

Leveraging Apache Kafka for Web Crawling and Data Processing

OpenStack Cork, 2017-11-21



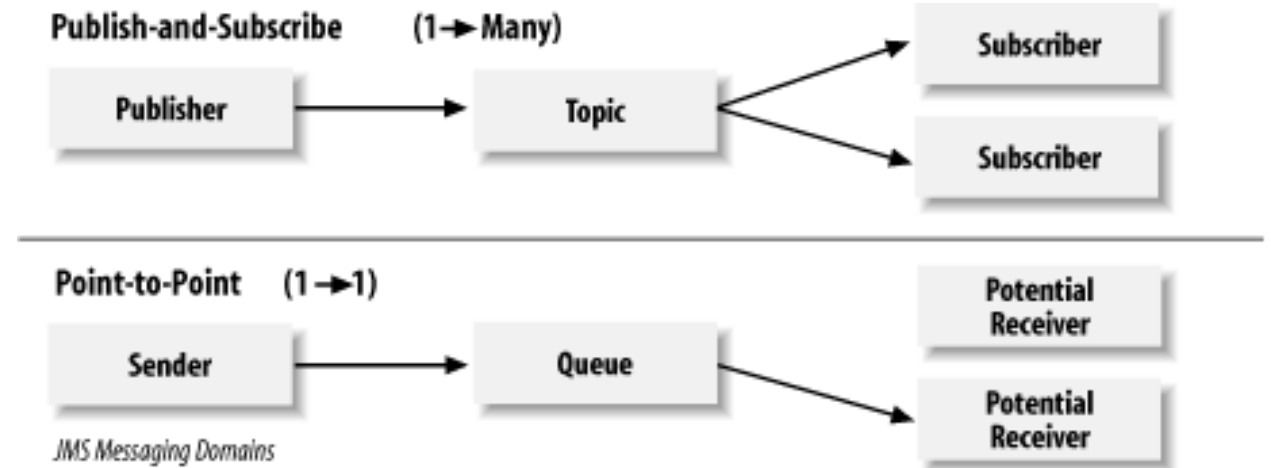
About Me

- Johannes Ahlmann
- fluquid.com
 - Sales & Client Intelligence
 - Intelligent Lead Generation
 - Large-scale web crawls
 - Gathering and Enriching Web Data
- webdata.org
 - Share Libraries and Best Practices
 - Bring Data Scientists and SME Companies together
 - [ForDevelopers](#)
 - [AwesomeAvailableDatasets](#)
- Contact:
johannes@fluquid.com



