



Real Time Data Analytics @ Uber

Ankur Bansal

November 14, 2016

About Me

- Sr. Software Engineer, Streaming Team @ Uber
 - Streaming team supports platform for real time data analytics: Kafka, Samza, Flink, Pinot.. and plenty more
 - Focused on scaling Kafka at Uber's pace
- Staff software Engineer @ Ebay
 - Build & scale Ebay's cloud using openstack
- Apache Kylin: Committer, Emeritus PMC

Agenda

- Real time Use Cases
- Kafka Infrastructure Deep Dive
- Our own Development:
 - Rest Proxy & Clients
 - Local Agent
 - uReplicator (Mirrormaker)
 - Chaperone (Auditing)
- Operations/Tooling

Important Use Cases

Real-time Price Surging



Rider eyeballs



Open car information

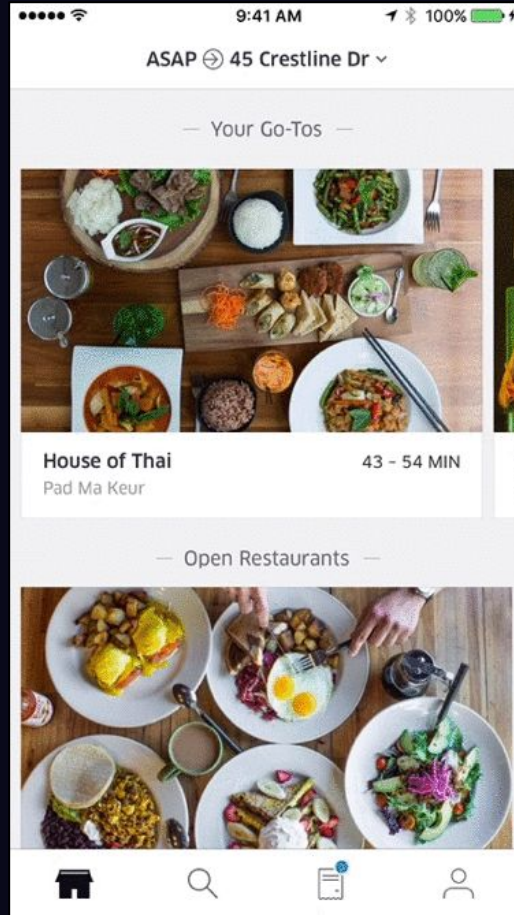
KAFKA

**Stream
Processing**



**SURGE
MULTIPLIERS**

Real-time Machine Learning - UberEats ETD

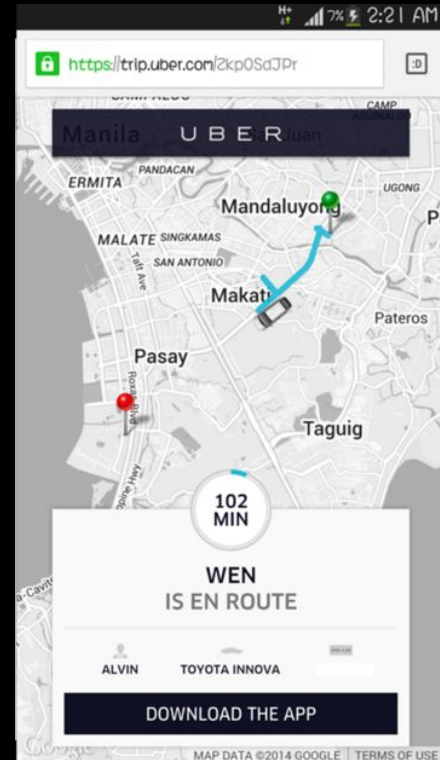


Experimentation Platform



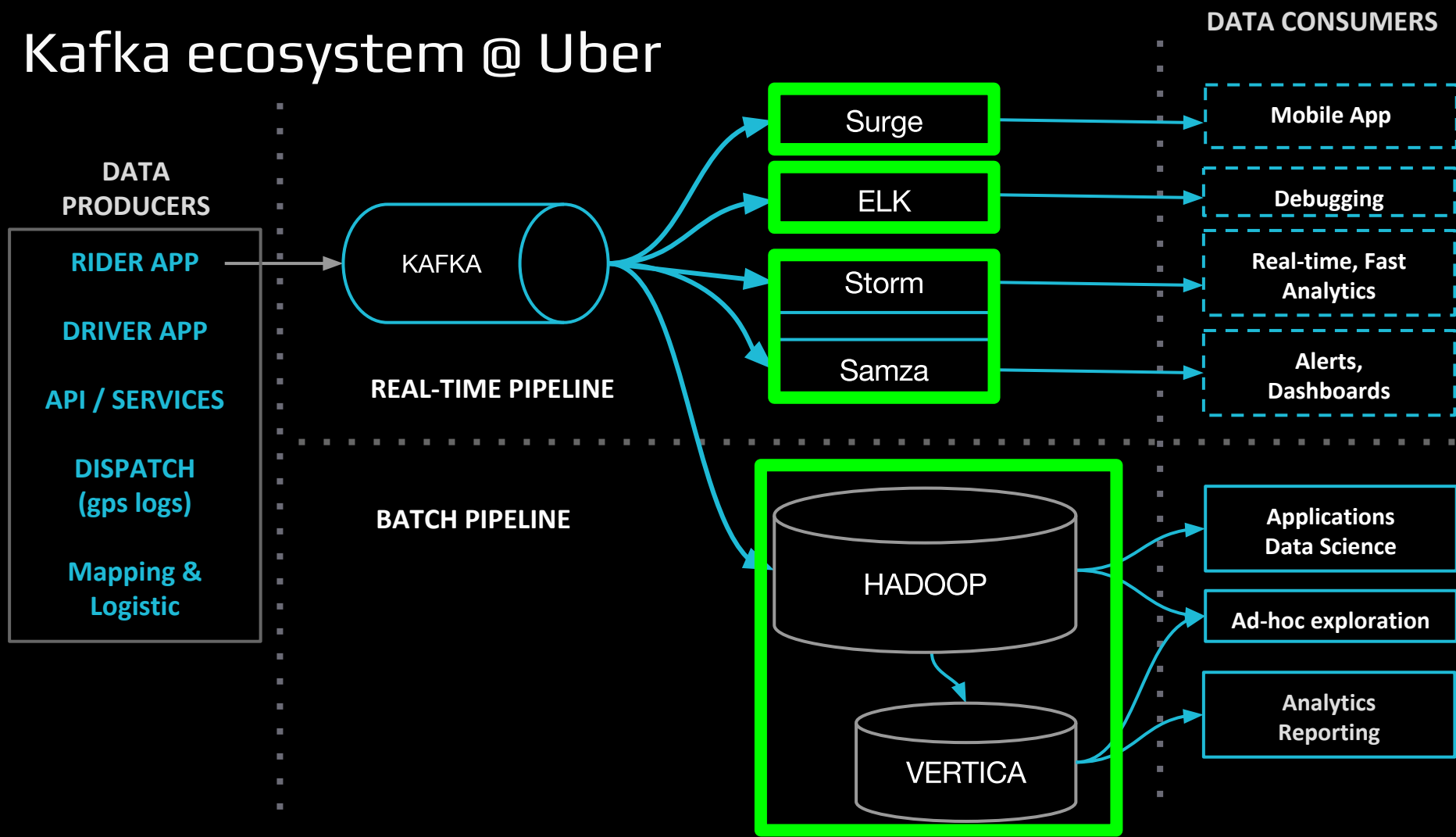
- Fraud detection
- Share my ETA

And many more ...



