



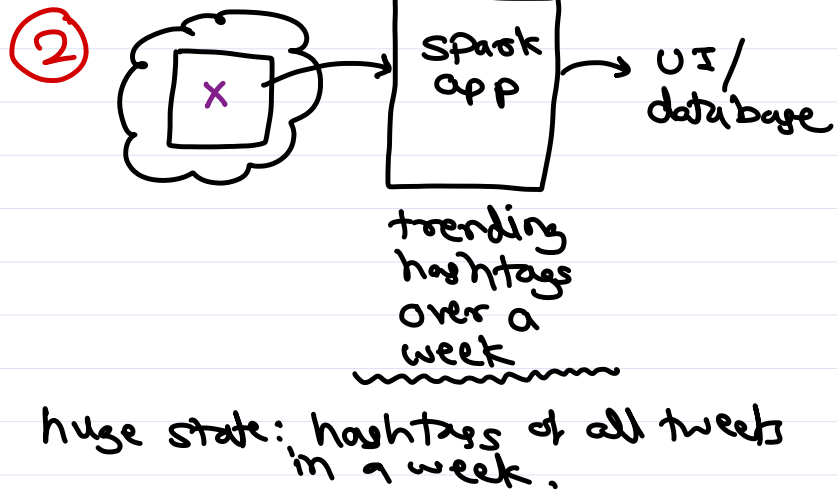
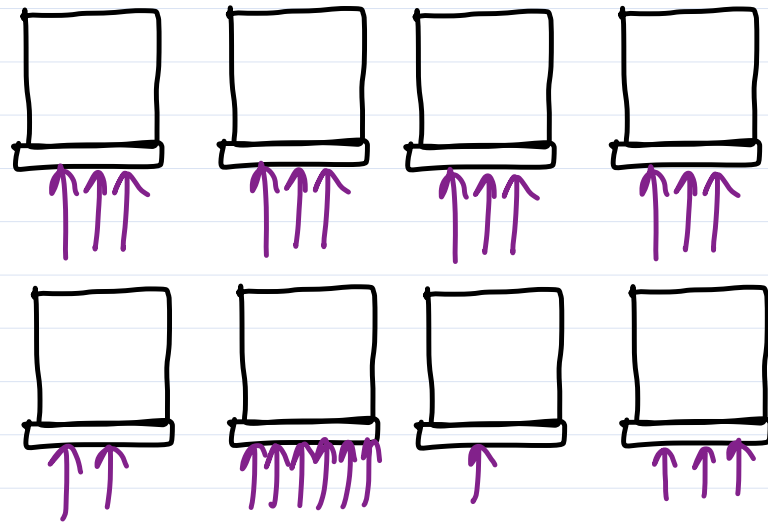
Big Data Technologies

Trainer: Mr. Nilesh Ghule.



