



Apache Kafka® *Internal Architecture*

Contents

1. The Fundamentals
2. Hands On: Tuning the Apache Kafka® Producer Client
3. Inside the Apache Kafka® Broker
4. Data Plane: Replication Protocol
5. The Apache Kafka® Control Plane
6. Consumer Group Protocol
7. Hands On: Consumer Group Protocol
8. Configuring Durability, Availability and Ordering Guarantees
9. Transactions
10. Hands On: Transactions
11. Topic Compaction
12. Tiered Storage
13. Cluster Elasticity
14. Geo-Replication
15. Hands On: Cluster Linking



The Fundamentals

What's Covered



1 Inside the Broker

Control Panel, Replication Protocols, Topic Compaction, etc.

2 Durability, Availability and Ordering Guarantees

3 Consumer Group Protocol

4 Transactions

5 Tiered Storage

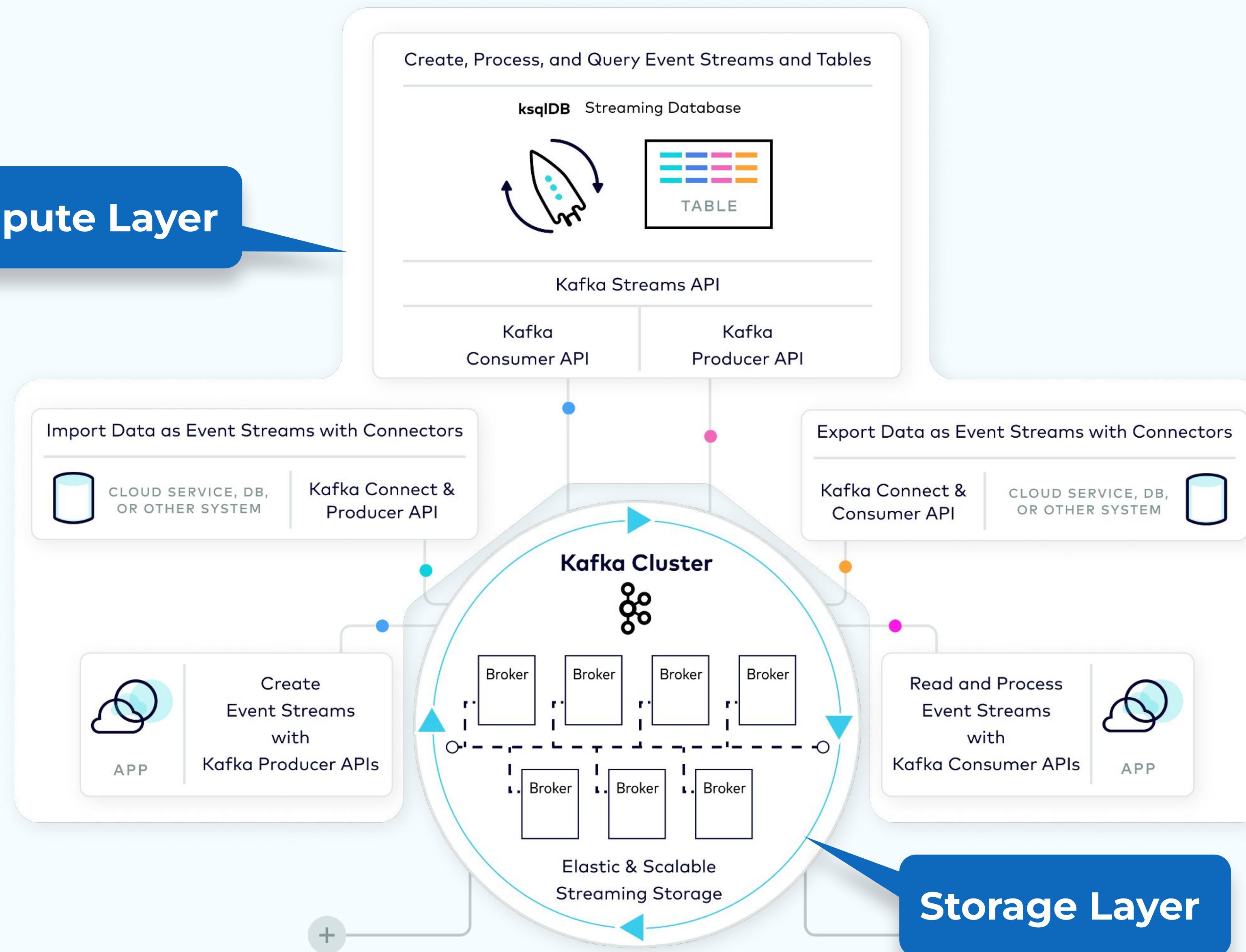
6 Elasticity

7 Geo-Replication

Overview of Kafka Architecture



Compute Layer



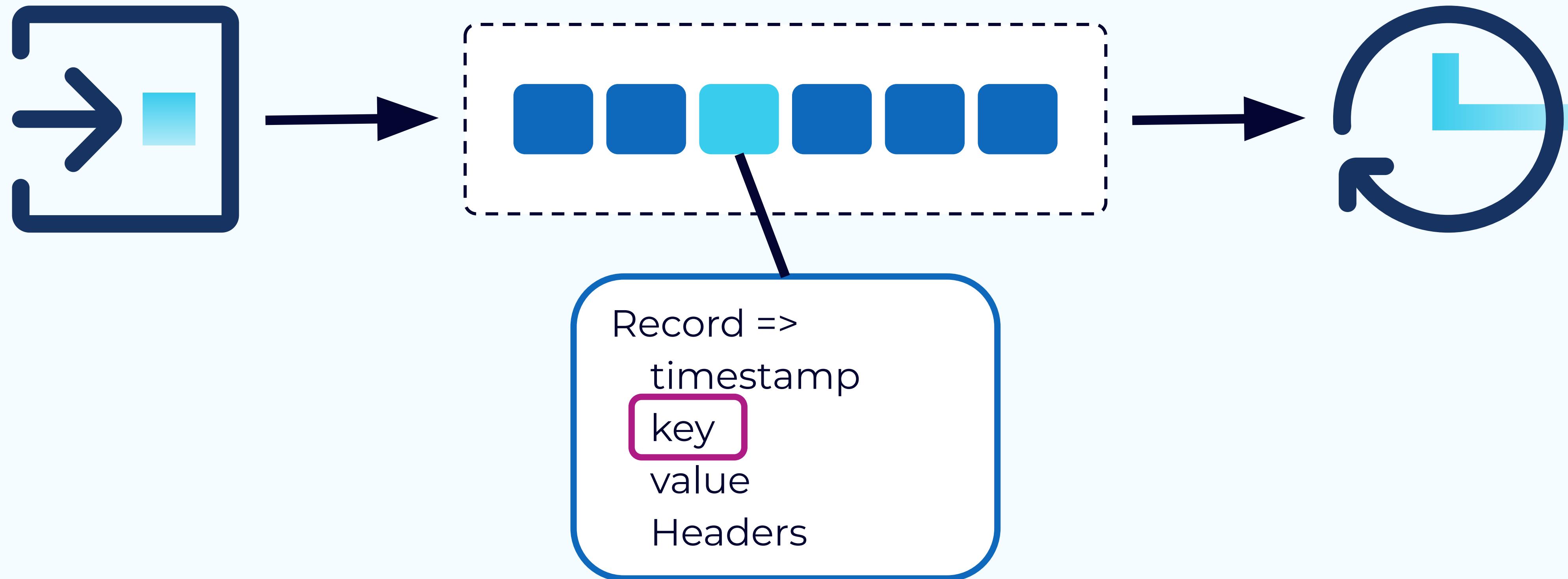
Events



Event Source

Event Stream

Event Processing
Application



Record Schema



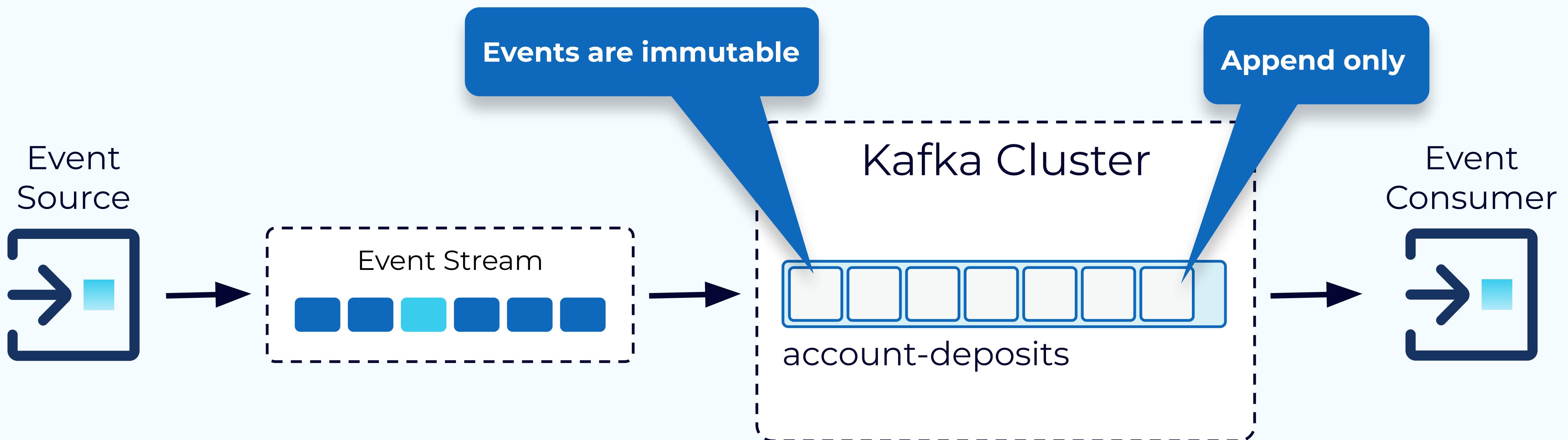
Record =>
timestamp
key
value
Headers

key/ value Bytes	Area	Description
0	Magic Byte	Confluent serialization format version number; currently always 0 .
1-4	Schema ID	4-byte schema ID as returned by Schema Registry.
5...	Data	Serialized data for the specified schema format.



