STORE

[Store streams of data safely in a distributed, replicated, fault-tolerant cluster.](#)

[Learn more »](#)



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

History/Facts/Figures

- Originally developed at LinkedIn in 2010
- Open sourced in 2011
- Graduation from the Apache Incubator in 2012
- Jun Rao, Jay Kreps, and Neha Narkhede (worked on Kafka at LinkedIn)
 - created a new company named **Confluent**
- Written in Java and Scala
- Kafka adoption exploded in 2017
 - Used by 1/3 of Fortune 500 companies
 - 6 of top ten travel companies, 7 of top ten banks, 8 of top ten insurance companies, 9 of top ten telecom companies uses Kafka in production
 - Version – 0.10.x
- Current version – 2.1.x
- Spring-kafka – 2.2.x (based on kafka-clients jar – Official client)
- Client are available for many programming languages

Basic concepts

3 main components of Kafka

Component	API	Description
Publish/Subscribe (Pub/Sub)	Producer API Consumer API	Sending and receiving messages efficiently and reliably at scale.
Kafka Connect	Connect API	Integrating Kafka with external data sources
Kafka Streams	Streams API	Processing messages in real time

Broker, record and topic

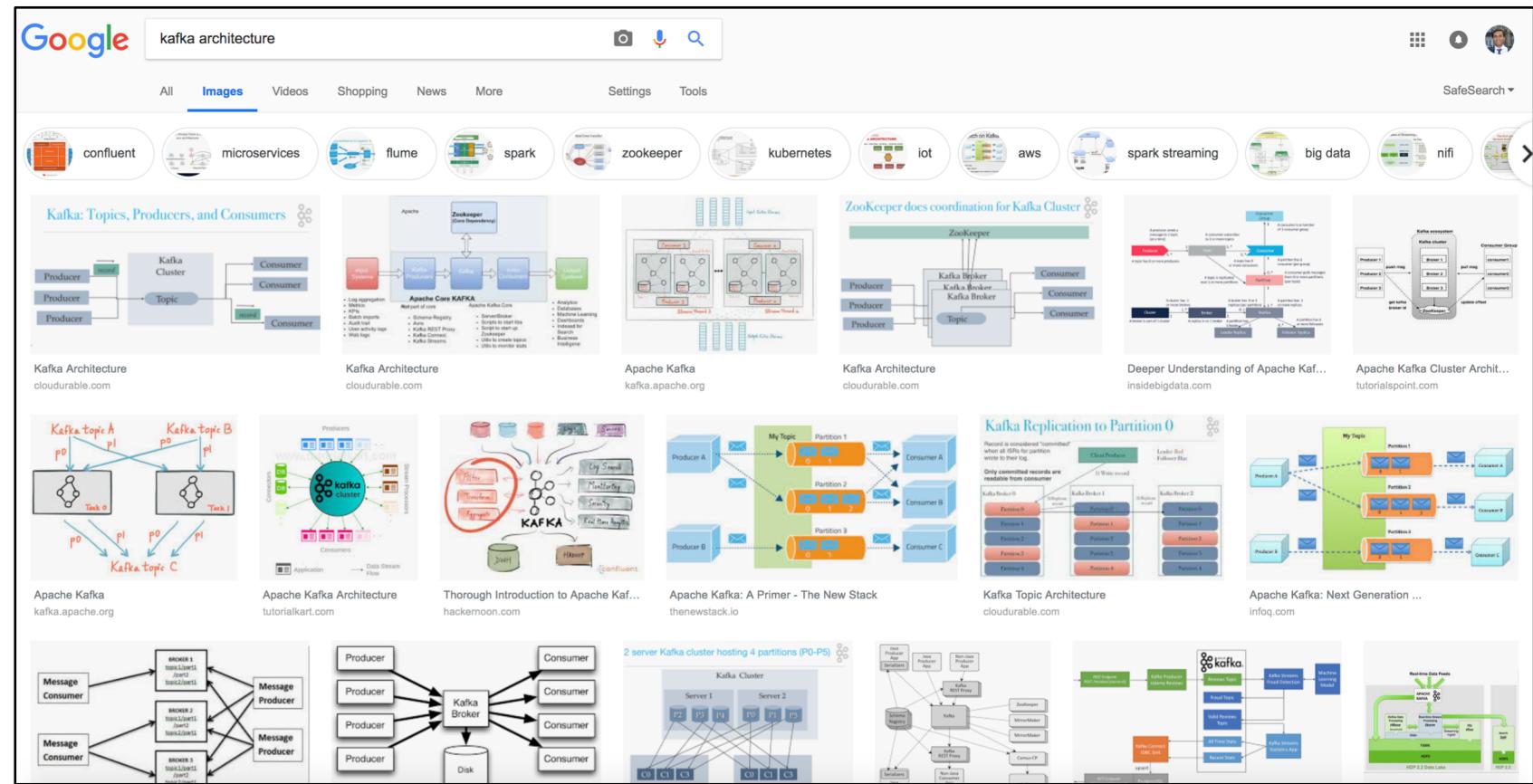
Broker / Server / Bootstrap server	Servers where data are stored. Multiple brokers can be used to form a cluster.
Message / Record	Data stored in Kafka. A record contains a key, a value, a timestamp and a list of headers
Topic / Event	Logical name for producing, consuming or separating records