Stream Processing

① Load Balancing

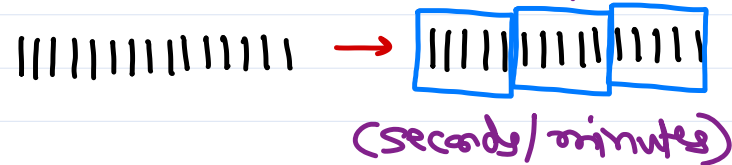


③ exactly once processing

- Ⓐ at least once (may repeat)
- Ⓑ at most once (may skip)
- Ⓒ exactly once

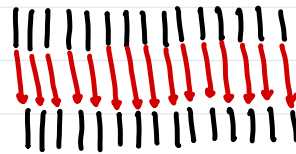
④ Low latency processing

@ micro-batch approach

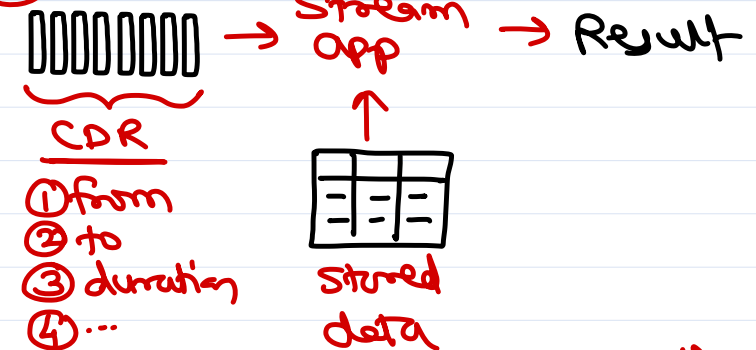


ⓑ Continuous processing

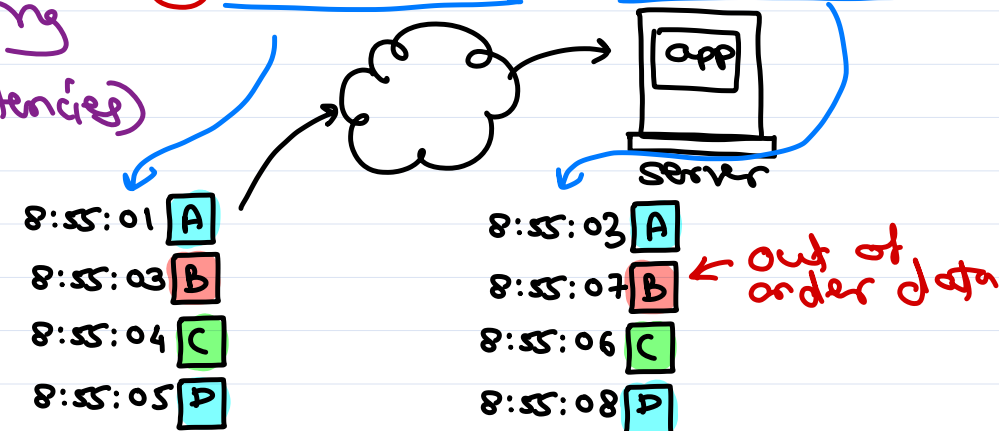
(lower latencies)