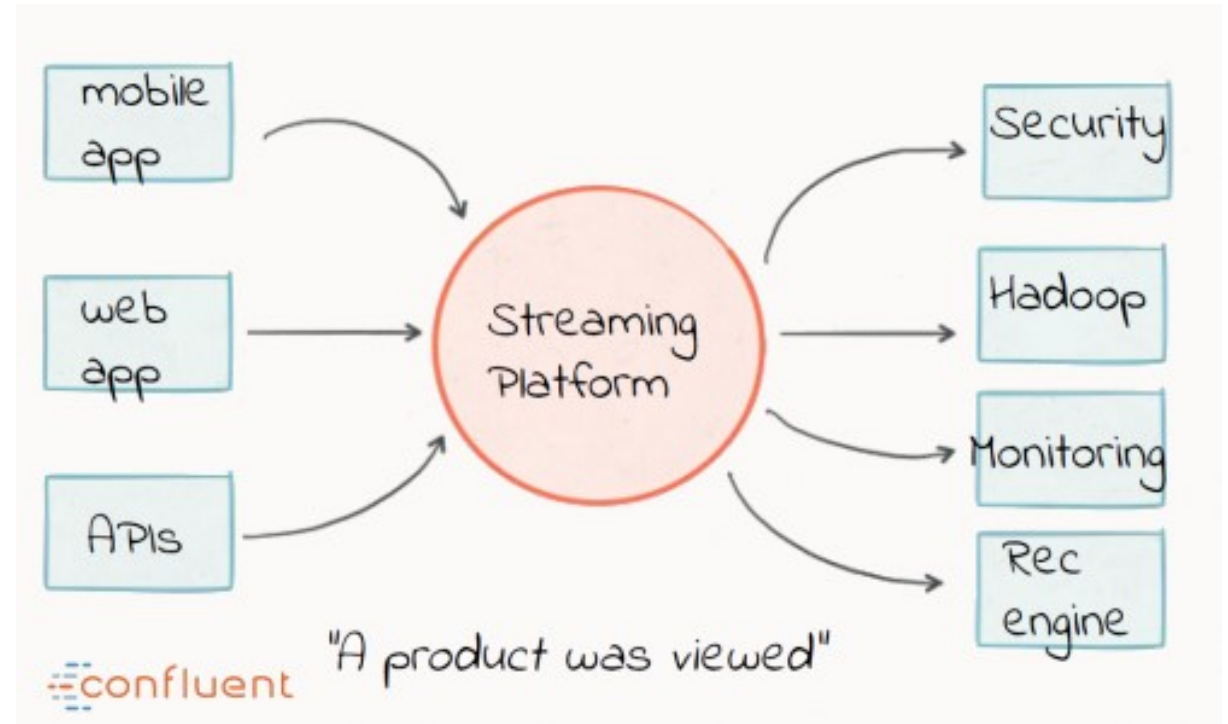
Background: Queues & PubSub

- PubSub
 - multiple subscribers
 - no scaling
- Queues
 - multiple consumers
 - single-subscriber
 - scaling, load balancing



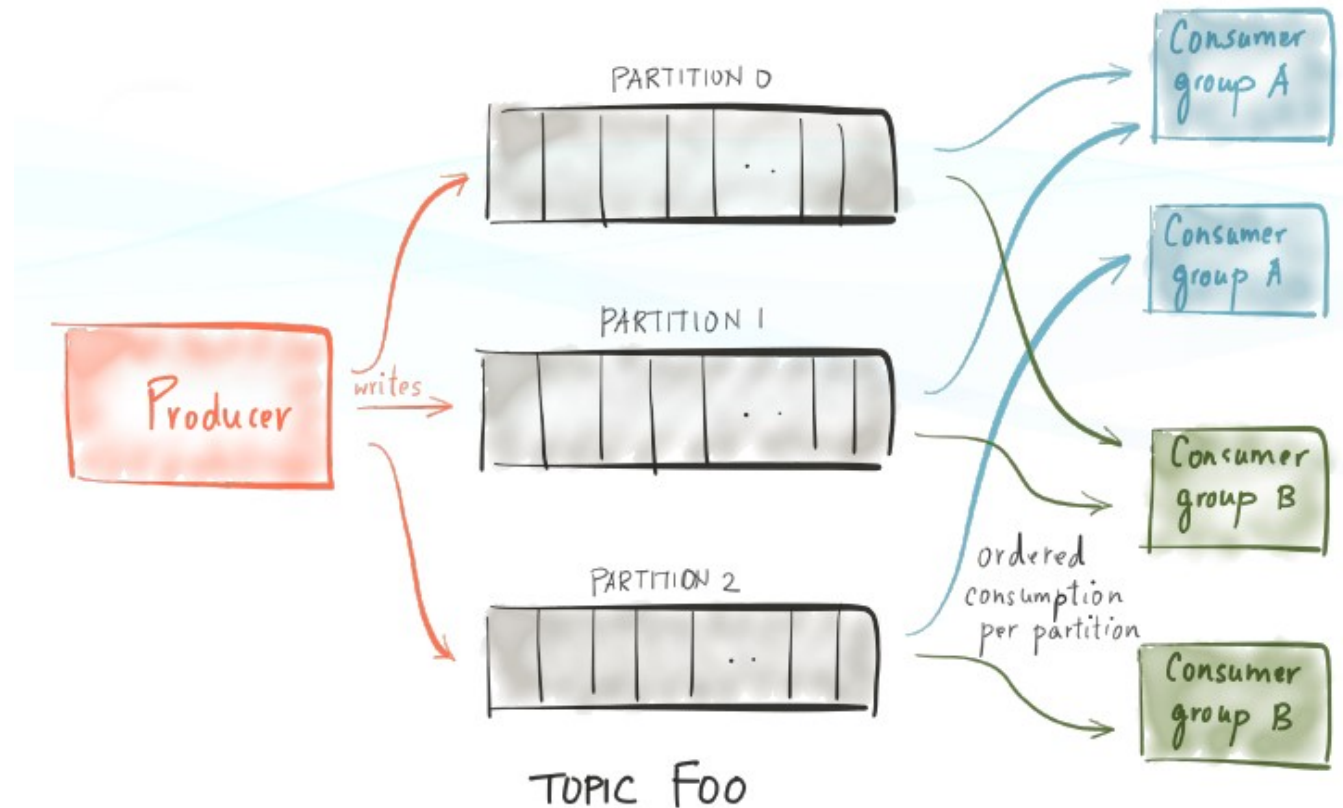
A high-throughput distributed messaging system