Zookeeper

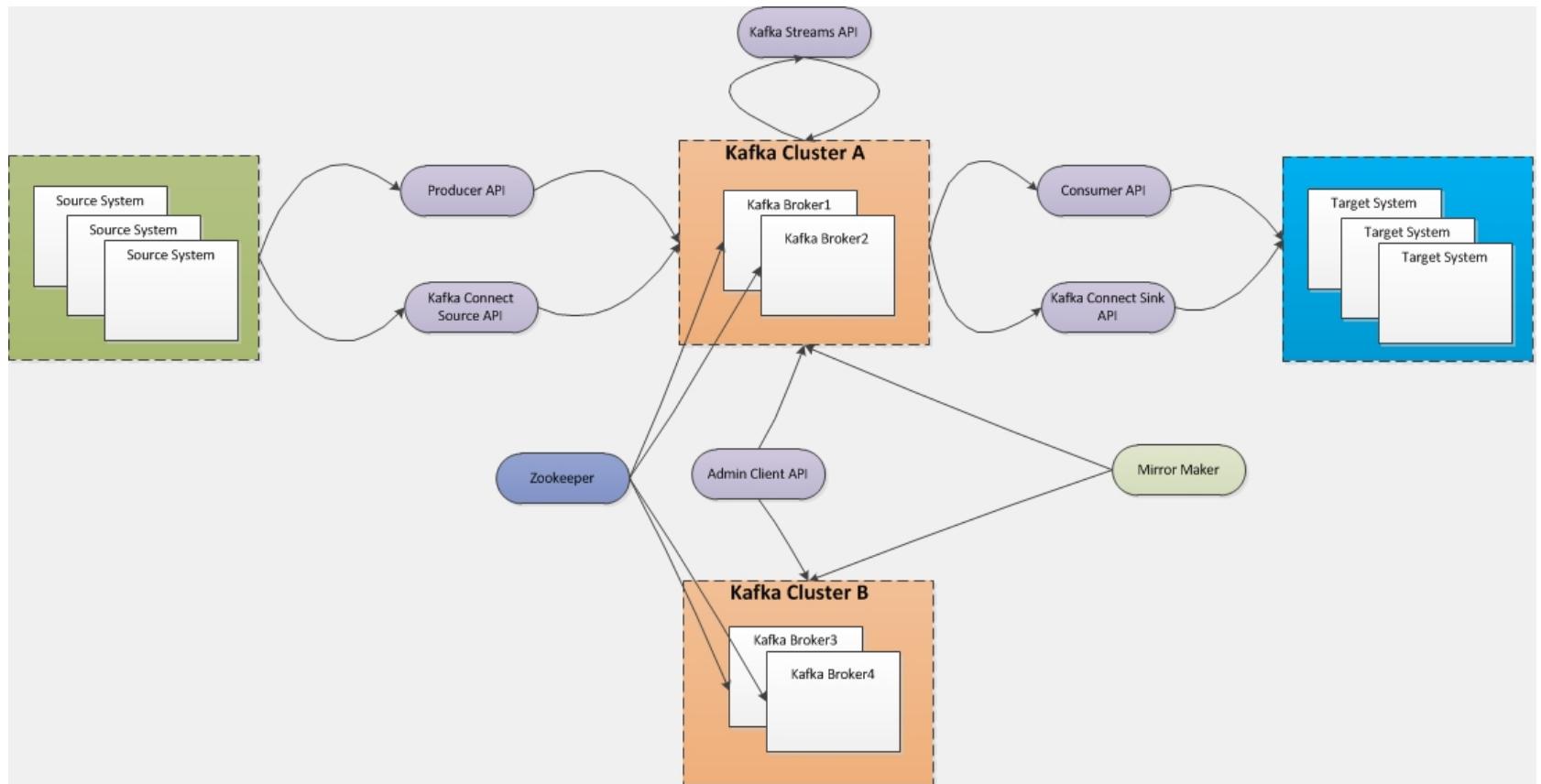
- Zookeeper is an open source server that enable highly reliable distributed co-ordination
- We just have to start Zookeeper and don't worry about it at all
- It helps in maintaining the cluster membership. For example, when a new broker is added, a broker is removed, a new topic is added or a topic is deleted, when a broker goes down or comes up etc, Zookeeper manages such situations, informs Kafka.

Kafka Architecture

Too many varieties...