⑤ Join with external data



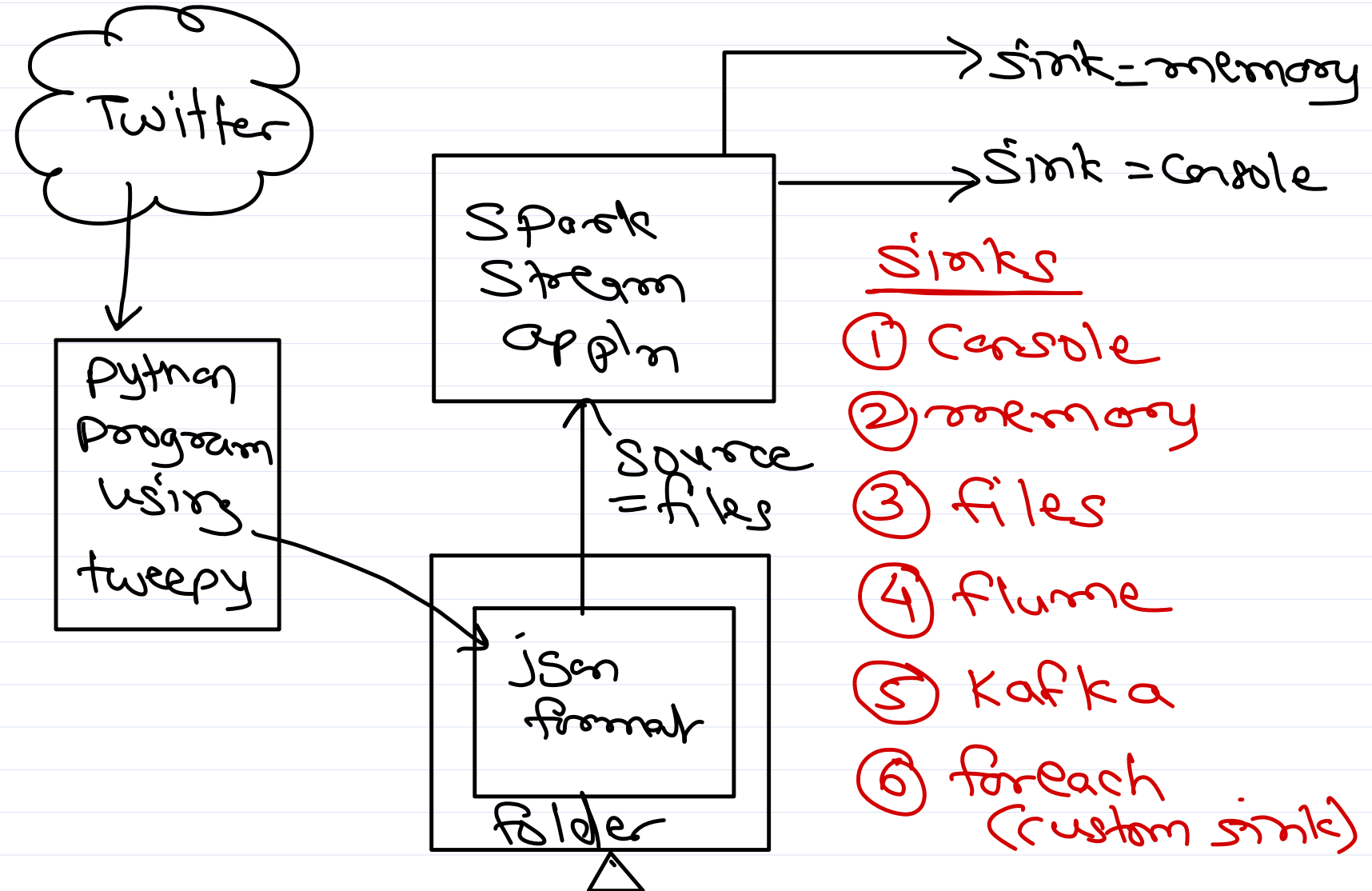
⑥ Event time vs processing time



Spark Streaming

Sources

- ① socket
- ② files (csv, json, ...)
