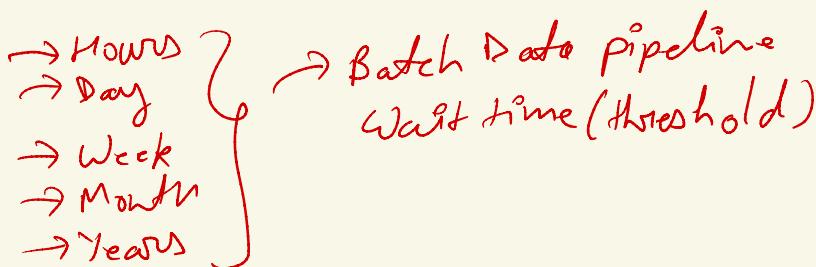
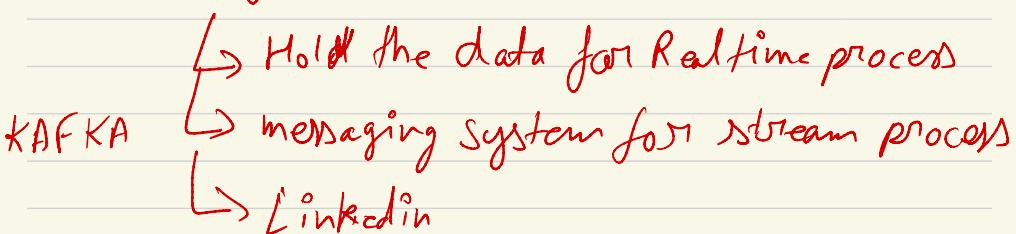



Types of Datapipelines

- ① Batch Data processing → Wait for sometime
Hadoop, Hive, Apache Spark
- ② Near Real Time (NRT) × 30 Sec, 1 Min, 2 min
- ③ Real Time Data Processing × Instantly



* Messaging Queues for Real Time Data Process



Real Time processes

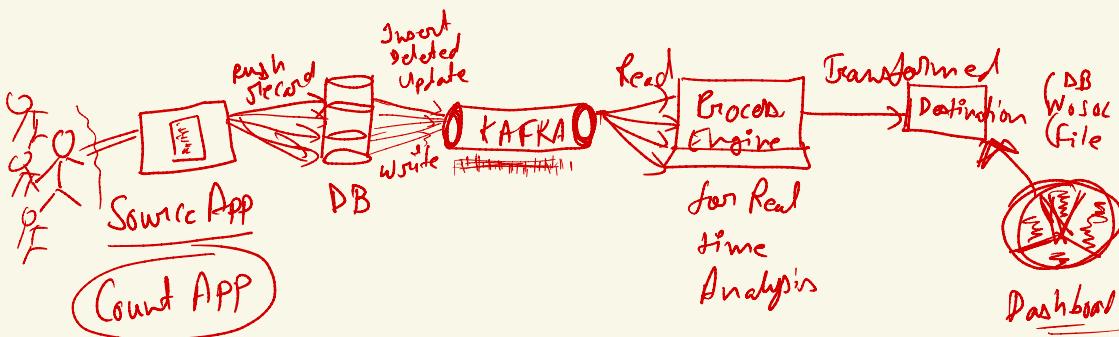
- ① Bank Transaction
- ② Application for Live Cricket scores and aggregations
- ③ Gaming Application (live tournaments)

④ Trending Twitter Hashtags

⑤ Any type of real time dashboards

↳ Stock Exchange Apps

↳ Live Counter in a mall



KAFKA Basic Actions

↳ Producers *

↳ Consumers *

↳ Connectors (Practical)



- Producers are publishing Data
- Consumers basically Subscribe

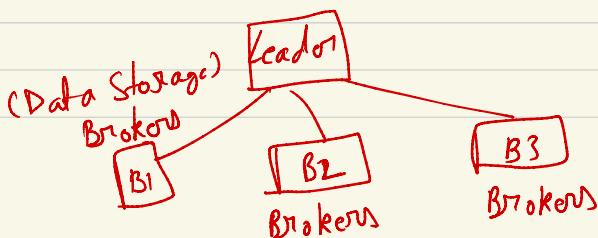
KAFKA Components

- ↳ Broker
- ↳ Topics
- ↳ Partitions
- ↳ Replicas
- ↳ Offset
- ↳ CommitOffset
 - Sync
 - Async
- ↳ Consumer Groups
- ↳ Parallelism in KAFKA
- ↳ Back pressure