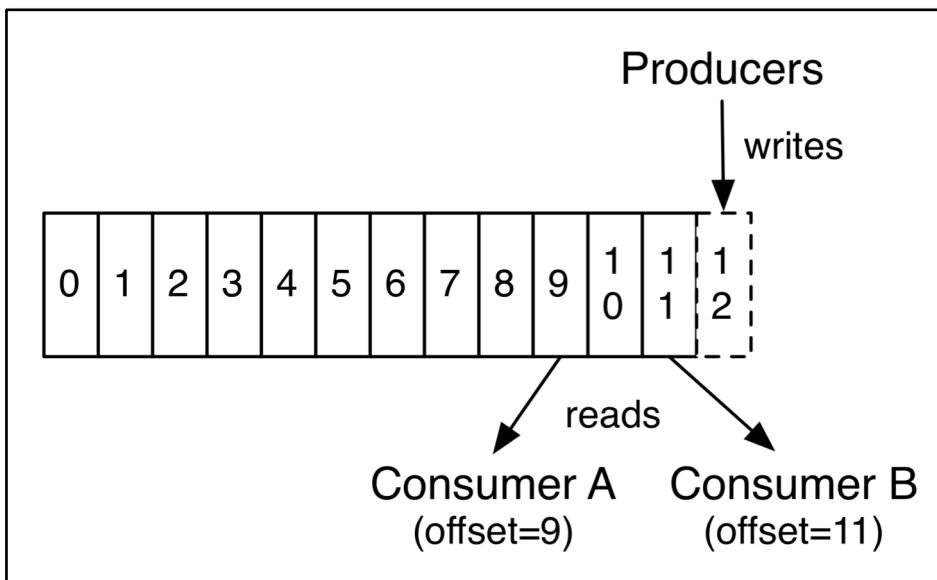


Kafka Architecture / Ecosystem



<https://iteritory.com/beginners-guide-apache-kafka-basic-architecture-components-concepts/>

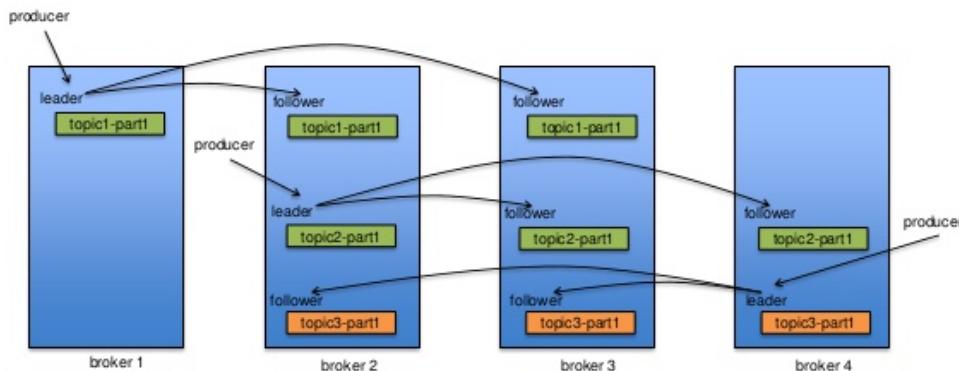
Commit log, retention policy and offset



- Commit log
 - ordered, immutable sequence of records
 - continually appended
 - storage mechanism in Kafka
- Retention
 - Common policies
 - ‘log.retention.hours’ - time
 - ‘log.retention.bytes’ - size
- Offset
 - sequence id assigned to a record in a commit log
 - Uniquely identifies a record
 - Maintained as metadata per consumer

