

Introduction to Kafka

⇒ Apache Spark
Apache Airflow
Apache Kafka

⇒ What is Apache Kafka?

① Kafka is open-source, distributed, event-streaming platform.

② Kafka is designed to handle large quantities of data, & to be scalable to various situations.

③ Event-streaming varies based on context.

↓
in Kafka, it means obtaining information from various sources, reliably storing that data, & providing it to any users or systems needing access.

④ Kafka has many uses & can fit into many different scenarios.

1) E-commerce:- allowing different systems to monitor sales data in up to real-time timeframes.

2) Kafka works well with various types of order tracking, maintaining the status of various entities such as new orders, processing a purchase, a shipment & delivery.

3) building on order tracking, Kafka is often used by ride-share or food delivery services. adding geo-graphical location information & providing live status update info for all drivers, riders & customers.

4) Kafka works well for sensor data networks, such as closely monitoring temperature data in a warehouse or various information from the safety systems in a car (or even a roller coaster).

5) Use for cyber security purposes, like spam filtering, patch monitoring & more.

⇒ Components of Kafka:-

① Topics :- A topic is a common message type stored within Kafka.

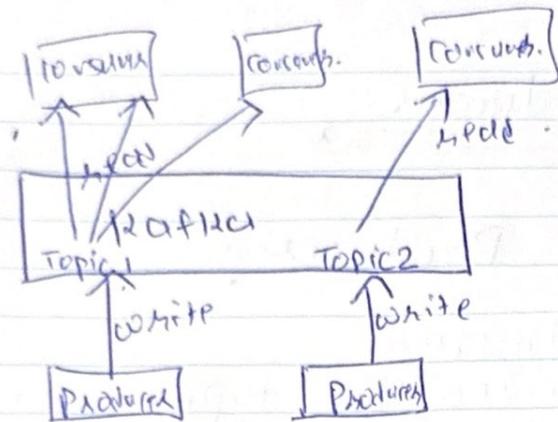
(eg:- like notes on where Kafka stores the events or messages it receives. There can be any # of topics in a Kafka system)

② Kafka Producer :- (writers)

which write events to topics. Producers can write to single or multiple topics.

③ Kafka Consumer :- (readers) which reads information for topics. There can be any # of consumers.

(2)



①

`bin/kafka-topic.sh --help`

↓
will output

↳ give
all the
options
available
on the
command.

Options

↓

:

description

↓

:

①

You want to see a list of topics stored within `kafka`. Then, use command,

`bin/kafka-topic.sh --bootstrap-server localhost:9092 --list.`

↓

give output

cleaned-devices

compromised-systems.

know-exploits

phising-sites.

⇒ Kafka Producers:-

① What is a Producer?

- refers to publisher
- Write messages to topics (Kafka topics)
- Messages sent via producers stored on Kafka for later use. (goal is to retain data on Kafka ASAP.)
- Many producers can write on Kafka at a time & each producer can write on single or multiple topics as required.

② Types of Producers:-

- Many types of producers are available, with varying features depending on needs.
- Most basic producer is Command-line producer.
- This is included in Kafka-software installation & is called Kafka-console-producer.sh
- for command-line tool API libraries
are available for Python, Java & many other languages.
↳ primary producers.
- Besides API libraries & Kafka command-line tools, there is another option called Kafka-connect.
↳ is designed to interact within existing systems, such as relational databases & automatically transfer data between them.

(3)

①

Kafka-console-producer.sh (tool)

→ found in `bin` folder, with other Kafka tools.

→ `bin/kafka-console-producer.sh`

→ Several command-line options available, including some required, others are not. This means they must always be present when using this command.

→ Required options:

①

`--bootstrap-server` ← specifies which server to use.
(Kafka-system)

`--bootstrap-servers localhost:9092` for communication with & defaults to `localhost:9092`.

②

`--topic` ← which requires the topic to write to, such as `topic phising-sites`.

