# Week 6

## Streaming - Kafka
### (Messaging System)
### [ Trasportation System ]

. Integration between <u>producer</u> & <u>consumer</u>

- Kafka <u>Handel</u> :
  - protocol - How data transported (TCP, HTTP, REST..)
  - Data format (Binary, CSV, JSON)
  - Data schema How data is shaped and might change.
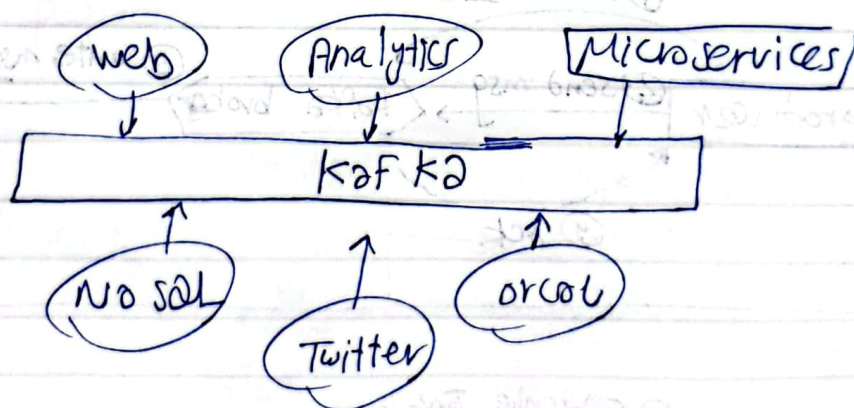
- latency less <u>10ms</u> [ Real Time ]

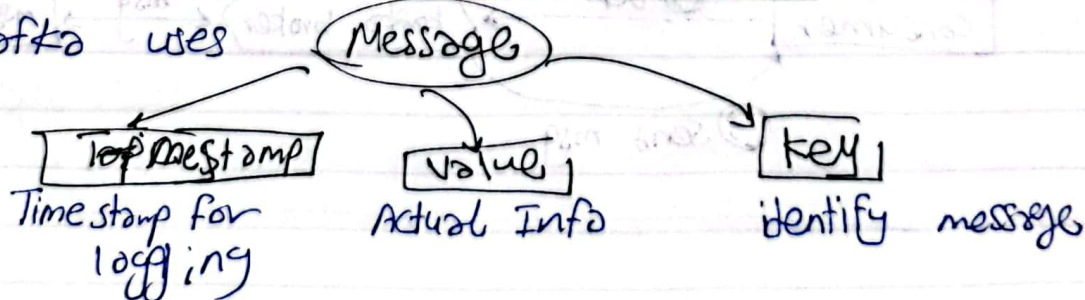- Consumers: Those that consume data web pages, micro services.

- produces: who supply data to consumes

- Connecting consumes direct to producers can lead to amorphous and (hard) to maintain architecture.
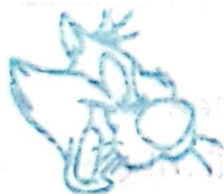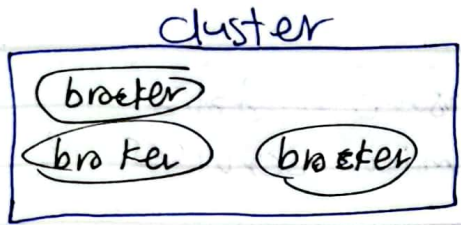
. kafka is intermediary component

```
  ( web )      ( Analytics )    | Microservices |
     |              |                   |
     v              v                   v
  +---------------------------------------------+
  |                 kafka                       |
  +---------------------------------------------+
     ^              ^                   ^
     |              |                   |
  ( No sql )     ( Twitter )        ( orcol )
```

. kafka uses ( Message )

| Topmestamp | value | key |

Time stamp for logging     Actual Info     identify message

Topic

Topic which set a group of msgs.

cluster

```
┌─────────────────────────┐
│  ⬭ brocker              │
│  ⬭ bro ker   ⬭ broeker  │
└─────────────────────────┘
```

Broker : pyshicol or virtual Machine
برسمار    where kofto work.