- ③ flume
- ④ kinesis
- ⑤ rate (testing)
- ⑥ kafka

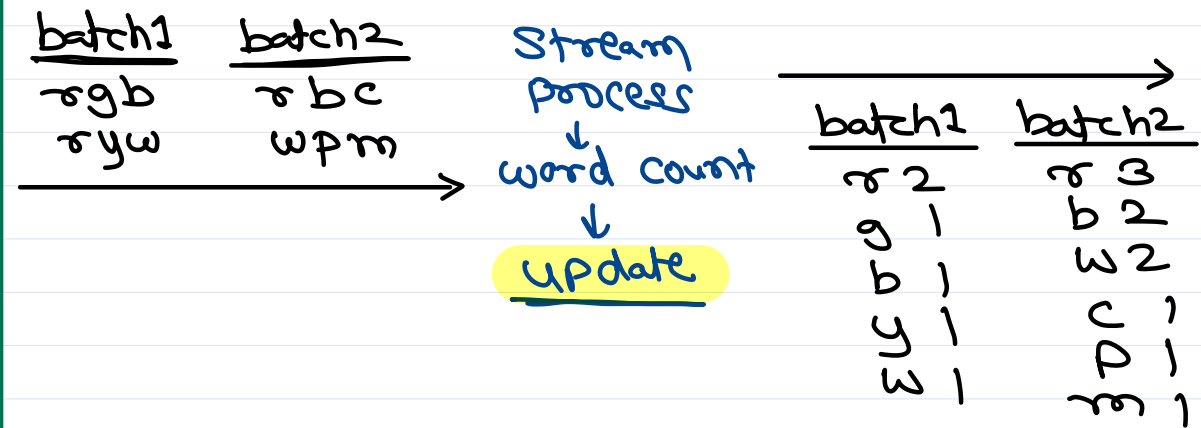
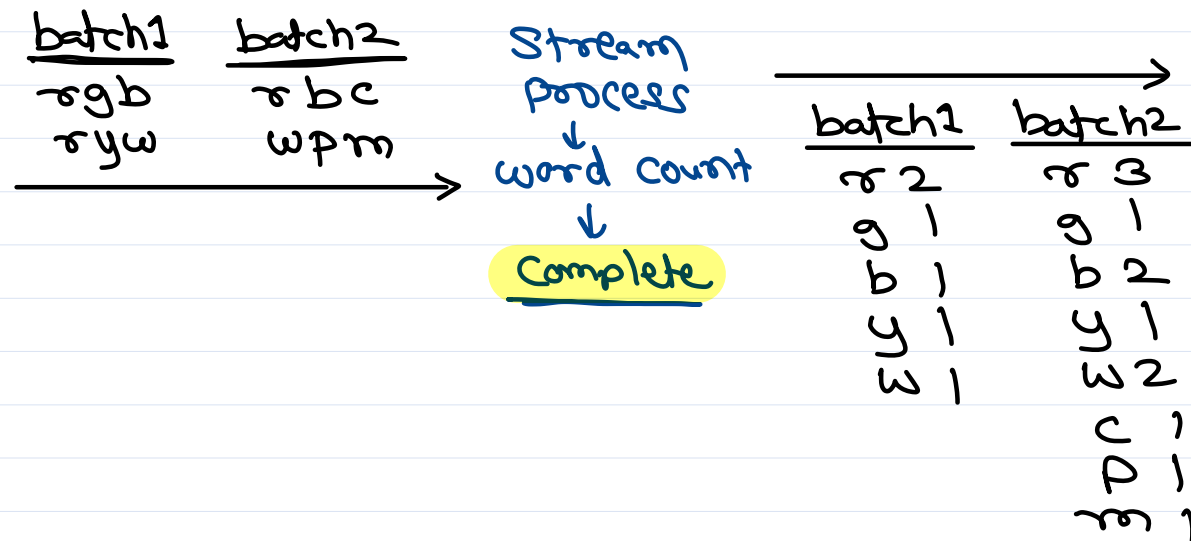
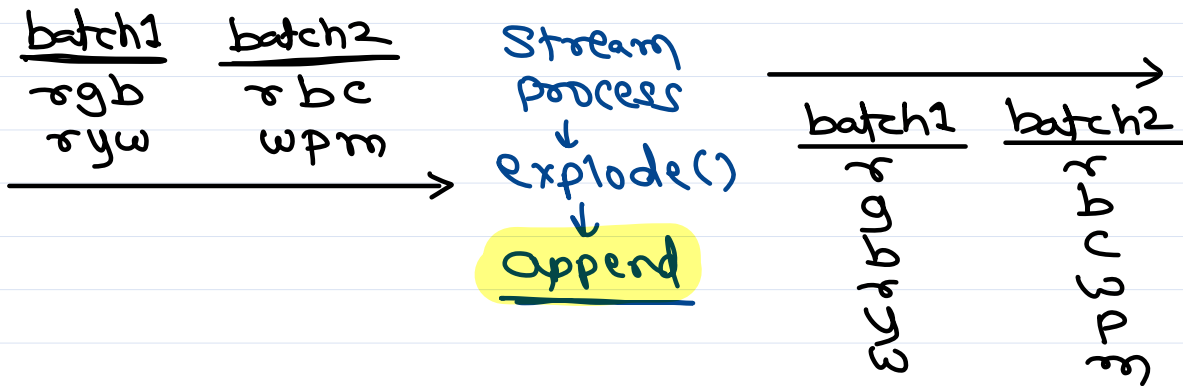