Apache Kafka is Uber's Lifeline

Kafka ecosystem @ Uber



Kafka cluster stats

100s of billion

Messages/day

100s TB

bytes/day

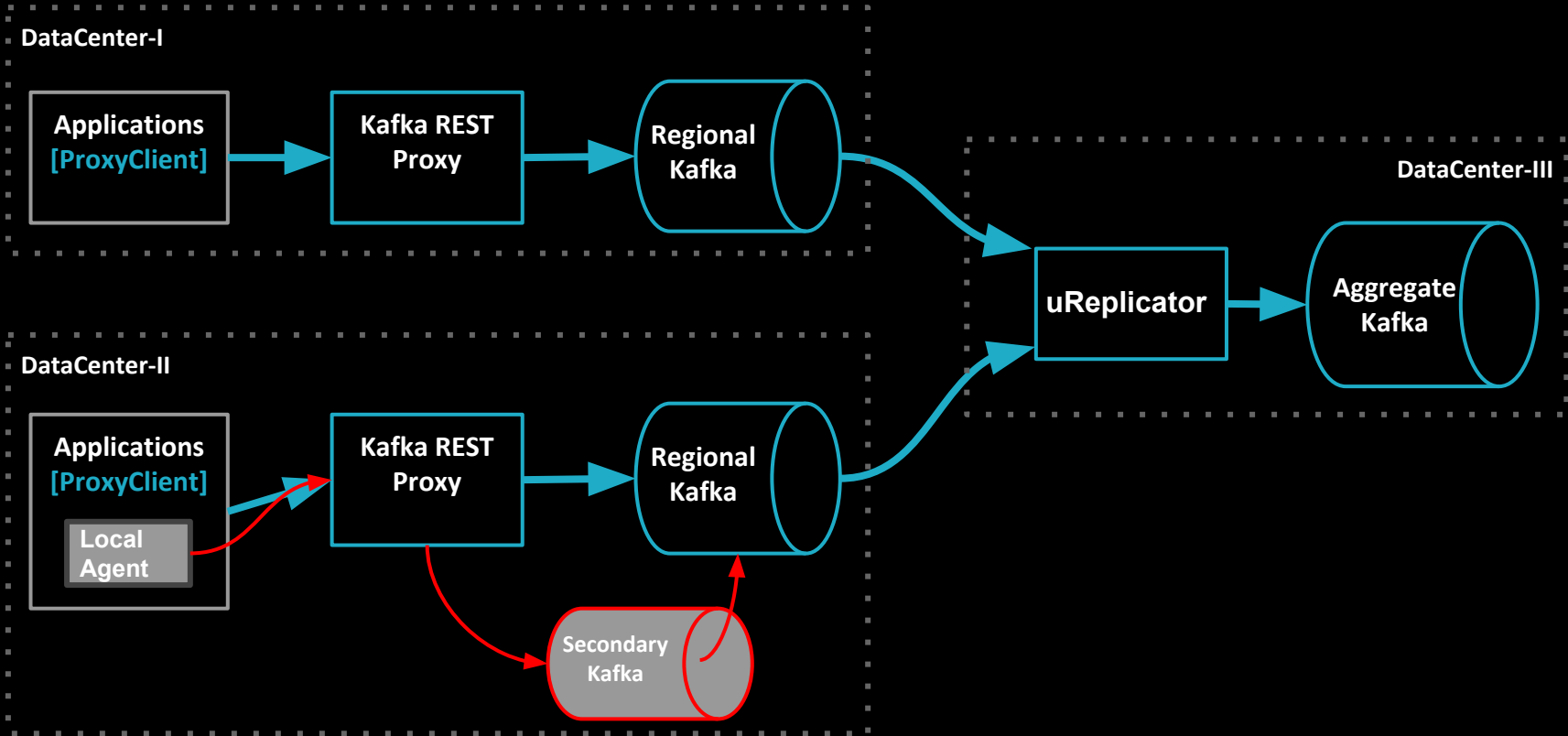
Multiple data centers

Kafka Infrastructure Deep Dive

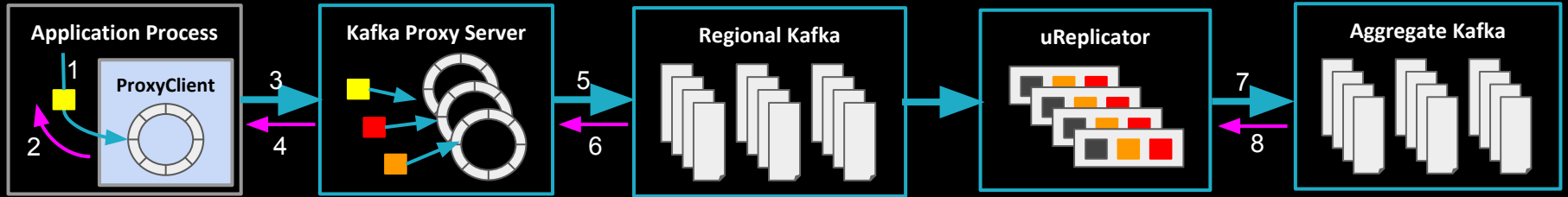
Requirements

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

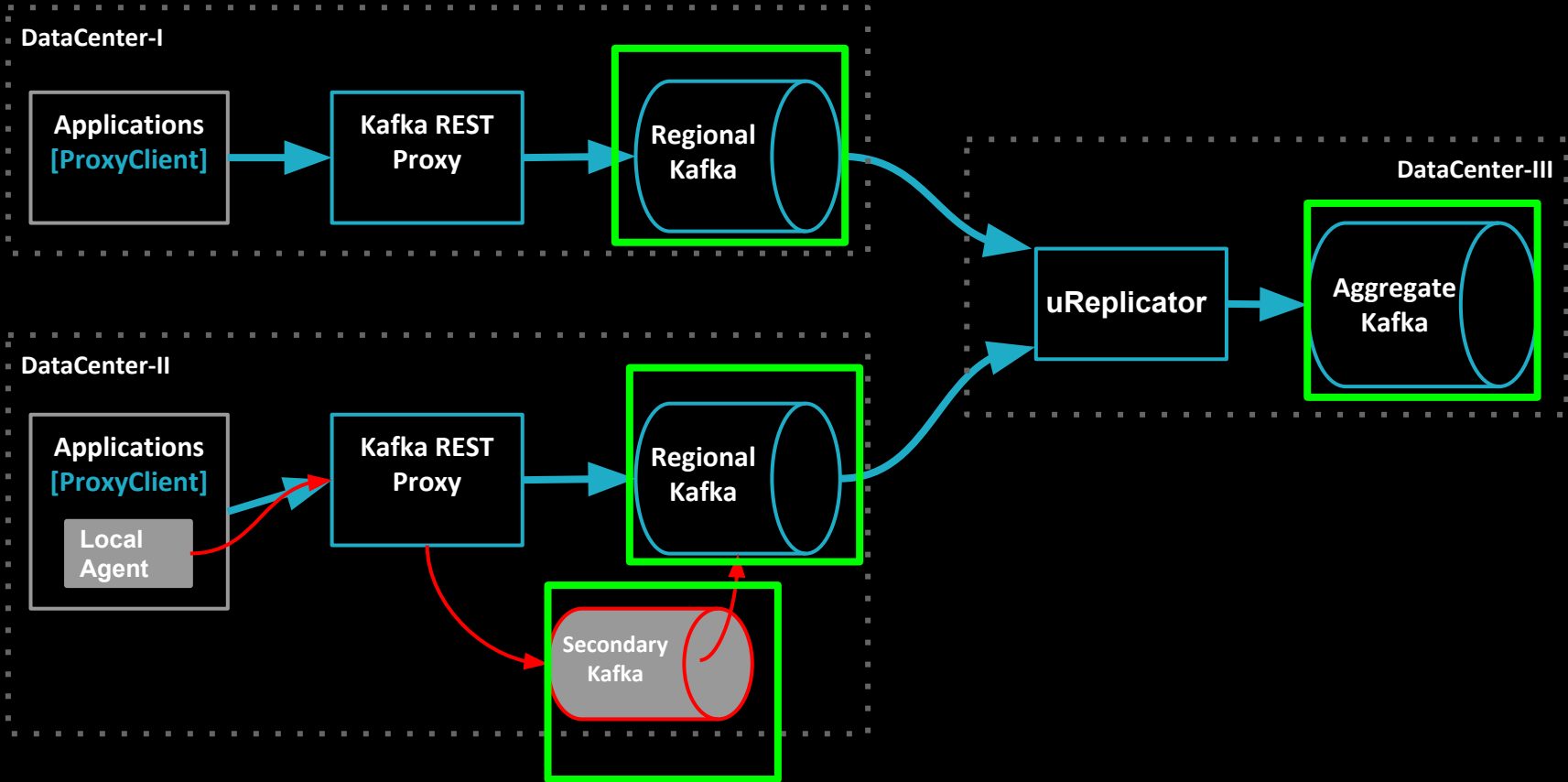
Kafka Pipeline



Kafka Pipeline: Data Flow



Kafka Clusters



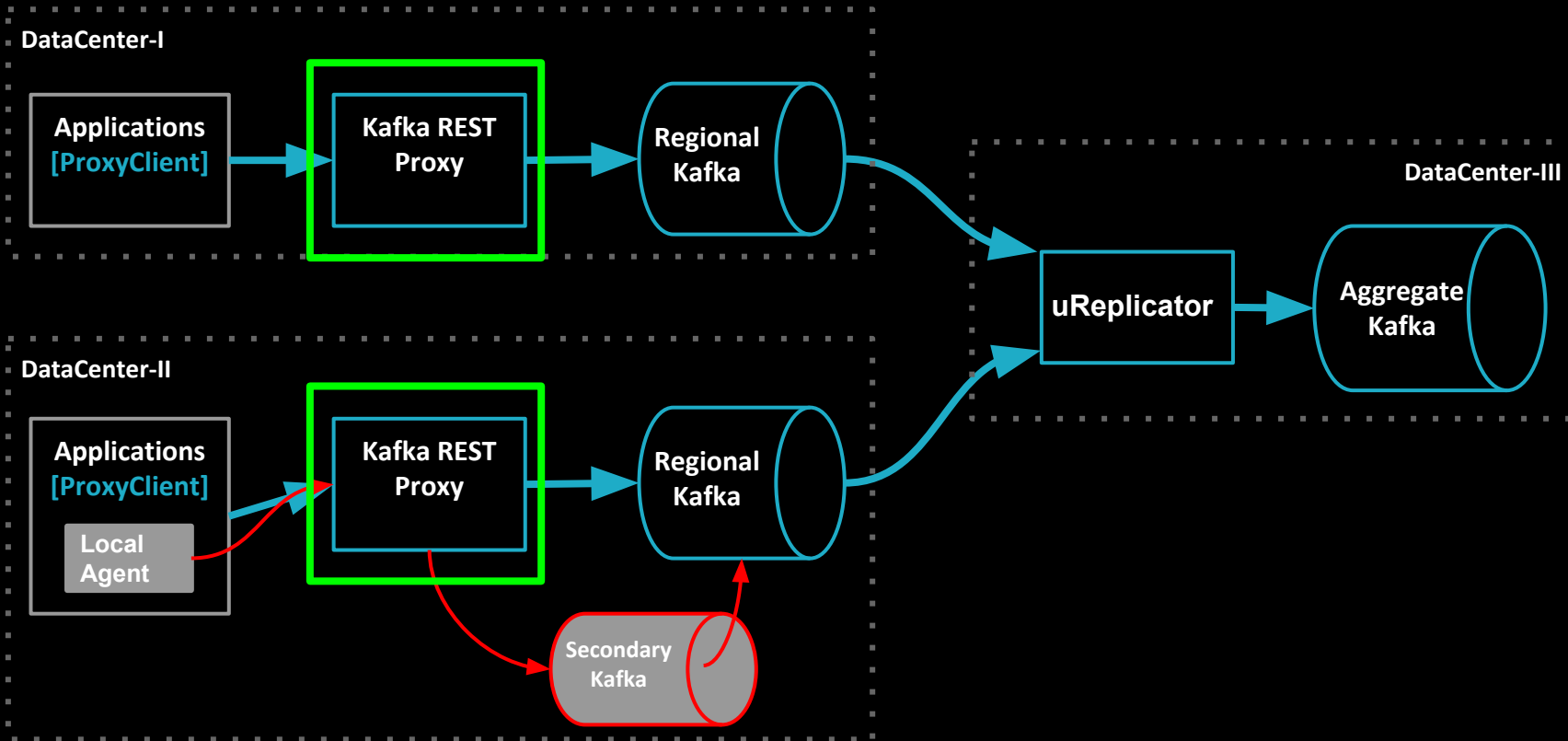
Kafka Clusters

- Use case based clusters
 - Data (async, reliable)
 - Logging (High throughput)
 - Time Sensitive (Low Latency e.g. Surge, Push notifications)
 - High Value Data (At-least once, Sync e.g. Payments)
- Secondary cluster as fallback
- Aggregate clusters for all data topics.

Kafka Clusters

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

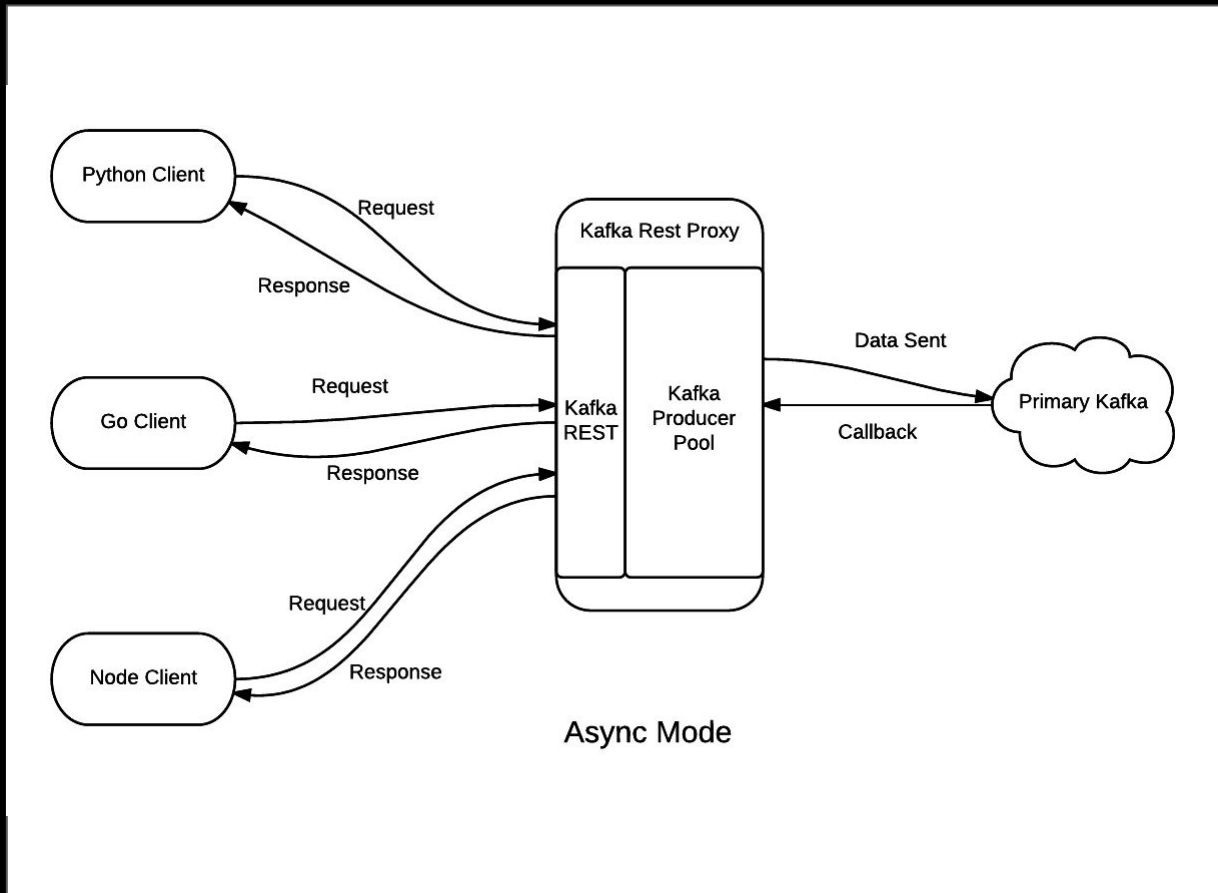
Kafka Rest Proxy



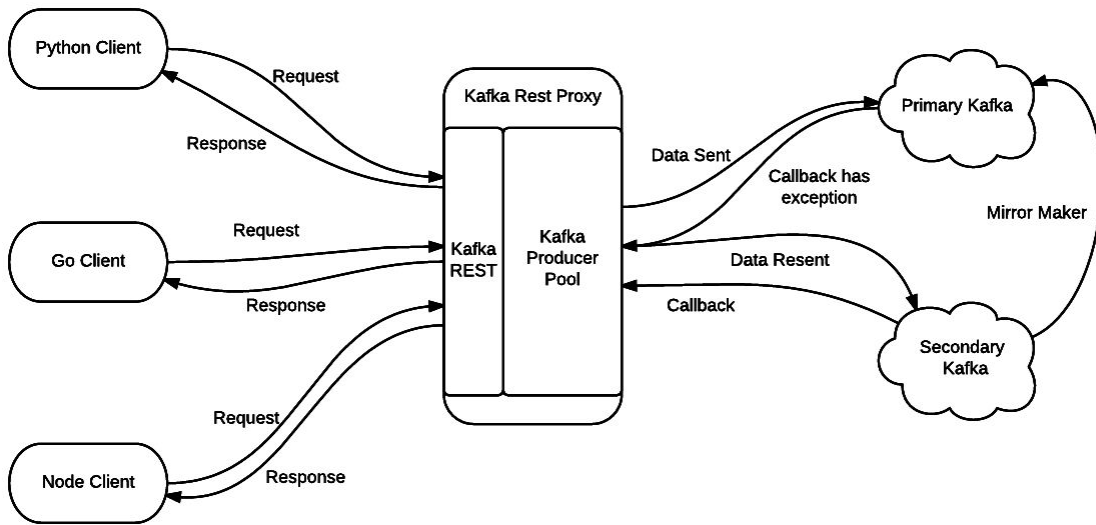
Why Kafka Rest Proxy ?

- Simplified Client API
- Multi-lang support (Java, NodeJs, Python, Golang)
- Decouple client from Kafka broker
 - Thin clients = operational ease
 - Less connections to Kafka brokers
 - Future kafka upgrade
- Enhanced Reliability
 - Primary & Secondary Kafka Clusters

Kafka Rest Proxy: Internals



Kafka Rest Proxy: Internals



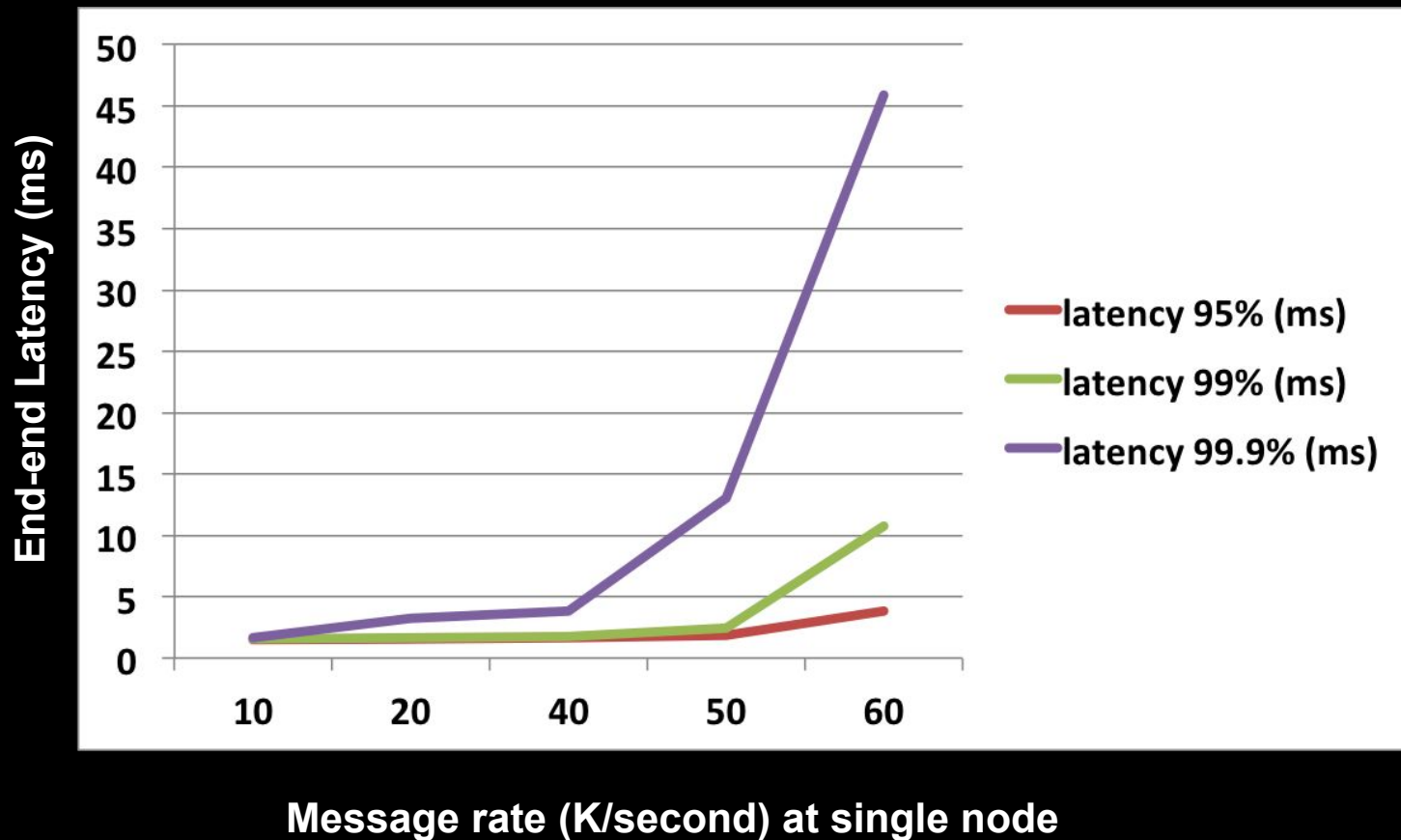
Async Mode

Kafka Rest Proxy: Internals

- Based on Confluent's open sourced Rest Proxy
- Performance enhancements
 - Simple http servlets on jetty instead of Jersey
 - Optimized for binary payloads.
 - Performance increase from 7K* to 45-50K QPS/box
- Caching of topic metadata.
- Reliability improvements*
 - Support for Fallback cluster
 - Support for multiple Producers (SLA based segregation)
- Plan to contribute back to community