`--topic phising-sites`

→ By default, this opens a connection where you can type messages.

In this case, you can write any messages & then use `[Ctrl + C]` to exit the tool.

→ You can also use the pipe(|) symbol to pass any message data directly into the kafka-console-producer tool.

①

Examples:-

①

→ Interactive examples

bin/kafka-console-producer.sh

--bootstrap-server localhost:9092

--topic testing.



When this command will run we can enter the messages to console & use [Ctrl+C] to exit.

> This is the 1st message.

> II II II 2nd II.

> II II II 3rd II.

(Ctrl+C to exit).

→ Interactive means program waits for our input before continuing.

(4)

(2)

→ Non - Interactive examples.

(1)

Pretty similar to interactive pg,
but using echo command.

(2)

echo "This is the first message!" |
bin/kafka-console-producer.sh

--bootstrap-server localhost:9092
--topic testing.



To write our message & pass it directly
to the kafka-producer.sh, this
would automatically write the message to
kafka without having to use ctrl+c to
exit.



Producers are primary component of
kafka.

Hilary

⇒ Kafka consumer:-

- Reads message from Kafka topics.
- Sometimes called subscribers
- consumer read messages stored within Kafka topics.
- can have many consumers.
- Consumers can read from one or multiple topics.
- Many different kind of consumers.
- most popular is command-line tool.

Kafka-console-consumer.sh

- there are API library available for Python, Java & many other languages.
- There are also component within Kafka connect that allow users to read data from Kafka topics.
- Kafka-console-consumer.sh tool.

↓

normally, found in bin folder, with other Kafka tools.

- Many option for consumer, but most required ones are:
--topic
--bootstrap-servers.
- By default, opens a connection where messages appear (read-time) & maintain position.
- Use Ctrl+C for exit.

(5)

⇒ Optional Arguments :-

① --from-beginning

① Tells Kafka system to send all messages in topic, even ones that may have been already read.

② --max-messages <number>

② Used to read up to a maximum number of messages before exiting.

⇒ Interactive example:-

```
bin/kafka-console-consumer.sh --bootstrap-servers  
localhost:9092 --topic testing --from-beginning.
```

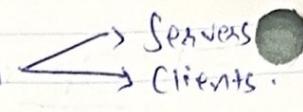
↓
can see all messages
[CTTC] to exit.

⇒ Non-interactive example:-

```
bin/kafka-console-consumer.sh --bootstrap-servers  
localhost:9092 --topic testing --from-beginning  
--max-messages 2
```

↓
return first 2 events & then exit
automatically.

-: Kafka - Architecture :-

⇒ Major component of kafka 

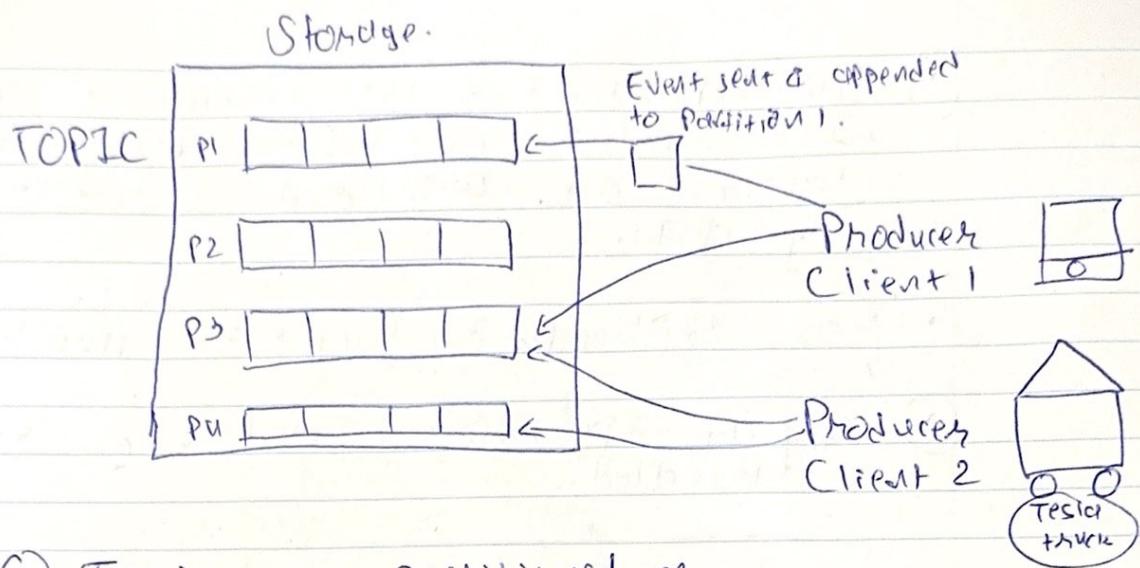
- ① In kafka, server is a cluster of one or more computers handle storing data & manage communications with kafka clients.
- ② It can also handle integration with other systems including databases, log files, & so on.

