

Simplifying Distributed Systems

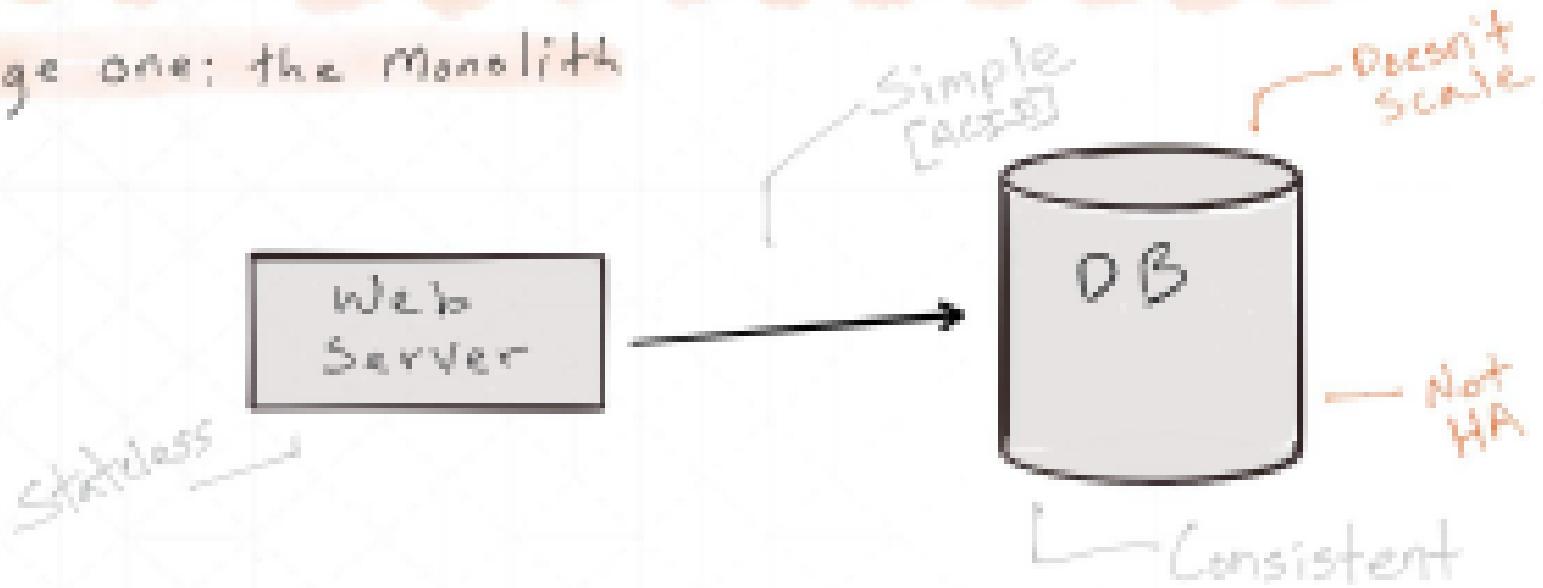
With
Apache
Kafka

Introduction

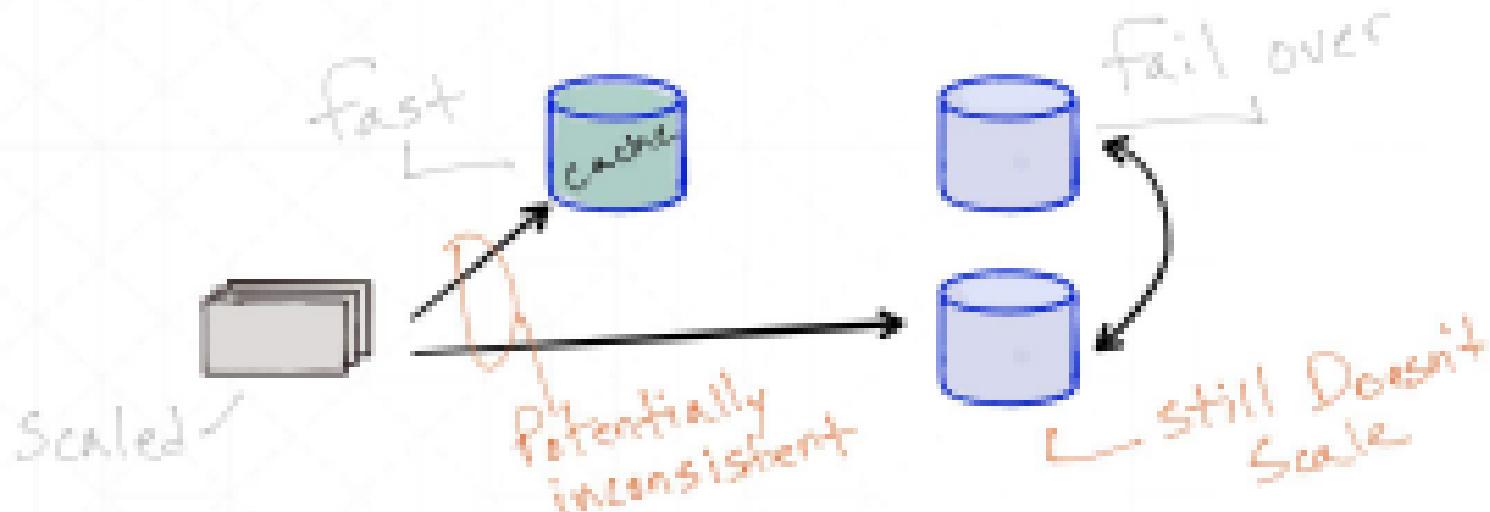
Why Distributed?