Event Stream

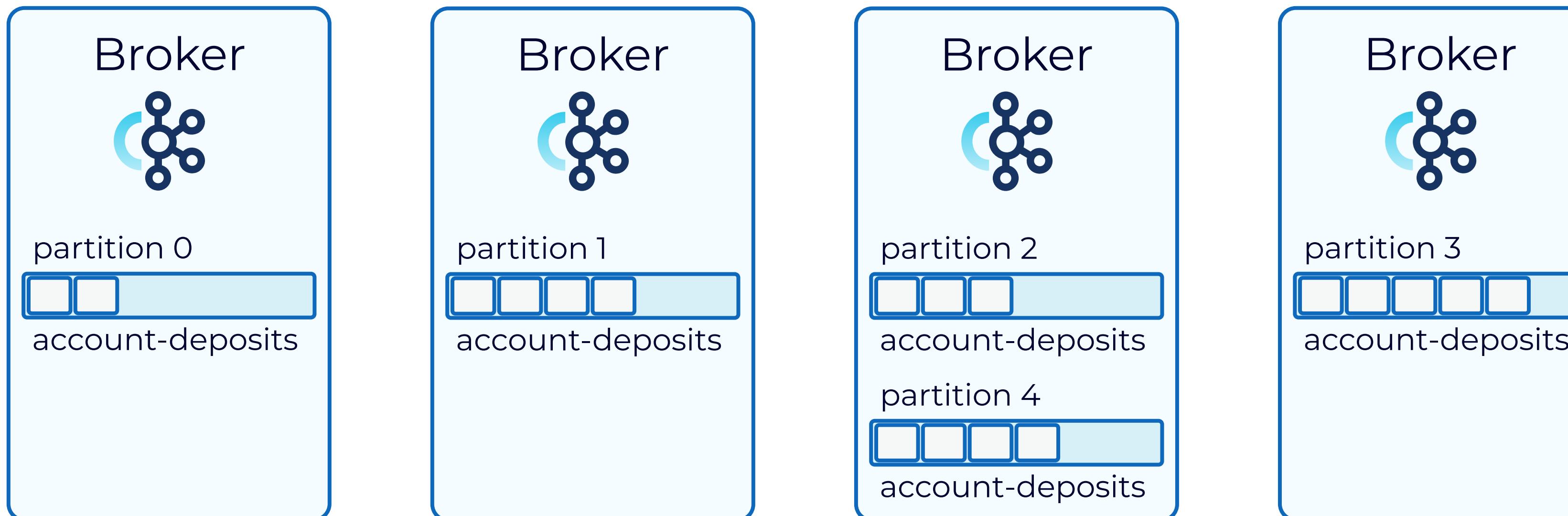
Kafka Topics



Topic Partitions

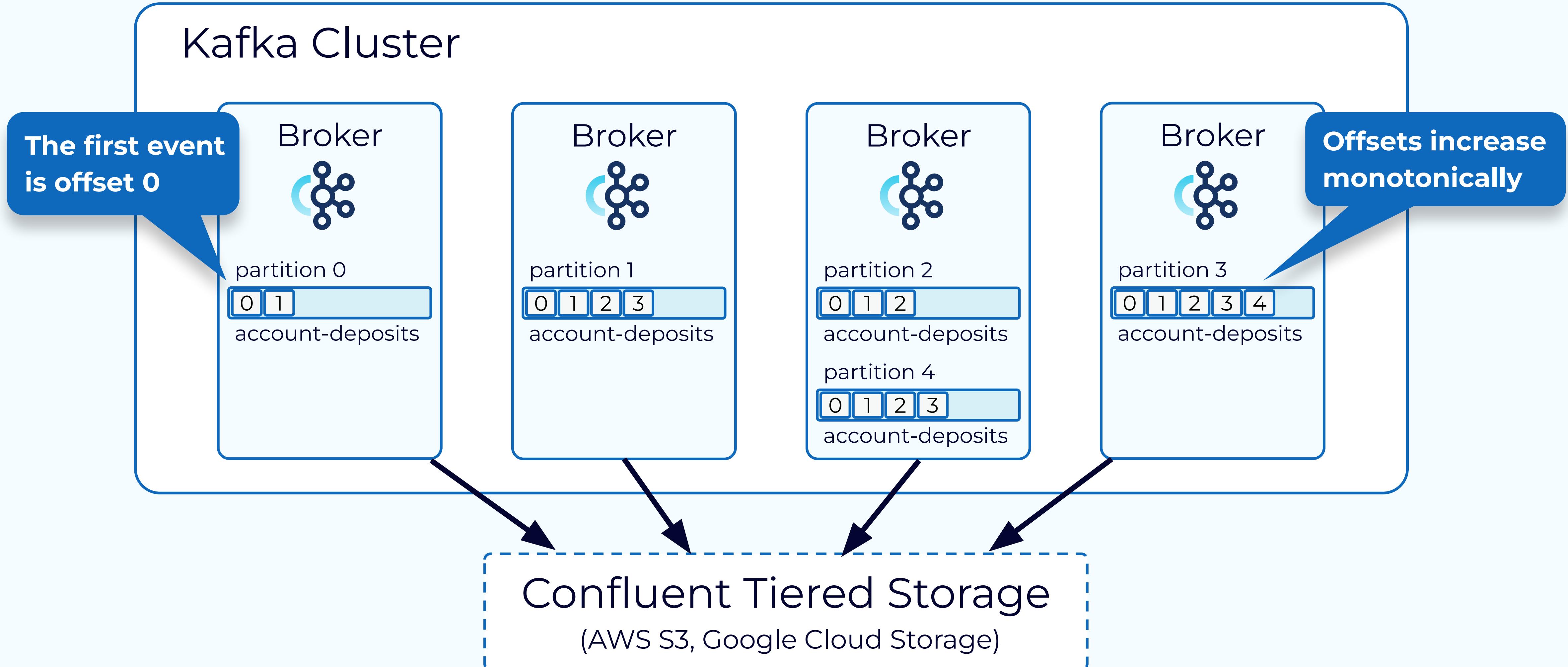


Kafka Cluster



Confluent Tiered Storage
(AWS S3, Google Cloud Storage)