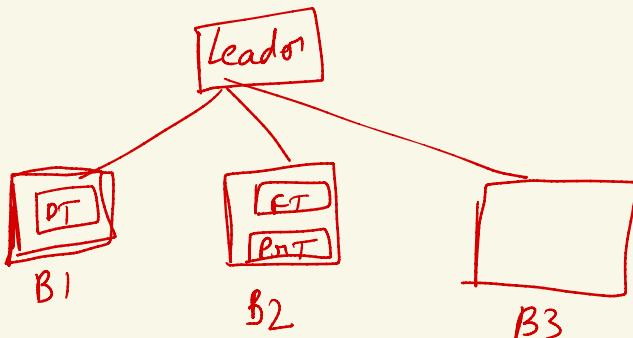
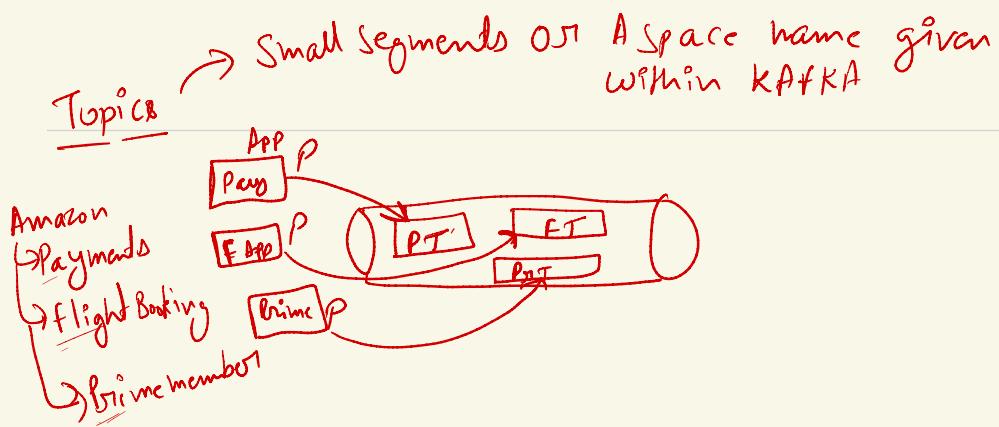
KAFKA Architecture:

- ↳ KAFKA Cluster

↳ Brokers means servers in a cluster



Topic



Partitions

Message = Record

KAFKA Topic

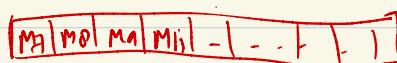
P0



P1



P2



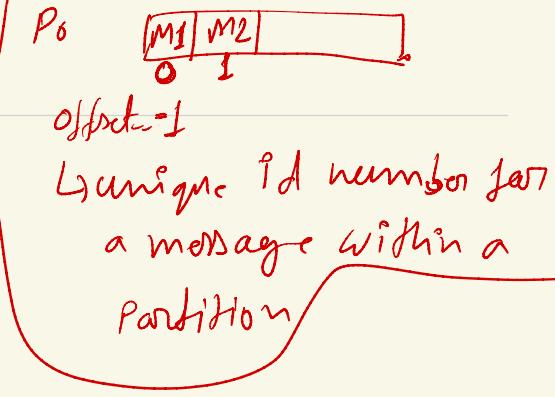
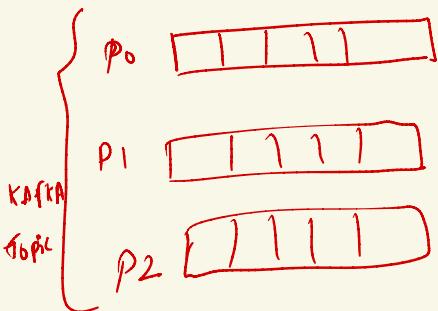
→ Current Partition Index = Offset