*Based on benchmarking & analysis done in Jun '2015

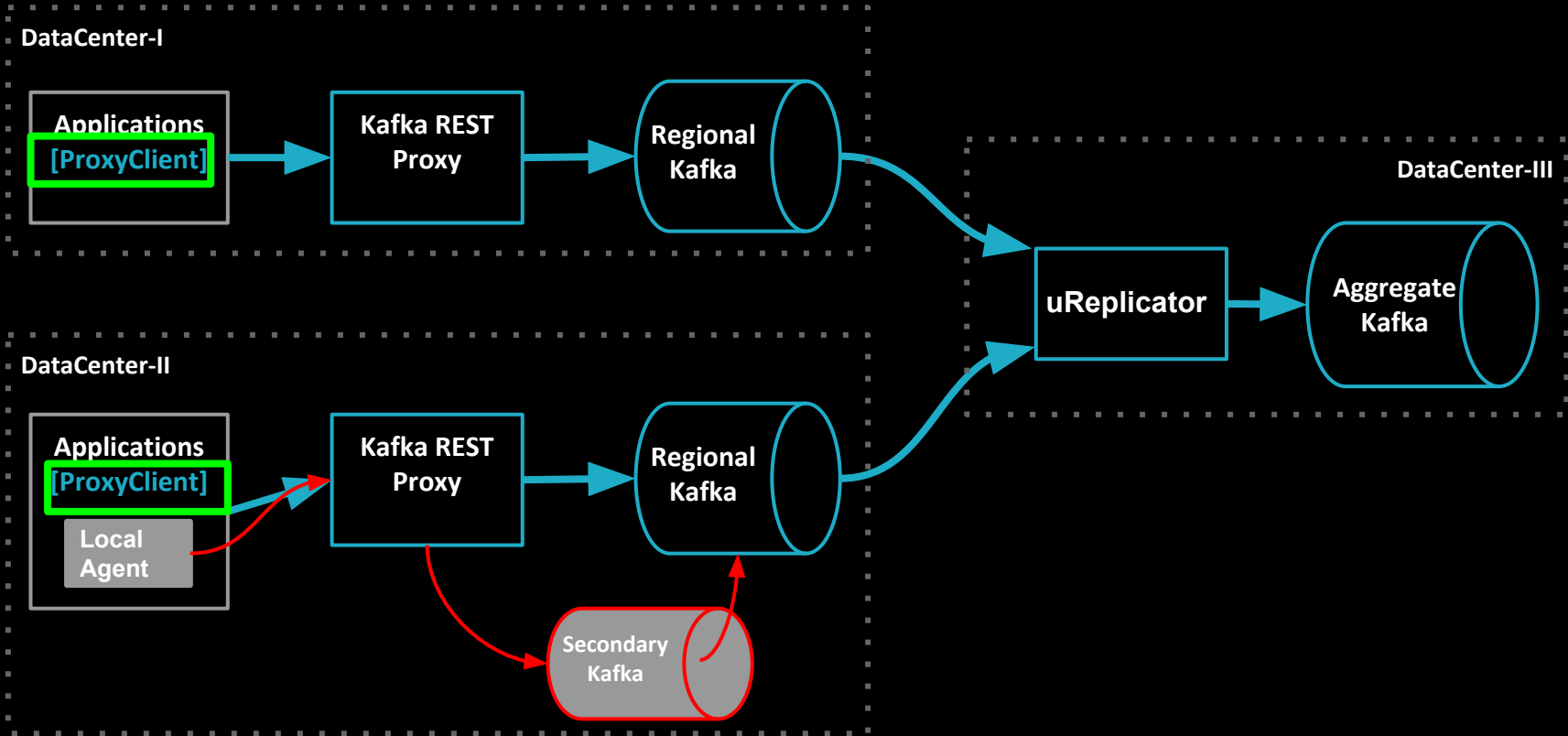
Rest Proxy: performance (1 box)



Kafka Clusters + Rest Proxy

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

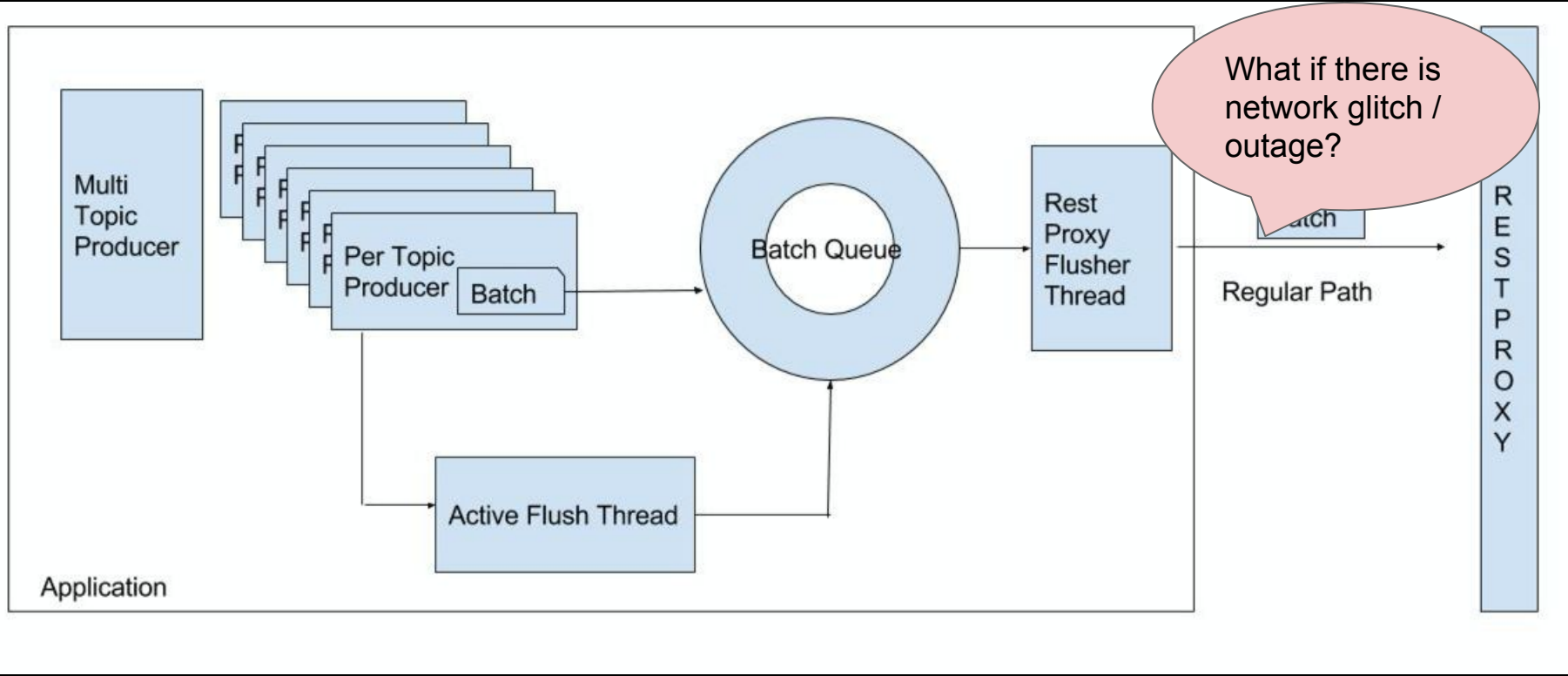
Kafka Clients



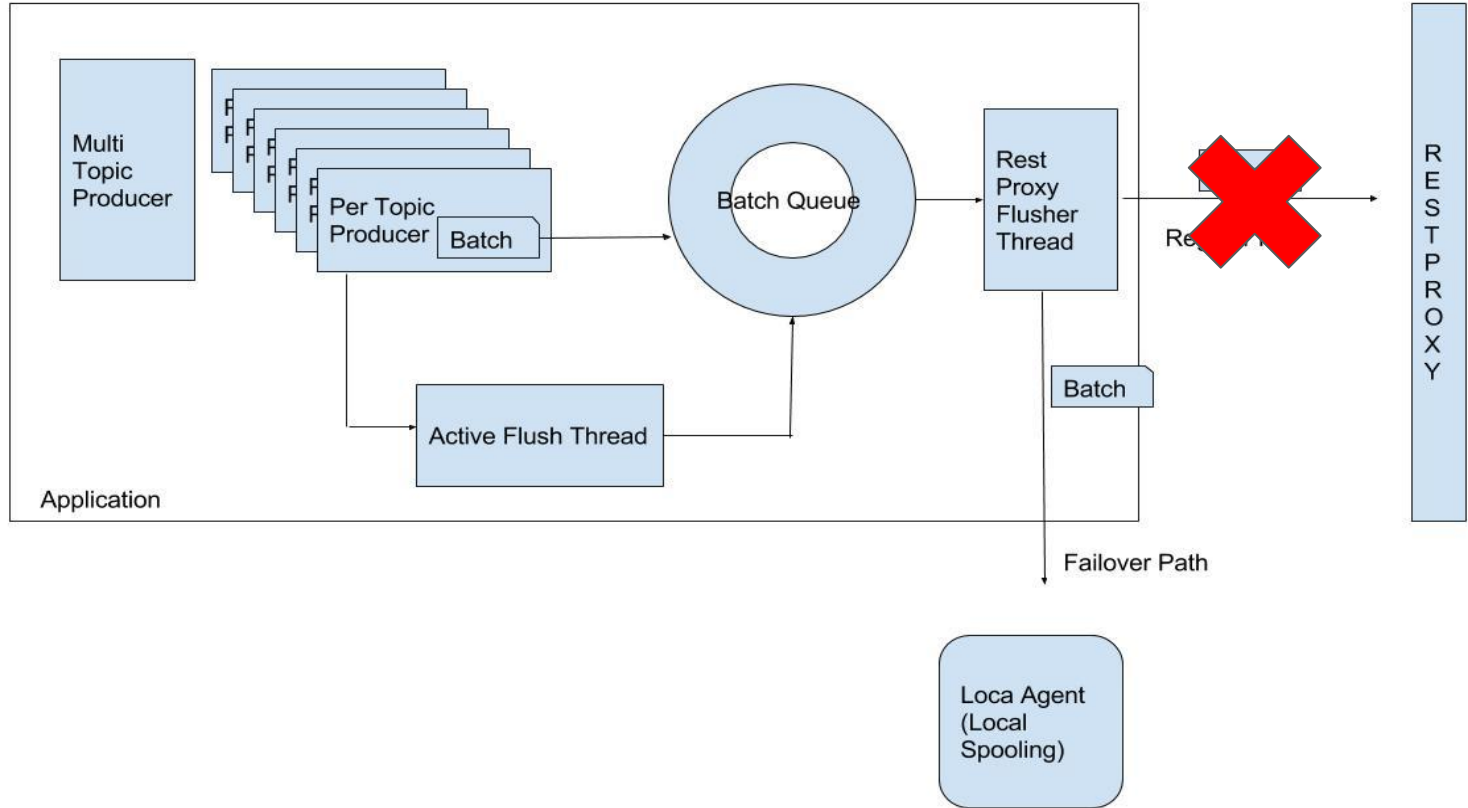
Client Libraries

- Support for multiple clusters.
- High Throughput
 - Non-blocking, async, batching
 - <1ms produce latency for clients
 - Handles Throttling/BackOff signals from Rest Proxy
- Topic Discovery
 - Discovers the kafka cluster a topic belongs
 - Able to multiplex to different kafka clusters
- Integration with Local Agent for critical data