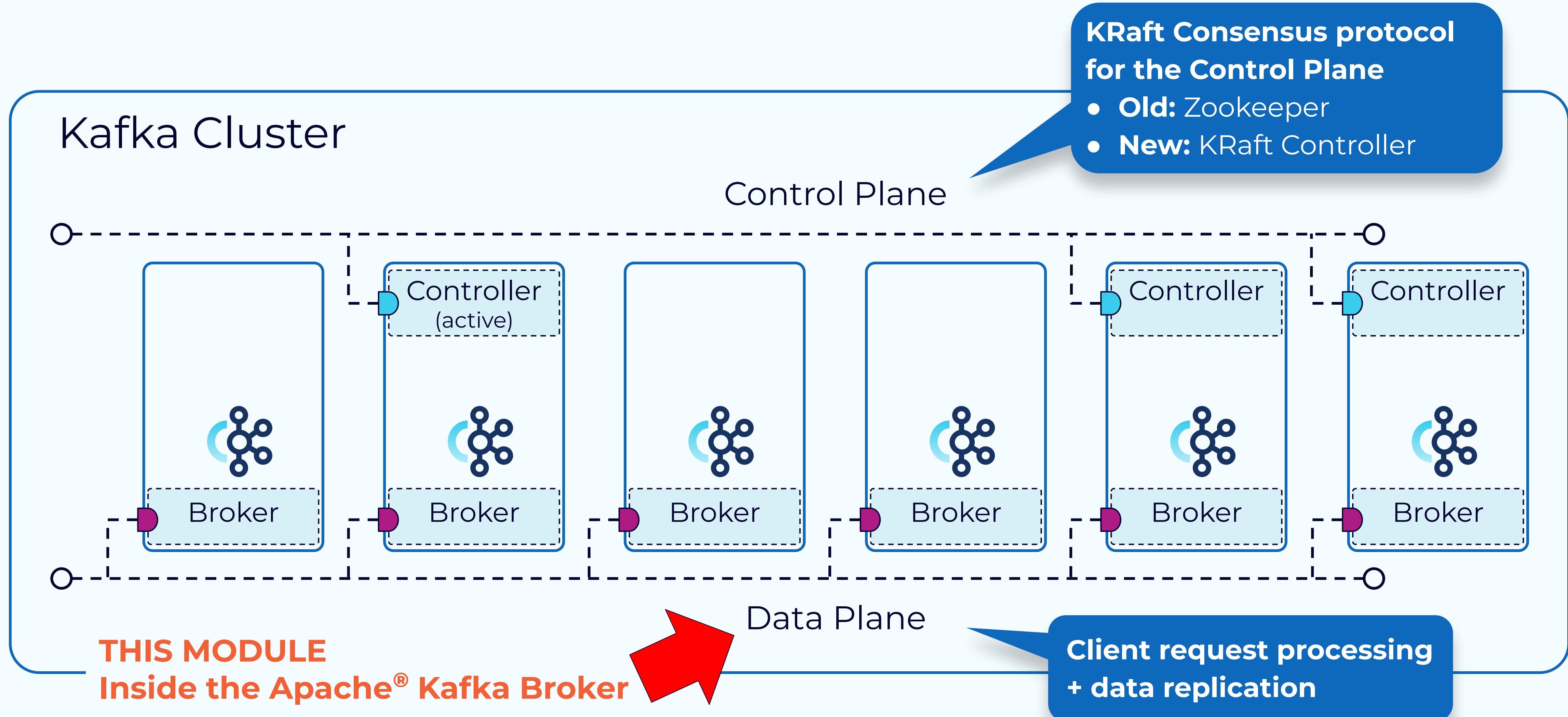
Partition Offsets



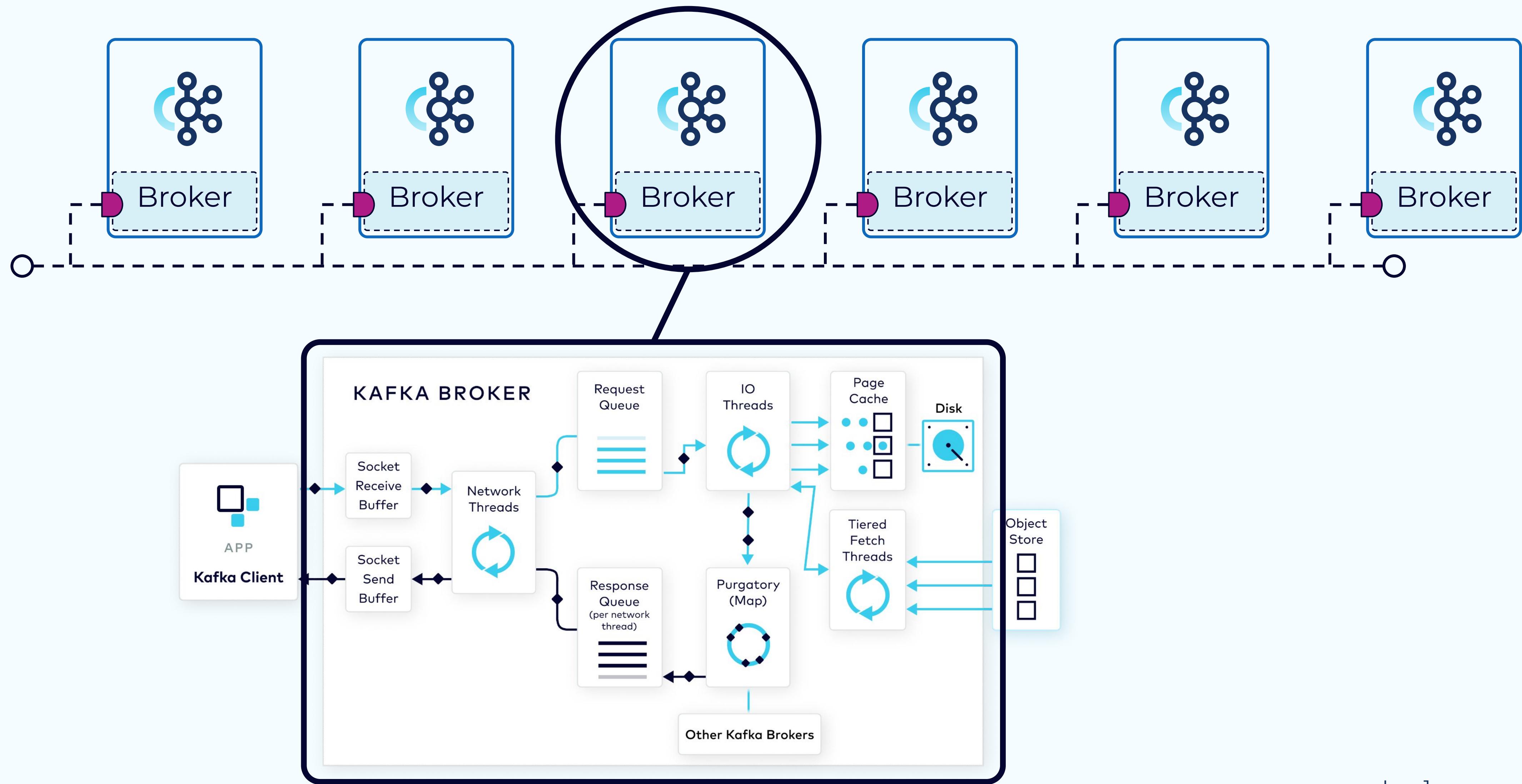


Inside the Apache Kafka® Broker

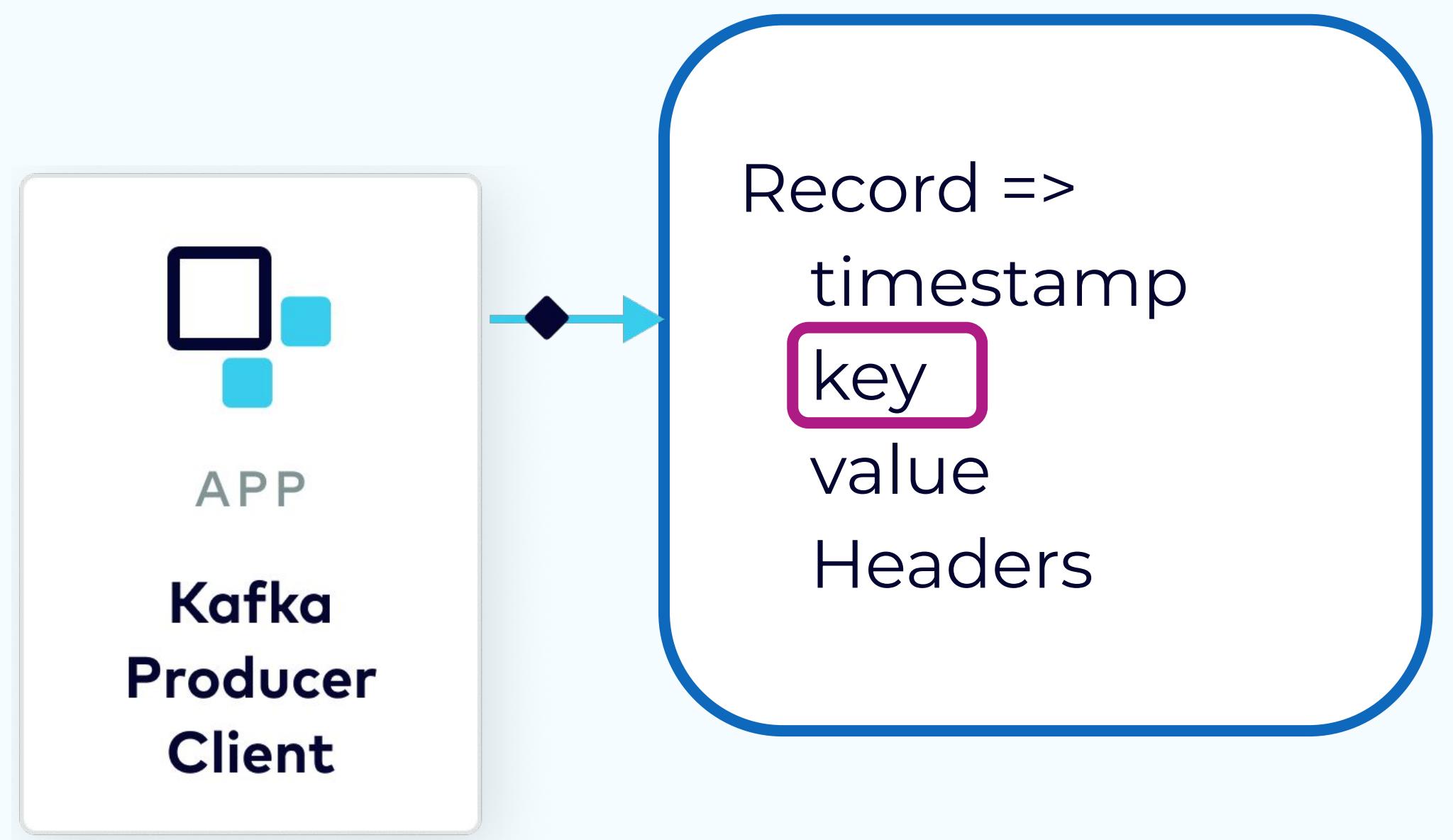
Kafka Manages Data & Metadata Separately



