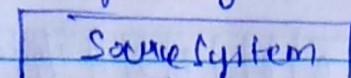


## Apache Kafka

Let's take a scenario we need, to exchange a data between source to target system.



Example for data streams:-

→ Website events.

→ Pricing data.

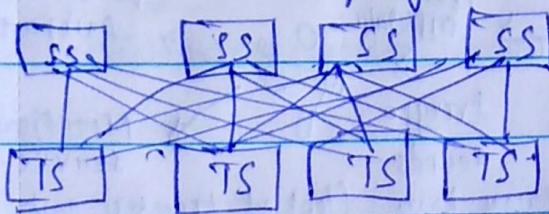
→ Financial transaction

→ User interaction.

→ CFC.

Once we have data in kafka, we can put it in any system we like.

After a while things get complicated.



→ So if we have 4 SS & 6 TS we need to write 24 integrations.

→ Each integration comes with difficulties around:

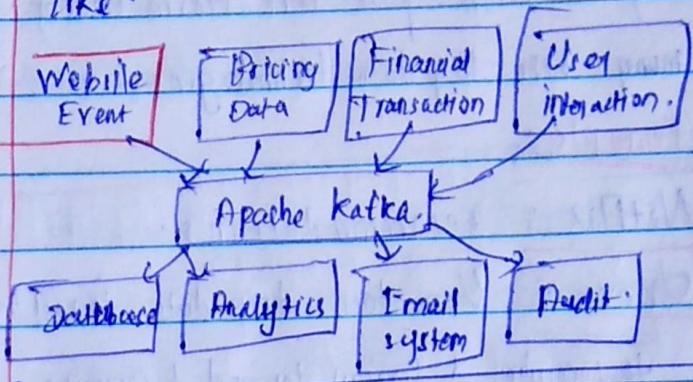
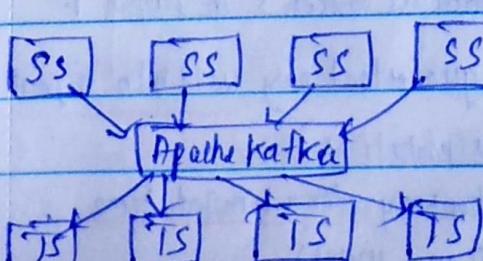
\* Protocol & how data is transported.  
(TCP, HTTP, REST, FTP...).

\* Data format - how data is parsed  
(Binary, CSV, JSON, Avro...).

\* Data schema & evolution - how data is shaped & may change.

→ Each source system will have an increased load from the connection.

Doupling of data streams & systems :-



### Why Apache Kafka

→ Created by LinkedIn, Open sourced project mainly maintained by Confluent.

→ Distributed, resilient architecture, fault tolerant.

→ Horizontal scalability.

\* can scale to 100 brokers.

\* can scale to millions of messages per second.

→ High performance (latency of less than 10ms) - real time.

→ Used by 2000+ firms, 35% of the Fortune 500.

## Apache Kafka Usecases

- Messaging System
- Activity Tracking
- Gather metrics from many different locations.
- Application log gathering
- Stream Processing (with the kafka Stream API (or Spark for example)).
- Decoupling of system dependencies.
- Integration with spark flink storm Hadoop and many other Big data technologies.

### For Examples:-

- Netflix: recommendations.
- Cab: User taxi trip data in real time to compute & forecast demand & compute surge pricing in real time.
- LinkedIn: prevent spam, collect user interactions to make better connection.
- Kafka is only used as a transportation mechanism!

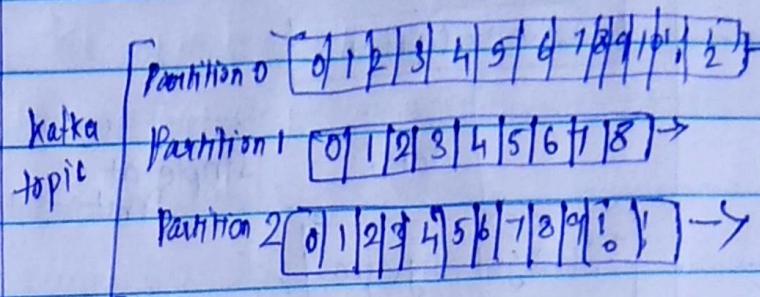
## Kafka Theory :

Topics, partitions & offsets:

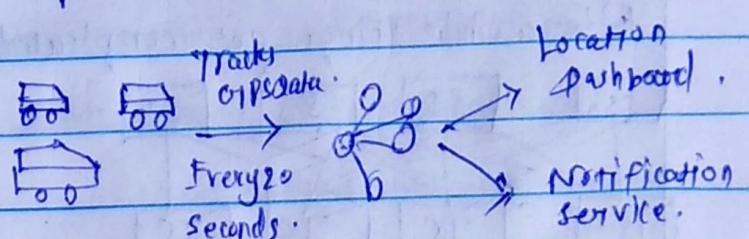
- Topics: a particular stream of data.
- Similar to a table in a database (without all the constraints)
  - You can have as many topics as you want
  - A topic is identified by its name.

Topics are split in partitions.

- Each partition is ordered.
- Each message written in a partition gets an incremental id called offset.



Topic Example : truck-gps.



→ Say you have a fleet of trucks, each truck reports its GPS position to Kafka.

→ You can have a topic truck-gps that contains the position of all trucks.

→ Each truck will send a message to Kafka every 20 seconds, each message will contain the truck ID and the truck position (latitude & longitude).

→ We choose to create that topic with 10 partitions.

### Offset:

→ Offset only have a meaning for a specific partitions.

Eg: offset 3 in partition 0 does not represent the same data as offset 3 in partition 1.