[How did this happen]

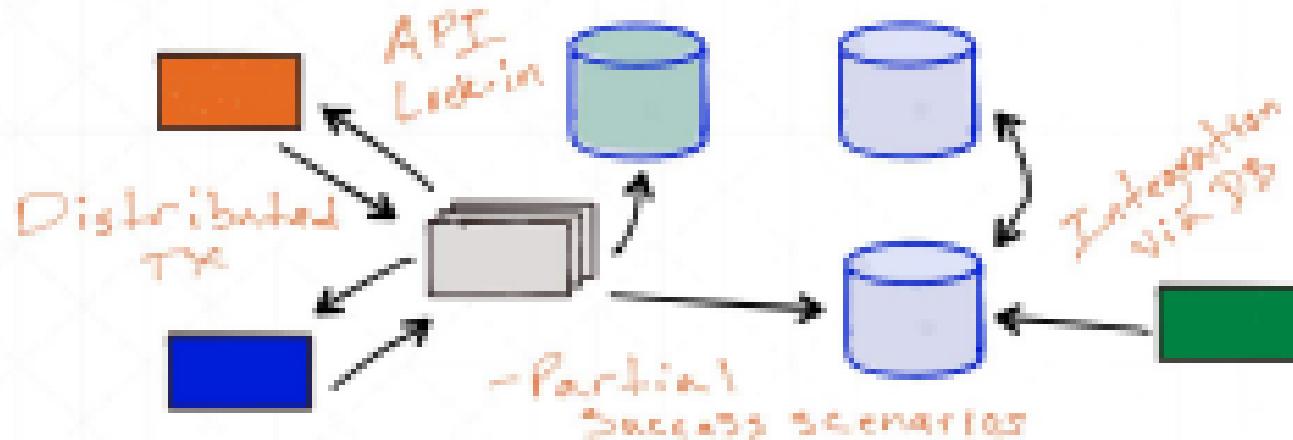
Stage one: the Monolith



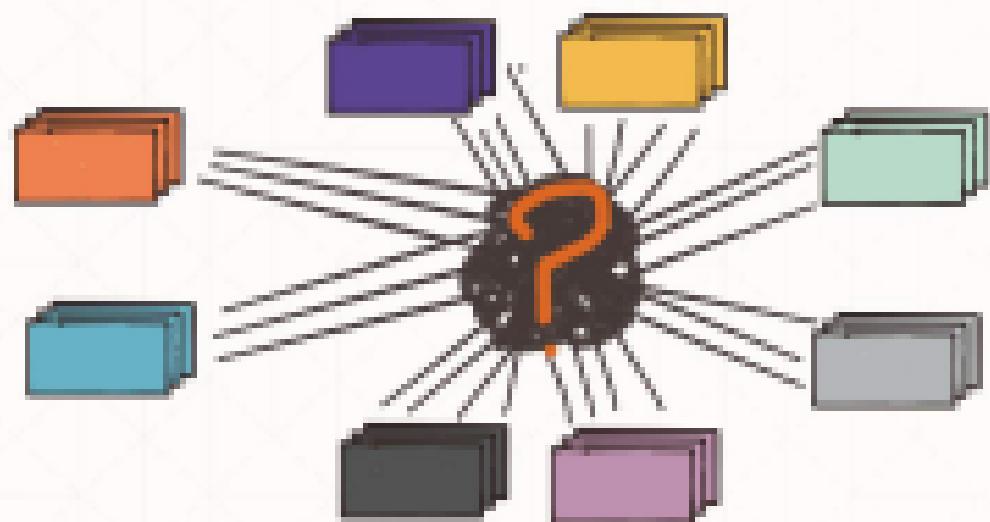
Stage 2 : HA/Scale



Stage 3: Integrations



So on...and so on



Microservices?