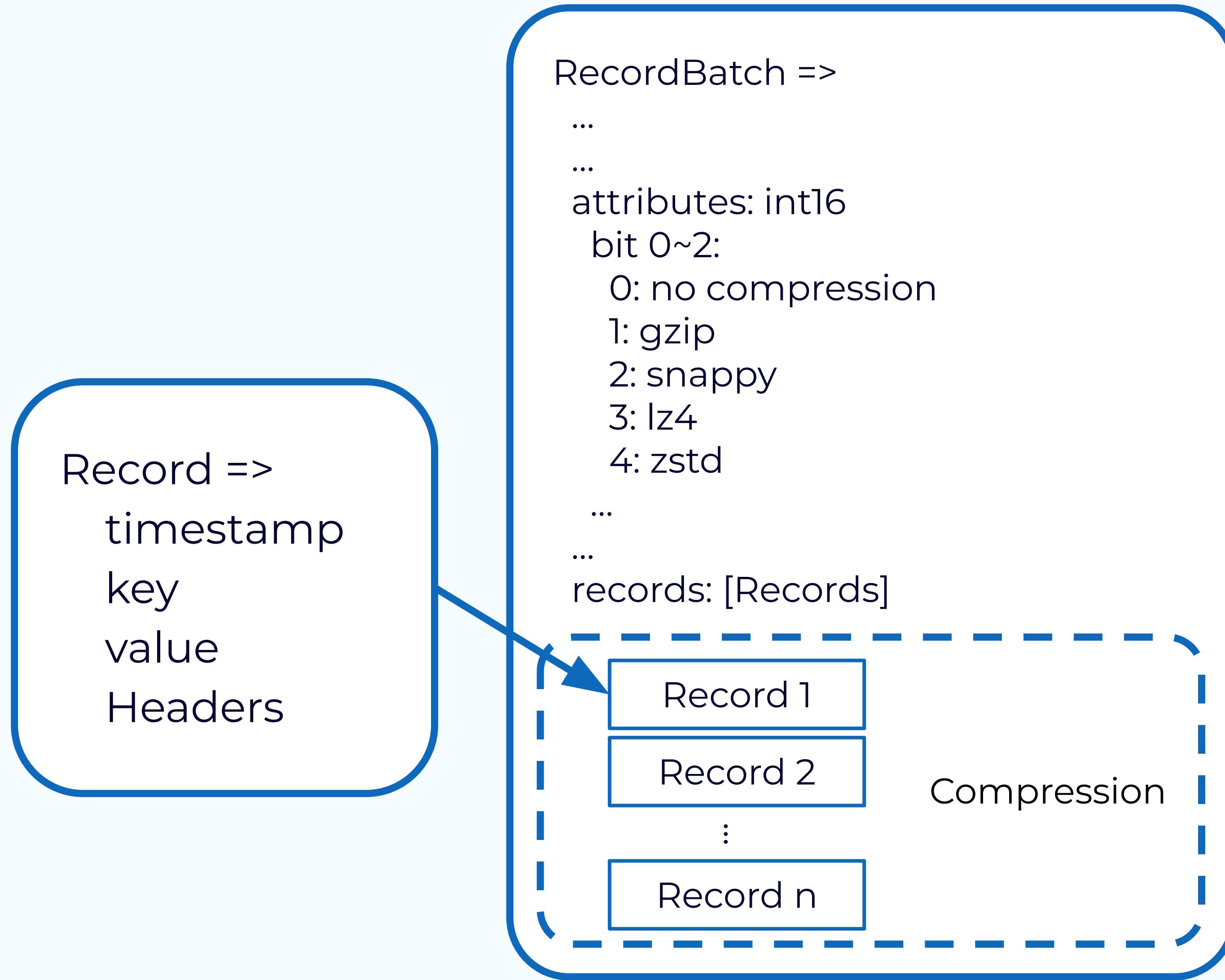
Inside The Apache Kafka Broker



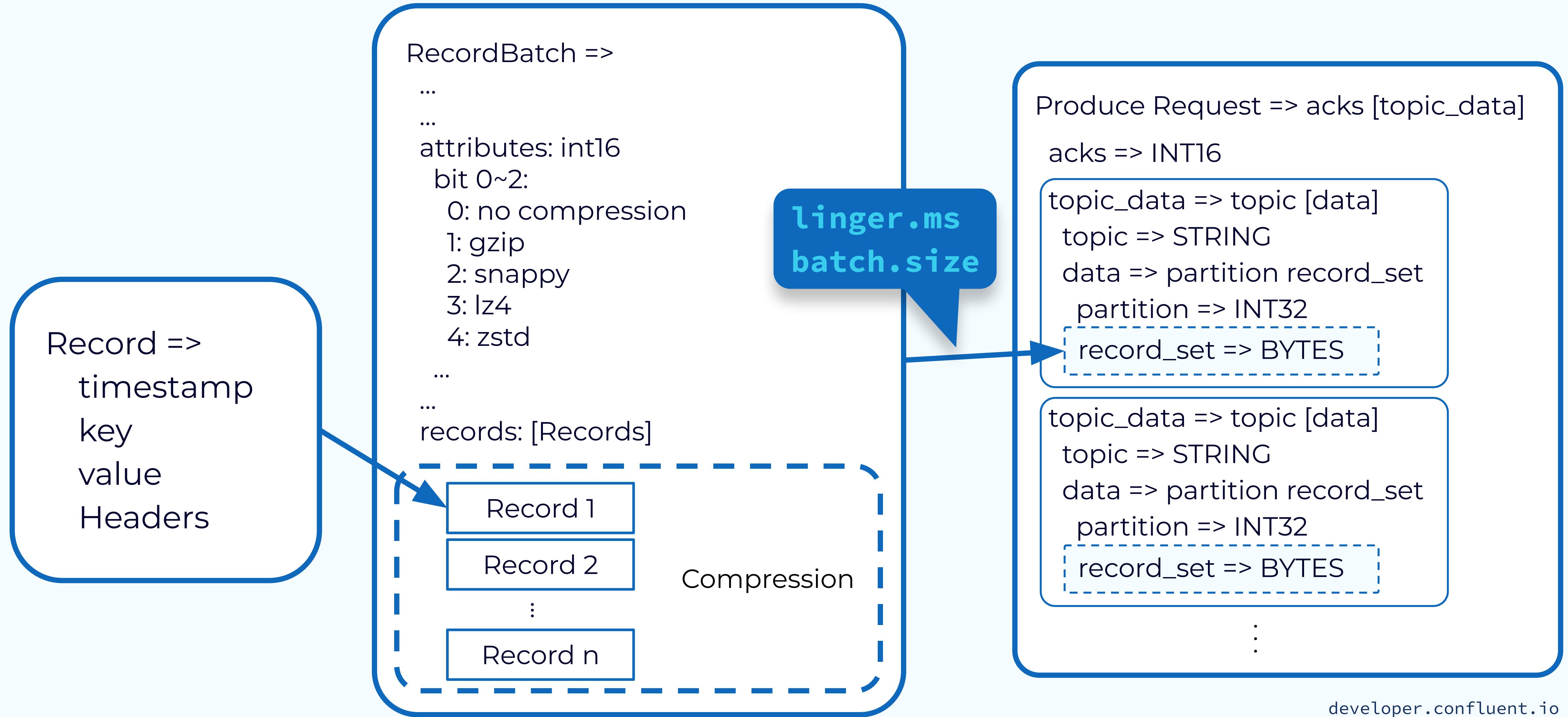
Assign Record to Topic Partition



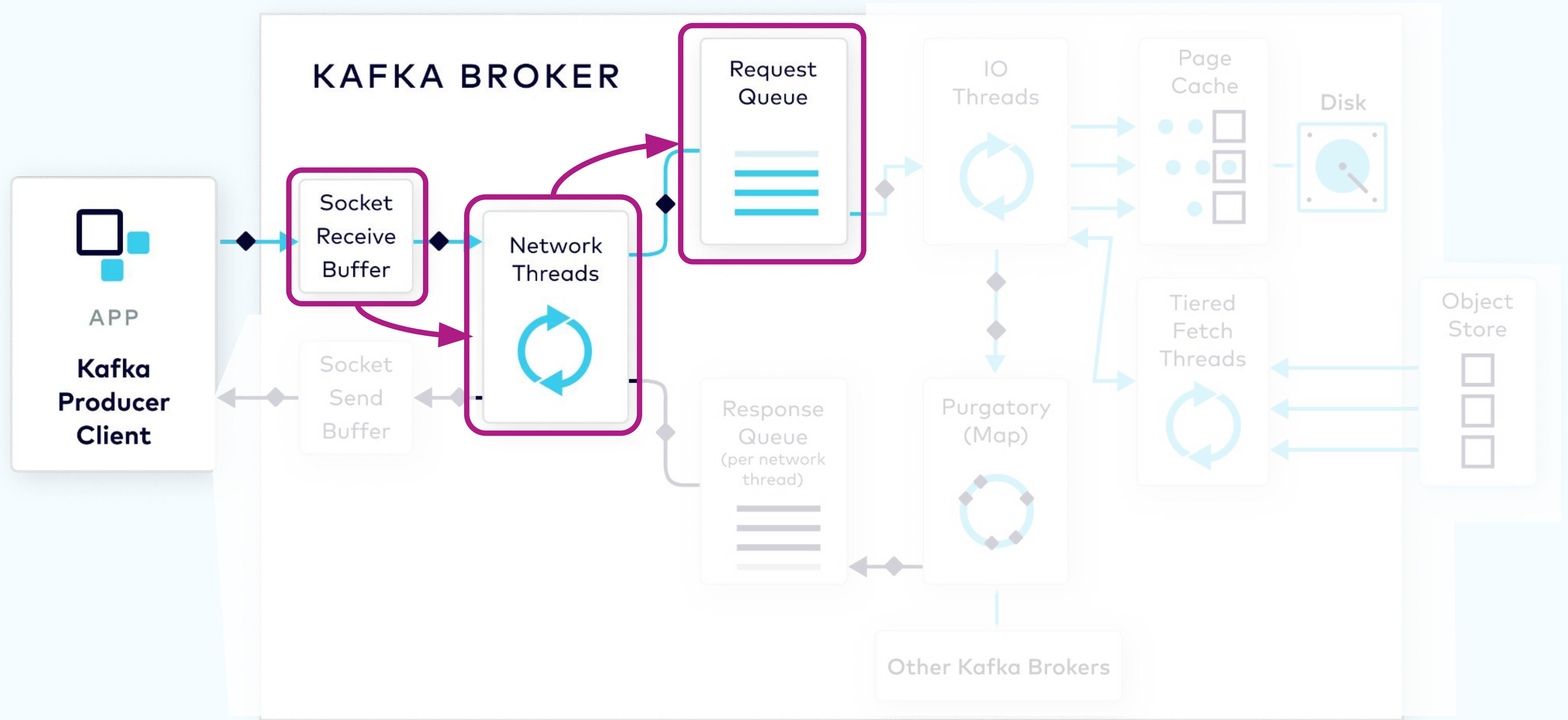
Records Accumulated into Record Batches