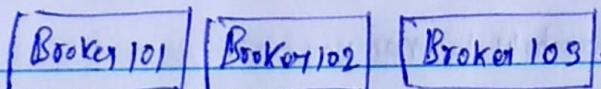
→ Order is guaranteed only within a partition (not across partitions).

→ Data is kept only for a limited time (Default is one week).

- ⇒ Once the data is written to a partition, it can't be changed (Immutability)
- ⇒ Data is assigned randomly to a partition unless a key is provided (more)

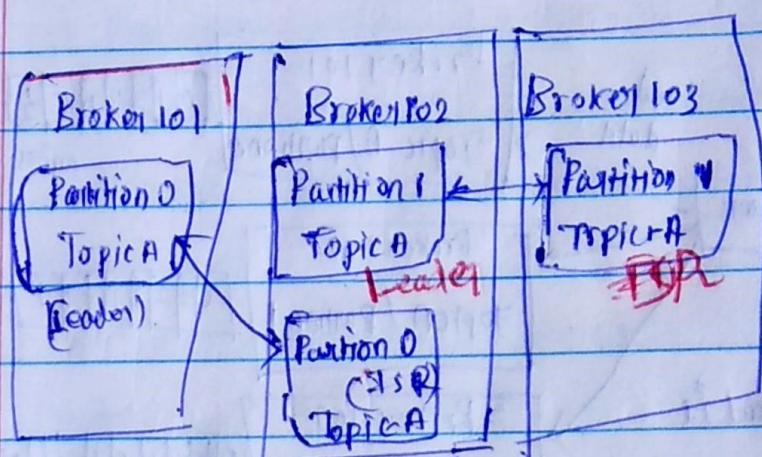
### Brokers :

- ⇒ It a holder for Topics
- ⇒ A kafka cluster is composed of multiple brokers (servers)
- ⇒ Each broker is identified with its ID (integer).
- ⇒ Each broker contains certain topic partitions.
- ⇒ After connecting to any broker (called a bootstrap broker) you will be connected to the entire cluster.
- ⇒ Big cluster have over 100 brokers.



### Topic replication factor

- ⇒ Kafka is a distributed system.
- In distributed system, we have replication.
- ⇒ If things go down application still works.
- ⇒ Topics should have a replication factor (usually 2 & 3)
- ⇒ This way if a broker is down, another broker can serve the data.
- ⇒ Example : Topic A with 2 partitions & replication factor of 2.



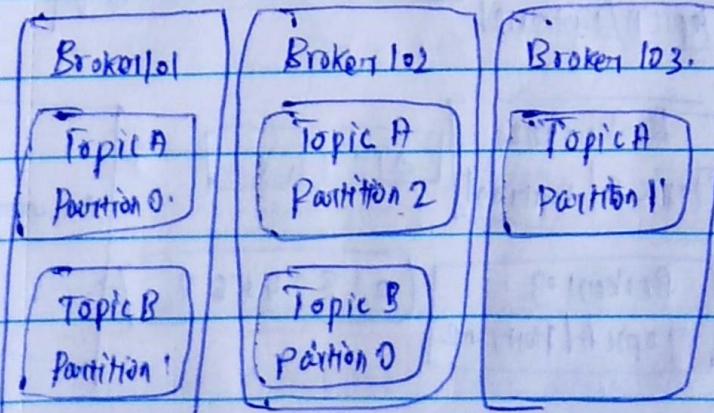
Example : we lost Broker 102.

Result : Broker 101 + 103 can still serve the data.

### Brokers & Topics :

- ⇒ Example of Topic-A with 3 partitions
- ⇒ Example of Topic-B with 2 partitions.

A kafka cluster with 3 Brokers.



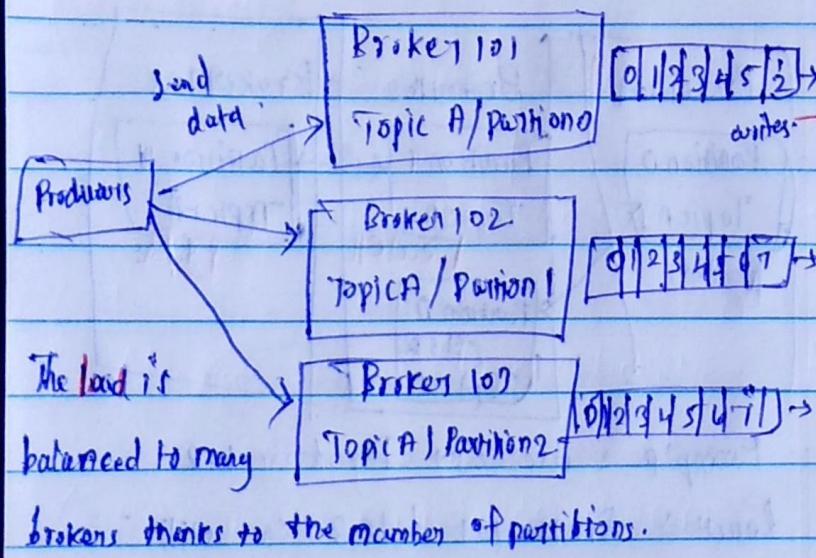
- ⇒ Data is distributed and Broker 103 does not have any Topic B data.

### Concept of Leader for a Partition :

- ⇒ At any time only one broker can be a leader for a given partition.
- ⇒ Only that leader can receive and serve data for a partition.
- ⇒ The other brokers will synchronize the data.
- ⇒ Therefore each partition has one leader and multiple ISR (in-sync replica).

## Producers:-

- How do we get data in kafka? That is the job of producers.
- Producers does the load balancing in round robin fashion.
- Producers write data to topics (which is made of partitions)
- Producers automatically know to which broker and partition to write to.
- In case of Broker failure, Producers will automatically reover.



acks=0: Producers won't wait for acknowledgement (possible data loss).