Sinks

- ① console
- ② memory
- ③ files
- ④ flume
- ⑤ kafka
- ⑥ foreach (custom sink)

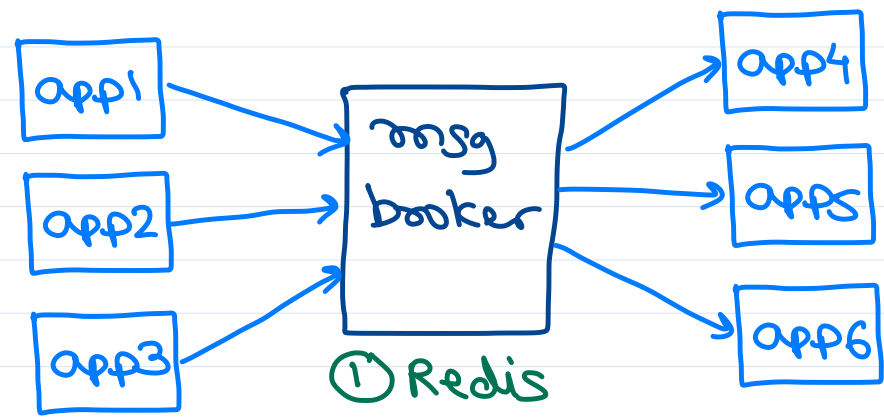


Stream Output Mode

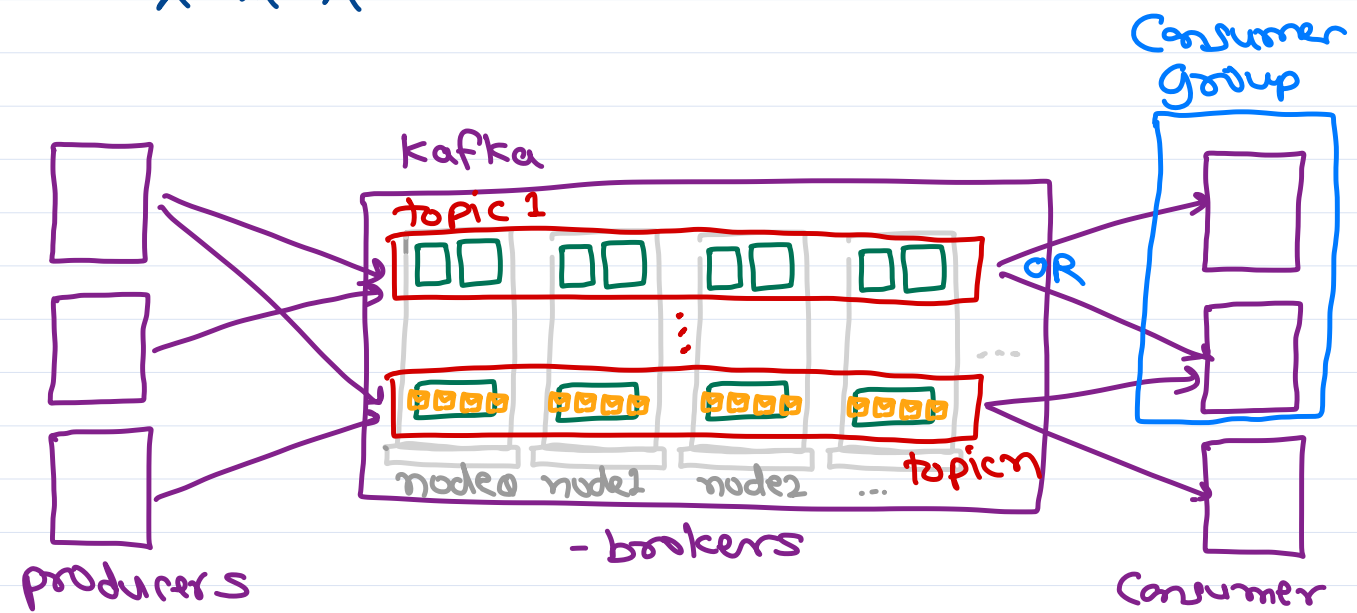
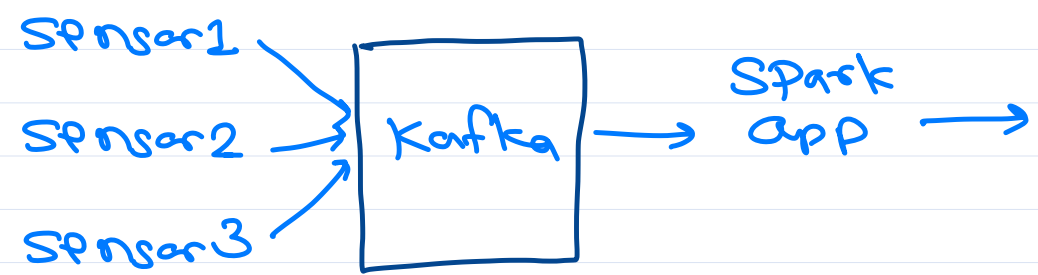
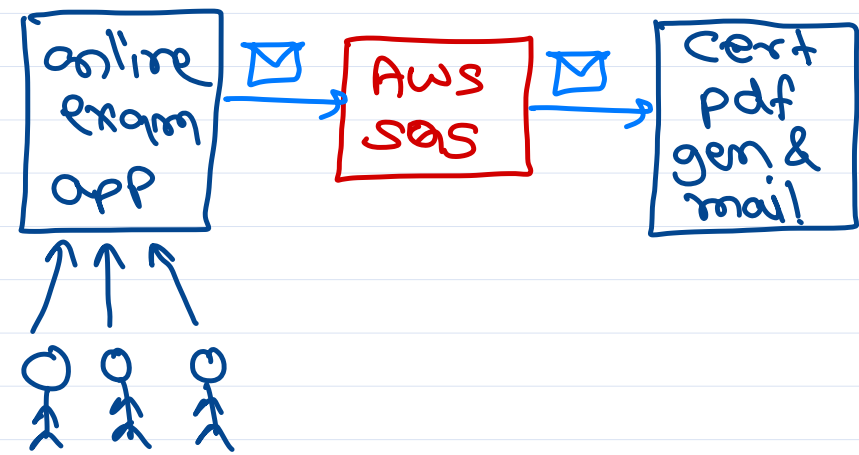