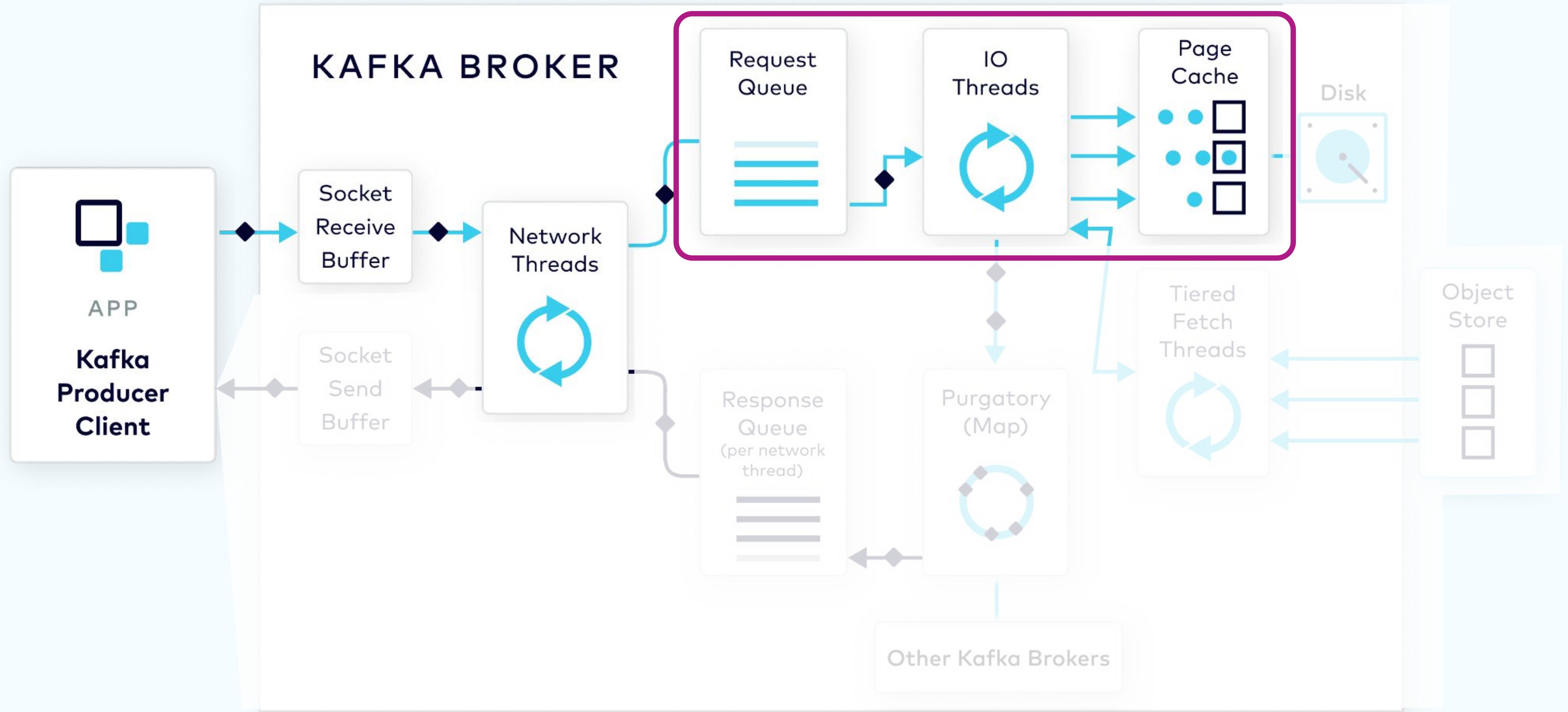
Record Batches Drained into Produce Requests



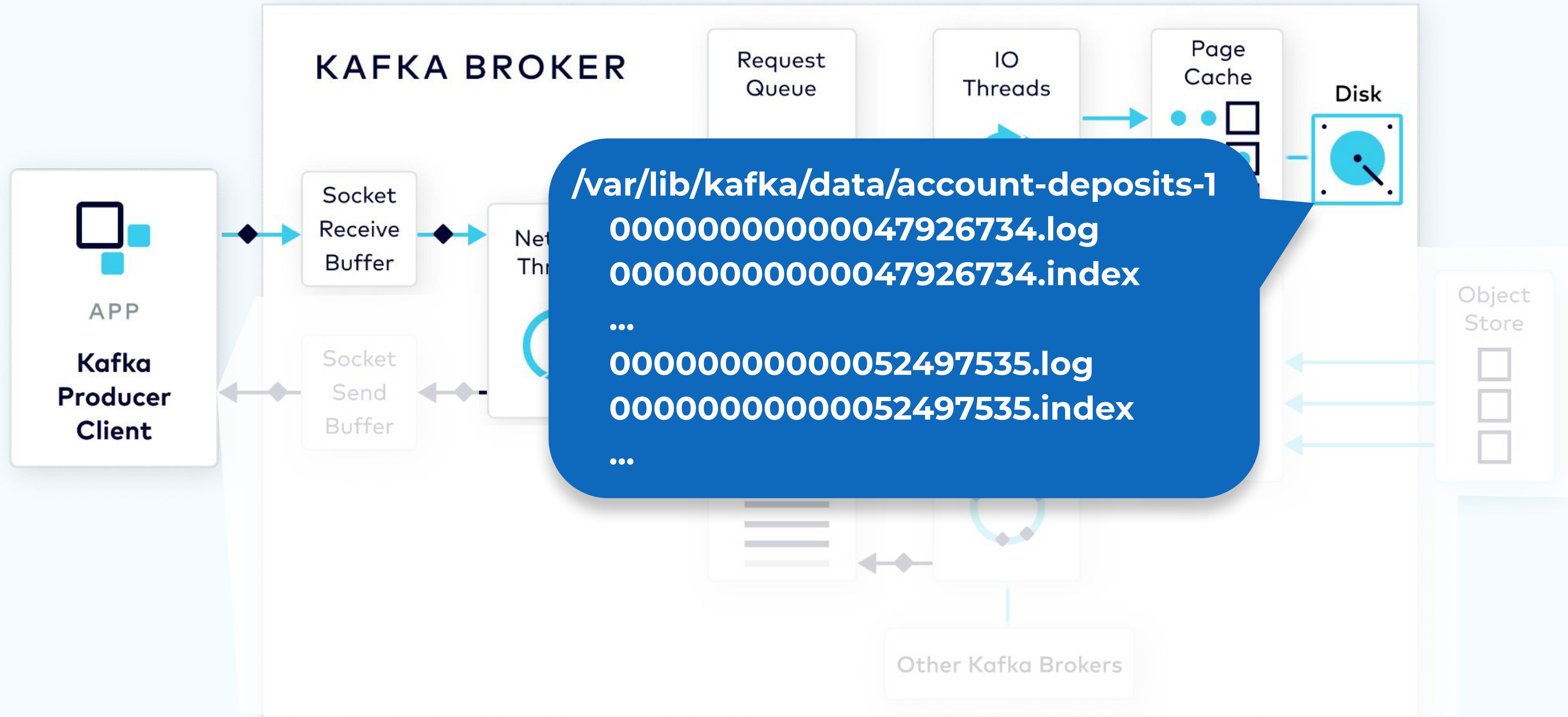
Network Thread Adds Request to Queue



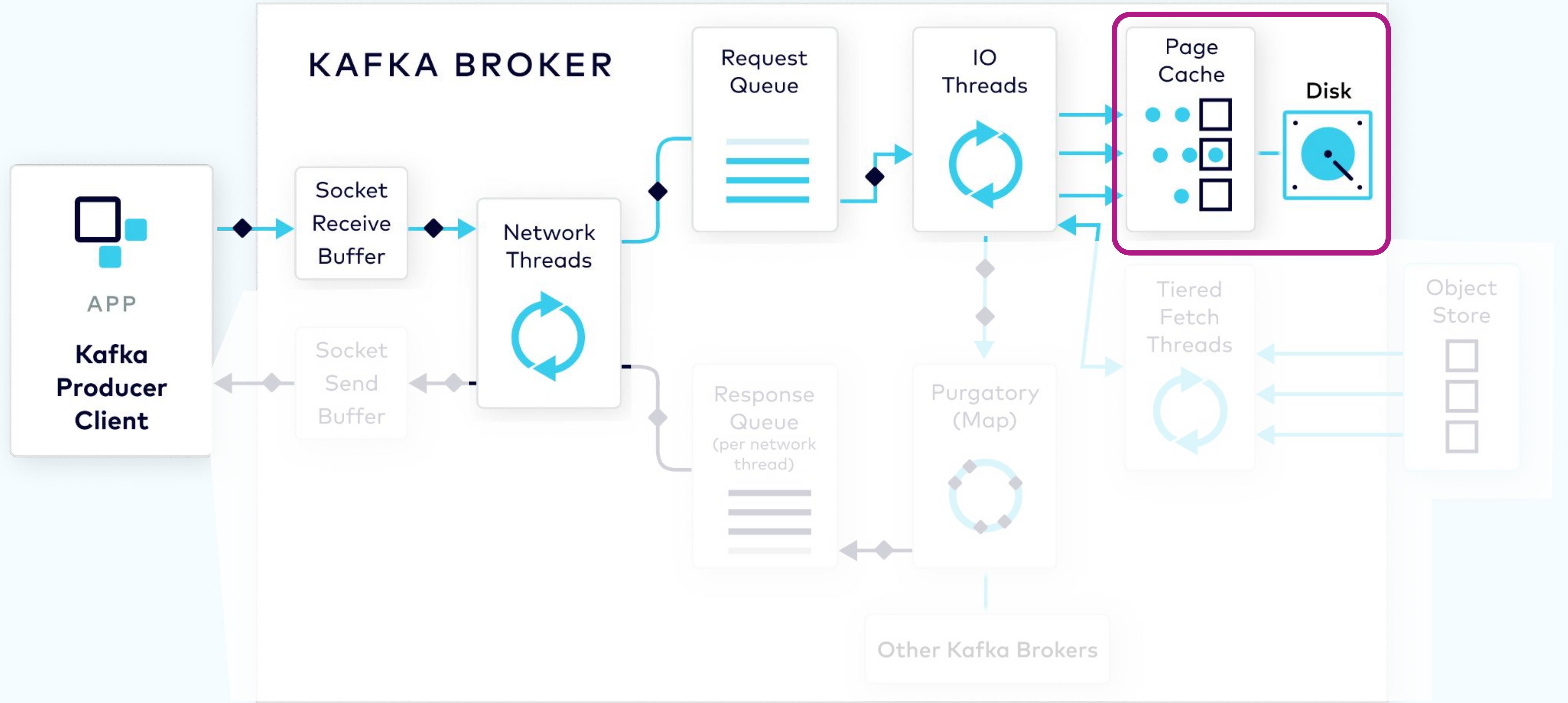
IO Thread Verifies Record Batch And Stores