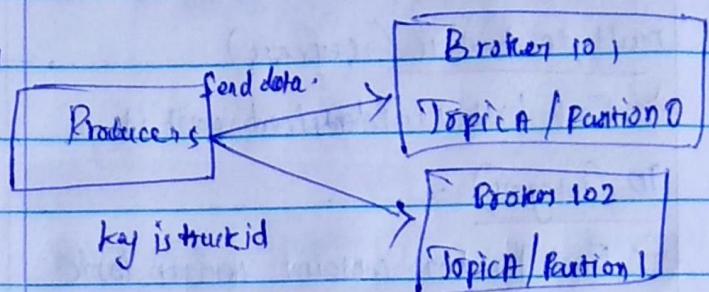
acks=1: Producers will wait for leader acknowledgement (limited data loss).

acks=all: Leader + replica acknowledgement (No data loss).

## Messages key:

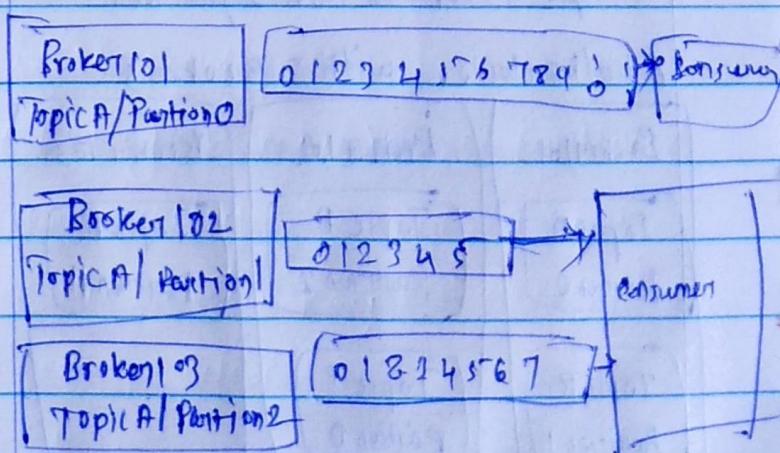
Producers can choose to send a key with message.. (string, number etc..)

- If key = null, data is sent round robin (Broker 101 then 102 then 103) -
- If key is sent, then all message for that key will go to same partition.
- A key is basically sent if you need message ordering for a specific field.  
Ex: track\_id



## Consumers:

- Consumers read data from a topic (identified by name)
- Consumer knows which broker to read from.
- In case of broker failures, consumers know how to recover.
- Data is read in order within each partition.
- Across partition offset order not matters.



## Consumer Groups:

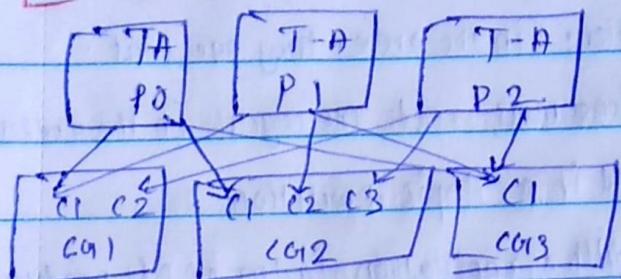
- Consumers read data in consumer group.
- Each consumer within a group read from exclusive partitions.
- If you have more consumers than partitions, some consumers will be inactive.

CG - Consumer group.

C - consumer.

T - Topic

P - Partition.

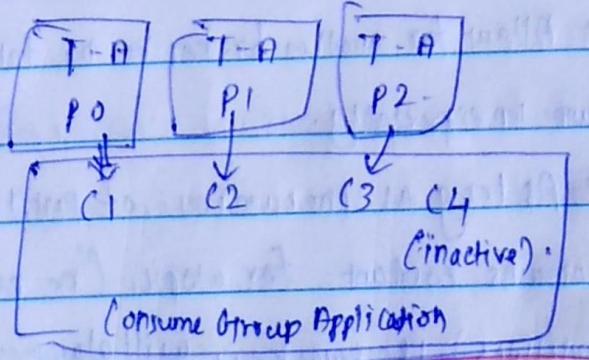


Consumers will automatically choose a Group Coordinator and a ConsumerCoordinator.

to assign a consumer to a partition.

What if too many consumers?

If you have more consumers than partitions, some consumer will be inactive.



## Consumer Offsets:

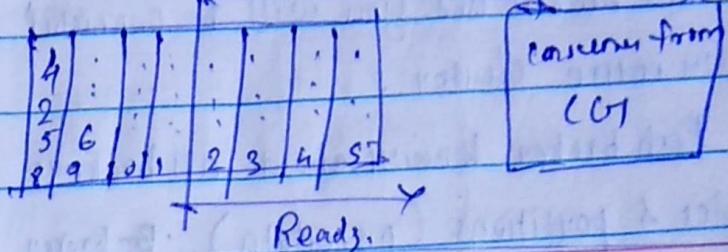
- Kafka stores the offset at which a consumer group has been reading.
- The offset committed live in a kafka workflow using kafka streams API.

topic named \_consumer\_offsets.

→ When a consumer in a group has processed data received from kafka, it should be committing the offset.

→ If a consumer dies, it will be able to read back from where it left off, thanks to the committed consumer offsets.

↓ committed offset.



## Delivery semantics for consumer:-

→ Committing offset implies it.

→ Consumer choose when to commit offsets.

→ There are 3 delivery semantics.

At most once:

→ offsets are committed as soon as the message is received.

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

At least once:

→ committed after message is processed.

→ If process goes wrong message will be read again.

→ This can result in duplicate processing of message. Make sure your processing is idempotent. (processing again the message won't impact your system).

Exactly once:

→ Can be achieved for kafka-to-kafka workflows using kafka streams API.