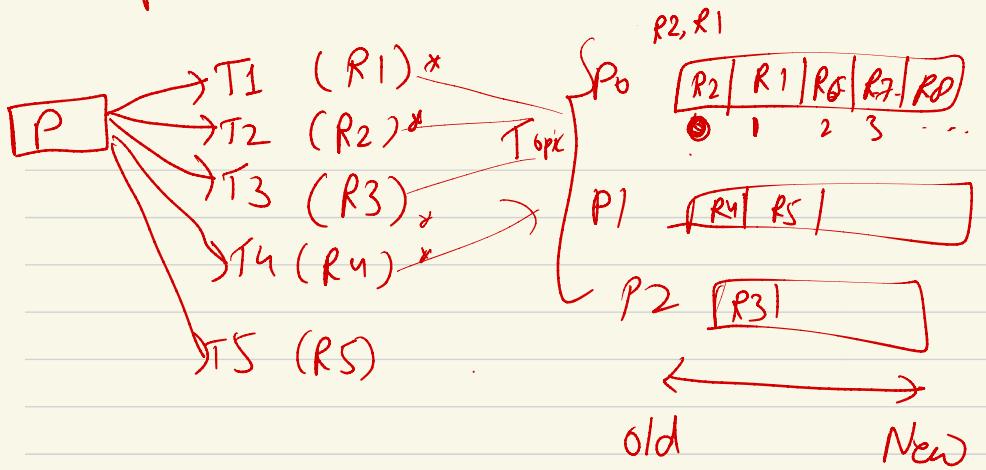
offset = 1

P0





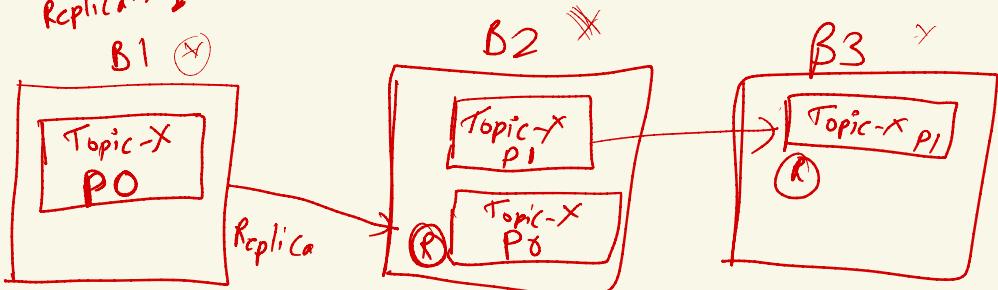
Order of messages will be maintained within a partition



Partitions
 ↗ fault tolerance
 ↗ Parallelism

Topic → X
Part → 2
Replica → 3

Lead ⚡



Message without key

P R
Message with key

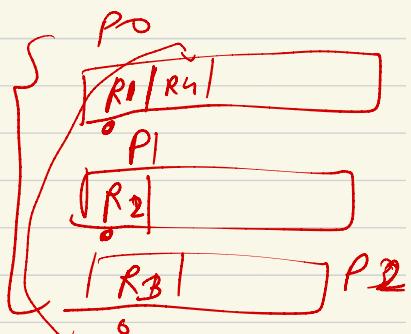
R1 {Key, R1y}
S R }

KAFKA Topic

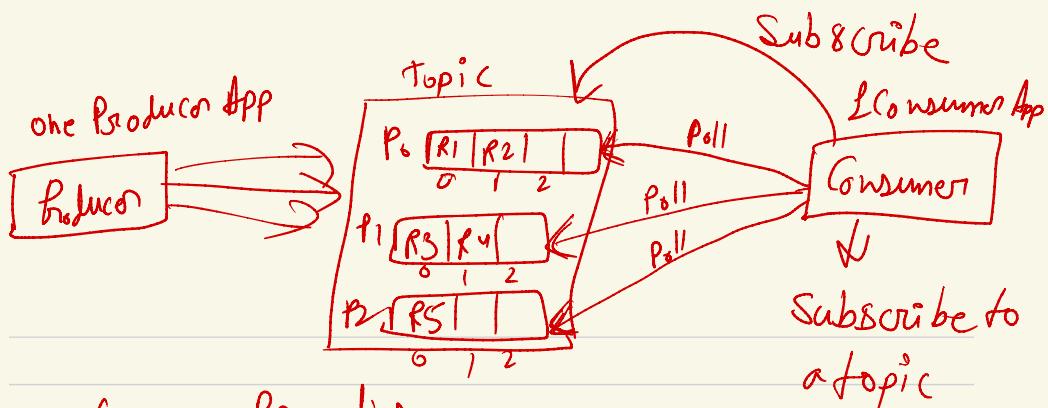
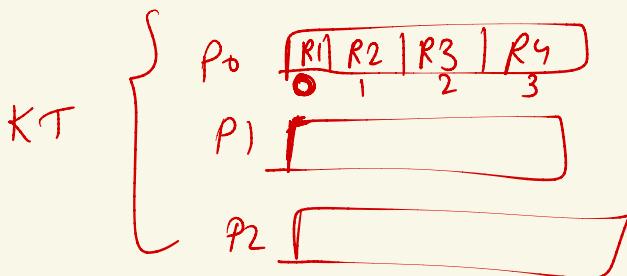
{6, R1y} → Key = 6
6 % 3 = 0