- Decouples Data Pipelines
- Scalable & Fault-Tolerant
- Kafka Functionalities
 - Messaging
 - Processing
 - Storing
- Performance (>100k/s)
 - Batching
 - Zero Copy I/O
 - Leverages OS Cache
- Durability



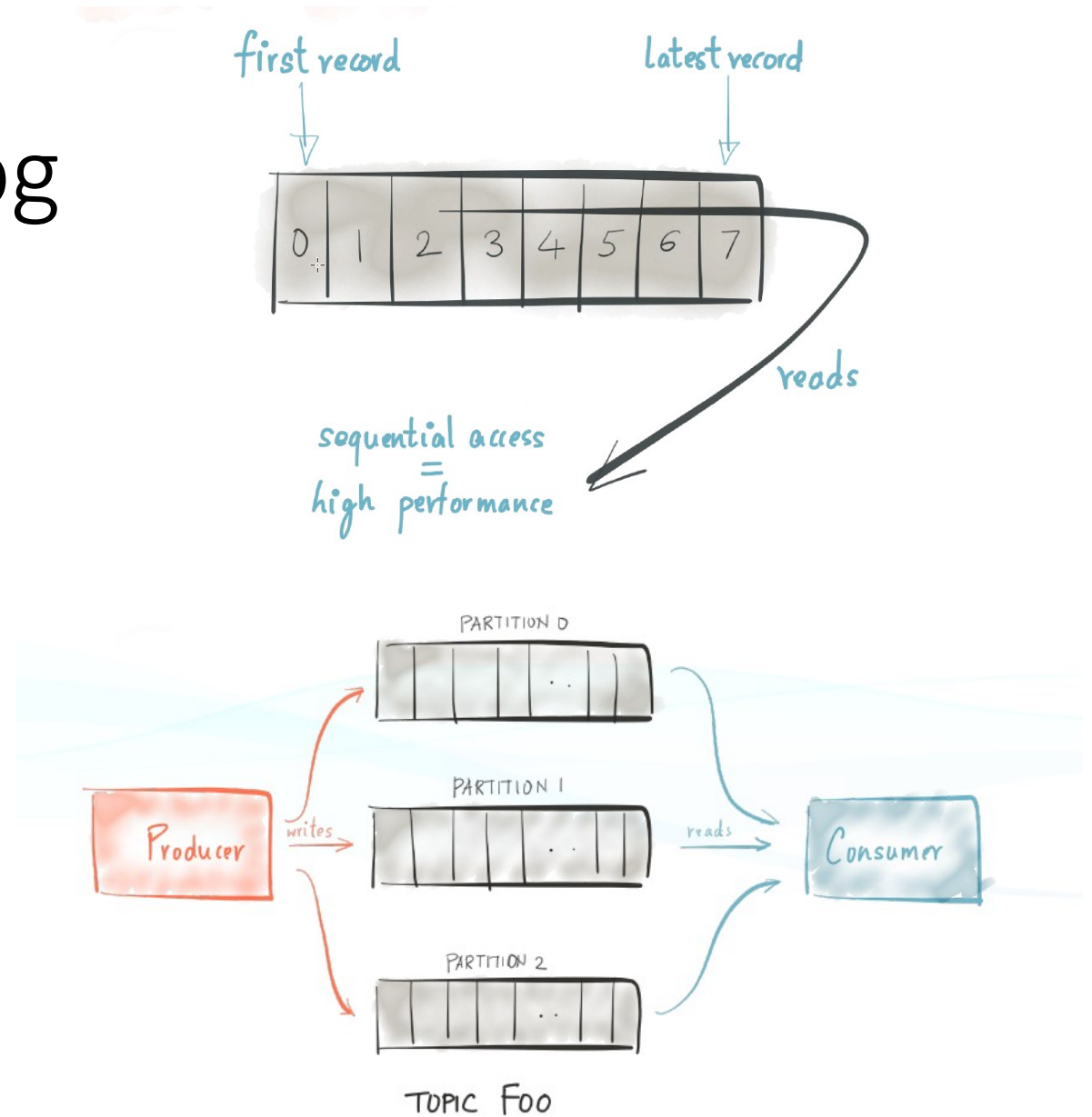
Core Concepts

- Producers
- Consumers
- Brokers
- Topics
 - Offsets
 - Broker Addresses



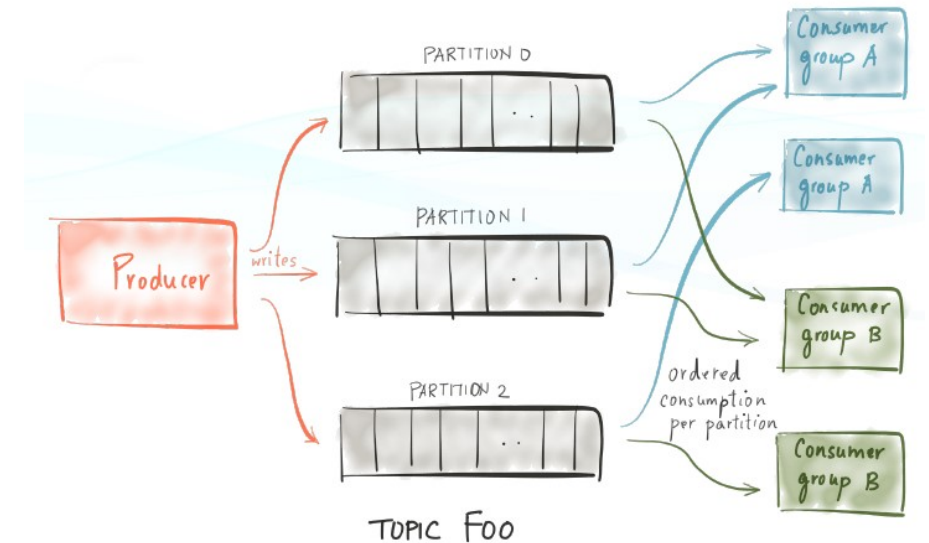
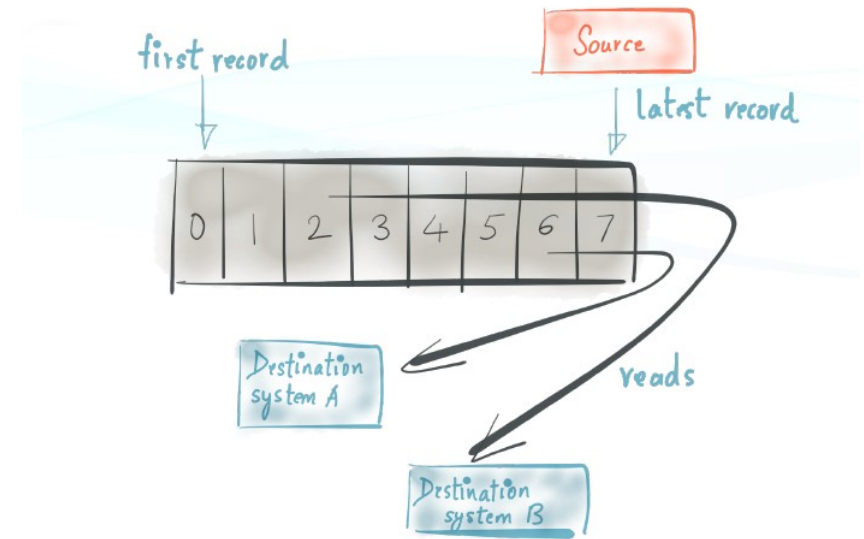
Key Idea: Partitioned Log