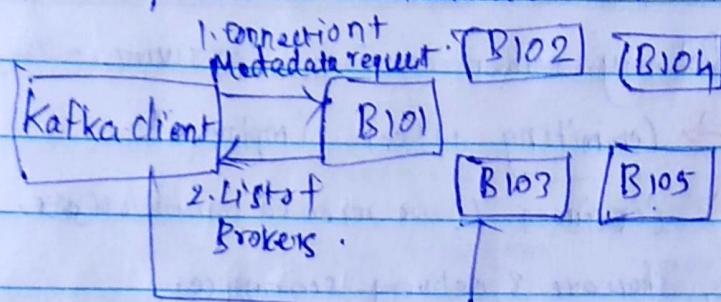
→ For Kafka → External system workflow use an idempotent consumer.

## Kafka Broker Discovery :-

Every kafka broker is also called a "bootstrap server".

That means that you <sup>only</sup> need to connect to one broker and you will be connected to the entire cluster.

Each broker knows about all brokers topics & partitions (metadata) • Br-Brokers



3. Can connect to the needed brokers.

## Zookeeper:

\* Manages brokers (keeps a list of them)

\* Helps in performing leader election for partitions.

\* Sends notification to Kafka in case of changes (e.g.: new topic, broker dies, broker comes up, delete topics, etc...)

\* Kafka can't work without Zookeeper.

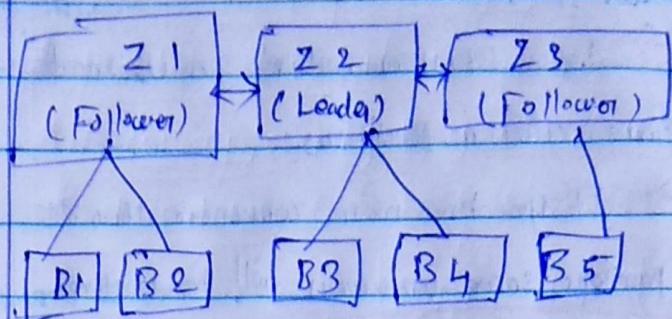
\* By Design, operates with an odd number of servers (3, 5, 7) (Non-leader)

\* Has a leader, the rest of the servers are followers (handle reads)

\* Zookeeper does not store consumer offsets with kafka Yugioh

Z - Zookeeper Server.

B - Kafka Broker



## Kafka Guarantees :-

→ Messages are appended to a topic partition in the order they are sent.

→ Consumers read messages in the order stored in a topic partition.

→ With a replication factor of N, producers and consumers can tolerate up to N-1 broker being down.

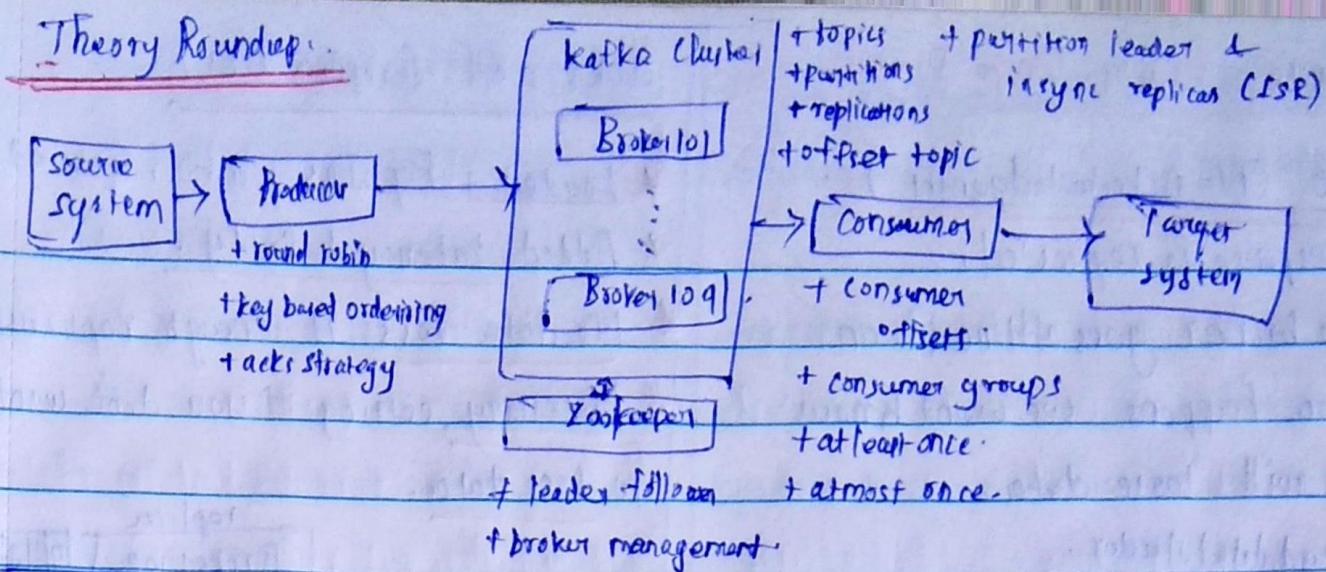
→ This is why a replication factor of 3 is good idea.

→ Allows for one broker to be taken down for maintenance.

→ Allow for another broker to be taken down unexpectedly.

→ As long as the number of partitions remains constant for a topic, (no new partitions) the same key will always go to the same partition.

## Theory Roundup



Command for starting zookeeper:

## ZooKeeper - server-start config //

11

# For Starting Zookeeper

## Starting kaffee

Kafka server - start config / server.properties

## Creating a topic

kafka topics -zookeeper 0.0.0.0:2181 -topic first topic

-- Create --partition 3 --replication factor 1

Please refer the word document for more commands.

Providing a key in producer.

Always goes to the same partition

Consumer is consumer group

automatically rebalanced to different partitions in topic.