Client Libraries



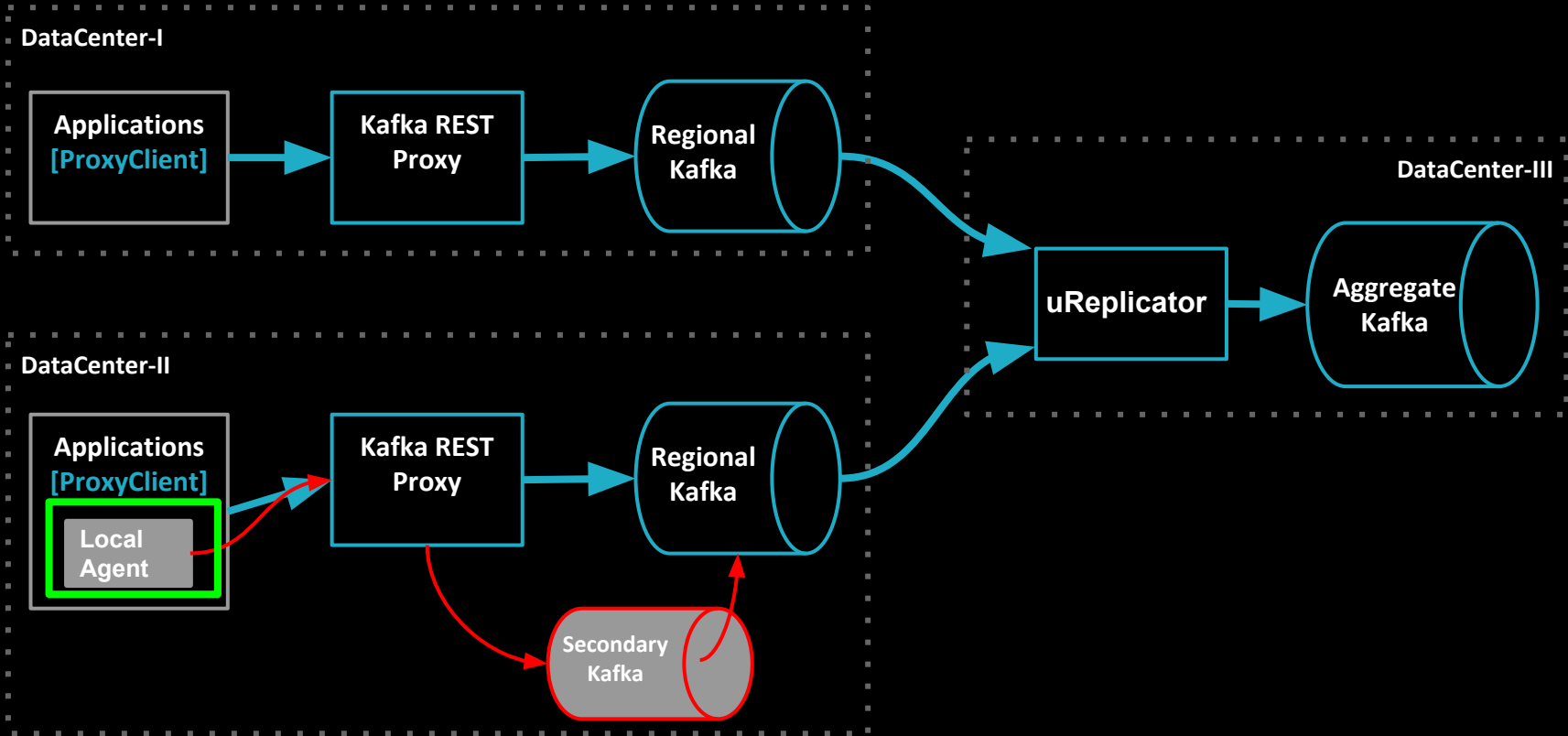
Client Libraries



Kafka Clusters + Rest Proxy + Clients

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

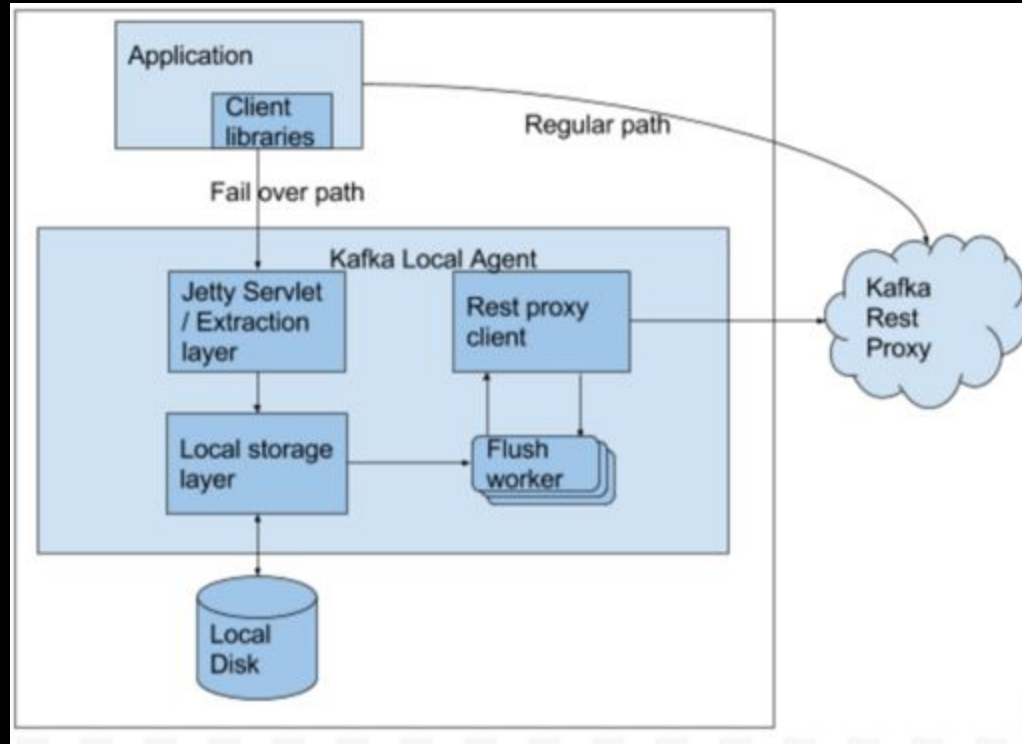
Local Agent



Local Agent

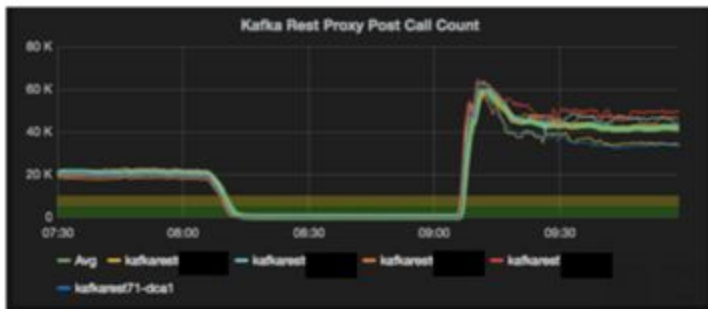
- Local spooling in case of downstream outage/backpressure
- Backfills at the controlled rate to avoid hammering infrastructure recovering from outage
- Implementation:
 - Reuses code from rest-proxy and kafka's log module.
 - Appends all topics to same file for high throughput.

Local Agent Architecture



Local Agent in Action

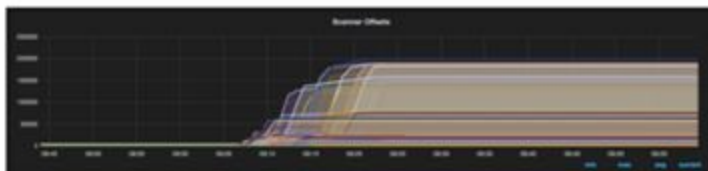
#1: Rest proxy data incoming traffic coming down to 0 around 8:10. Recovery around 9:10 with a spike from client libraries buffers



#2: Traffic starts moving to local agent in outage window



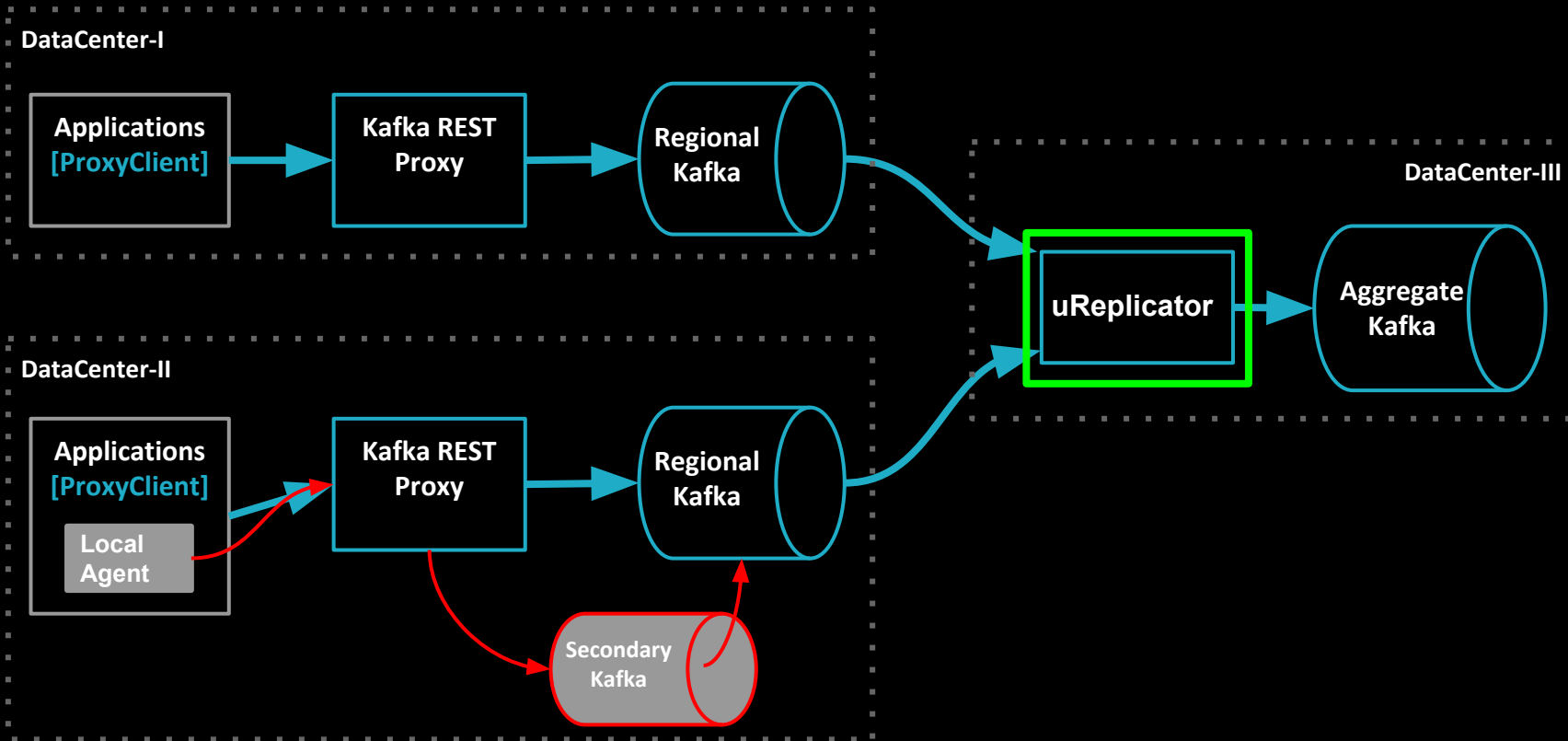
#3: Local agent producing to Rest proxy after recovery



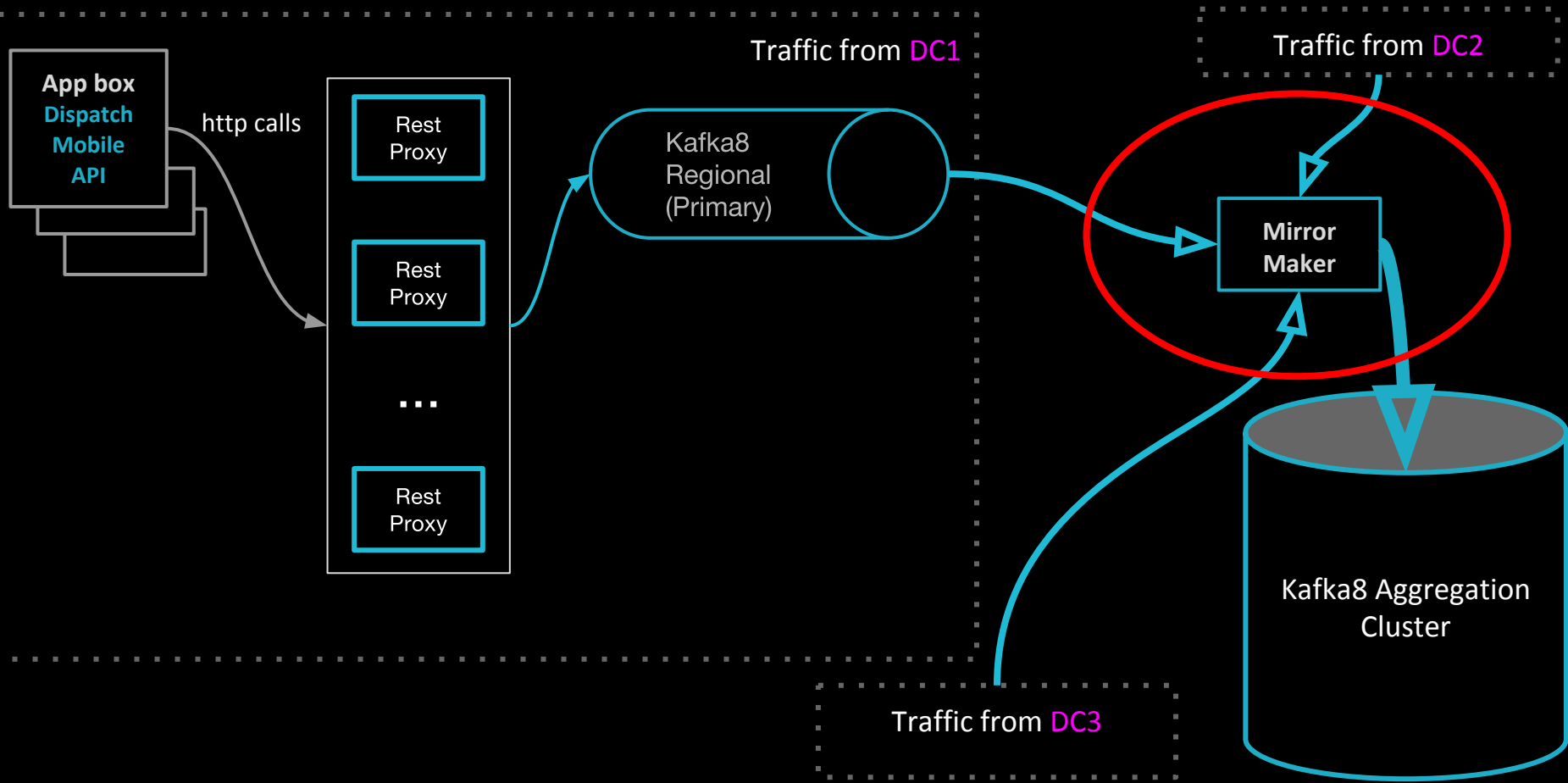
Kafka Clusters + Rest Proxy + Clients + Local Agent