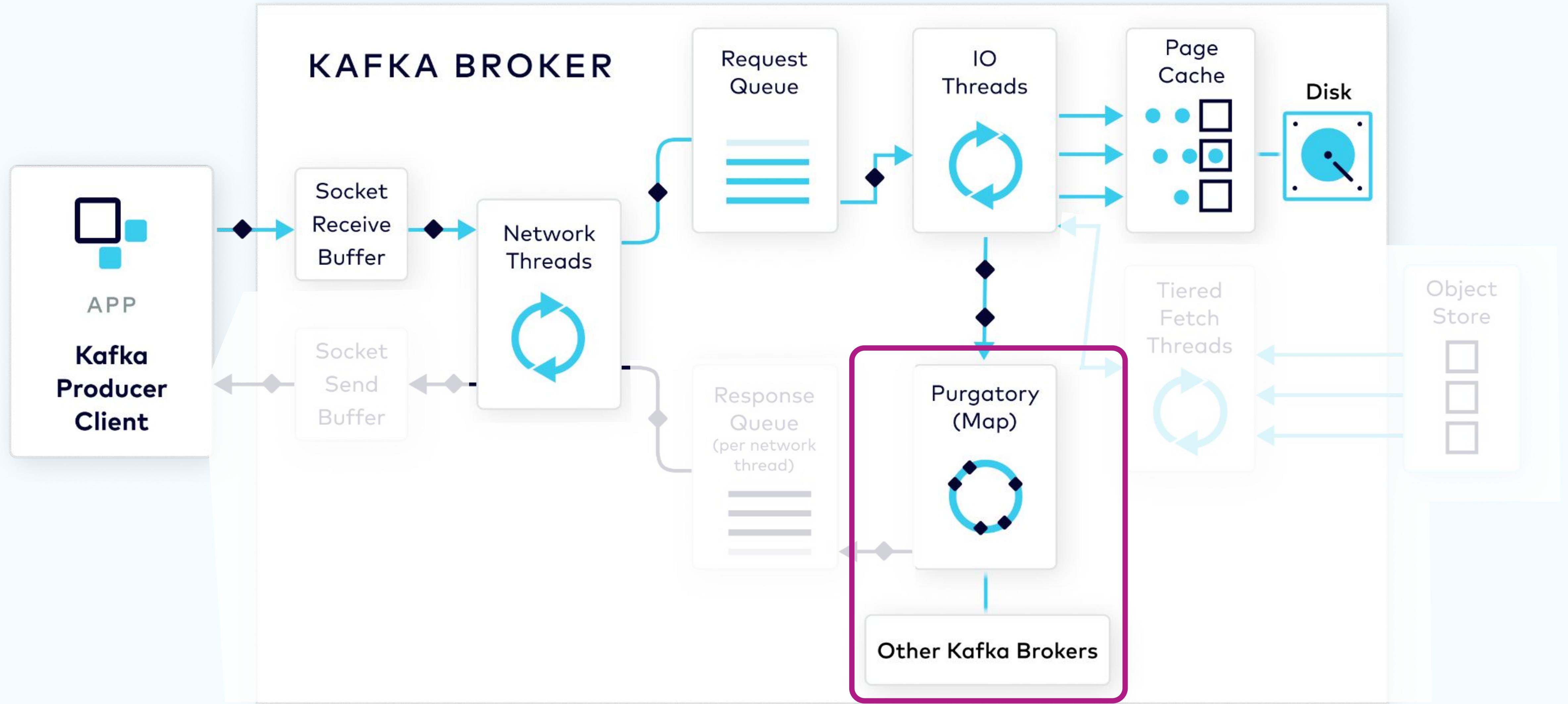
Kafka Physical Storage



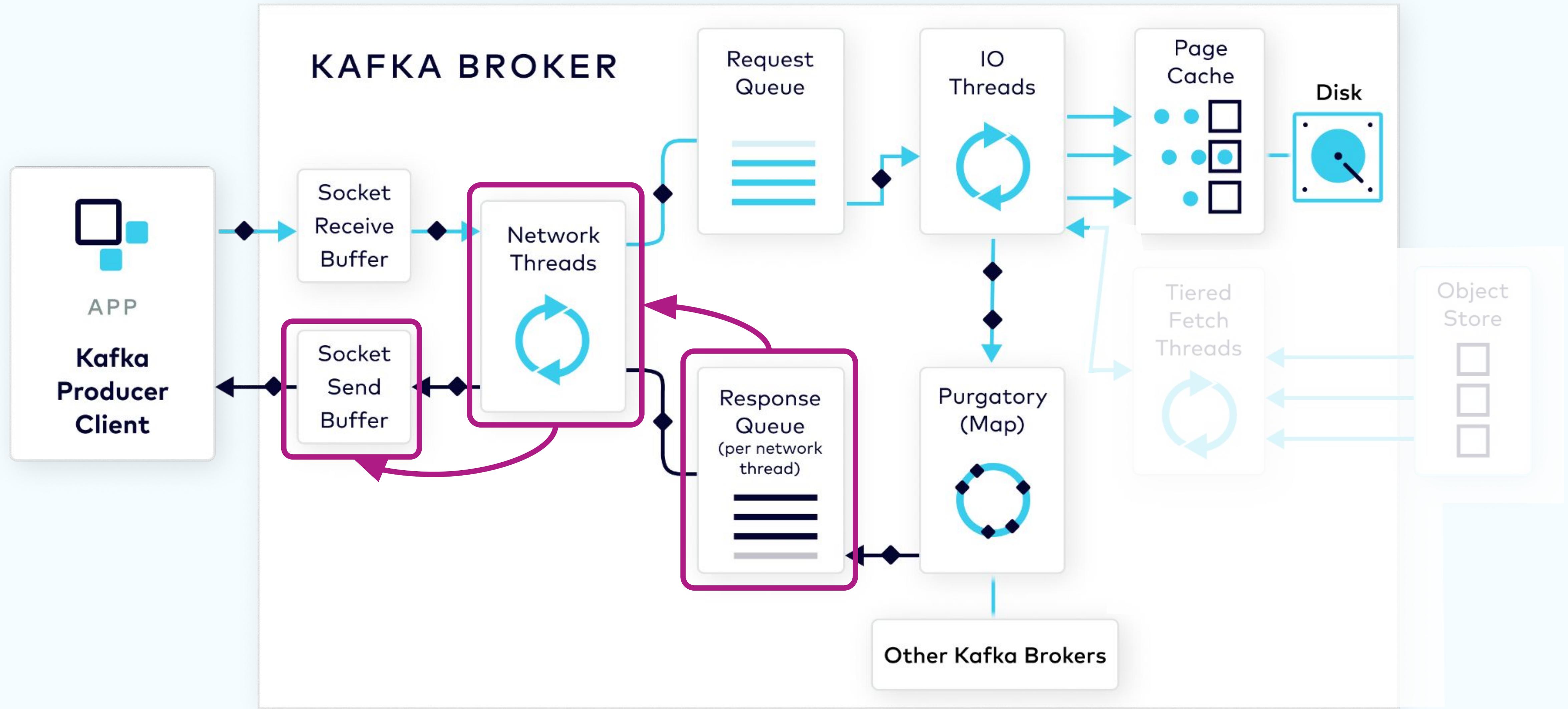
By Default, Data Is Written Asynch to Disk