## Producers Ack Dags Dive:

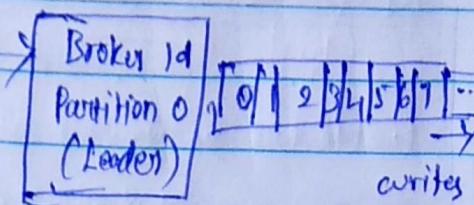
ack=0 (No Acknowledgement)

\* No response is requested.

\* If the broker goes offline or an exception happens we won't know and we will lose data.

Send data to leader.

Producer



\* Useful for data where it's okay to potentially lose messages.

Metric Collection

Log collection.

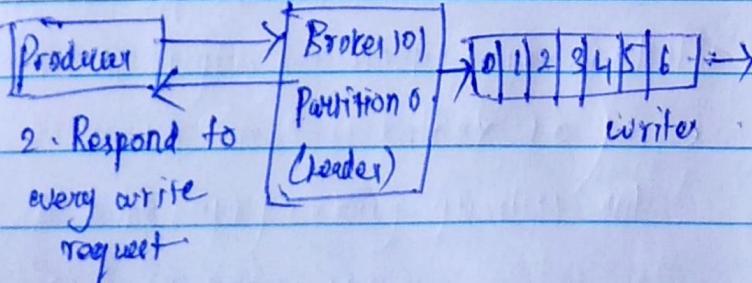
\* Good in performance.

ack=1 (Leader acks)

\* Leader response is requested, but replication is not guaranteed (happens in the background).

\* If an ack is not received, the producer may vary.

1. Send data to leader.



\* If the broker goes offline but replicas haven't replicated the data yet. we have a data loss.

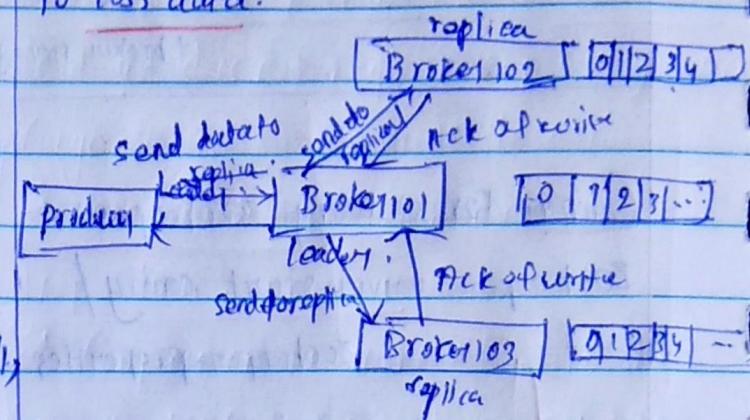
acks=all (replica acks)

\* Leader + Replicas ack requested.

\* Added latency & Safety.

\* No data loss if enough replicas.

\* Necessary setting if you don't want to lose data.



\* Acks-all must be used in conjunction with min.insync.replicas.

\* Min.insync.replicas can be set at the broker (or topic level) (override).

\* min.insync.replicas=2 implies that at least 2 brokers that are ISR (including leader) must respond that they have the data.

\* That means if you use replication.factor=3 min.insync=2, acks=all you can only tolerate one broker going down otherwise producer will require an expn in fct.

## Producer Retries

\* In case of transient failures, developers are expected to handle exceptions, otherwise the data will be lost.

\* Example of transient failures:

→ NotEnoughReplicasException.

\* There is a retries setting

→ default to 0 for Kafka <= 2.0

→ default to 2147483647 for Kafka >= 2.1

\* The retry.backoff.ms setting is by default 100 ms.

\* In Kafka = 1.0.0 there's better solution with idempotent producers.

### Idempotent Producers:

Problem: - The producer can introduce duplicate messages in kafka due to network errors.

### Producer Timouts:

\* If retries > 0, for example,

retries = 2147483647;

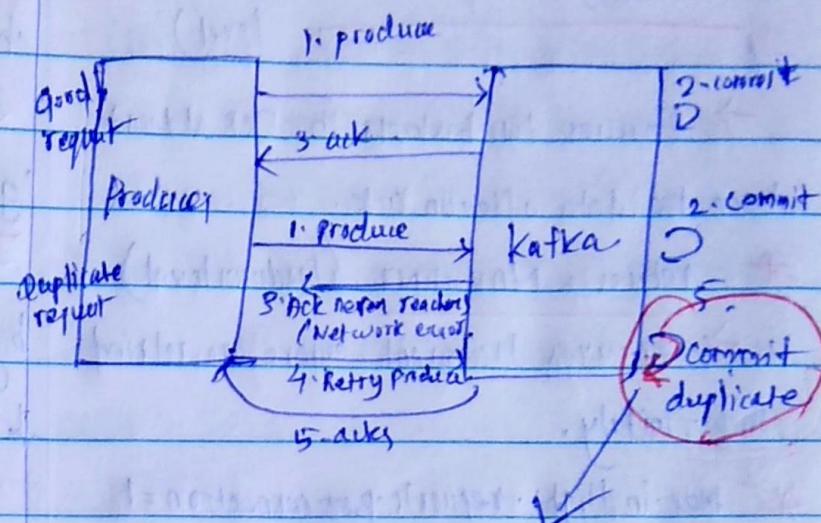
\* the producer won't try the request for ever, it's bounded by a timeout

\* For this, you can set an intuitive

Producer Timeout (KIP-91-kafka2.1).

\* delivery.timeout.ms = 10000 ms = 2 min

\* Records will be failed if they can't be acknowledged in delivery.timeout.ms.



where all in idempotent

producer, won't introduce duplicates on network error.

\* If detect the duplicate, don't commit twice.

\* In Kafka v=0.1.1 we can define idempotent producer.