The Problems

- Point-to-Point Communication
- Multi Dispatch Updates
- Distributed Transactions
- Integration Via Database
- API Coupling

Apache Kafka

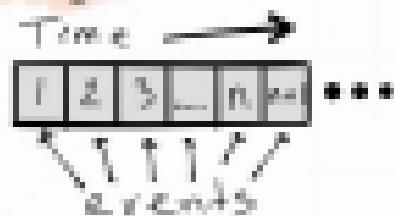
- Originally built by LinkedIn
- Apache top level project since 2012
- Used by many notable companies
- Modeled as a distributed transaction log

Design

Append only log

- Commonly used in databases [no move]
- Ordered
- Analog to "Stream"

The log



More Design

A topic has 1-N partitions

Topics can be replicated

Partitions can reside
on different brokers

Scales based on
partition count

Topic



Clusters

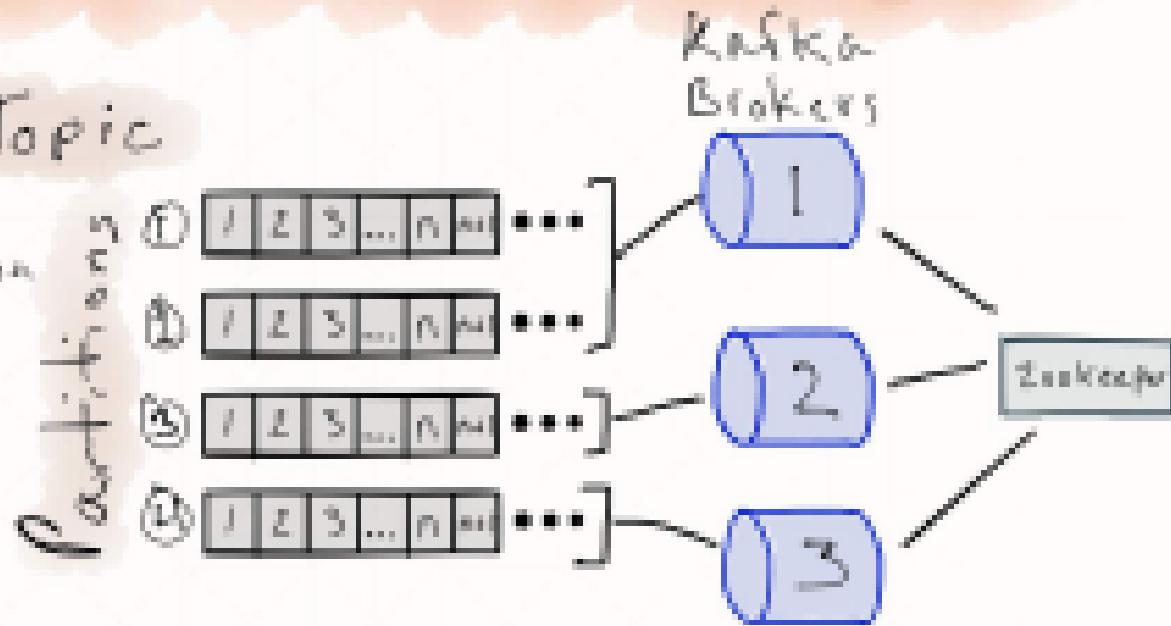
Zookeeper Topic

used for:

