## → Kafka
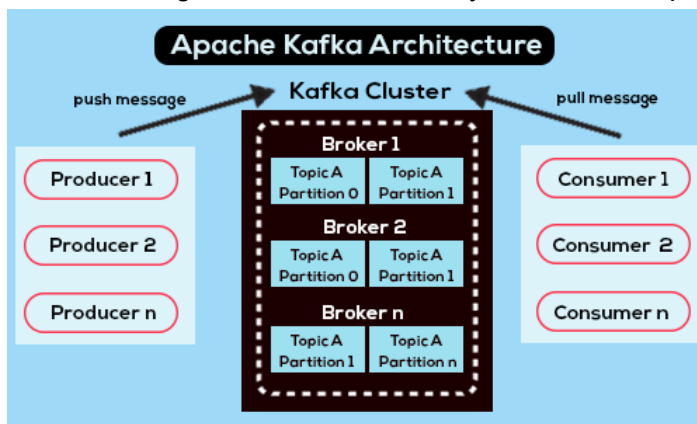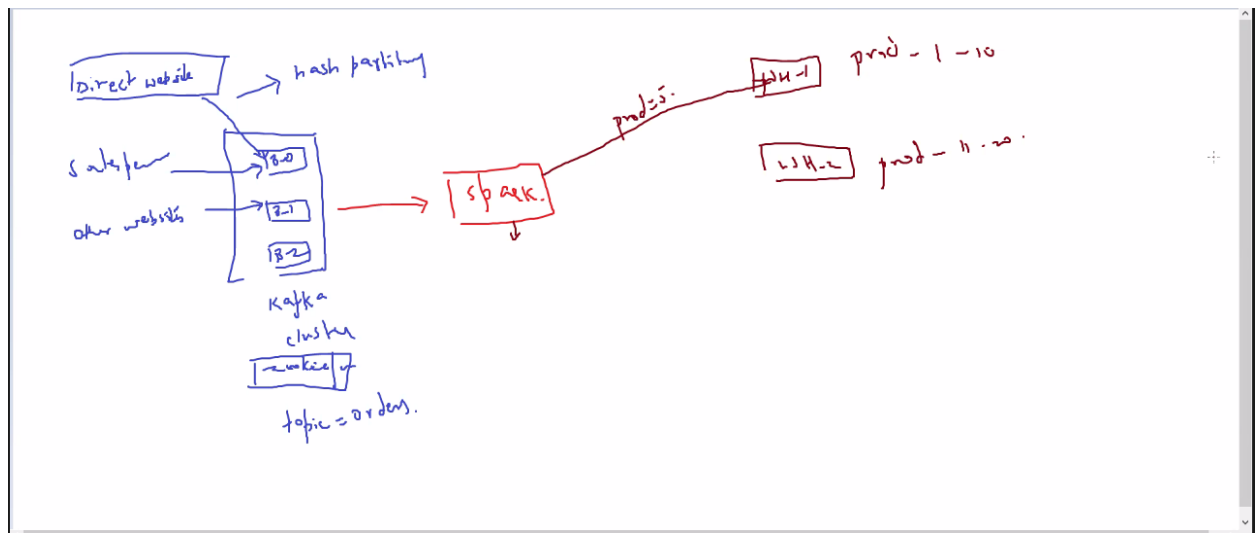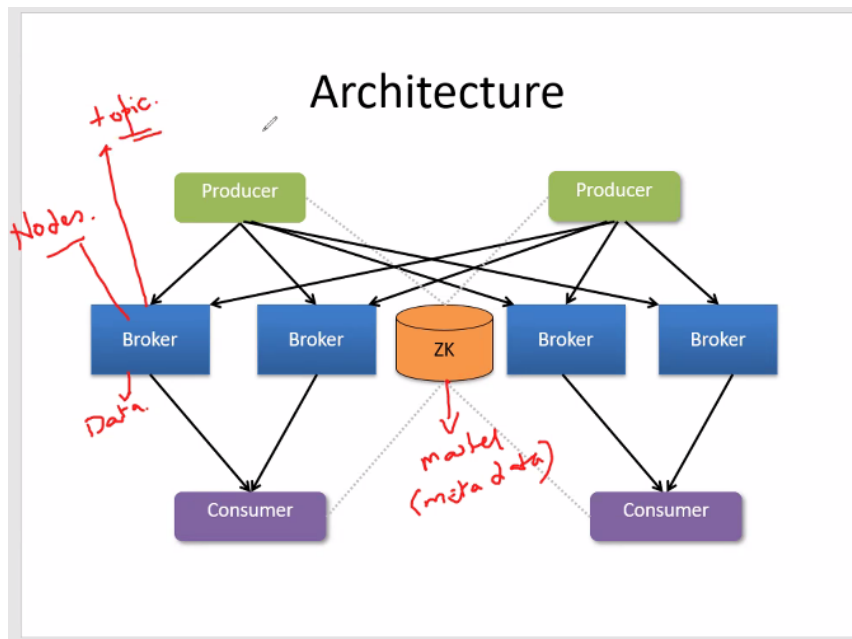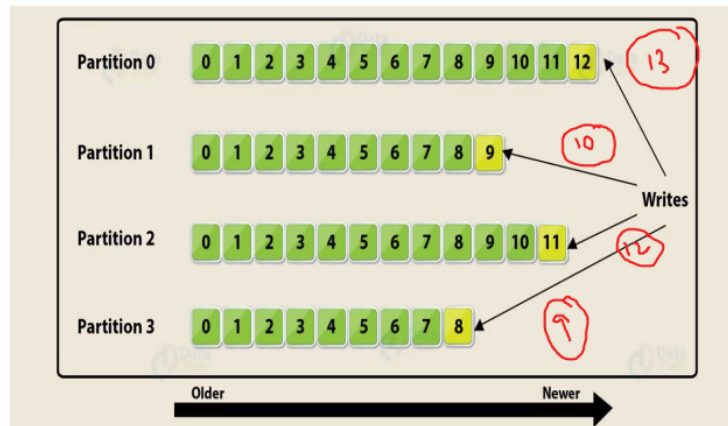