\* If provider stable & safe pipeline

\* If comes with:

\* retries = Integer.MAX\_VALUE (2^31-1)

\* max.in.flight.requests = 1

\* max.in.flight.requests = 5

\* These settings are applied automatically after your producer has started if you don't set them ~~automatically manually~~.

properties.set("enable.idempotence", "true");

\* See it to 1 if you need to ensure ordering.

## Safe producer summary:-

Kafka >= 0.11.

- \* acks = all (producer level) .  
→ Ensures data is properly replicated before an ack is received.
- \* min.insync.replicas = 2 (broker/topic level)  
→ Ensures two brokers in ISR at least have the data after an ack.
- \* retries = MAX\\_INT (producer level) .  
→ Ensures transient errors are retried indefinitely.
- \* max.in.flight.requests.per.connection = 1 (producer level) .

→ Ensures only one request is tried at any time, preventing message reordering in case of retries .

## Kafka >= 0.11

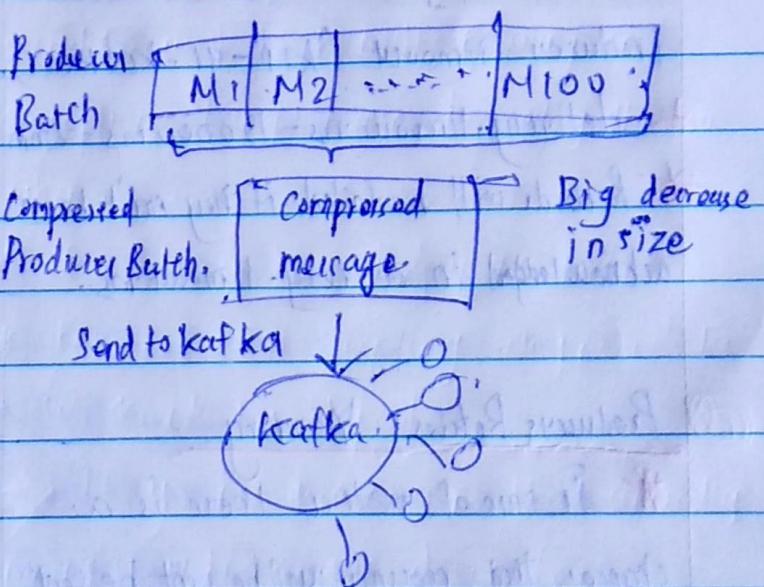
- \* enable.idempotence = true (producer level) + min.insync.replicas = 2 (broker/topic level)
- \* implies (acks = all, retries = MAX\\_INT, max.in.flight.requests.per.connection = 1) if kafka >= 0.11 (or 5 if kafka >= 1.0)
- \* while keeping ordering guarantees and improving performance.

→ Running a safe producer might impact throughput & latency always test for your use case.

## Message compression:-

- \* Producers usually send data that is text based for example with JSON data.
- \* In this case, it is important to apply compression to the producer.
- \* Compression is enabled at the producer level and does not require any configuration change in the Broker or in the consumers.
- \* compression.type can be none (default) 'gzip', 'lz4', 'snappy'.
- \* Compression is more effective the bigger the batch of message being sent to Kafka.

### M - Message



### Advantages:-

- \* Much smaller producer request size (compression ratio upto 4x !)
- \* Faster to transfer data over the network less latency
- \* Better throughput
- \* Better disk utilisation in kafka (stored messages on disk are smaller).

## Disadvantages:

- Producers must commit some CPU cycles to compression.
- Consumers must commit some CPU cycles to decompression.

## Overall:

→ Consider testing snappy (or LZ4) for optimal speed / compression ratios.

## Recommendation:

- Test all algorithms and decide.
  - Use it in production & especially if you have high throughput.
    - Consider tweaking linger.ms & batch.size to have bigger batches, therefore, more compression & higher throughput.
- By introducing some lag (for example  $\text{linger.ms} = 5$ ) we increase the chance of messages being sent together in a batch.
- So at the expense of introducing a small delay, we can increase throughput, compression & efficiency of our producer.
- If a batch is full before the end of the linger.ms period, it will be sent to Kafka right away.

## Producer Batching:

- By default, Kafka tries to send record as good as possible.
- It will have up to 5 requests in flight, meaning up to 5 messages individually sent at the same time.
- After this, if more messages have to be sent while others are in flight, Kafka is smart & will start batching them while they wait to send them all at once.

→ This smart batching allows Kafka to increase throughput while maintaining very low latency.

→ Batches have higher compression ratio, so better efficiency.

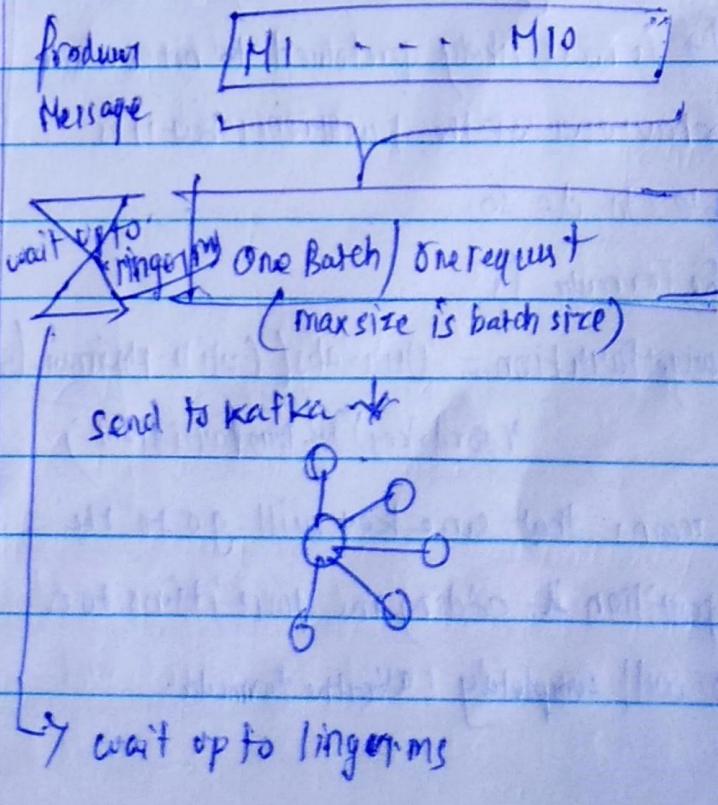
## Linger.ms & Batch.size:

→ Linger.ms: Number of millis a producer is willing to wait before sending a batch or not (Default 0).