Partitions, replicas, ISRs (In-Sync Replicas), leader and followers

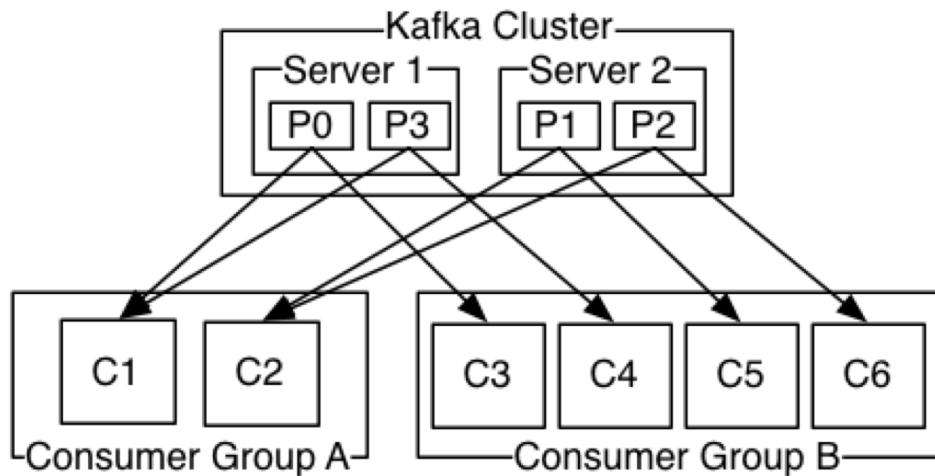
Extend to Multiple Partitions



- Leaders are evenly spread among brokers

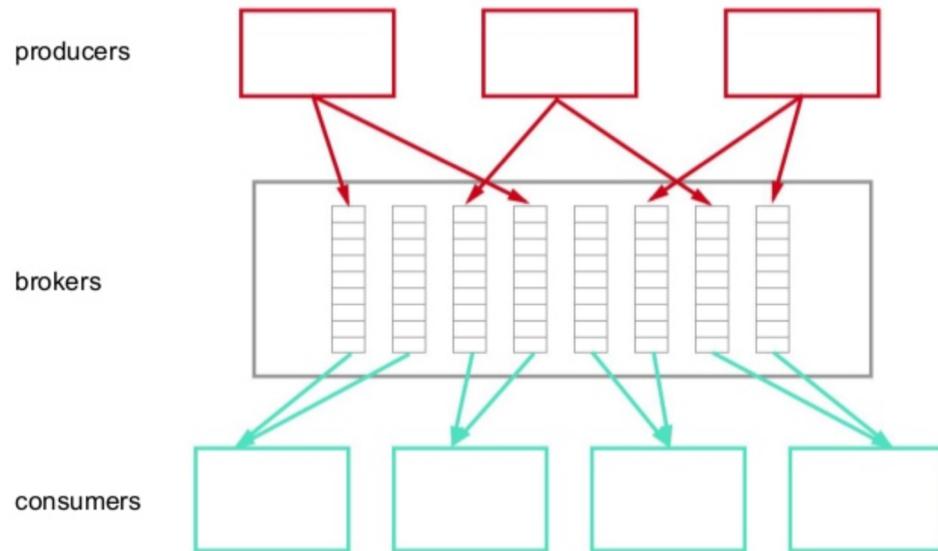
- Each topic should have atleast 1 partition and 1 replica
- Commit log = Partition
- 1 Partition = 1 Replica
- Record is said to be committed when all replicas per partition writes to their log

Consumer group



- Consumer(s) of a topic should belong to a consumer group
- Consumer instances can be in separate processes or on separate machines
- Guaranteed processing of record by atleast one consumer instance
 - Acknowledgement
- No guarantee of order across multiple consumer instances

Core Kafka Pt. 5: How Scalable is Kafka?



- **No bottleneck!**
 - Many brokers
 - Many producers
 - Many consumers
- **Limits?**
 - Internet giants are driving the limits higher - you won't need to worry.
 - e.g. LinkedIn > 1 trillion messages / day through Kafka clusters.
 - 100 brokers / 2 billion messages a day is "straightforward" to operate
 - Don't over partition
~< 100k partitions

Quick start - Kafka command line tool

Hand-on – Command line

- Download and install Kafka from -
https://www.apache.org/dyn/closer.cgi?path=/kafka/2.0.1/kafka_2.11-2.0.1.tgz
- Extract the contents of zip file to a directory in your filesystem
- Follow instructions in <http://kafka.apache.org/quickstart>
 - Start Zookeeper
 - Start Kafka server/broker



Topic commands

- Create topic

```
bin/kafka-topics.sh --create --zookeeper localhost:2181 --replication-factor 2 \  
--partitions 4 --topic my-replicated-topic
```

- Describe topic

```
bin/kafka-topics.sh --describe --zookeeper localhost:2181 --topic my-replicated-topic
```

- Delete topic

```
bin/kafka-topics.sh --delete --zookeeper localhost:2181 --topic my-replicated-topic
```

Pub/Sub – Kafka Producer/Consumer

Hands-on – Pub/Sub

Base project - <https://github.com/shaunthomas999/kafka-workshop-base>

Complete project - <https://github.com/shaunthomas999/kafka-workshop>



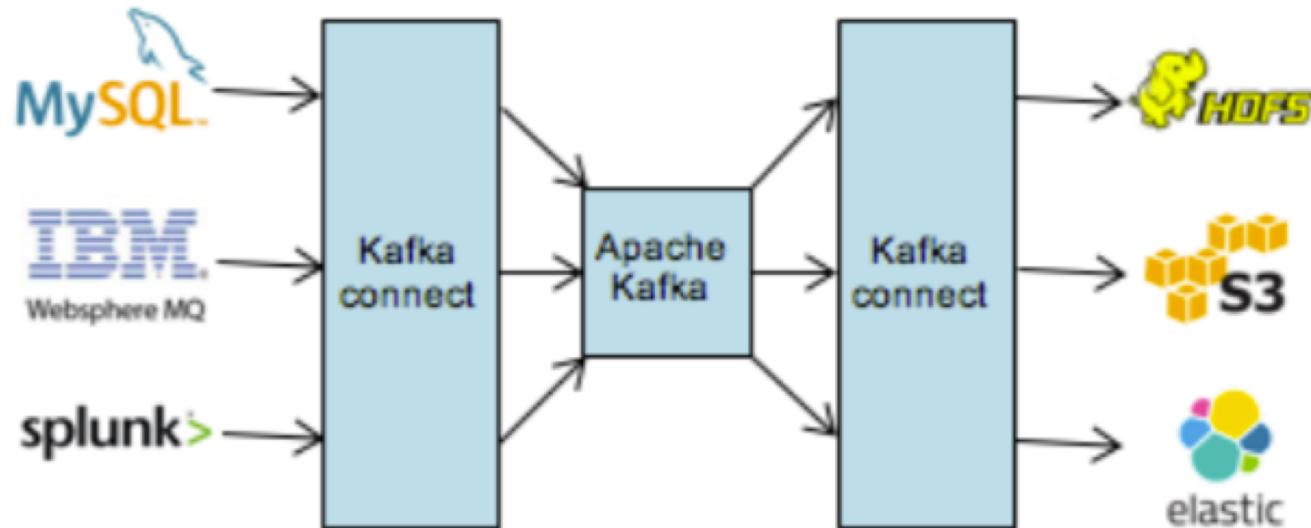
```
StockPrice {  
    String name;  
    String symbol;  
    String type;  
    String datetime;  
    BigDecimal price;  
    String currency;  
}
```