- ① Kafka clients, which read data via kafka consumers or write data via producers.
- ② Kafka clients can also process data locally and then store that information elsewhere or back to kafka itself.

⇒ Kafka Server :-

- ① Acting as a kafka broker.
- ② Storage, also known as the kafka broker.
 - ↓
 - server also handles data storage.
 - handles communication b/w consumers & producers.
- ③ Data written by producers is stored & organized via topics.

(6)

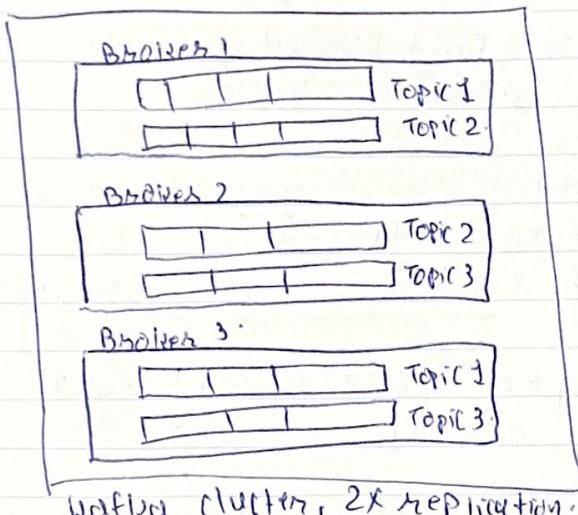


- Topics are partitioned, or stored in separate topics.
- The individual messages are stored in a given partition based on an event id, such as customerid or topicid.
- Each messages are retrieved in some order as written.

\Rightarrow Partitions & replication:-

- Kafka is a fault-tolerant system, meaning that if a system within cluster goes offline, the others can provide the data.
- If one system goes offline, others can provide the requested data.
- For each Kafka cluster, there is a replication factor. \rightarrow [Max # of failures - 1 = replication factor]

- The replication factor minus the cluster number of system failures the cluster can withstand without losing any data.
- Max replication factor = number of servers
- Copying partitions are how replication is handled.



consider an example cluster with 3 brokers.

3 topics defined & set with 2x replication.

2 copy of each partition within the cluster.

$2 \times \text{replication} = \text{can lose 1 system. If before data failure.}$

if Broker 2 fails

\downarrow
replication factor 2x

$2 - 1 = \text{can handle 1 failed system.}$

We are ok, because each topic is exists even Broker 2 fails within cluster.

if Broker 1, Broker 2 fails

\downarrow
More failure than supported (1 system).

Topic 1 & 3
are available so
not a complete
failure
Topic 2 is no longer available in system.

(7)



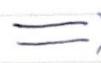
=>

I have 1replica cluster-brkersh.

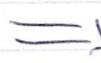
find o # of
brkersh in the
cluster.