---

Logs

kofto store and assige each msg ID
then store it in logs

① declor Topic abc

② send msg    ③ write msg    Logs of
                              TopiC 'abc'
producer → kofto broker → | msg1 , msg2 |

④ ack

---

① subscribe Toplc Abc

④ ack                     ② check    logs
                            msg
Consumer → kofto broker ←───→ | msg1 , msg2 |

③ send msg

---

Consumer - offset

special topic keep track of each consumer
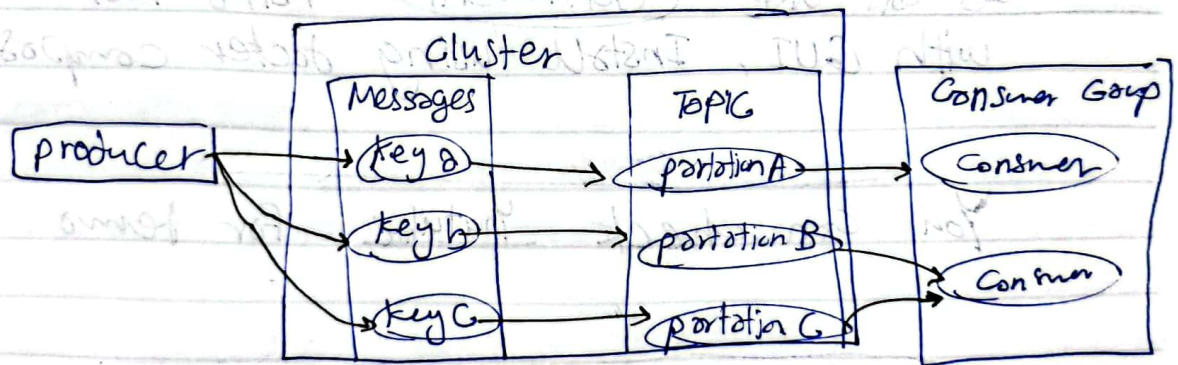and topic.
keep track of what consumer do.

e.X   If consumers dies or what ever happen
      and get back Consumer_offset will continue
      from post statue.

## Consumer Group

- a set of consumers with an ID for Group & ID for each Consumer But that help with:
  - Duplicated  - redundant msgs.

---

## Partations

Logs contains _msgs ID_, we call _Topic logs_ as _partations_.



## Replication

partations are replicated across multiple brokers in kofka. fault tolerance.

## Kafka Configurations

retention (ms): Time specific topic will be available before deleting.

replication : Replication factor, Number of times partitions will be replicated.

consumer configuration
offset: msgs ID which have been read
by consumers.

producer configuration
acks :
0 ⟶ fire forget
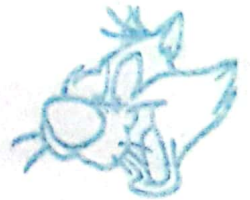1 ⟶ wait for leader

## Installing kafka:

we will use (confluent) kafka tool
with GUI, Install using docker compose.

You can check Youtube for demo.

## Avro & schema Registry:

msgs can be any thing plain txt --- binary
this give Flexibility to kafka.

But we must take care of compatibility

PHP producer ──PHP object──> kafka ──WTF ??──> CSV Consumer

## Avro

Data serialization System, unlike JSON
Avro stores schema separated from record.

- Smaller record size.
- Schema evolution : Change schema over time.
- Automatic validation

## Schema registry

Contains info about schema, we validate
from schema reg.

## Schema Evolution

for ex Backward check against
last version, upgrade consumers first.

forward upgrade producers first.

## Kafka Streams

- Kafka streams used when we work with
  input & output stored in Kafka Clusters.
  stream & produce data From Kafka to Kafka

- streams process events with latency of
  milliseconds, deal with msgs ASA they arrive

. streans works with DSL
( Domain specific Language ) , which simplify
creation of stream.

. kafka stream is powerful and simple
Spark and flink more powerful , But (hord.)

> kafka streams work with kafka produce
and consumer only.

_____

~~state~~ . Stream vs state
- ~~stream changelog~~ - individual msg
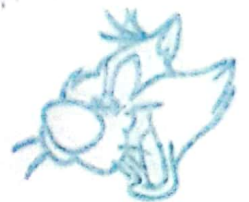- ~~kTable~~ - read sequential

. stream ( k stheams)
- individual msg read sequentially.

. state ( k-state)
- stream changelog
- view of stream at certin time.

_____

Topology
How input data transfoned to output.

stream processor
Transform data.

- Native language for kafka streams is _scala_. We will use _faust python libirary_ in `JVM`

## Windowing
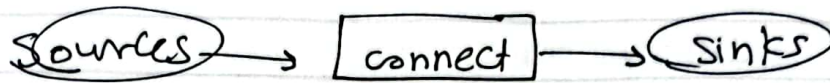Time reference in which series of events happen.

- Time based window
  - window has predetermined size "s"

- Session - based
  - window based on key to start and end not Absolute time.

- Kafka support _Threads_, And control number of threads.

- _Global K-Table_
  - Brodcast variable for
    - more convenient
    - efective Joints
    - But incresse local storage

- Interactive queries
  Allows External apps to query your Stream apps directly.

- **Kafka connect**
  - stream data between **external** application

Sources → [ connect ] → Sinks

## KSQL

KSQL offers consumers such as Data science a tool for analyzing Kafka streams.