Pub/Sub - Advanced

- Messages types
 - Json
 - Protobuf
 - Avro
- Schema Management
 - <https://docs.confluent.io/current/schema-management.html>

Kafka Connect

Kafka Connect Overview



Discover Connectors

<https://www.confluent.io/hub/>

Mind: Sources/Sink

The screenshot shows the Confluent Connector Hub interface. At the top, there's a navigation bar with the Confluent logo, a search bar labeled "Search connectors", and links for Product, Cloud, Developers, Blog, Docs, Download, and a magnifying glass icon. Below the navigation is a filter bar with tabs for All, Verified, Sources, Sinks, and Community. The main area displays a 4x4 grid of connector cards. Each card includes the Confluent logo, the connector name, the provider (Confluent, Inc. or otherwise), and download/installation statistics.

Category	Connector Name	Provider	Downloads	Installations	Rating (%)
Sources	Kafka Connect Salesforce	Confluent, Inc.	3	0%	-
	Kafka Connect Syslog	Confluent, Inc.	19	0%	-
	Kafka Connect Datagen	Confluent, Inc.	99	0%	-
	Kafka Connect Neo4j	Confluent, Inc.	61	88%	-
Sinks	Kafka Connect S3	Confluent, Inc.	161	86%	-
	Kafka Connect JDBC	Confluent, Inc.	487	91%	-
	Kafka Connect Avro Converter	Confluent, Inc.	86	83%	-
	Kafka Connect JMS	Confluent, Inc.	34	73%	-
Other	Kafka Connect IBM MQ	Confluent, Inc.	-	-	-
	Kafka Connect ActiveMQ	Confluent, Inc.	-	-	-
	Kafka Connect Elasticsearch	Confluent, Inc.	-	-	-
	Kafka Connect HDFS	Confluent, Inc.	-	-	-

Connect REST API

CONNECT REST API	MEANING
GET /connectors	Return a list of active connectors
POST /connectors	Create a new connector
GET /connectors/{name}	Get the information of a specific connector
GET /connectors/{name}/config	Get configuration parameters for a specific connector
PUT /connectors/{name}/config	Update configuration parameters for a specific connector
GET /connectors/{name}/status	Get the current status of the connector

- Typically runs in distributed mode
- Connect can be managed through REST API

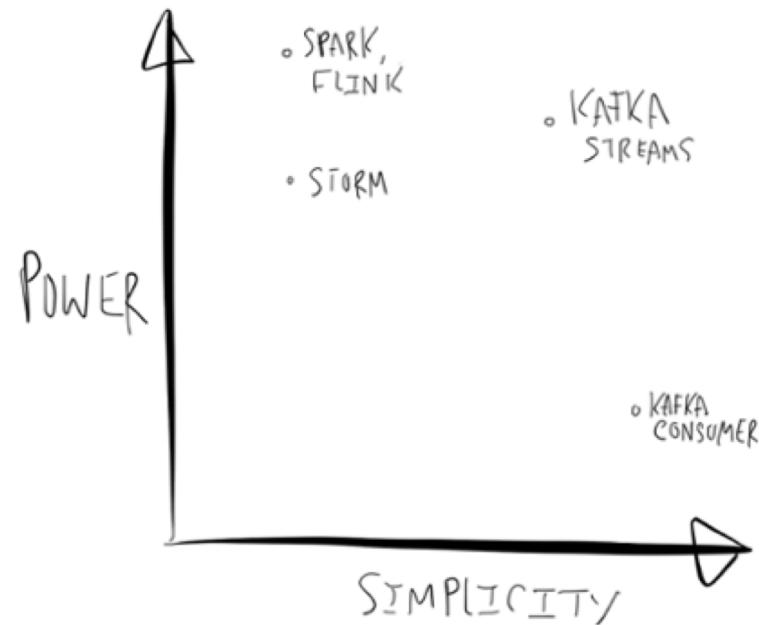
Kafka Streams

Stream Processing Introduction

- Stream Processing is a Big data technology
- Query continuous data stream and detect conditions fast within a small time period from the time of receiving the data
- It is also called by many names: real-time analytics, streaming analytics, complex event processing, real-time streaming analytics, and event processing

<https://medium.com/stream-processing/what-is-stream-processing-1eadfca11b97>

Streaming platforms comparison



<https://www.confluent.io/blog/introducing-kafka-streams-stream-processing-made-simple/>

Kafka Streams Types

- KStreams and KTables (high level)
 - Stream processing DSL
 - Offers high level operators like filter, map, grouping, aggregation, joins, windowing etc.
- Stream Processors (low level)
 - Can create custom operators
 - Low level development approach
- KSQL
 - Streaming SQL for Apache Kafka
 - <https://www.confluent.io/product/ksql/>

KStreams vs KTables

	KSTREAM	KTABLE
CONCEPT	Each record is treated as an <u>append</u> to the stream.	Each record is treated as an <u>update</u> to an existing key
USAGE	Model append-only data such as click streams.	Model updatable reference data such as user profiles.