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

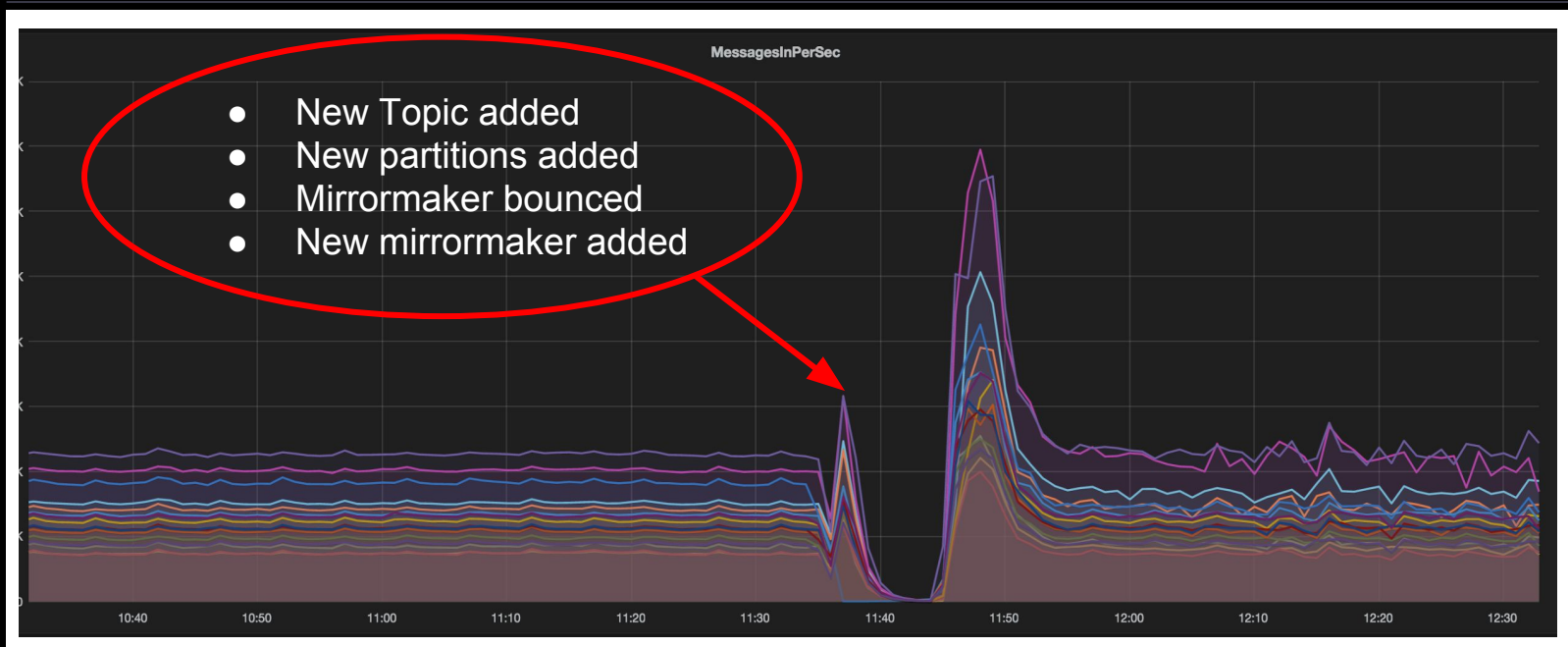
uReplicator



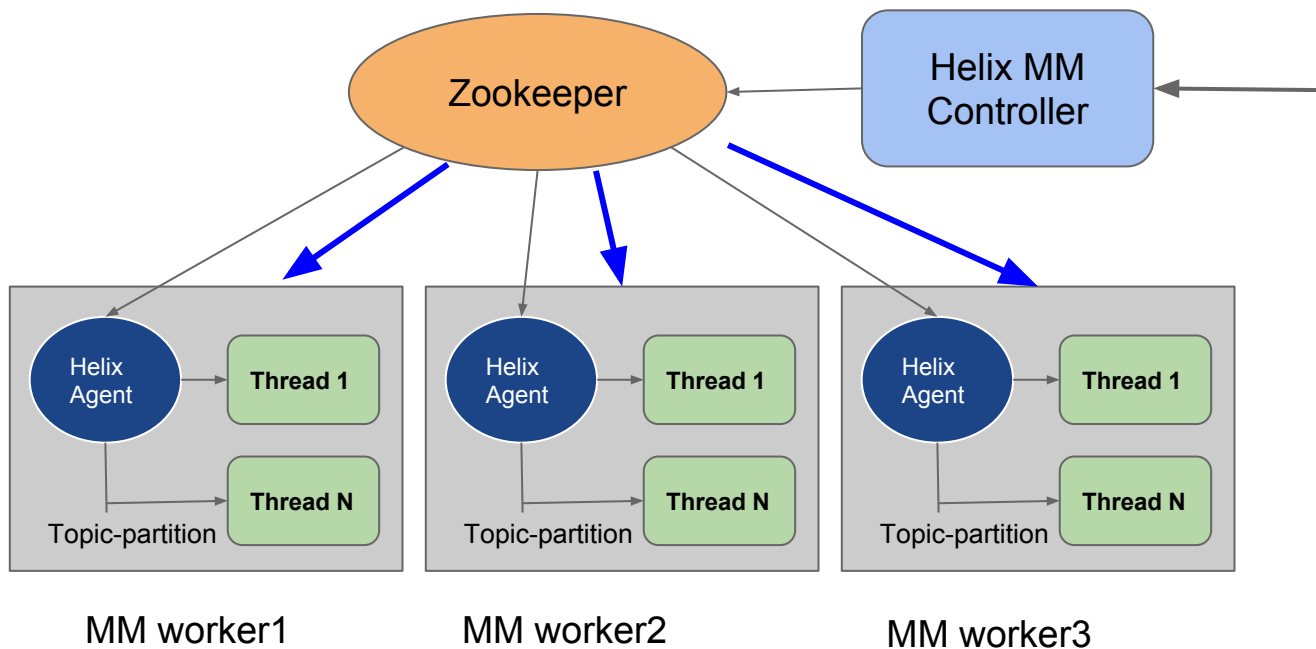
Multi-DC data flow



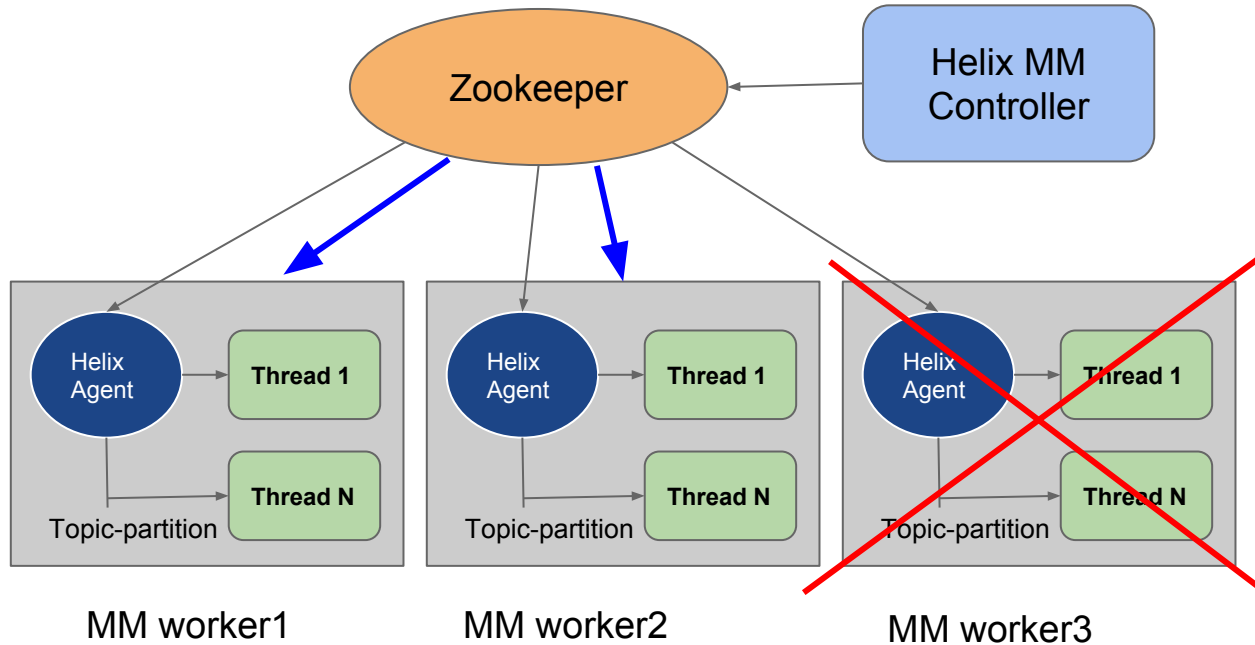
Mirrormaker : existing problems



uReplicator: In-house solution



uReplicator



Kafka Clusters + Rest Proxy + Clients + Local Agent

- Scale to 100s Billions/day → 1 Trillion/day
- High Throughput (Scale: 100s TB → PB)
- Low Latency for most use cases(<5ms)
- Reliability - 99.99% (#Msgs Available /#Msgs Produced)
- Multi-Language Support
- Tens of thousands of simultaneous clients.
- Reliable data replication across DC

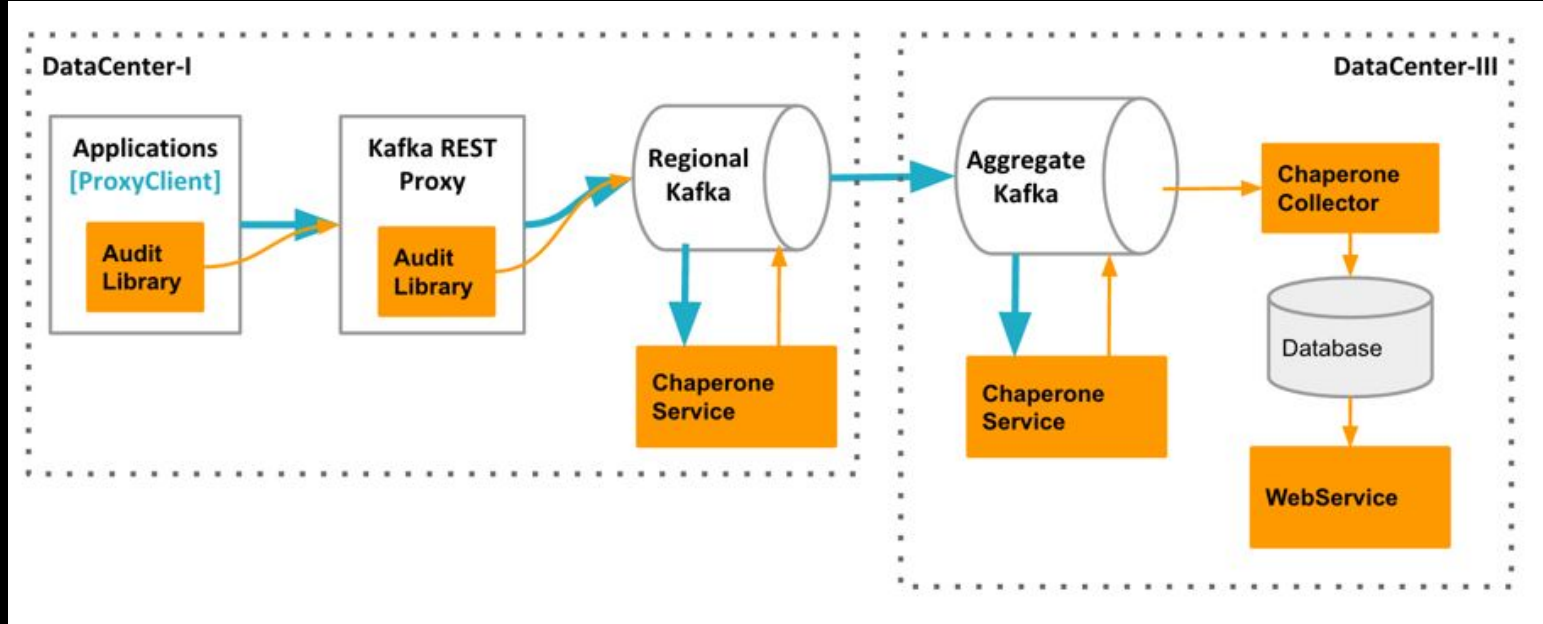
uReplicator

- Running in production for 1+ year
- Open sourced: <https://github.com/uber/uReplicator>
- Blog: <https://eng.uber.com/ureplicator/>

Chaperone - E2E Auditing



Chaperone Architecture



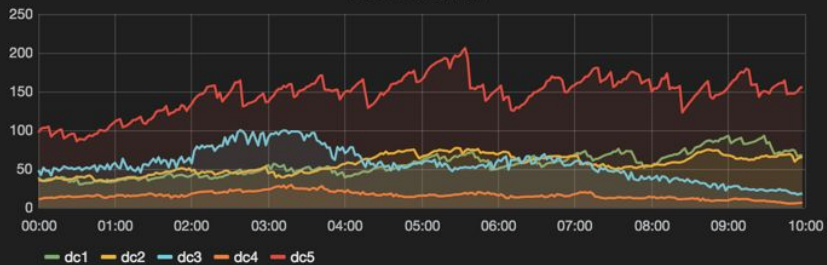
Chaperone : Track counts



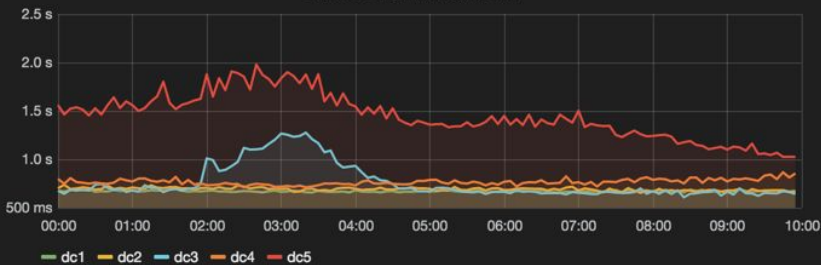
Chaperone : Track Latency

MSGRATE AND LATENCY

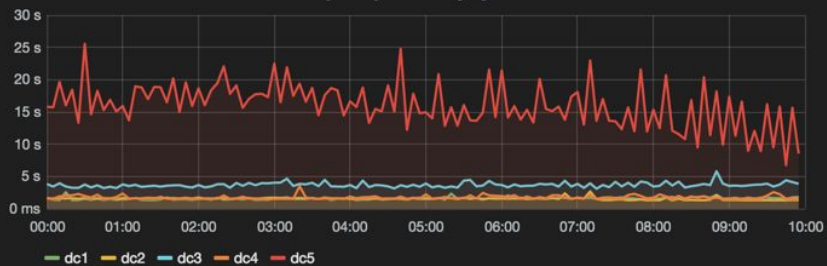
topic01 msg rate



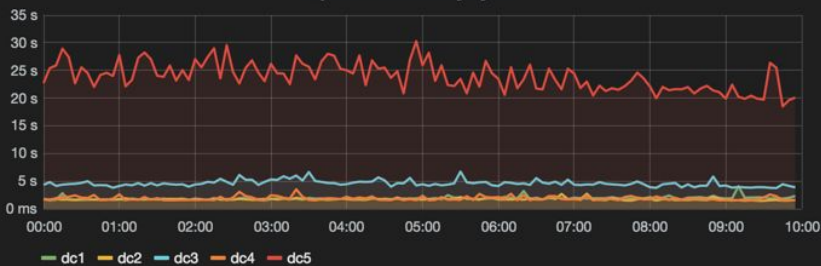
topic01 avg latency by tier



topic01 p99 latency by tier



topic01 max latency by tier



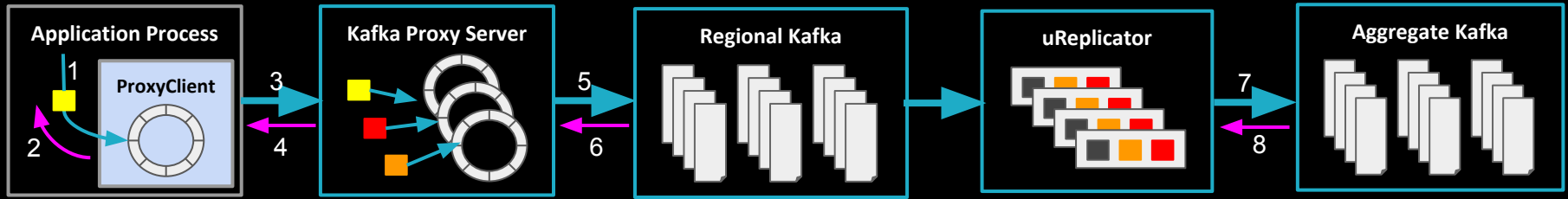
Chaperone

- Running in production for 1+ year
- Planning to open source in ~2 Weeks

At-least Once Kafka



Why do we need it?