Creating & Managing Hadoop Clusters



=> It can possible by Apache Zookeeper.



=> What is Zookeeper?



• Zookeeper is a framework to manage information & provide services necessary for running distributed systems.

• Zookeeper is primarily used by developers to create distributed applications, but users do interact with it to manage those applications.

• e.g. Of applications using Zookeeper include
① HDFS → distributed processing framework.
② Hadoop → distributed processing framework.
③ Neo4j → Graph db.



=> What Zookeeper actually does?



• It provides various services necessary to run distributed applications. This includes handling any configuration information & the naming of systems to prevent conflicts.

- It also provides the ability to synchronize across systems, such as determining what systems are available, when they should start, how services can reach them, & so forth.
- Zookeeper also provides any other basic service that would be required for a group of systems to communicate.

⇒ Zookeeper is designed as a framework to prevent individual distributed applications from implementing custom versions of services.

for eg:- how power plug or water hose nozzle implements common standards vs each having its own.

This allow for easier configuration, implementation & interaction.

⇒ So, Kafka use Zookeeper for cluster management.

- Newer cluster management tool available called **KRaft**.

- Kafka uses two primary configuration files for server / cluster. (two files)

config/zookeeper.properties.

config/server.properties

(8)

zookeeper.properties

The directory where the
snapshot is stored.
dataDir = /tmp/zookeeper

The port at which the
clients will connect
clientPort = 2181

handles information
needed for a
basic zookeepers
settings setup
including where to
store zookeeper
data & what
network port to
run on.

Portion of zookeeper.properties.

⇒ The primary Kafka server configuration file is the config/server.properties file.

① This contains considerably more information than the zookeeper.properties file & defines information specific to Kafka installation.

② Includes info regarding Kafka brokers, details, any network configuration, where to store events, & basic topic configuration including any replication details.

...
A comma separated list of directories
under which to store log files.
logs.dirs = /tmp/kafka-~~logs~~logs.

The default number of log partitions per topic.
num.~~partitions~~ partitions = 1.

⇒ Kafka servers are started in two steps.

* Starting a Kafka cluster:-

→ 1) Using command

bin/zookeeper-server-start.sh config/zookeeper.properties

This starts up the basic Zookeeper server & the configuration details found in the zookeeper.properties file.

There's an extensive amount of output that will vary based on your system.

→ 2) bin/kafka-server-start.sh config/server.properties

This will start Kafka services as defined in the server.properties file.

→ To stop a Kafka cluster, we do the reverse using kafka-server-stop.sh & zookeeper-server-stop.sh

NOTE:- because Kafka must cleanly shutdown any open connections, & because it uses Zookeeper services, the shutdown order is reversed.

(a).

① Kafka topics:-

1	2	3	4	5
---	---	---	---	---

E-commerce orders

- Introduction to Apache Kafka:-

⇒ Create & manage Kafka topics.

① Kafka topics are a logical grouping of events, where each event or message is referring to the same type of thing, such as orders, downloads, virus detections, etc.

② Topic is similar to a table in a relational database, where each event has a similar structure and meaning.

③ Topic sometimes refer to as "Event log." & messages written to the log are immutable. (means messages in a topic can be read or created but can't modified.).

④ Messages can be removed based on age, but otherwise, it is available to be read at any point in the future, given enough storage.

⇒ multiple ways to create a topics in Kafka, but we'll primarily use the bin/kafka-topics.sh script.

Topic
new topic
on topic
search.

bin/kafka-topics.sh --bootstrap-server localhost
--topic topicname --create.

NOTE:-

Still we need to add the
--bootstrap-server option.

e.g:-

```
$ bin/ kafka-topics.sh --bootstrap-server  
localhost:9092 --topic orders --create
```

System will return text "Created
topic orders."

⇒ Other Variations:-



Optional arguments for

```
bin/kafka-topics.sh  
--create
```

```
--replication-factor <x>
```



Define the replication factor of the topic.