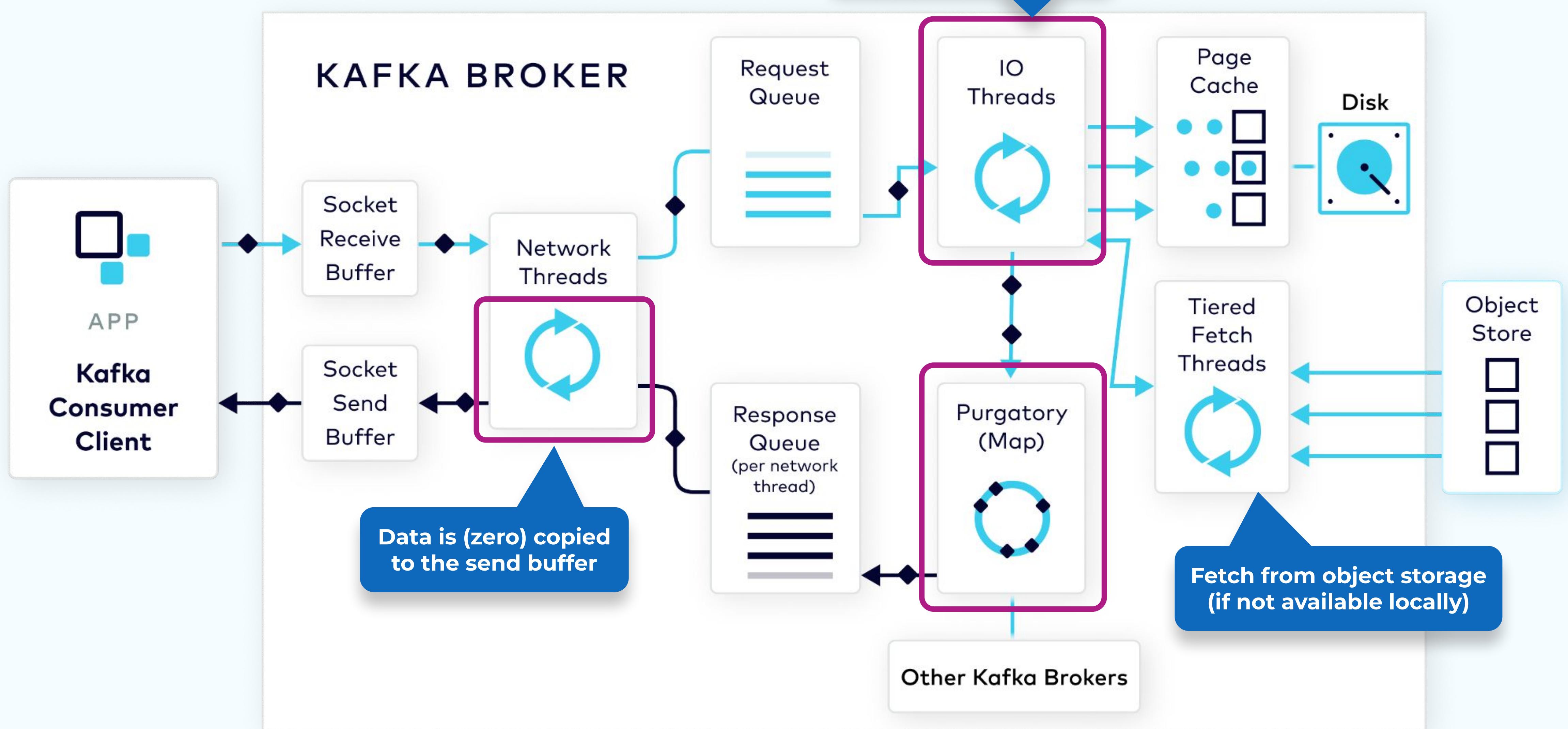
Purgatory Holds Requests Being Replicated



Response Added to Socket Send Buffer



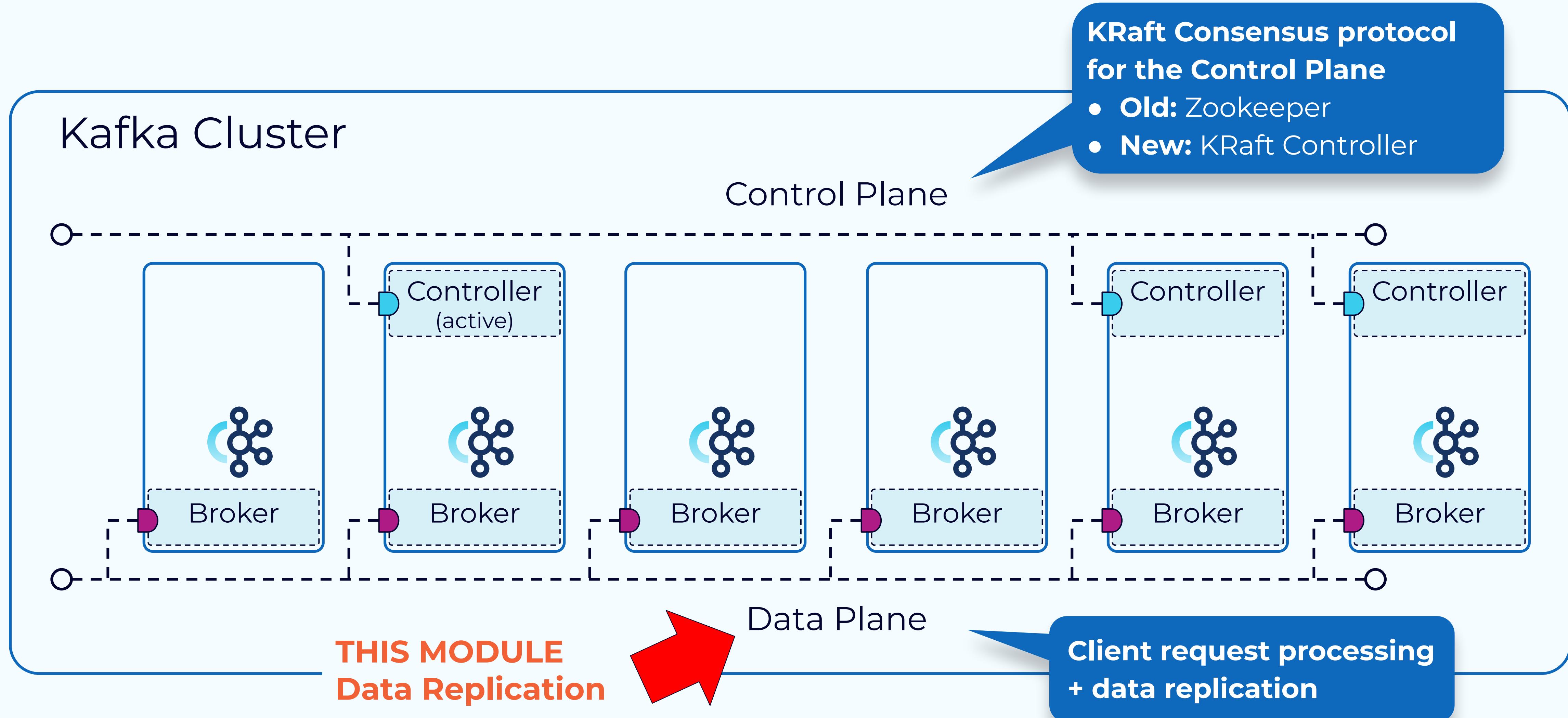
Fetch Requests





Data Plane: Replication Protocol

Kafka Manages Data & Metadata Separately

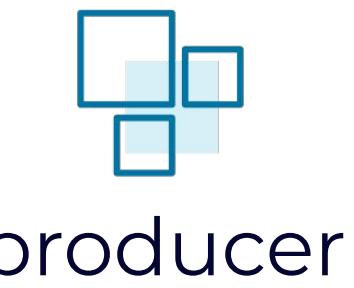
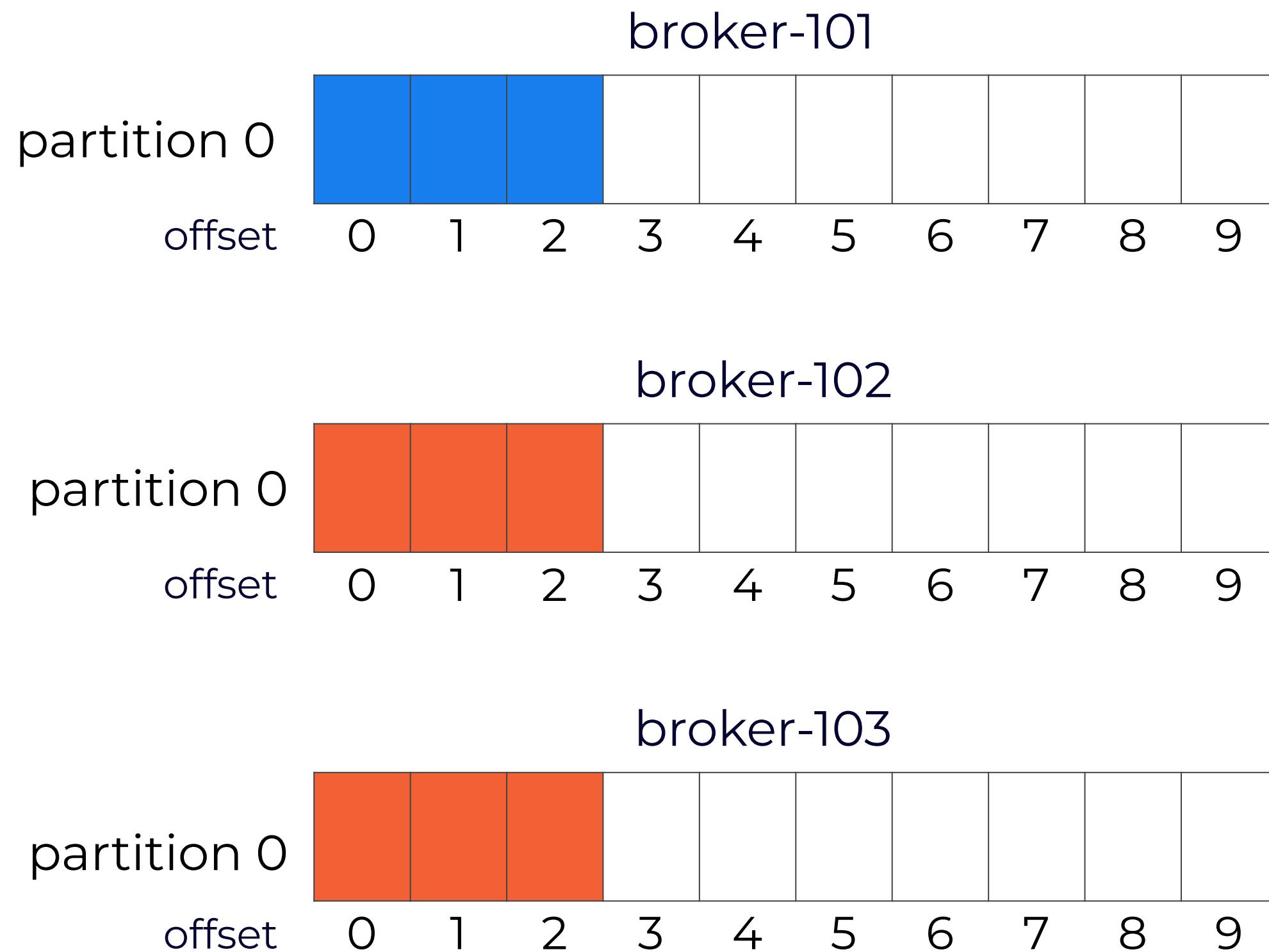