- Most of infrastructure tuned for high throughput
 - Batching at each stage
 - Ack before produce (ack'ed != committed)
- Single node failure in any stage leads to data loss
- Need a reliable pipeline for High Value Data e.g. Payments

How did we achieve it?

- Brokers:
 - `min.insync.replicas=2`, can only tolerate one node failure
 - `unclean.leader.election=false`, need to wait until the old leader comes back
- Rest Proxy:
 - Partition Failover
- Improved Operations:
 - Replication throttling, to reduce impact of node bootstrap
 - Prevent catching up nodes to become ISR

Operations/Tooling



Partition Rebalancing

Cluster Broker Summary

Number of brokers: 28
Number of brokers by state:
State: ONLINE Count: 28

Membership bytes per broker:

MAX:	MIN:	DIFF:
AVG:	STD:	STD% of AVG: 9.4949%
P50:	P75:	P90:
P95:	P99:	

Total bytes per broker:

MAX:	MIN:	DIFF:
AVG:	STD:	STD% of AVG: 3.8469%
P50:	P75:	P90:
P95:	P99:	

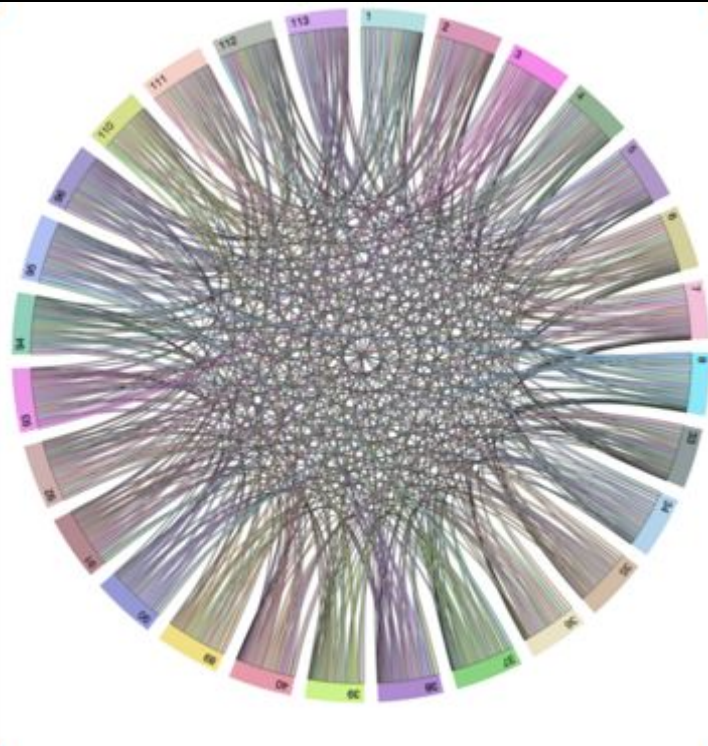
Number of replicas per broker:

MAX:	MIN:	DIFF:
AVG:	STD:	STD% of AVG: 58.7987%
P50:	P75:	P90:
P95:	P99:	

Broker Adjacency (The size of the partitions replicas intersecting any two brokers)

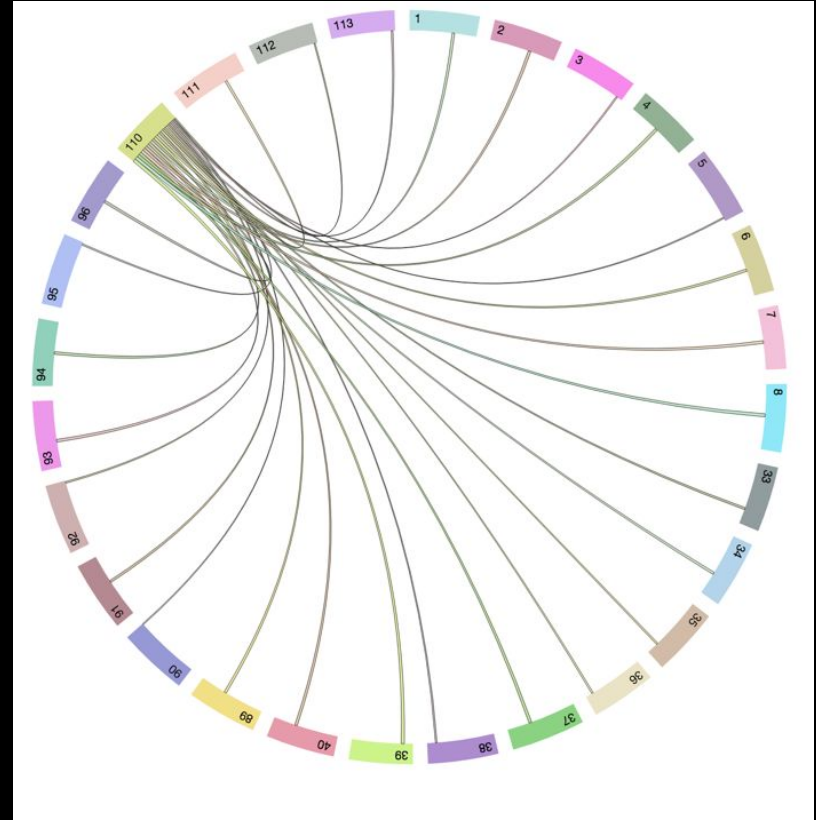
MAX:	MIN:	DIFF:
AVG:	STD:	STD% of AVG: 38.8983%
P50:	P75:	P90:
P95:	P99:	

Broker ID	Disk Bytes	Num Parts	State	Hostname
1		3842	ONLINE	localhost
2		2817	ONLINE	localhost
3		2486	ONLINE	localhost
4		2738	ONLINE	localhost
5		2874	ONLINE	localhost
6		2712	ONLINE	localhost
7		2911	ONLINE	localhost
8		2668	ONLINE	localhost
33		2714	ONLINE	localhost
34		3862	ONLINE	localhost
35		3999	ONLINE	localhost
36		3868	ONLINE	localhost
37		2821	ONLINE	localhost
38		2688	ONLINE	localhost
39		2683	ONLINE	localhost
48		2843	ONLINE	localhost
89		725	ONLINE	localhost
90		553	ONLINE	localhost
91		568	ONLINE	localhost
92		479	ONLINE	localhost
93		951	ONLINE	localhost
94		397	ONLINE	localhost
95		635	ONLINE	localhost
96		454	ONLINE	localhost
118		518	ONLINE	localhost



Partition Rebalancing

- Calculates partition imbalance and inter-broker dependency.
- Generates & Executes Rebalance Plan.
- Rebalance plans are incremental, can be stopped and resumed.
- Currently on-demand, Automated in the future.



XFS vs EXT4



Summary: Scale

- Kafka Brokers:
 - Multiple Clusters per DC
 - Use case based tuning
- Rest Proxy to reduce connections and better batching
- Rest Proxy & Clients
 - Batch everywhere, Async produce
 - Replace Jersey with Jetty
- XFS

Summary: Reliability

- Local Agent
- Secondary Clusters
- Multi Producer support in Rest Proxy
- uReplicator
- Auditing via Chaperone

Future Work

- Open source contribution
 - Chaperone
 - Toolkit
- Data Lineage
- Active Active Kafka
- Chargeback
- Exactly once mirroring via uReplicator

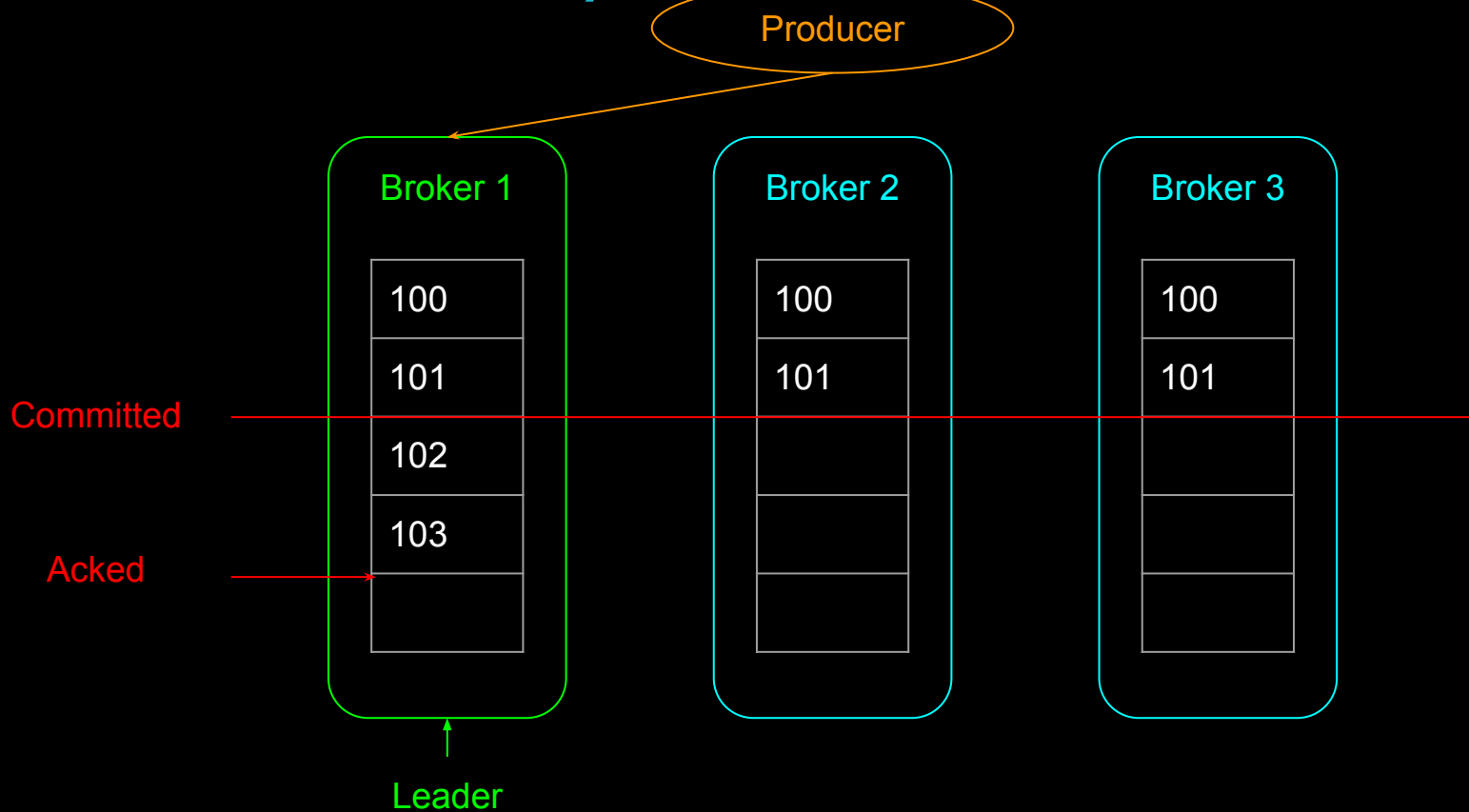
Questions ?