- Very fast, due to zero copy I/O and batching
- Uses sendfile and OS buffer cache
- Sequential writes to FS
- Order guaranteed within partition
- Scaling



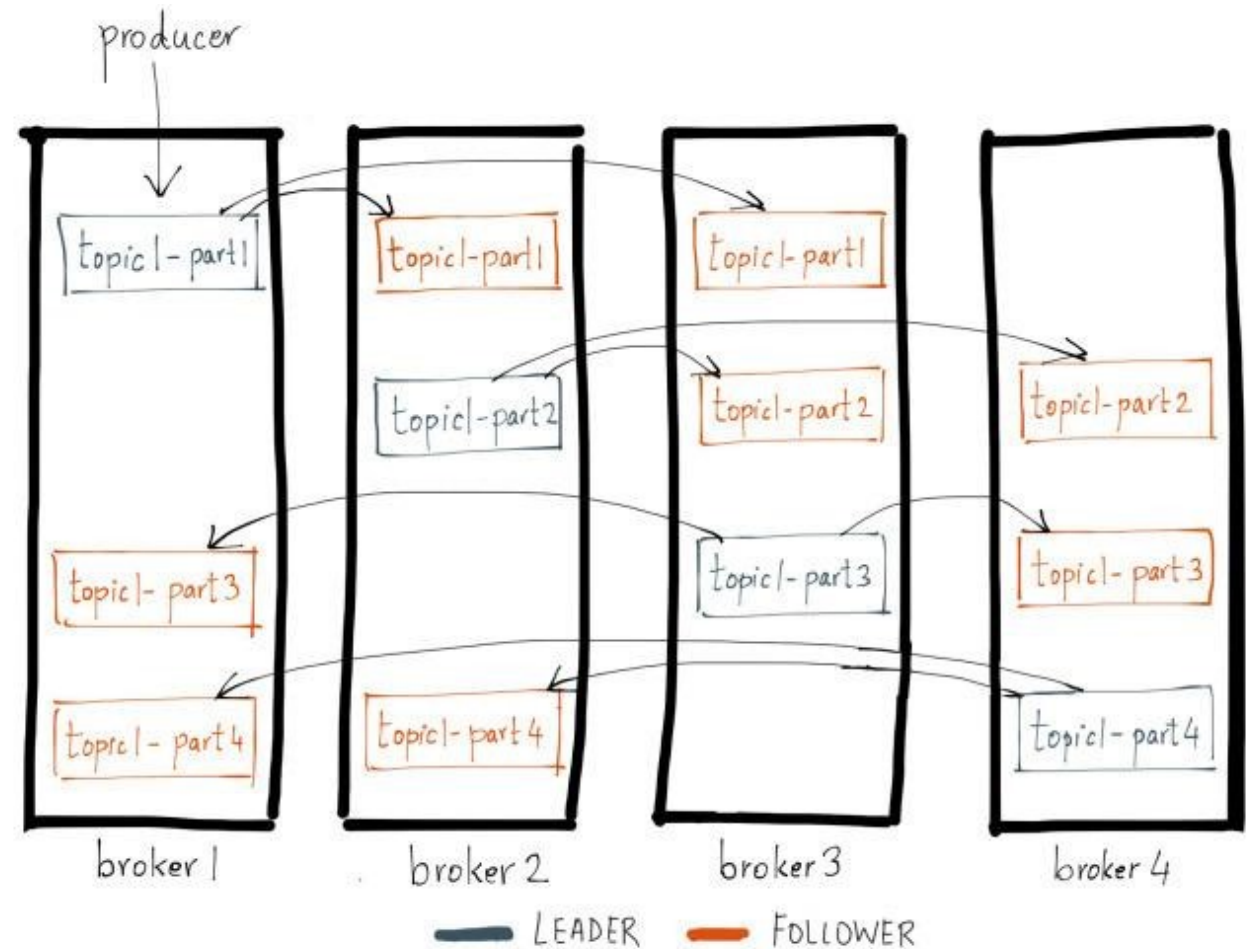
Logs & PubSub

- Consumers can be transient
- Consumer Groups
- Delivery Semantics
 - at least once (default)
 - at most once
 - exactly once
- Retention Policy
- Reprocessing