(key , value) records	Sum of values As KStream	Sum of values As KTable
(“k1”, 2) (“k1”, 5)	7	5

Hands-on – KStream

Method Summary

All Methods	Instance Methods	Abstract Methods
<code>KStream<K,V>[]</code>	<code>branch(Predicate<? super K,>? super V>... predicates)</code>	Creates an array of <code>KStream</code> from this stream by branching the records in the original stream based on the supplied predicates.
<code>KStream<K,V></code>	<code>filter(Predicate<? super K,>? super V> predicate)</code>	Create a new <code>KStream</code> that consists of all records of this stream which satisfy the given predicate.
<code>KStream<K,V></code>	<code>filterNot(Predicate<? super K,>? super V> predicate)</code>	Create a new <code>KStream</code> that consists all records of this stream which do <i>not</i> satisfy the given predicate.
<code><KR,VR> KStream<KR,VR></code>	<code>flatMap(KeyValueMapper<? super K,>? super V,> extends java.lang.Iterable<? extends KeyValue<? extends KR,>? extends VR>>> mapper)</code>	Transform each record of the input stream into zero or more records in the output stream (both key and value type can be altered arbitrarily).
<code><VR> KStream<K,VR></code>	<code>flatMapValues(ValueMapper<? super V,>? extends java.lang.Iterable<? extends VR>>> mapper)</code>	Create a new <code>KStream</code> by transforming the value of each record in this stream into zero or more values with the same key in the new stream.
<code><VR> KStream<K,VR></code>	<code>flatMapValues(ValueMapperWithKey<? super K,>? super V,>? extends java.lang.Iterable<? extends VR>>> mapper)</code>	Create a new <code>KStream</code> by transforming the value of each record in this stream into zero or more values with the same key in the new stream.
<code>void</code>	<code>foreach(ForeachAction<? super K,>? super V> action)</code>	Perform an action on each record.
<code><KR> KGroupedStream<KR,V></code>	<code>groupBy(KeyValueMapper<? super K,>? super V,KR> selector)</code>	Group the records of this <code>KStream</code> by their current key using the provided <code>KeyValueMapper</code> and default serialization.
<code><KR> KGroupedStream<KR,V></code>	<code>groupBy(KeyValueMapper<? super K,>? super V,KR> selector, Serialized<KR,V> serialized)</code>	Group the records of this <code>KStream</code> by their current key using the provided <code>KeyValueMapper</code> and specified serialization.
<code>KGroupedStream<K,V></code>	<code>groupByKey()</code>	Group the records by their current key using the original values and default serialization.
<code>KGroupedStream<K,V></code>	<code>groupByKey(Serialized<K,V> serialized)</code>	Group the records by their current key using the original values and using the specified serialization.

<https://github.com/shaunthomas999/kafka-workshop>

<https://kafka.apache.org/20/javadoc/org/apache/kafka/streams/kstream/KStream.html>



Architectures

STREAM PROCESSING

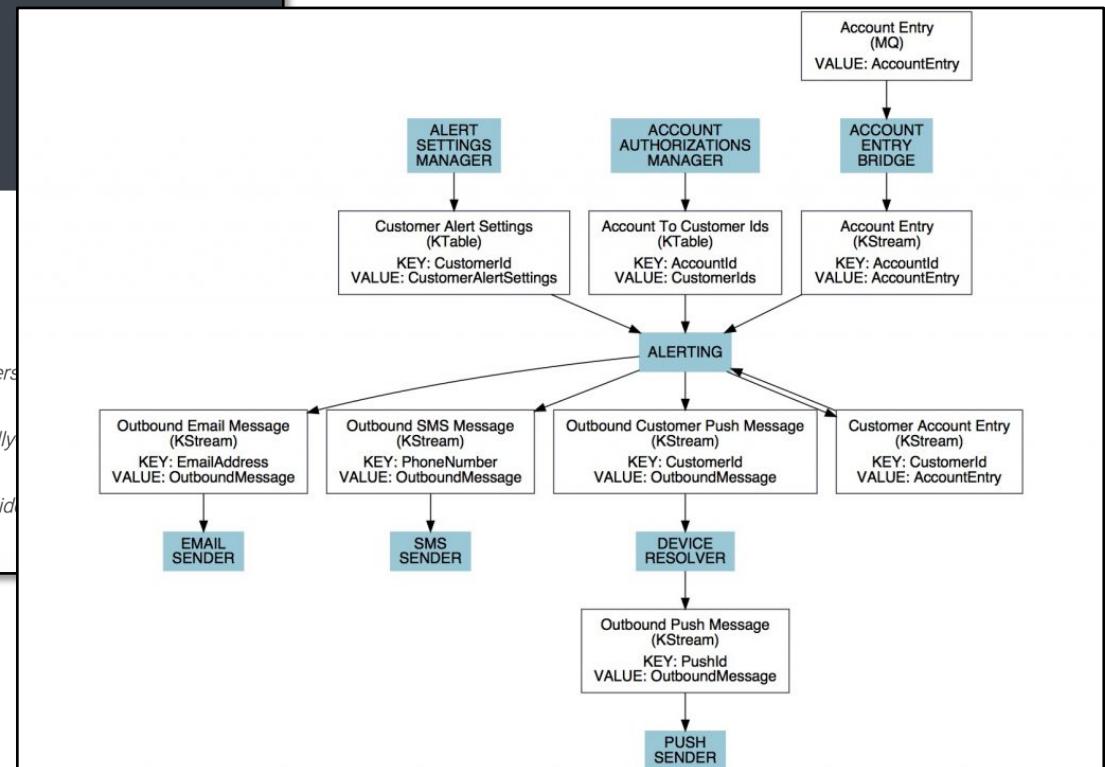
Real-time Financial Alerts at Rabobank with Apache Kafka's Streams API