Kafka

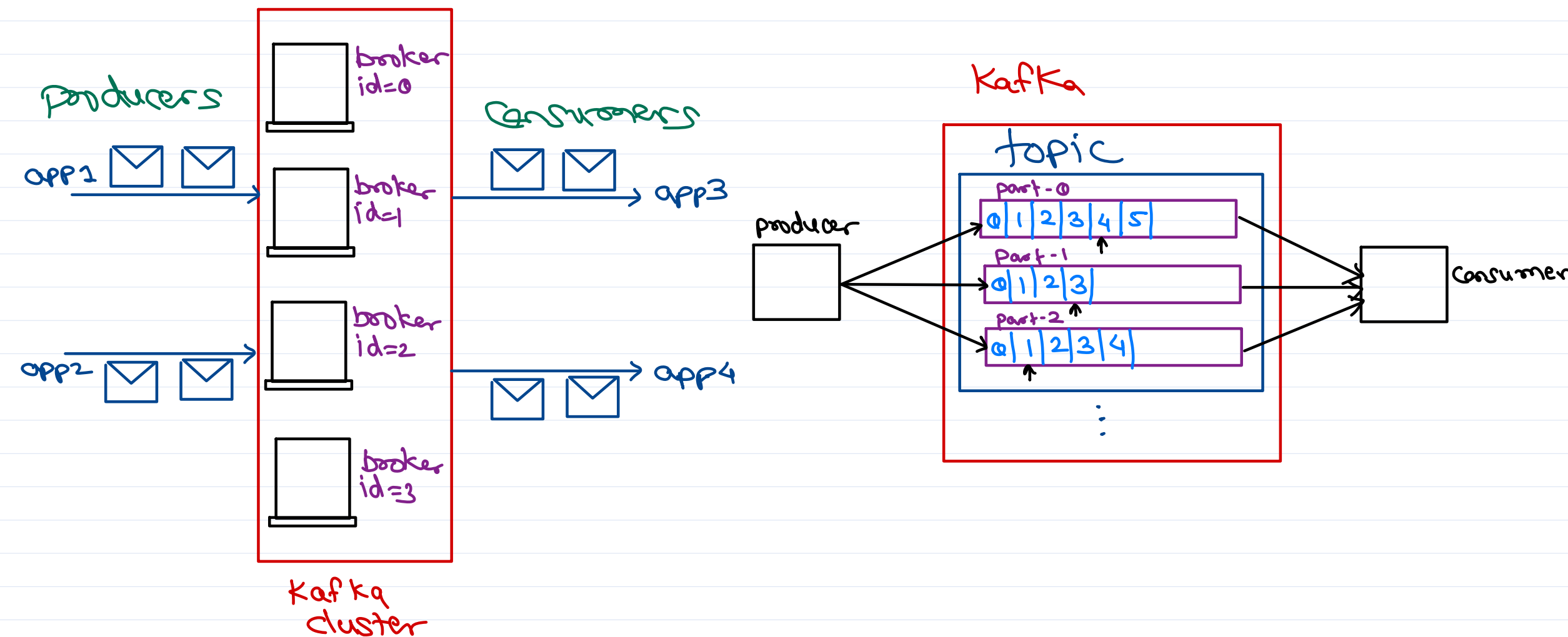
* Distributed message broker



- ① Redis
- ② Kafka
- ③ RabbitMQ
- ④ AWS SQS



Kafka cluster and topics



Kafka Architecture

Kafka producers follow
 ① Round Robin or ② Hashing
 to balance load across parts
 for each topic.

prod 1 →

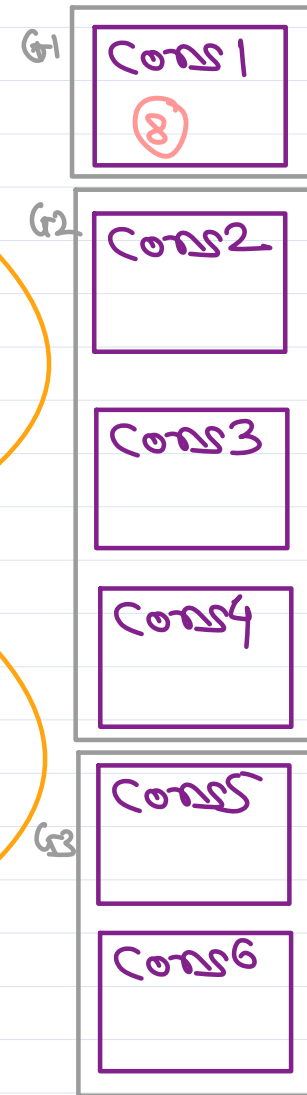
prod 2 → ✓

prod 3 →

prod 4 →

prod 5 →

broker 1
broker 2
broker 3





Thank you!

Nilesh Ghule <nilesh@sunbeaminfo.com>