ankur@uber.com

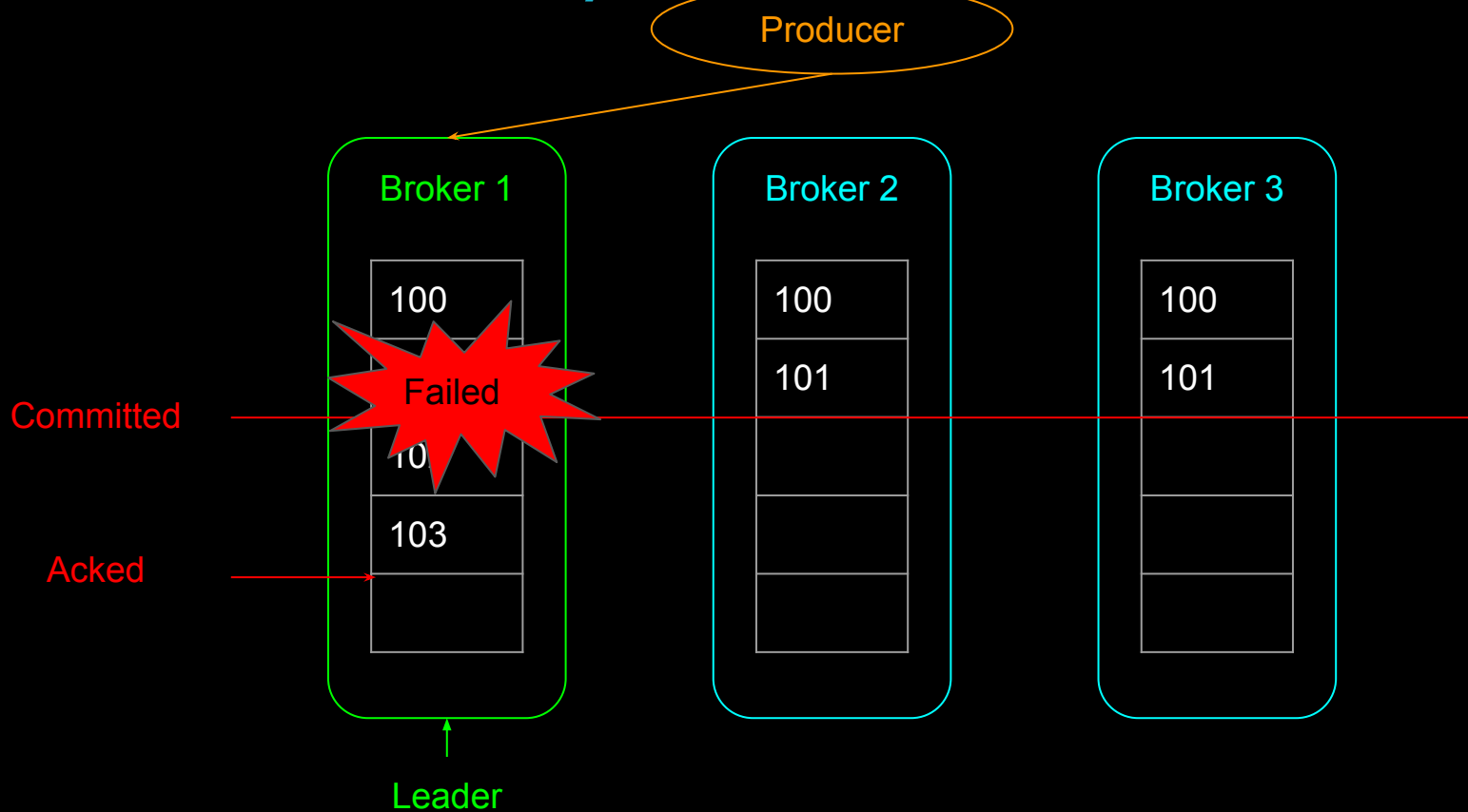
Extra Slides



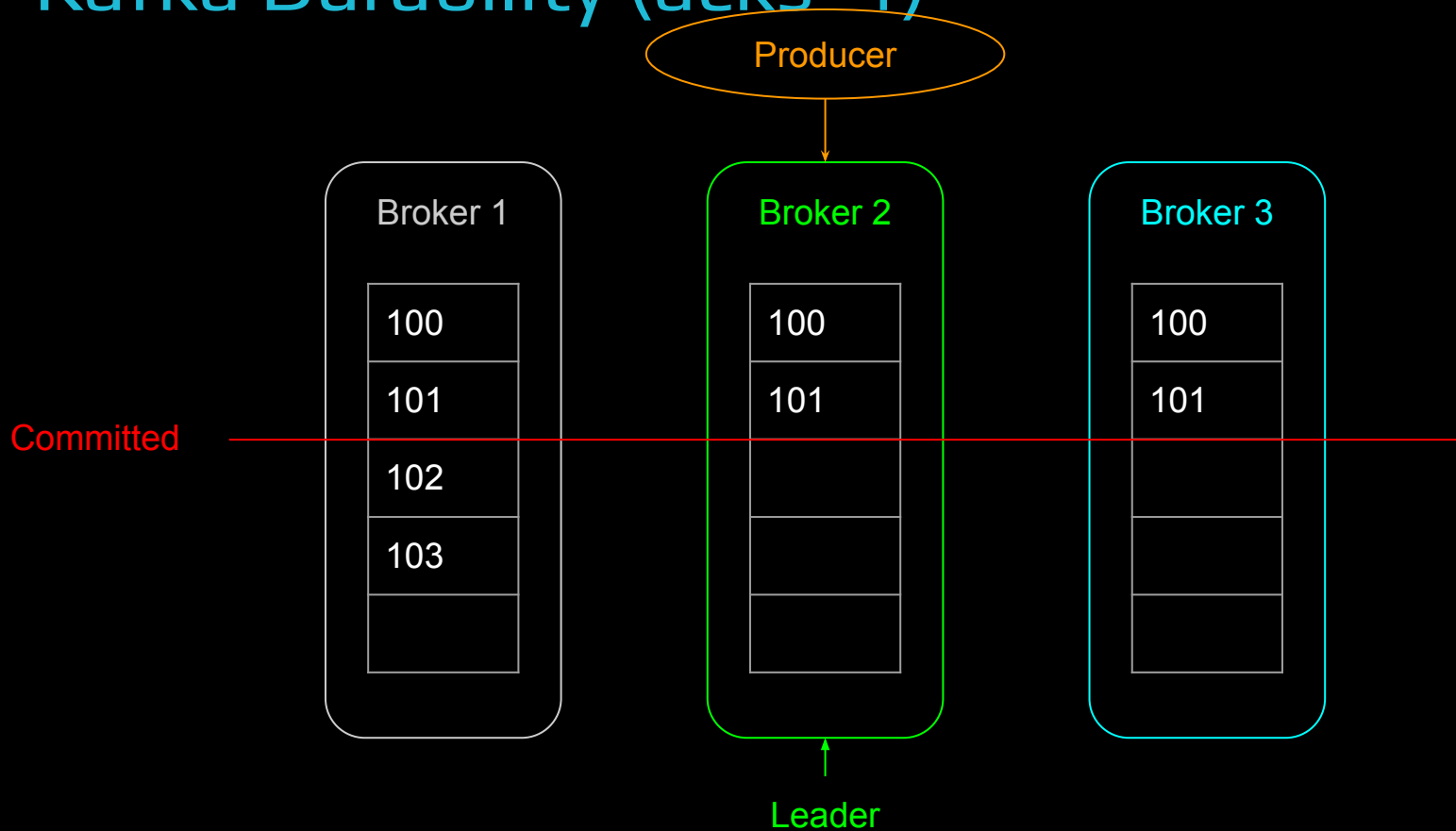
Kafka Durability (acks=1)



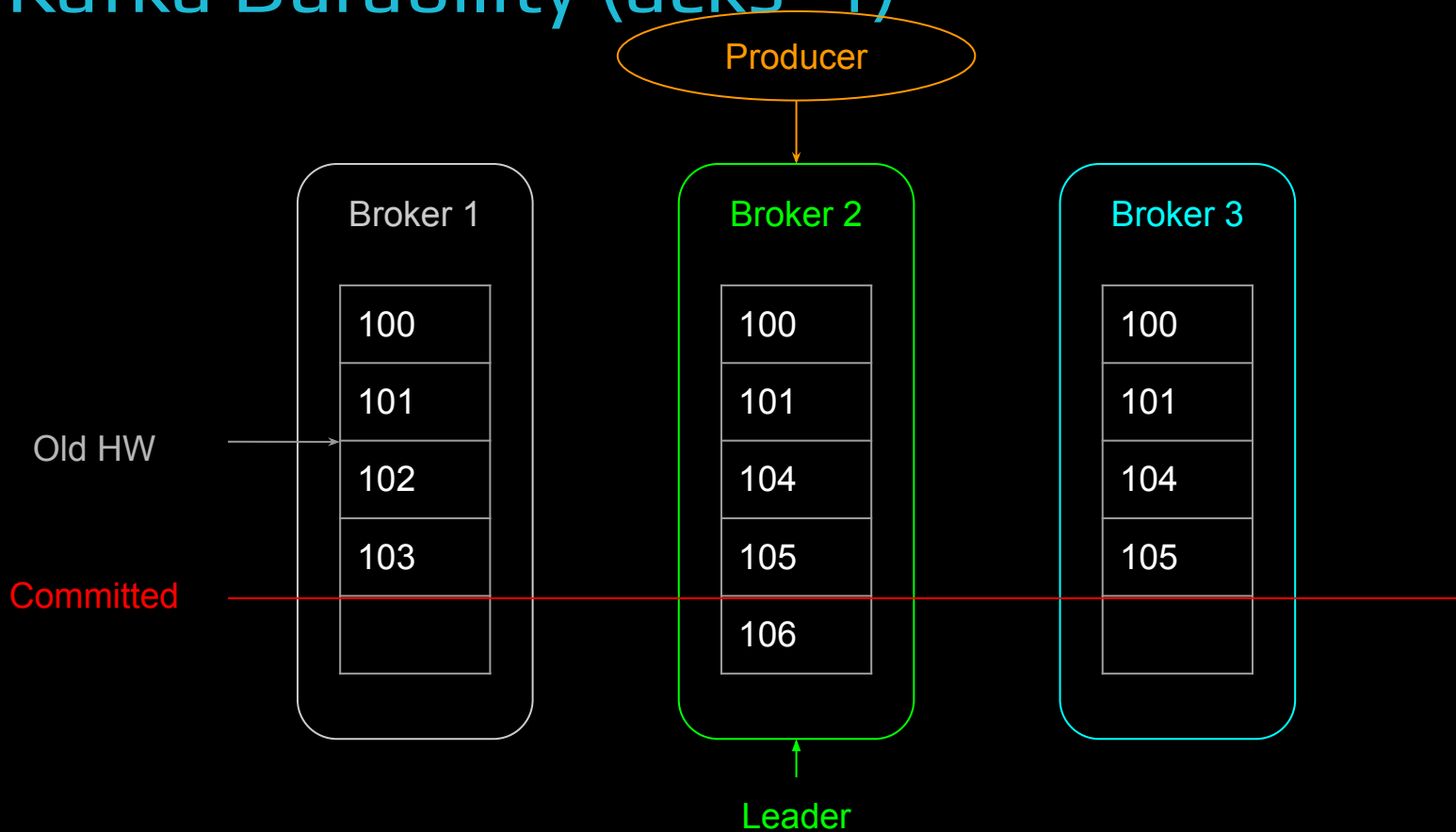
Kafka Durability (acks=1)



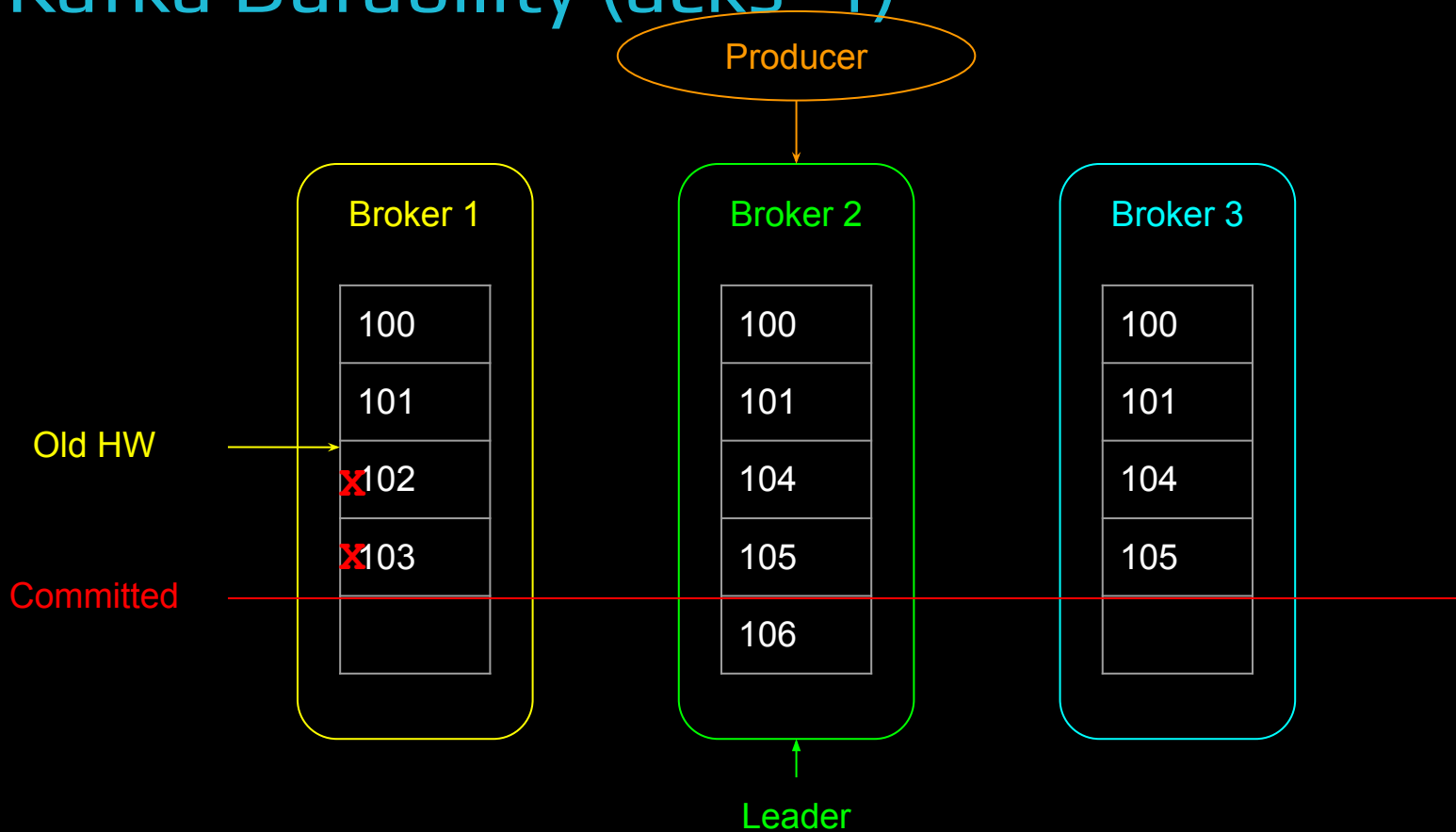
Kafka Durability (acks=1)