Jeroen van Disseldorp

July 20, 2017

This article discusses the use of Apache Kafka's Streams API for sending out alerts to customers. Rabobank is based in the Netherlands with over 900 locations worldwide, 48,000 employees, and €681B in assets. It is a bank by and for customers, a cooperative bank, a socially responsible bank. It aims to be market leader across all financial markets in the Netherlands. Rabobank is also committed to being a leading bank in the field of food and agriculture worldwide. Rabobank provides financial products and services to millions of customers around the world.



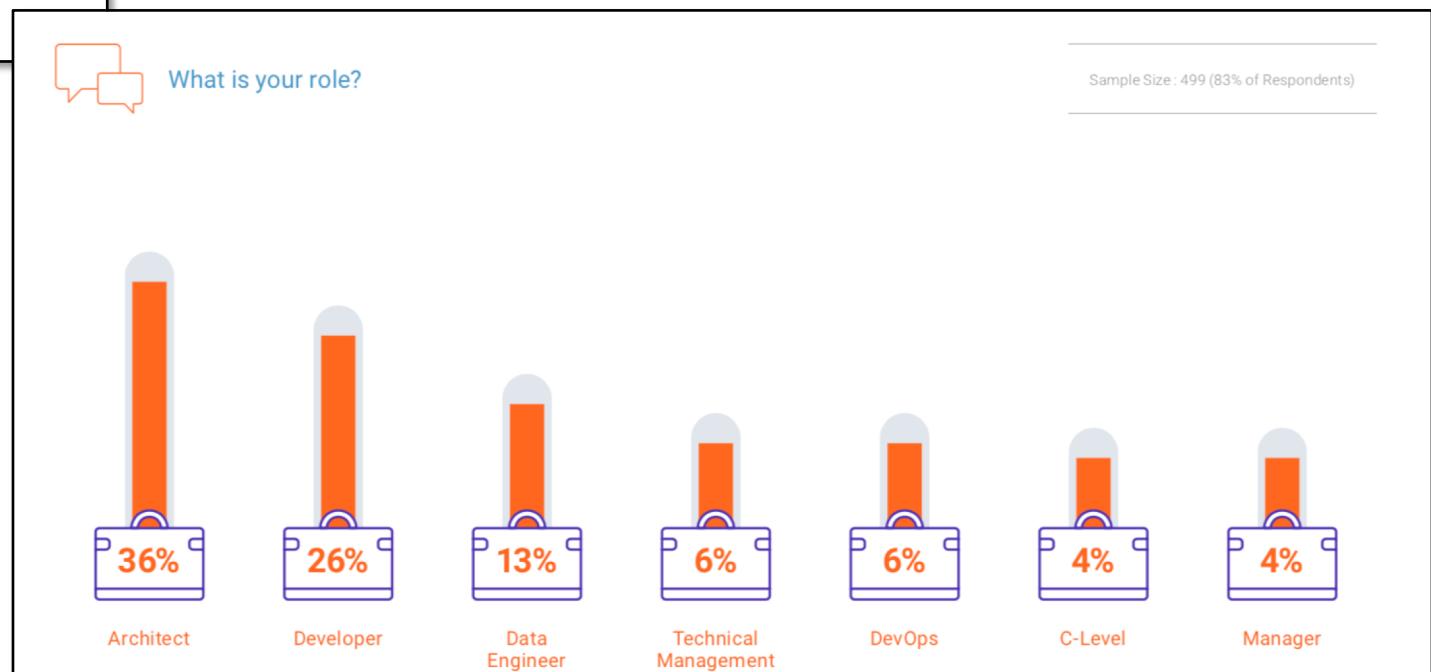
<https://www.confluent.io/blog/real-time-financial-alerts-rabobank-apache-kafkas-streams-api/>