a. <u>Kafka cluster</u>, called cluster because it has multiple brokers
b. Using kafka, we collect real-time data from different publishers
c. distributed publish-subscribe messaging system
d. <u>Publisher</u> is sender and <u>Consumer/Subscriber</u> is receiver
e. <u>Topic</u>:
    i.    an object where data is stored representing similar types of data
    ii.    Can have many K, V paired messages
f. Once order is collected from publishers, then it is sent to spark app which is a subscriber, and then spark app sends to different data warehouses
g. One <u>Partition</u> for each Broker, but each partition can have different number of messages
h. A kafka server is called a <u>Broker</u>, a bridge between producers & consumers
i. Kafka does not process data, it's just a storage layer
j. <u>Zookeeper</u> is used to store information about Kafka Cluster & details of consumer clients, manages brokers, mandatory to run zookeeper server



k.
l.

# Architecture

topic.

Nodes.

Producer    Producer

Broker    Broker    ZK    Broker    Broker

Data.

master
(metadata)

Consumer    Consumer

---

Direct website    → hash partition

salesforce

other websites

B-0
B-1
B-2

Kafka
cluster

zookeeper

topic = orders.

spark.

process.

WH-1    prod - 1 - 10

WH-2    prod - 11 - 20.

## Partitions for one Topic



Handwritten annotations on the slide:
- Circled numbers: 13 (pointing to Partition 0), 10 (Partition 1), 12 (Partition 2), 9 (Partition 3)
- = 44 messages in my topic!

## Consumer



Handwritten annotations:
- V.J.D - Ric
- Amazon;
- millions;
- vendor; → despatch.
- s*b.
- sup
- sup.
- Common.

# → Oozie
    a. Java Web App to schedule hadoop jobs
    b. Combines multiple jobs sequentially into one logical unit of work
    c. Three types of oozie jobs
        i. Oozie Workflow
            1. Directed Acyclic Graphs, specifying a sequence of actions to execute
        ii. Oozie Coordinator
            1. Jobs are recurrent Oozie Workflow jobs that are triggered by time and data availability
        iii. Oozie Bundle

1. Provides a way to package multiple coordinator and workflow jobs and to manage the lifecycle of those jobs