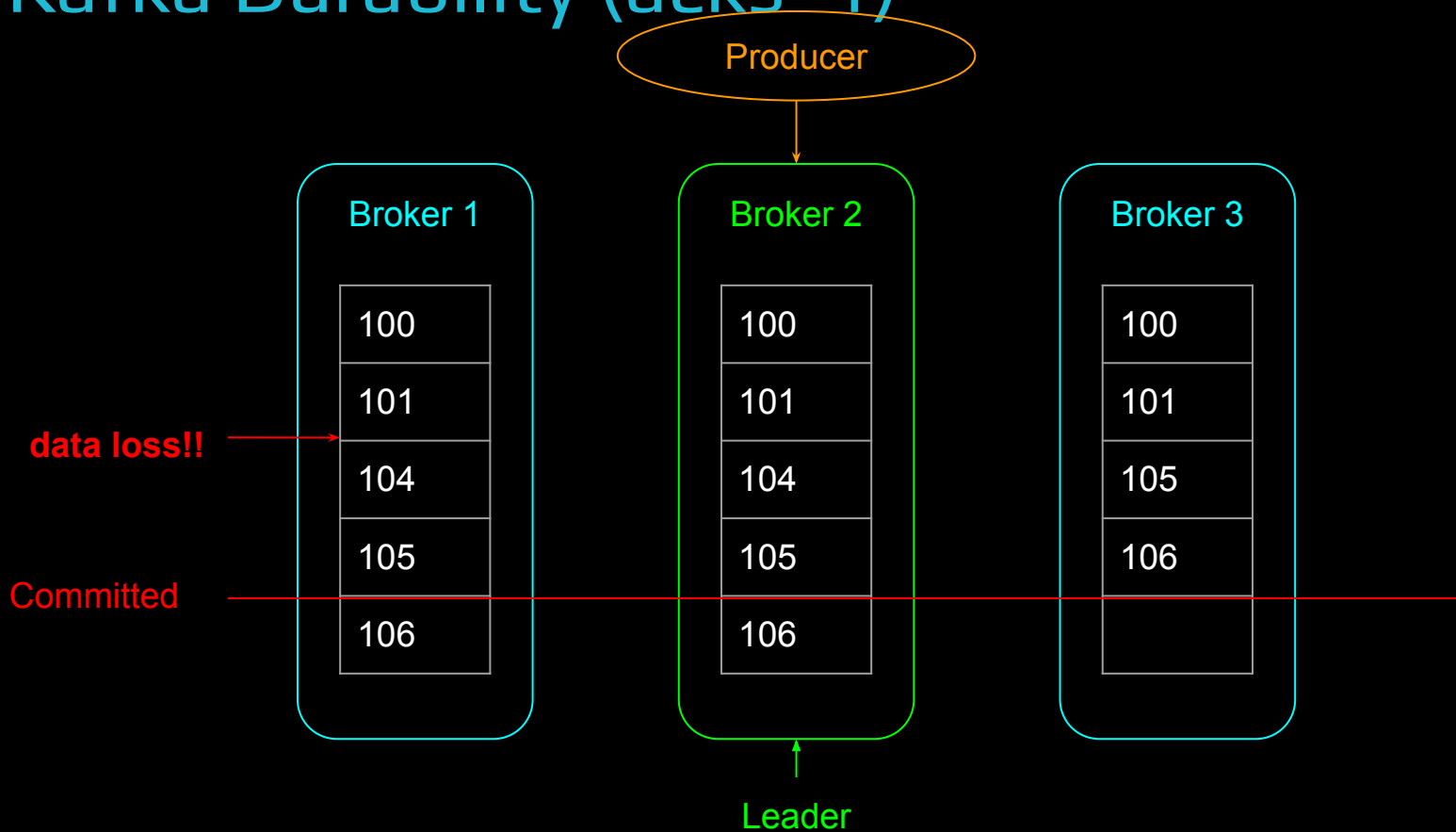
Kafka Durability (acks=1)



Kafka Durability (acks=1)

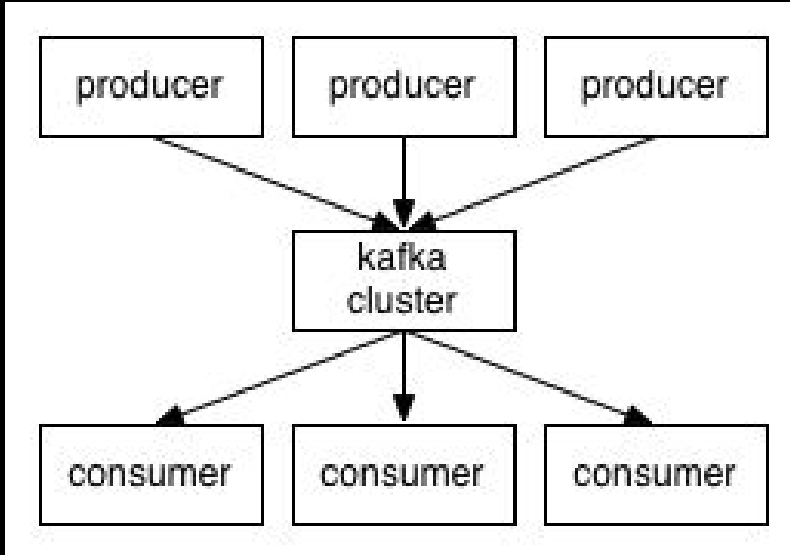


Kafka Durability (acks=1)





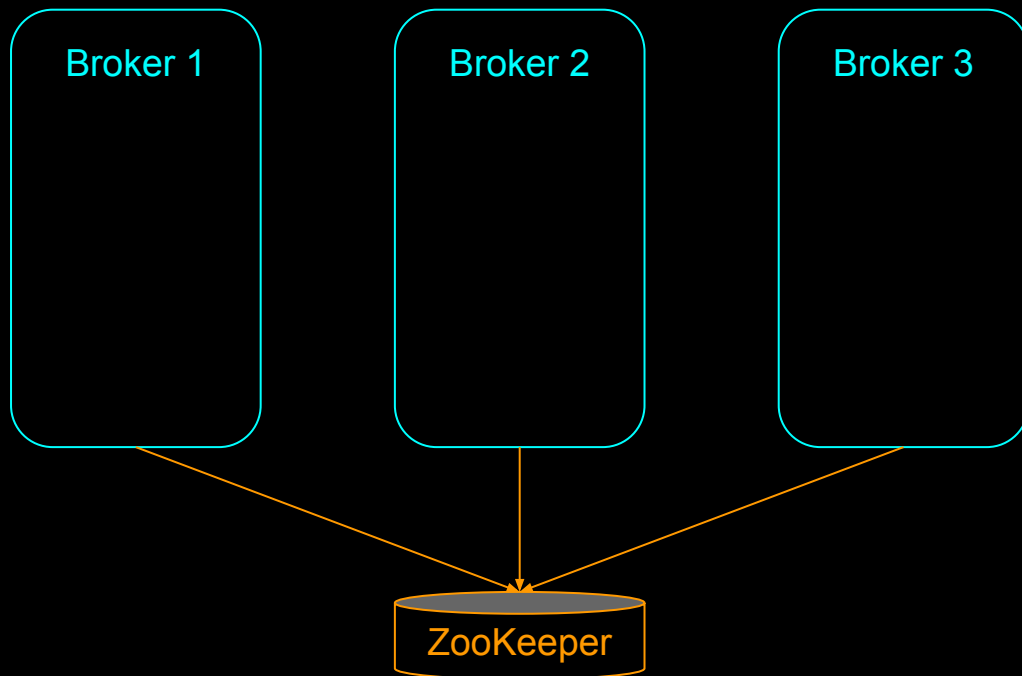
Distributed Messaging system



- High throughput
- Low latency
- Scalable
- Centralized
- Real-time

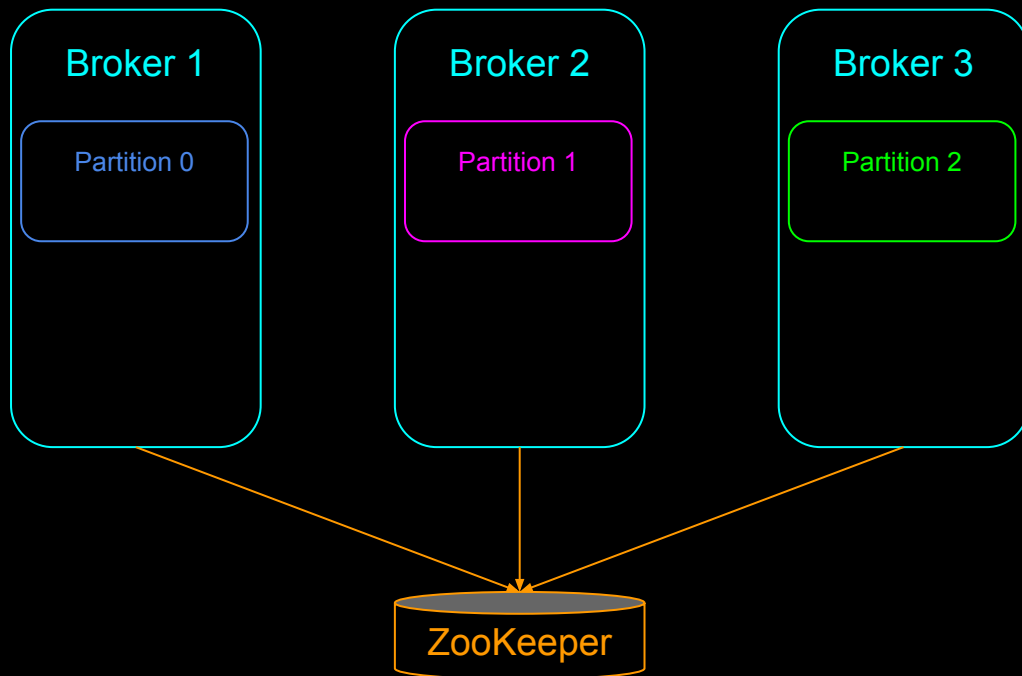
What is Kafka?

- **Distributed**
- Partitioned
- Replicated
- Commit Log



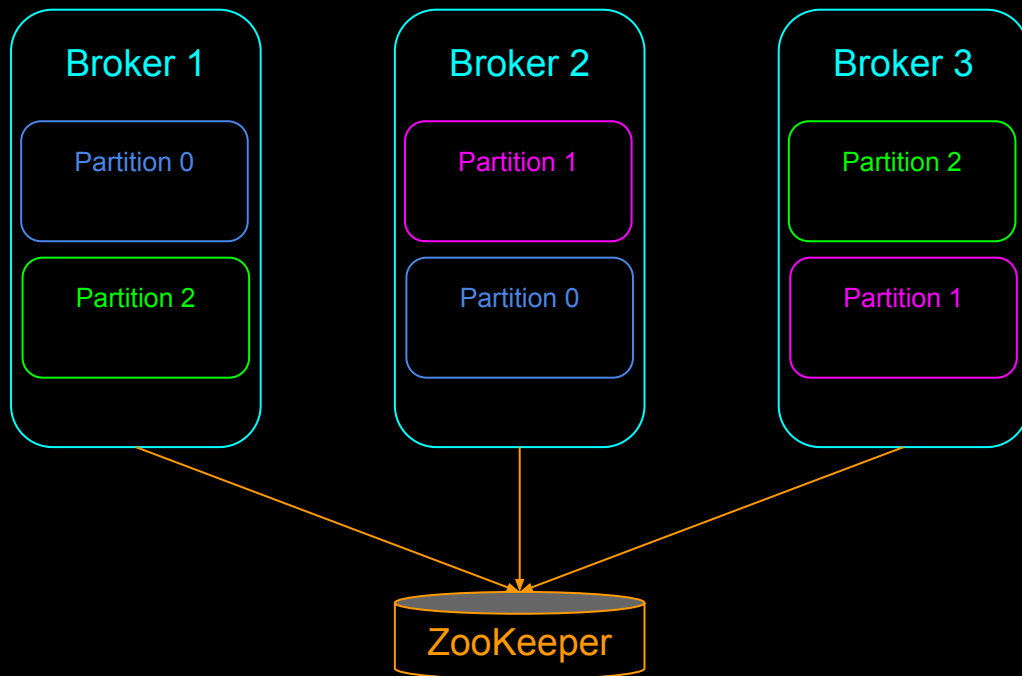
What is Kafka?

- Distributed
- **Partitioned**
- Replicated
- Commit Log



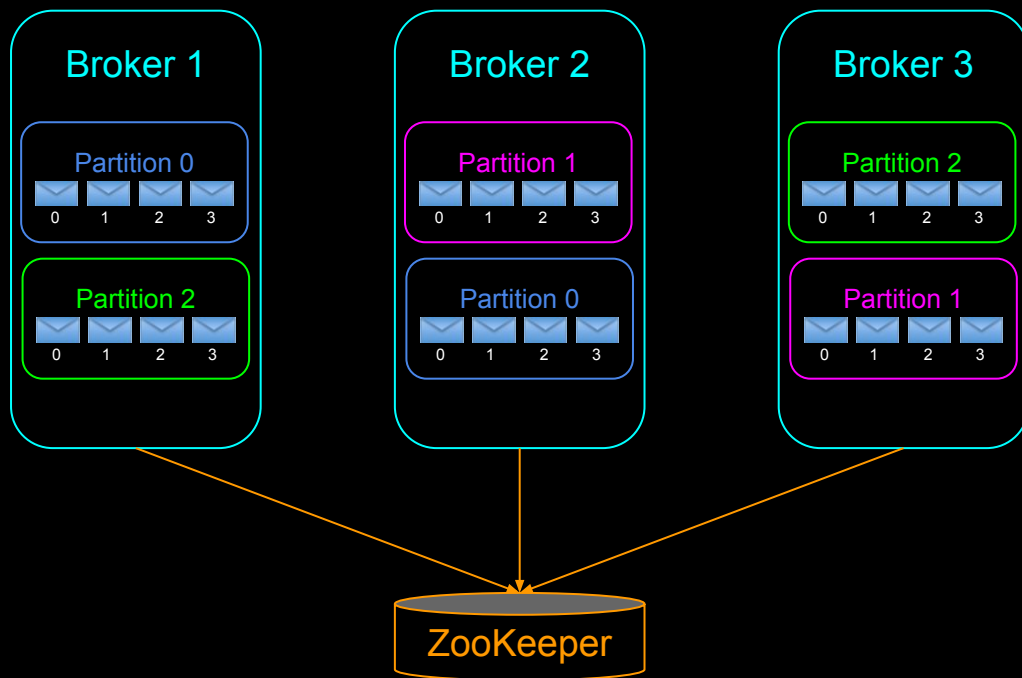
What is Kafka?

- Distributed
- Partitioned
- **Replicated**
- Commit Log



What is Kafka?

- Distributed
- Partitioned
- Replicated
- Commit Log



Kafka Concepts