- Leader election
- partitions

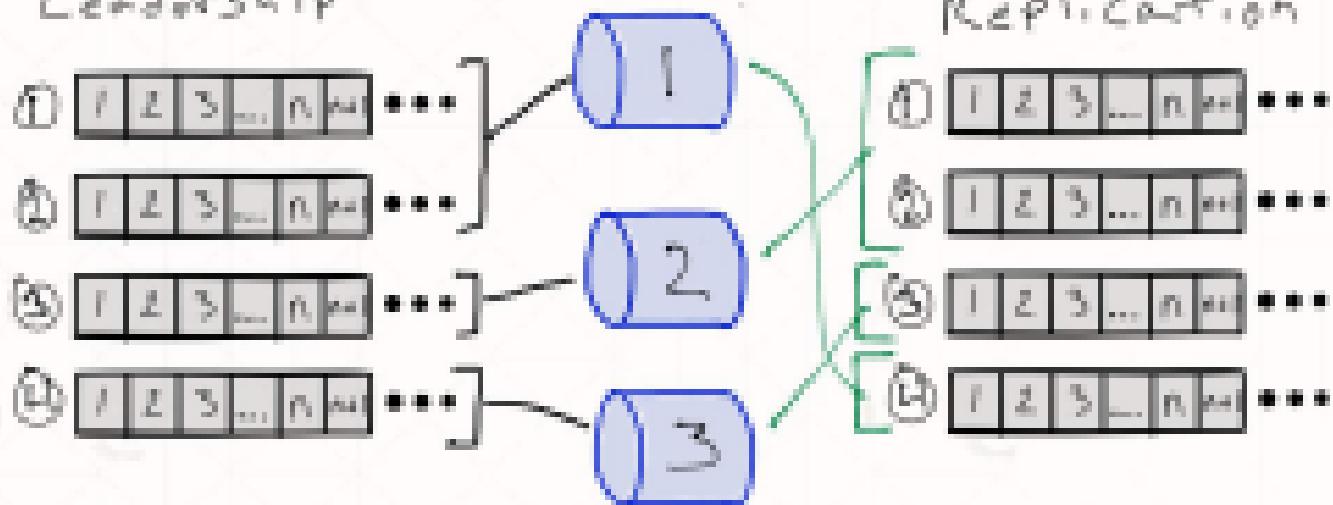
- Configuration

- Broker Topics
- Consumers



Clusters, Replication

Leadership

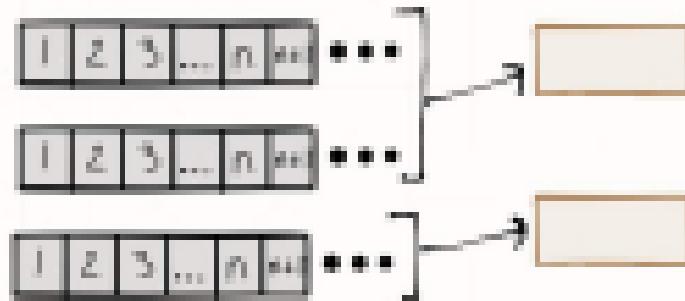


Replication

Consumers

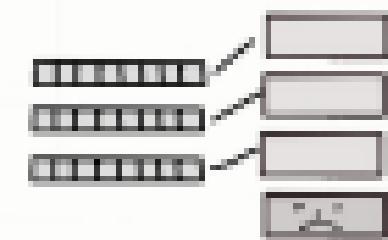
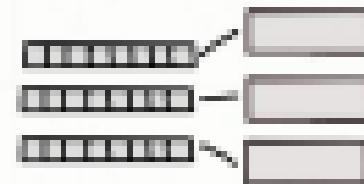
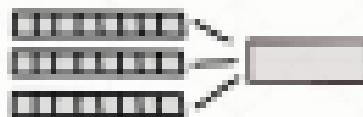
Grouped by group id

- Many consumers per group
- Consumer count <= # of partitions
 - Extra will be idle
- Each consumer in group assigned 1 or more partitions until all are assigned.
- Automatic partition reassignment on consumer count change

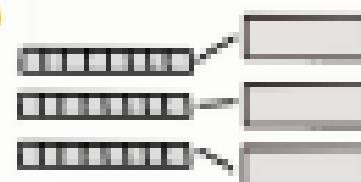
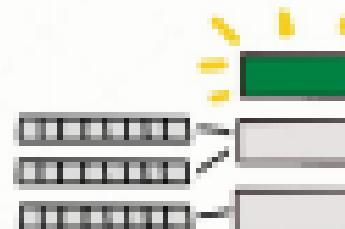
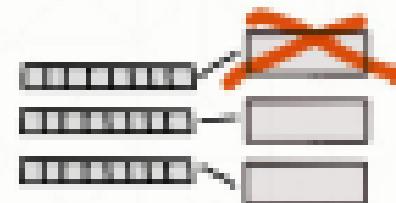
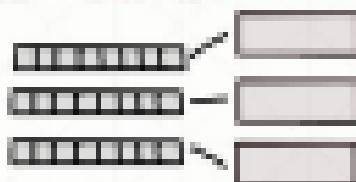


Consumers

Partitions to consumers



fail over



Consumers

Data is never
"consumed" if
remains on the
broker according
to retention policy

Consumer group

Offset

... after next
commit



Producers

Send unkeyed messages

- Same key always goes to same partition

Configurable delivery guarantees

- ACK 0 - no guarantee
- ACK all - all partitions received
- Bulk ACK - no replications required