→ By introducing some lag (for example  $\text{linger.ms} = 5$ ) we increase the chance of messages being sent together in a batch.

→ So at the expense of introducing a small delay, we can increase throughput, compression & efficiency of our producer.

→ If a batch is full before the end of the linger.ms period, it will be sent to Kafka right away.



## Batch size: (batch-size)

- Maximum number of bytes that will be included in a batch. The default is 16KB.
- Increasing a batch size to something like 32KB or 64KB can help increasing the compression, throughput efficiency of requests.
- Any message that is bigger than the batch size will not be batched.
- A batch is allocated per partition; so make sure that you don't set it to a number that's too high otherwise you'll run out of memory.
- (You can monitor the average batch size metric using Kafka Producer Metrics).

Max-block-ms & Buffer-memory.

- If the producer producer faster than the broker can take, the record will be buffered in memory.
- $\text{buffer\_memory} = 33554432 / 32\text{mb}$  The size of the send buffer.
- That buffer will fill up over time & fill back down when the throughput to the broker increases.
- If that buffer is full (all 32MB), then the `.send()` method will start to block (won't return right away).

max-block-ms = 60000: the time the `.send()` will block until throwing an exception. Exception case basically thrown when,

- \* The producer has filled up its buffer.
- \* The broker is not accepting any new data.
- \* 60 seconds have elapsed.
- If you hit an exception that usually means your brokers are down (or) overloaded as they can't respond to requests.

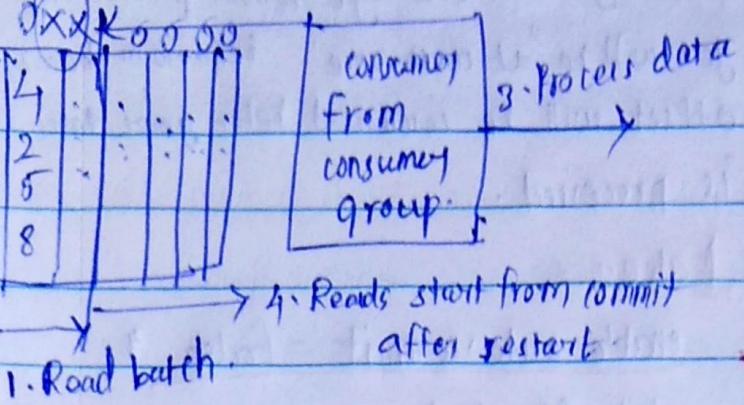
$$\text{targetPartition} = \text{Util} - \text{abs}((\text{Util}) \cdot \text{partition} - \text{recordKey})) / \text{numPartitions}$$

This means that same key will go to the same partition & adding the partitions to a topic will completely alter the formula.

## Delivery Semantics:-

At most once: offsets are committed as soon as the message batch is received. If the processing goes wrong, the message will be lost (it won't be read again).

There two duras are lost:  
committed offset

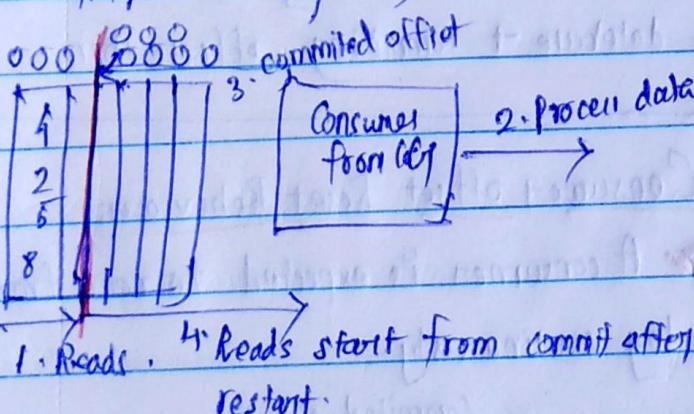


## Bottom Line:-

for most application you should use at least once processing & ensure transformation / Processing one idempotent.

## 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.
- Make sure your processing is idempotent. (Processing again won't impact the system).



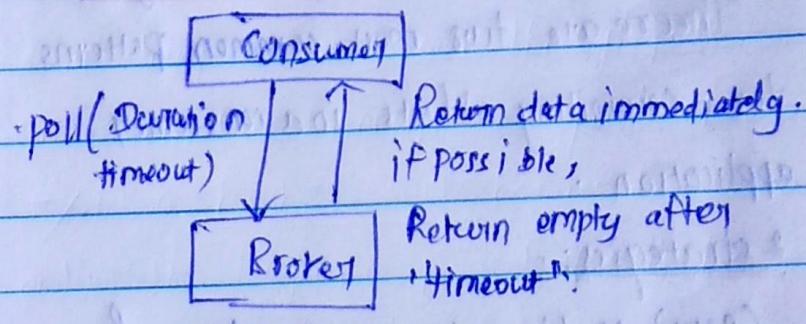
Exactly once: Can be achieved for kafka  $\Rightarrow$  kafka workflows using kafka streams API.

For kafka sink workflow use an idempotent consumer.

## Consumer Poll Behaviour:

$\Rightarrow$  Kafka consumer have a 'poll' model, while many other messaging bus in enterprises have a push model.