{8, R1y} → 8 % 3 = 2

{7, R1y} → 7 % 3 = 1



$\{6, R1\}$, $\{8, R2\}$, $\{6, R3\}$, $\{6, R4\}$

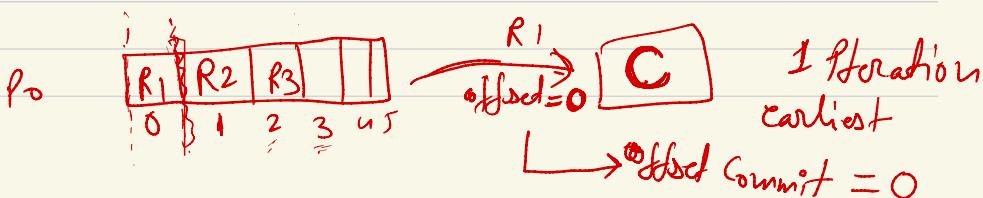


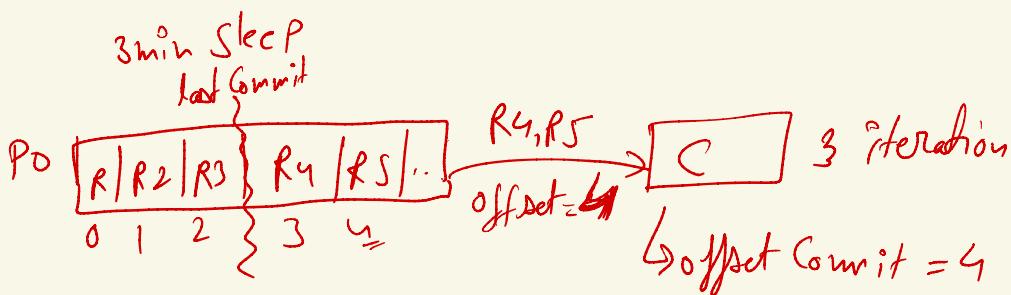
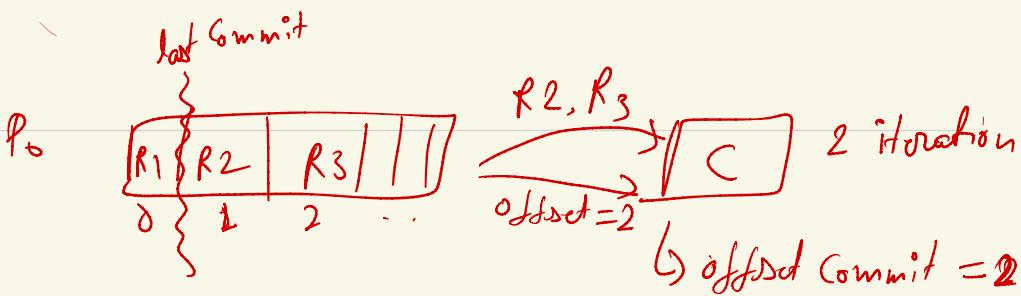
Consumer Properties

→ Latest (Read data from each partition from current offset)