Survey results

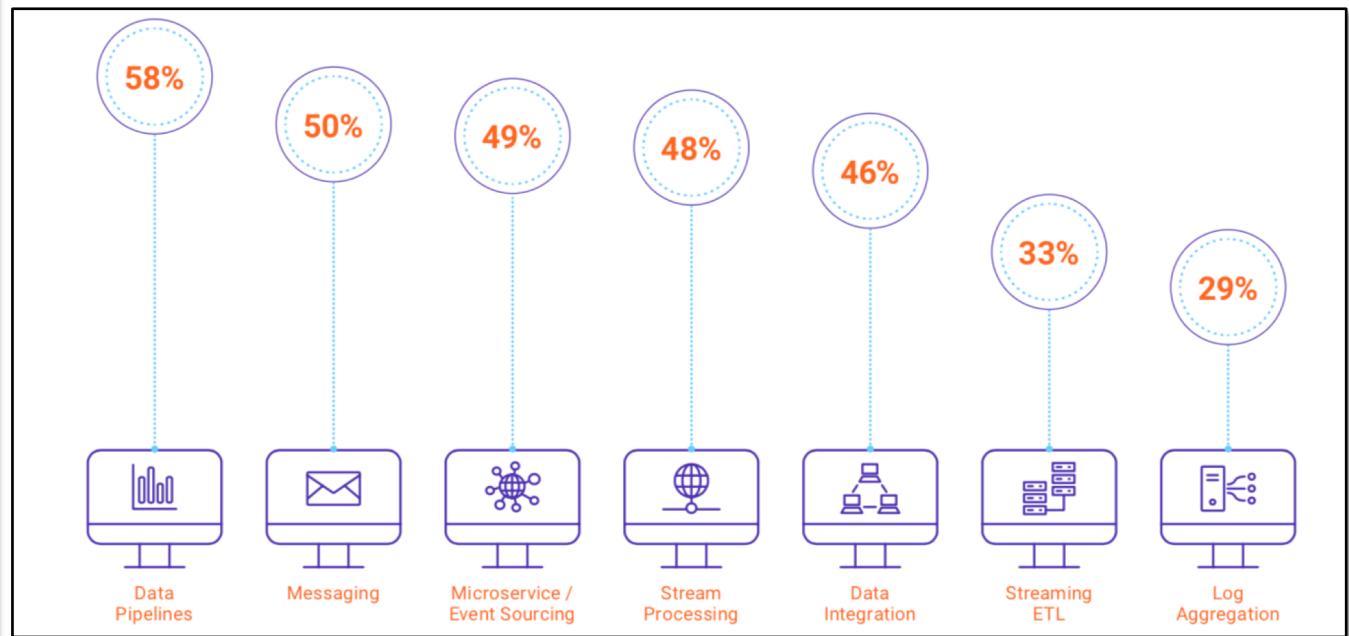
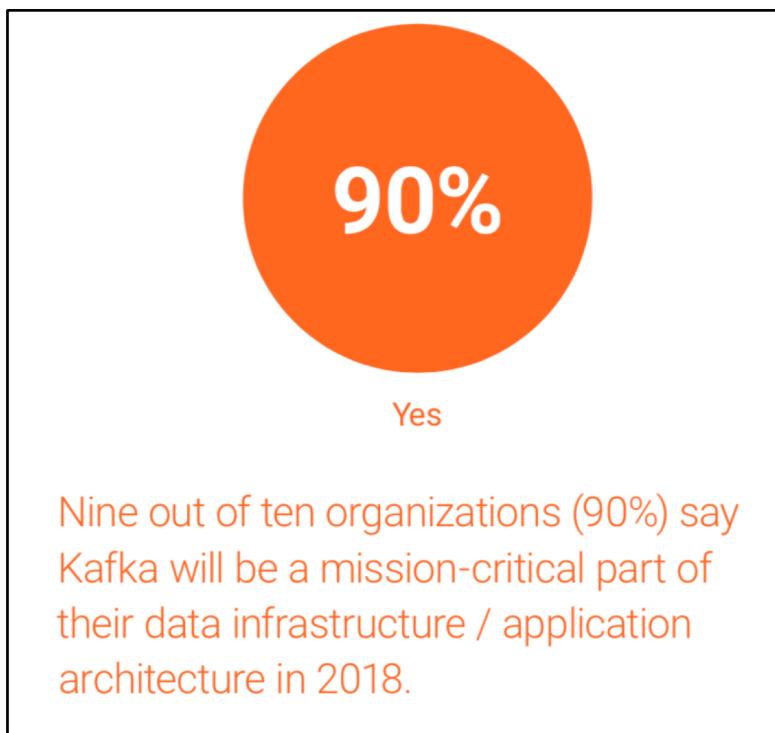
2018 Apache Kafka Report

Confluent surveyed over 600 Apache Kafka users from 59 countries in order to better understand the behavior, usage and attitudes of Kafka's growing user base.

<https://www.confluent.io/apache-kafka-report/>



Important results/trends



Contact

shaunthomas999@gmail.com

www.shaunthomas999.com

Twitter / LinkedIn / Facebook: shaunthomas999