$\Rightarrow$  This allow consumers to control where in the log they want to consume, how fast & gives them ability to replay events.



## Consumer Poll Behaviour:

$\Rightarrow$  Fetch min bytes (default 1):  
★ Controls how much data you want to pull at least on each request.

★ Helps improving the throughput & decreasing request number.

★ At the cost of latency.

Max poll records (default 500):

★ Controls how many records to receive per poll request.

$\Rightarrow$  Increase if your message are very small & have a lot of available RAM.

Max partitions fetch bytes (default 1MB):

\* Maximum data returned by the broker per partition.

\* If you read from 100 partitions, you will need a lot of memory (RAM).

Fetch max bytes (default 50MB):

\* Maximum data returned for each fetch request (covers multiple partitions).

\* The consumer performs multiple fetches in parallel.

\* Change these settings only if your consumer maxes out on throughput already.

### Consumer Offset Commit Strategies:

There are two most common patterns for committing offsets in a consumer application.

#### 2 strategies:-

(easy) enable.auto.commit = true & synchronous processing of batches.

(medium) enable.auto.commit = false & manual commit of offsets.

#### Title:-

```
while(true){  
    -- consumer.poll( - )  
    doSomethingSynchronous( - );  
}
```

→ with auto commit, offset will be committed automatically, for you at regular interval (auto.commit.interval.ms = 5000 by default)

every time you call .poll().  
→ If you don't use synchronous processing, you will be "at most once" behaviour bcoz, offset will be committed before your data is processed.

#### Fakes:-

```
enable.auto.commit = false  
while(true){  
    batch = consumer.poll( - );  
    if(isReady(batch)){  
        doSomethingSync(batch);  
        consumer.commitSync();  
    }  
}
```

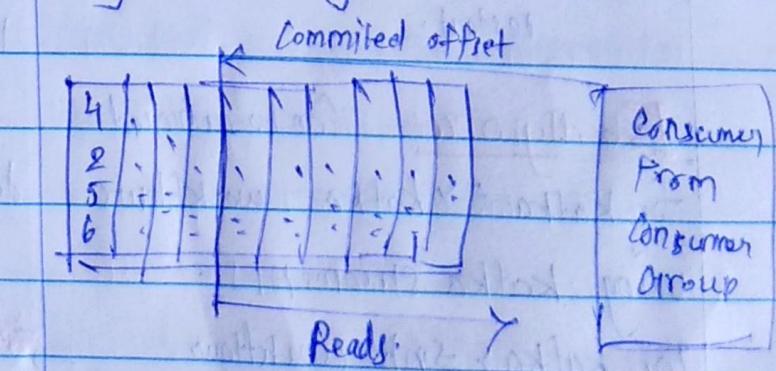
→ you control when you commit offsets.

& what condition for committing them?

→ Example: accumulating records into buffer and then flushing the buffer to a database + committing offsets then.

### Consumer offset Reset Behavior:-

→ A consumer is expected to read from a log continuously.



→ But if your application has a bug, your consumer can be down.

→ If kafka has a retention of 7 days & your consumer is down for more than 7 days, the offsets are invalid. → Use replay capability in case of unexpected behaviour.

### Auto-offset-reset value:

Latest → will read from the end of the log  
earliest → will read from the start of the log  
none → will throw an exception if no offset is found.

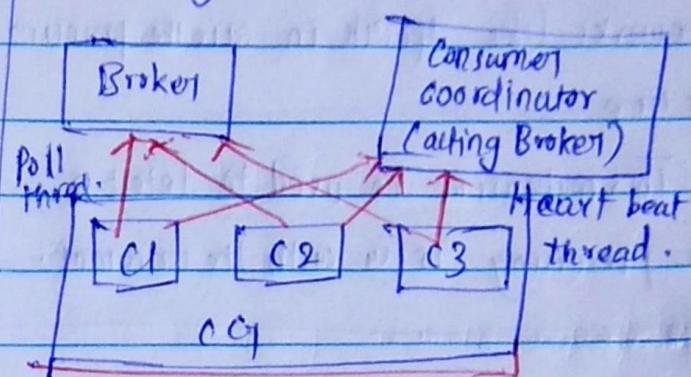
Additionally, consumer offsets can be lost.

→ If consumer has not read new data in 1 day (kafka < 2.0)

→ If consumer has not read new data in 7 days (kafka >= 2.0).

This can be controlled by the broker setting offset.retention.minutes.

### Controlling Consumer Lifelines:



→ Consumers in a group talk to a Consumer Group Coordinator.

→ To detect consumers that are down there is a heartbeat mechanism and a poll mechanism.

→ To avoid issues, consumers are encouraged to process data fast & poll often.

### Replaying data for a consumer group:

→ Take all the consumers from a specific group down.

→ Use 'kafka-consumer-group' command to set offset to what you want.

→ Restart consumers.

### Bottom line:

→ Set proper data retention period & offset retention period.

→ Ensure the auto offset reset behaviour is the one you expect/want.

### Consumer Heartbeat Thread:

session.timeout.ms (10 sec)

→ Heartbeats are sent periodically to the broker.

→ If no heartbeat is sent during that period, the consumer is considered dead.

→ Set even lower to faster consumer rebalances.

### HeartbeatIntervalMs (3 sec)

→ How often to send heartbeats.

→ Usually set to 1/3rd of session.timeout.

This mechanism is used to detect a consumer application being down.

## Consumer Poly Thread :-

max-poll-intrvalons (5 min) :-

→ Maximum amount of time b/w two poll() calls before declaring the consumer dead.

→ This particularly relevant for Big Data frameworks like Spark in case the processing take time.

→ This mechanism is used to detect a data processing issue with the consumer.

## Kafka Connect :-

→ Code & connectors re-use.

CQRS - Common Query Responsibility

Segregation.