Now many copies of a topic
exist within a Kafka
cluster.

```
--partitions <x>
```



Specifies the numbers of partitions to
split the topic into.

→ This is a performance optimization
but allows Kafka to handle higher
throughput.

(10)

for eg:-

```
$ bin/kafka-topics.sh --bootstrap-server localhost:9092 --topic orders --create --replication-factor 3 --partitions 3
```

↓
Created topic orders.

⇒ **--describe** → gives us details about the configurations of the kafka topics.
Argument in ↓ kafka-topics.sh

```
bin/kafka-topics.sh --bootstrap-server <server>  
--topic <topicname> --describe
```

for eg:-

```
$ bin/kafka-topics.sh --bootstrap-server localhost:  
9092 --topic orders --describe.
```

↓
Output:

```
Topic: orders-new TopicId: <topicid>  
PartitionCount: 3 ReplicationFactor: 1
```

Config's:

Topic: orders-new	Partition: 0	Leader: 0	Replicas: 2
Topic: " "	" : 1	" "	" "
Topic: " "	" : 2	" "	" "

⇒ Removing topics :-

```
$ bin/kafka-topics.sh --bootstrap-server  
<server> --topic <topicname> --delete.
```



- ① will remove the topic from a leaf kafka server (cluster).
- ② Removing topic also removes all messages in the topic.

for eg:-

```
$ bin/kafka-topics.sh --bootstrap-server  
localhost:9092 --topic Orders --delete
```

↓ simply returns to prompt.



Kafka - troubleshooting :-

→ One common issue you'll run into is that Kafka tries to be useful to users in its default configuration. Unfortunately, it sometimes causes problems instead.

Sounds confusing
but it
may lead to
significant
issue. →

One such issue is that, by default, a topic doesn't need to exist before a producer writes to it.

Consider what happens if producer misspelled a topic name, such as orders instead of orders.

(12)

→ The events would be written to the wrong topic, & consumers would not see these events.

→ list of topics

```
$ bin/kafka-topics.sh --bootstrap-server localhost:9092 --list.
```

Orders.

```
$ echo "Test message" | bin/kafka-console-producer --bootstrap-server localhost:9092 --topic test
```

Orders ← mispelled here.

```
$ bin/kafka-topics.sh --bootstrap-server localhost:9092 --list.
```

↓
Orders ↴ two topics.
Orders.

⇒ Connectivity Problems:-

- Kafka is a networked service, even on a single-node installation.
- Any networking issue can cause problems with Kafka communications.
- Best way to identify these problems is to look at the output of commands to determine the potential issue.

Let's consider trying to get
a list of topics from an installed
using the `kafka-topics.sh` script.

• `$ bin/kafka-topics.sh --bootstrap-server
localhost:9092 --list`

But error like (Output) →

WARN [AdminClient clientId=adminclient-
Connection to node-1(1) could not be
established.
Node may not be available.

connection
not may up
established.

⇒ To mitigate it, the first is to check
that the service is running using
something like ps & check if the
kafka process is present.

tlp.

(ps aux | grep kafka, netstat -tlnp | grep 9092)

listed for processes running on a system.

provide list of existing connection.

grep is basic search text tool.

⇒ We can also check for any firewall
issues & verify that we have the correct
IP port combo for our kafka installation.

• Verify correct port/IP.

⇒ Another Common Problem:-

① Consumers:-

① Use `--from-beginning` if you need older messages.

Otherwise, only new messages in the topic, since the consumer started, will be received.

① Remember `--max-messages` to only read a specific message quantity, otherwise the consumer script will never exit.

① All tools:-

① Remember to include `--bootstrap-server` options.

This is initial Kafka server that all programs communicate with to run Kafka commands.

① Most error messages are pretty clear & will point you to likely possible solutions.

① Tools also have a `-help` option to get further details when using these scripts.