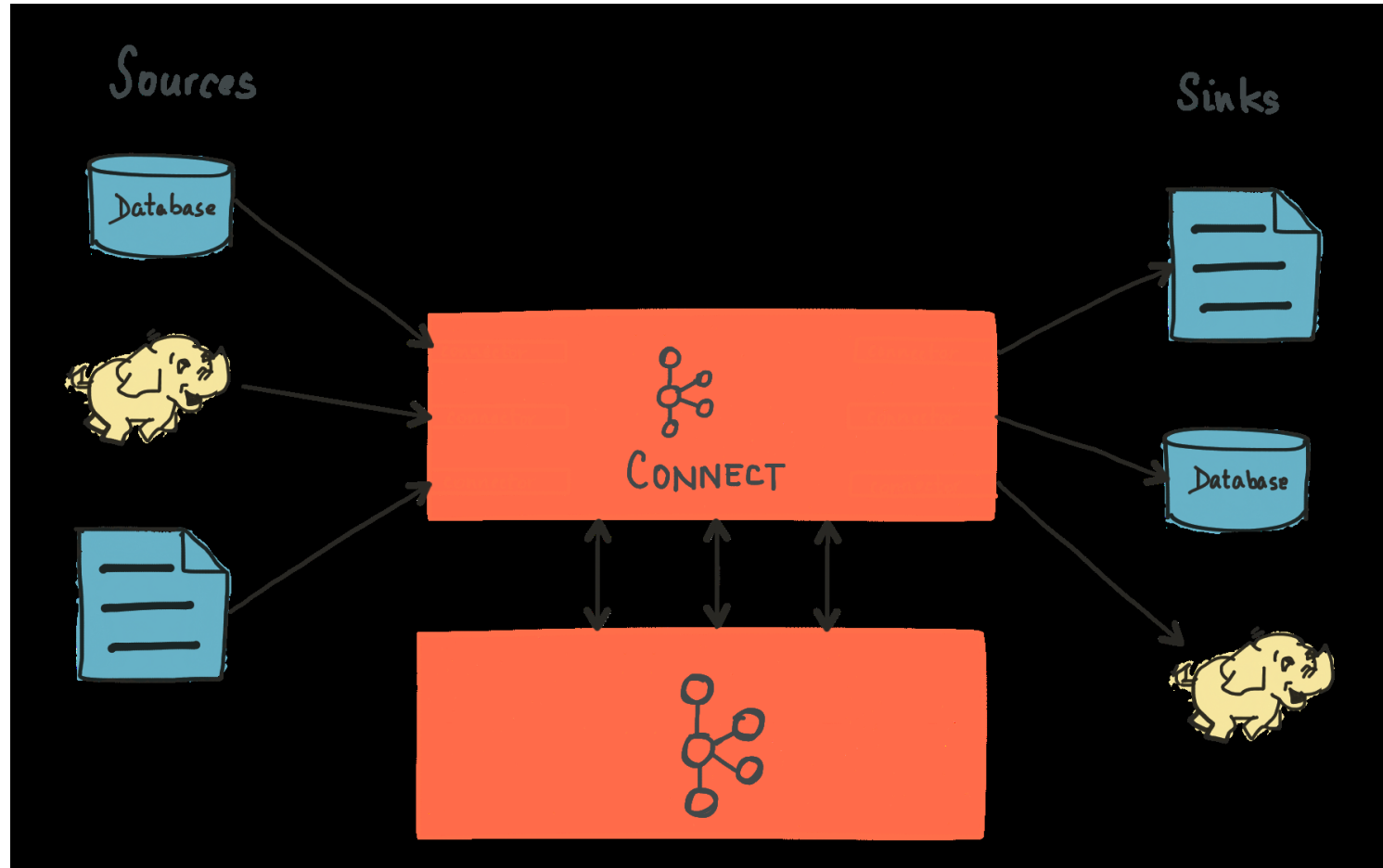


Partitions & Replication

- Partitions configurable
- Partition allocation
 - round-robin
 - semantic partition by key
- Replication
 - optional
 - 1 leader, 0 or more followers
 - sync or async
 - flush delay configurable