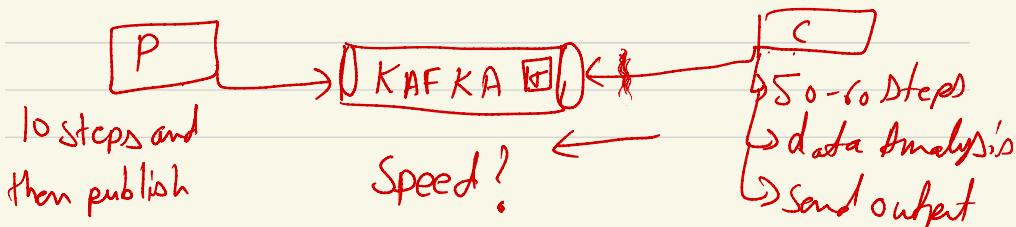
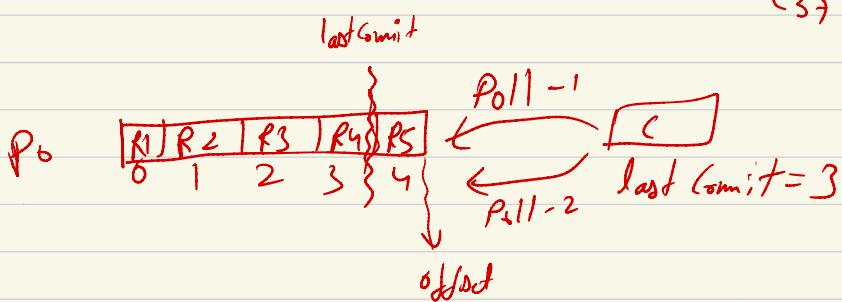
→ Earliest (Read for each partition from start offset)

Consumer needs to do the offset Commit

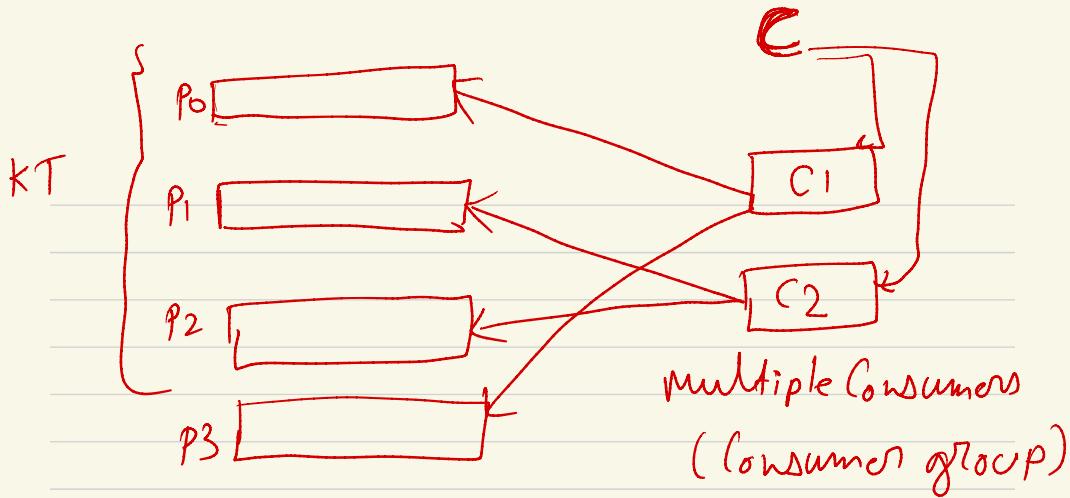
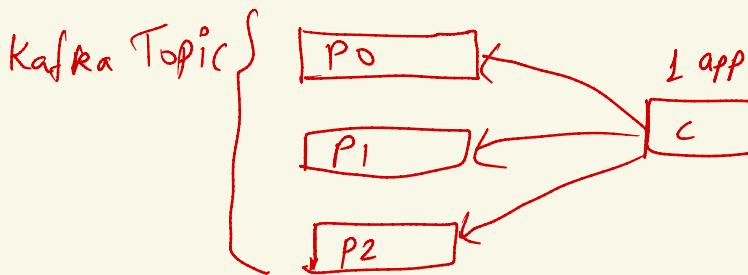




$$\text{consumer_0} = \begin{cases} 1 \mapsto 0 \\ 2 \mapsto 2 \\ 3 \mapsto 4 \end{cases}$$