Attempts to batch for increased throughput
+ configurable

Key

A B C

1	2	3	...	n	...
---	---	---	-----	---	-----

X Y Z

1	2	3	...	n	...
---	---	---	-----	---	-----

1 2 3

1	2	3	...	n	...
---	---	---	-----	---	-----

A B C

1	2	3	...	n	...
---	---	---	-----	---	-----

A B C

1	2	3	...	n	...
---	---	---	-----	---	-----

X Y Z

1	2	3	...	n	...
---	---	---	-----	---	-----

O P Q

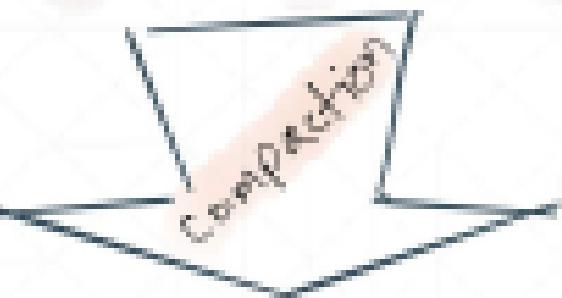
Log Compaction

- Based on keyed messages
- Always retains last message for a key
- Empty payload removes key

use as a
data store

Log Compaction

A * m A A X T



delete →
A

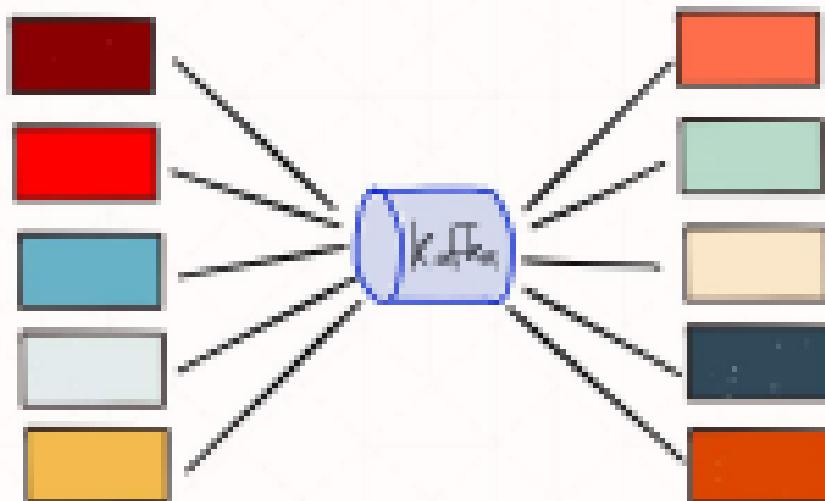
M A C S B M A

the delete

← Remaining in log for
a while, for trailing
consumers

Detangling

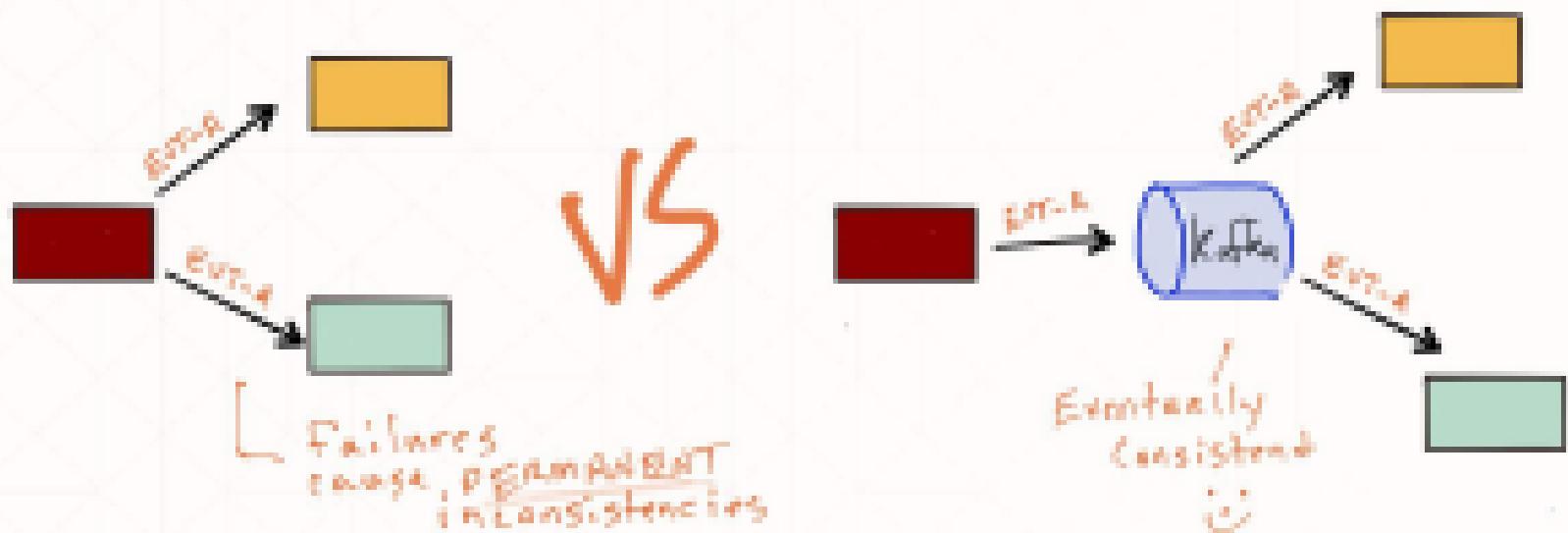
Replace P2P with an event hub



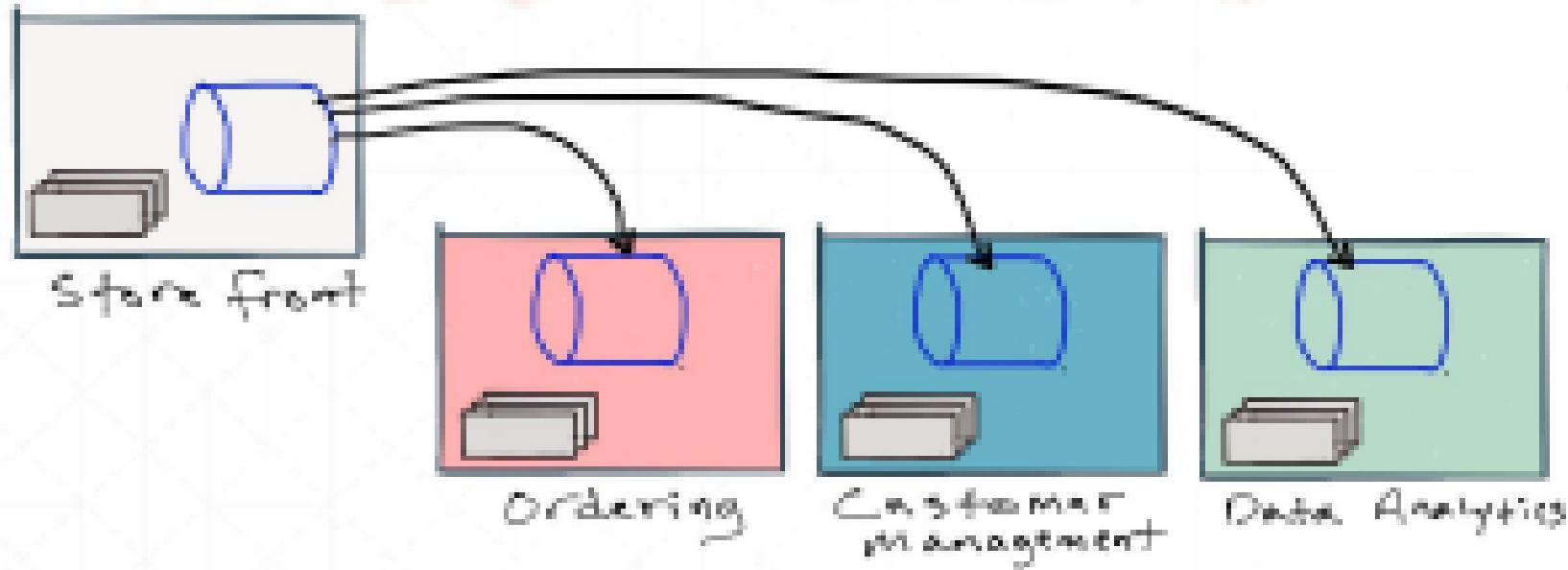
New services
can consume the
data they need.

Create custom
DBs and caches
suited to their
domain.