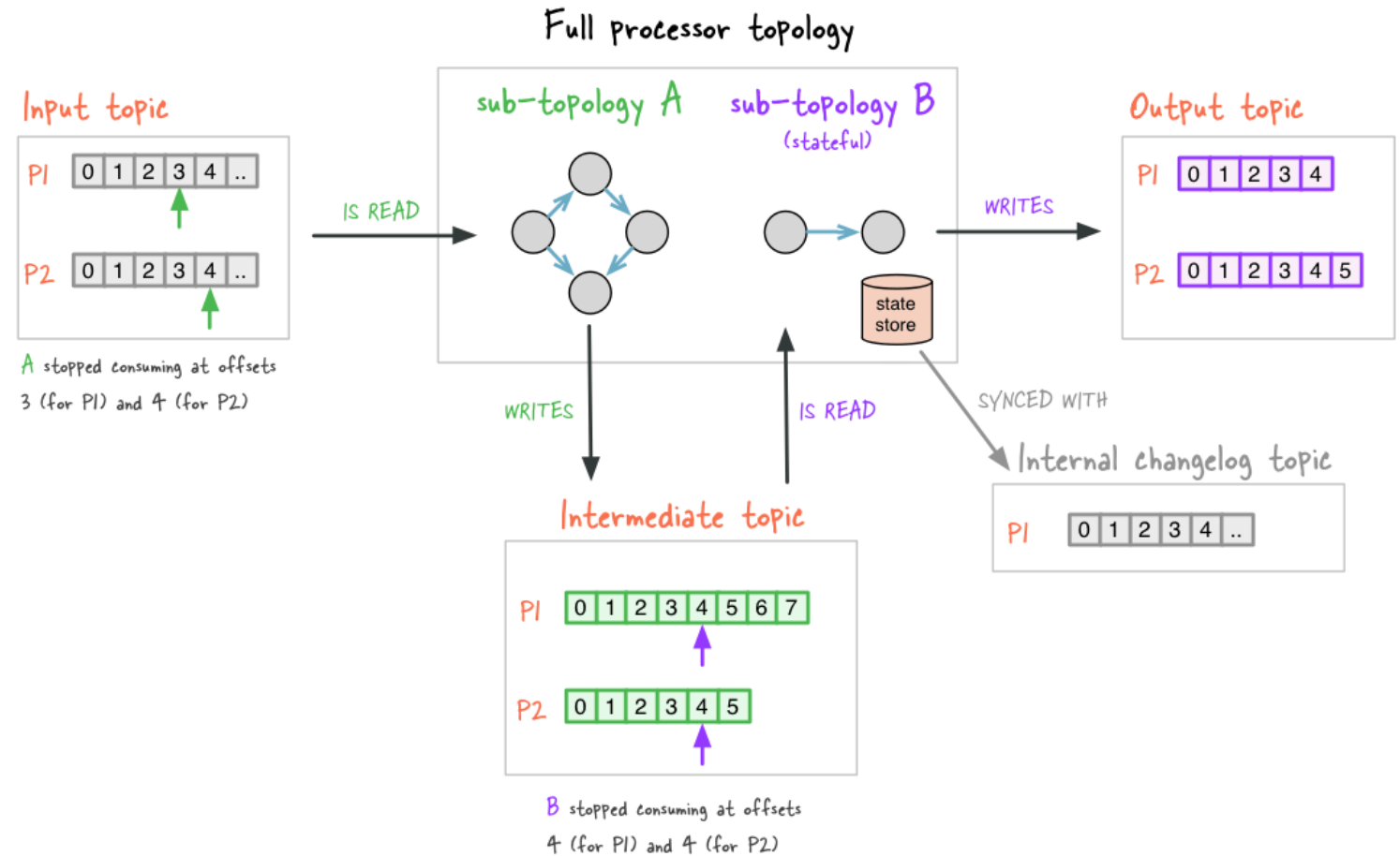
Kafka Connect



Kafka Streams

- Operations
 - filter
 - map
 - join
 - aggregate
- KStream
- KTable
 - manages local state
- Windows

} stateful



Summary

scalability of a filesystem

- hundreds of MB/s
- many TBs per server
- commodity hardware

guarantees of a database

- persistence
- ordering

distributed by design^I

- replication
- partitioning
- horizontal scalability
- fault tolerance