Backpressure } \Rightarrow because of
Speed mismatch

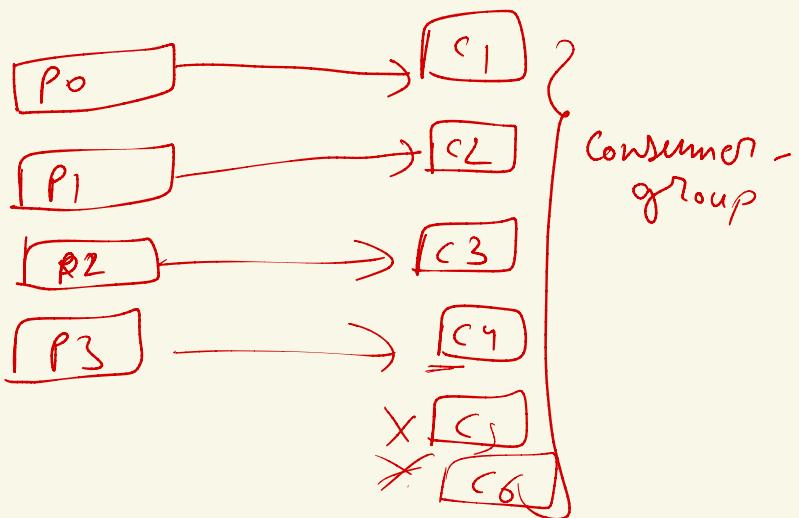


C = Consumer()

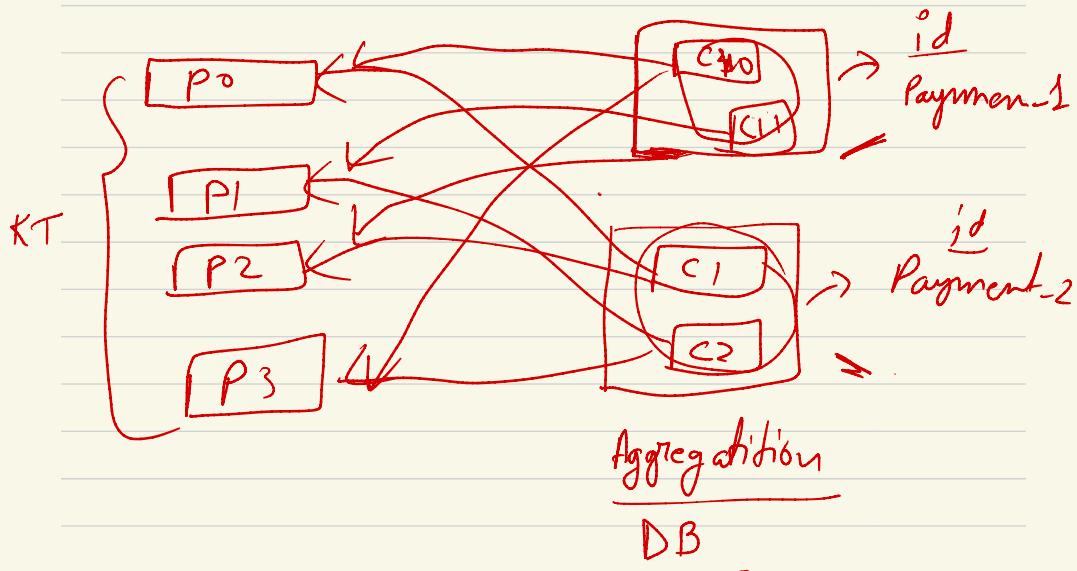
C = Consumer(groupI) ("Platfrom Consumer")

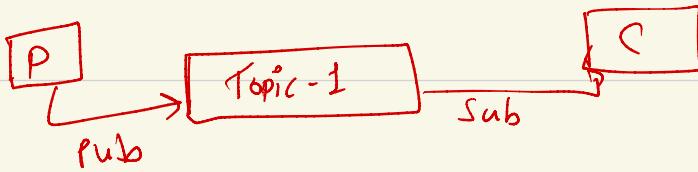
At a time only 1 consumer of a consumer group can read data from a specific

Partition



$\text{f}_1 \text{HOT}, \text{HDFJ}$

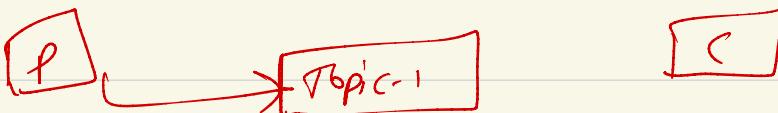




↳ Sync = Wait until the operation is complete

↳ async = Means fire & forget

↳ Don't wait for anything



step → Poll the data ^{to Record}

② → transform ✘

③ → Dump into destination ✘

④ → Commit offset ↳

⑤ → ① ↳ Sync Commit

⑥ → ↳ Async Commit

① Sync Commit

↳ Where we are not concerned latency

- Data Consistency
- Avoid data loss while processing

② Asynchronous Commit

- Where we need latency
- Inconsistent data
- data loss as well