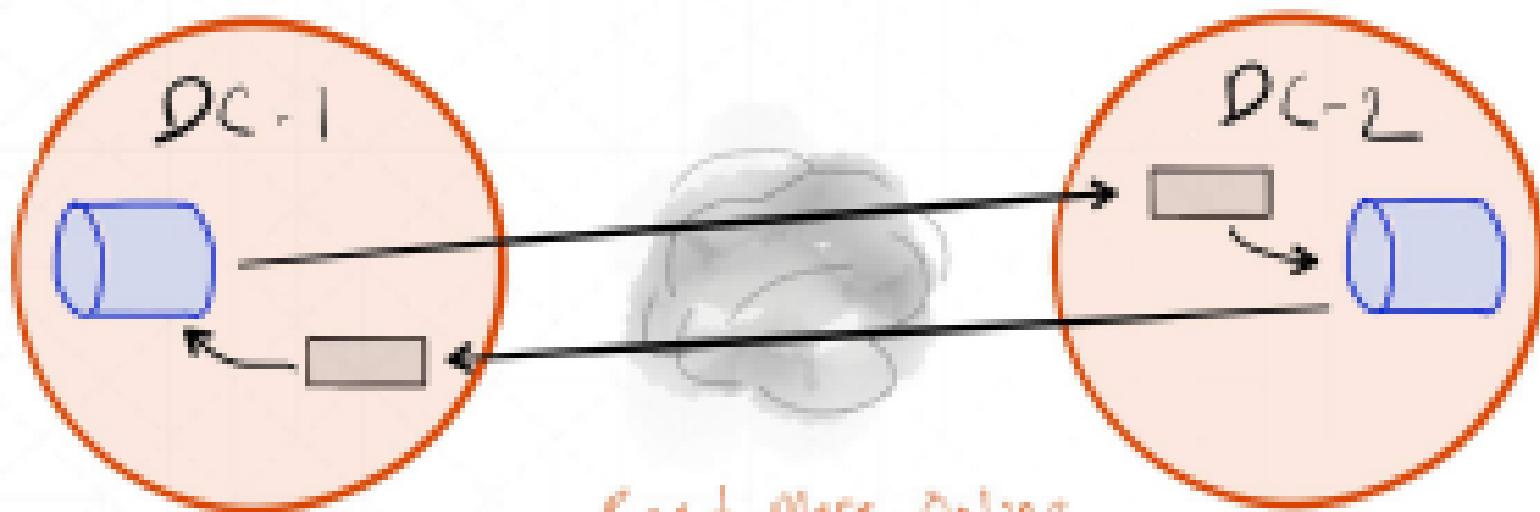
Detangling



System Integrations



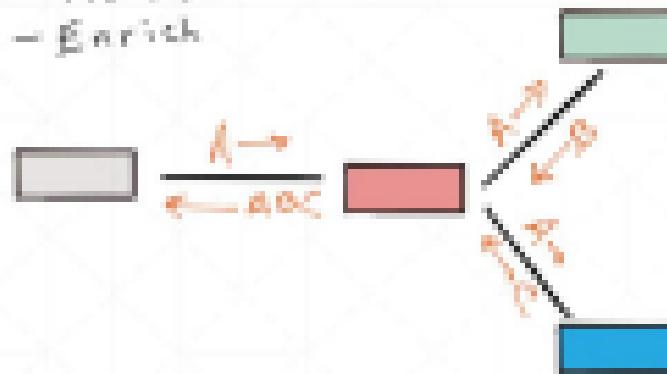
DC Replication



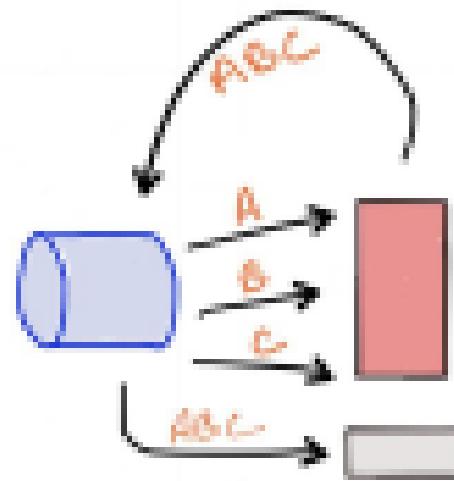
Lead More Online

Stream Processing

- Many services are essentially stream processors
- Transform
- Enrich



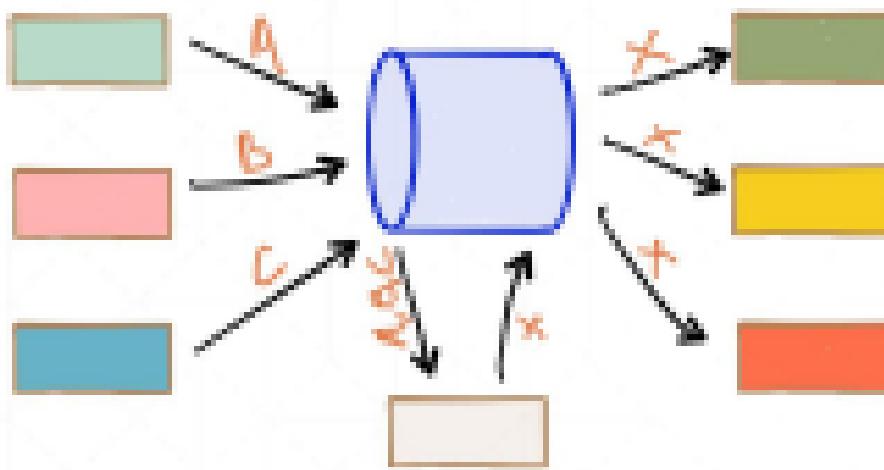
VS



Stream Processing

Transform
messages into
common formats
centrally.

Speed up
Integrations



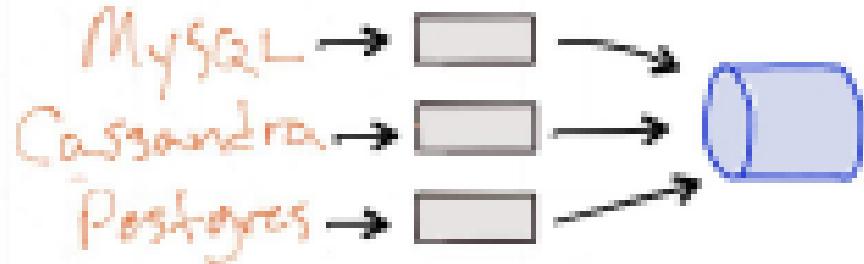
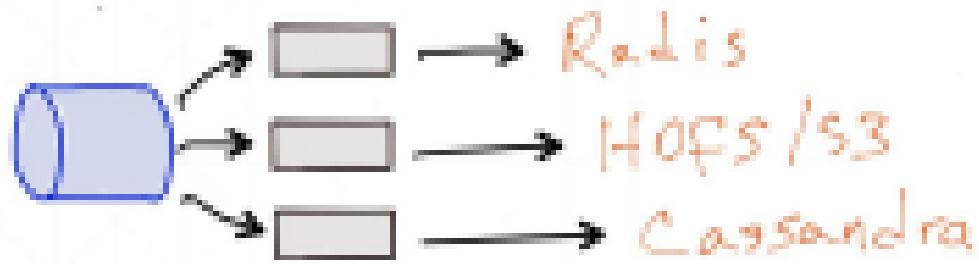
Stream Export/Import

— Export into stores

- Databases
- Caches
- Object Stores

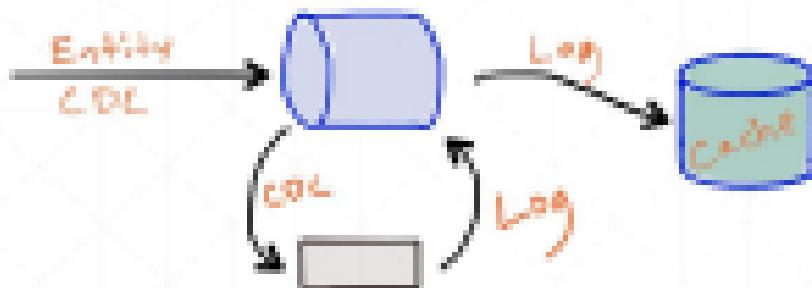
— Import from stores

- Databases
 - MySQL
 - PostgreSQL
 - etc.
- Filesystems



Streams to Caches

— Compacted log (all data)



- Cache always up to date
- Load new cache by reading from start of Log
- Same for new DB (no more schema migration)
- Small datasets can be in Process Memory

Message Replay

- Great when things go sideways
- Can replay all retained messages
Or from any offset



fix and Replay



Things to keep in mind

- Data formats
 - smaller better (bandwidth)
 - well defined (schemas)
 - Managed evolution
- Consumers
 - ; competent

Things to keep in mind

- Message guarantees
 - at most once
 - at least one
 - exactly once
 - requires consumer integration

you will
you of
oops

Things to keep in mind

- Dead letters / Retry
 - Don't skip messages
 - poison Bad Data
 - fix broken services
 - replay if necessary

Things to Keep in Mind

- Kafka / Zookeeper configuration
 - Kafka can be CP or AP (CAP)
 - over provision partitions
(Don't go nuts!)
 - repartitioning breaks ordering temporarily

Things to keep in mind

Monitoring is critical!

- Brokers
- Zoo keeper
- Consumers (lag)
 - Barrow can help

Questions?

john.holland@objectpartners.com
@JohnHolland9