Kafka Data Replication



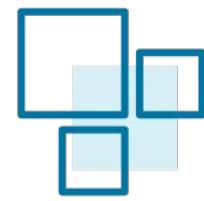
Each topic partition
is replicated to
multiple brokers

Given n replicas, no data
loss with n-1 replica failure

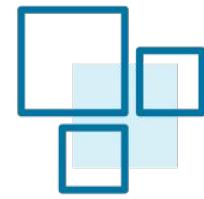


producer

Consumer
Group

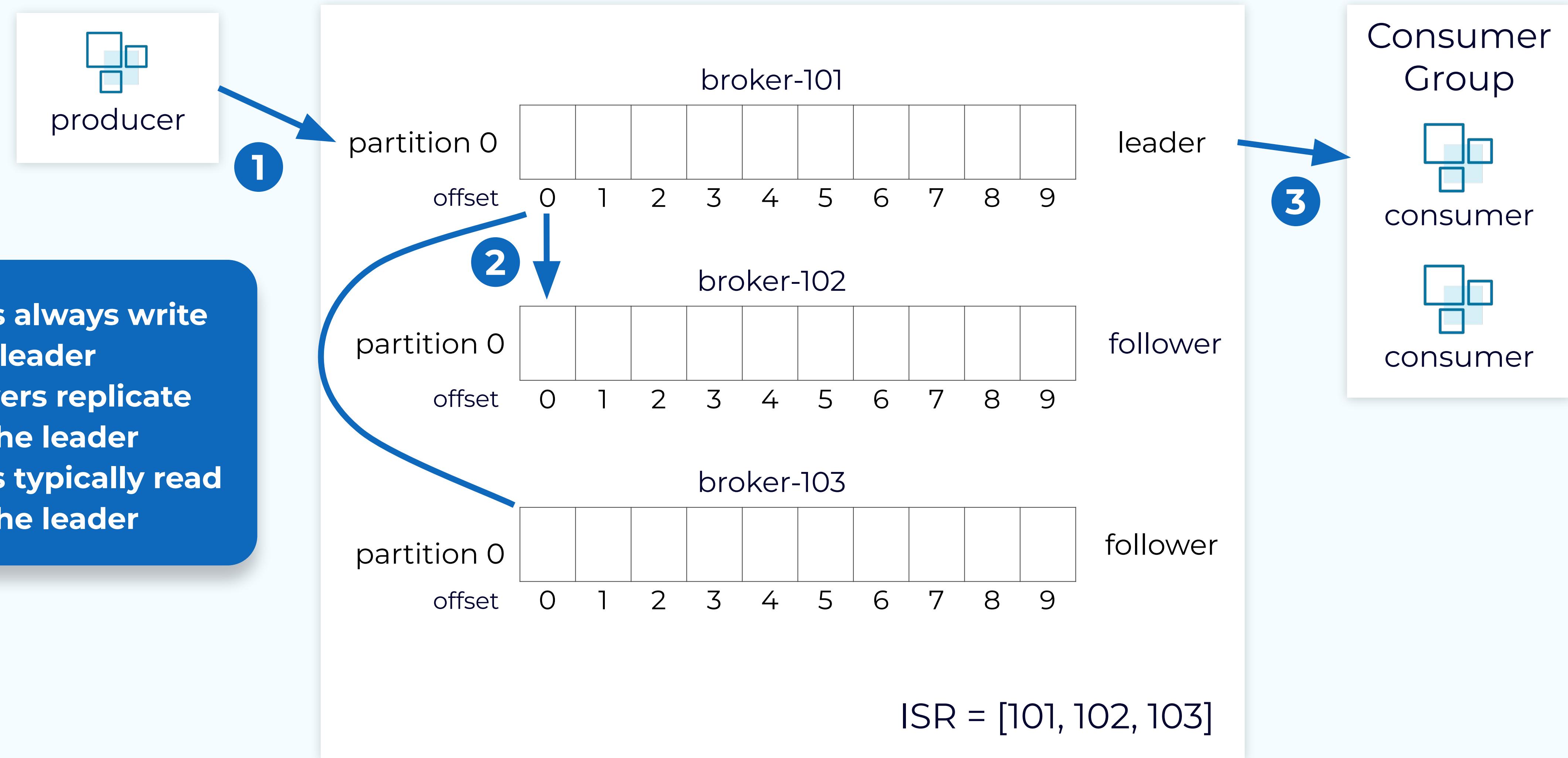


consumer

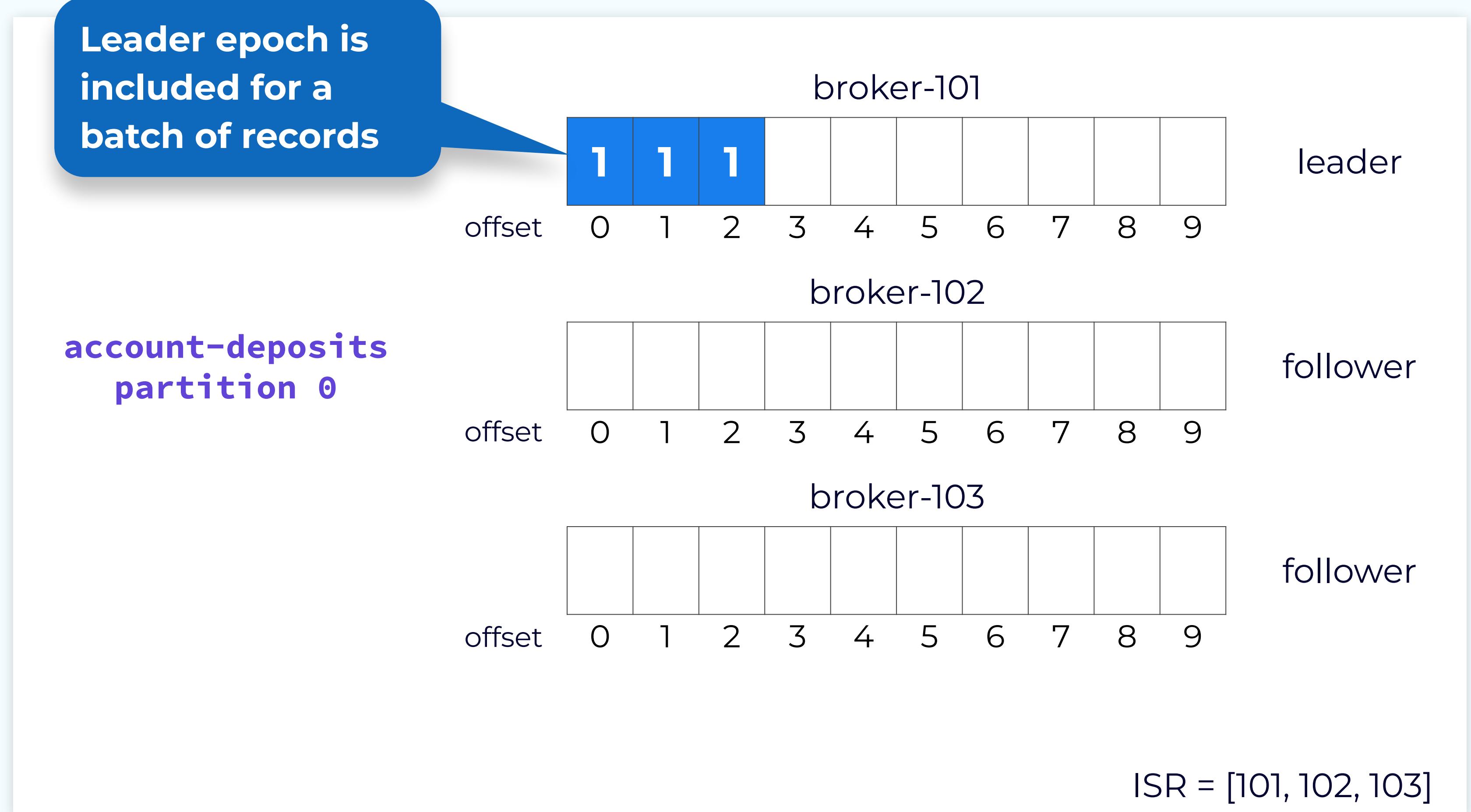


consumer

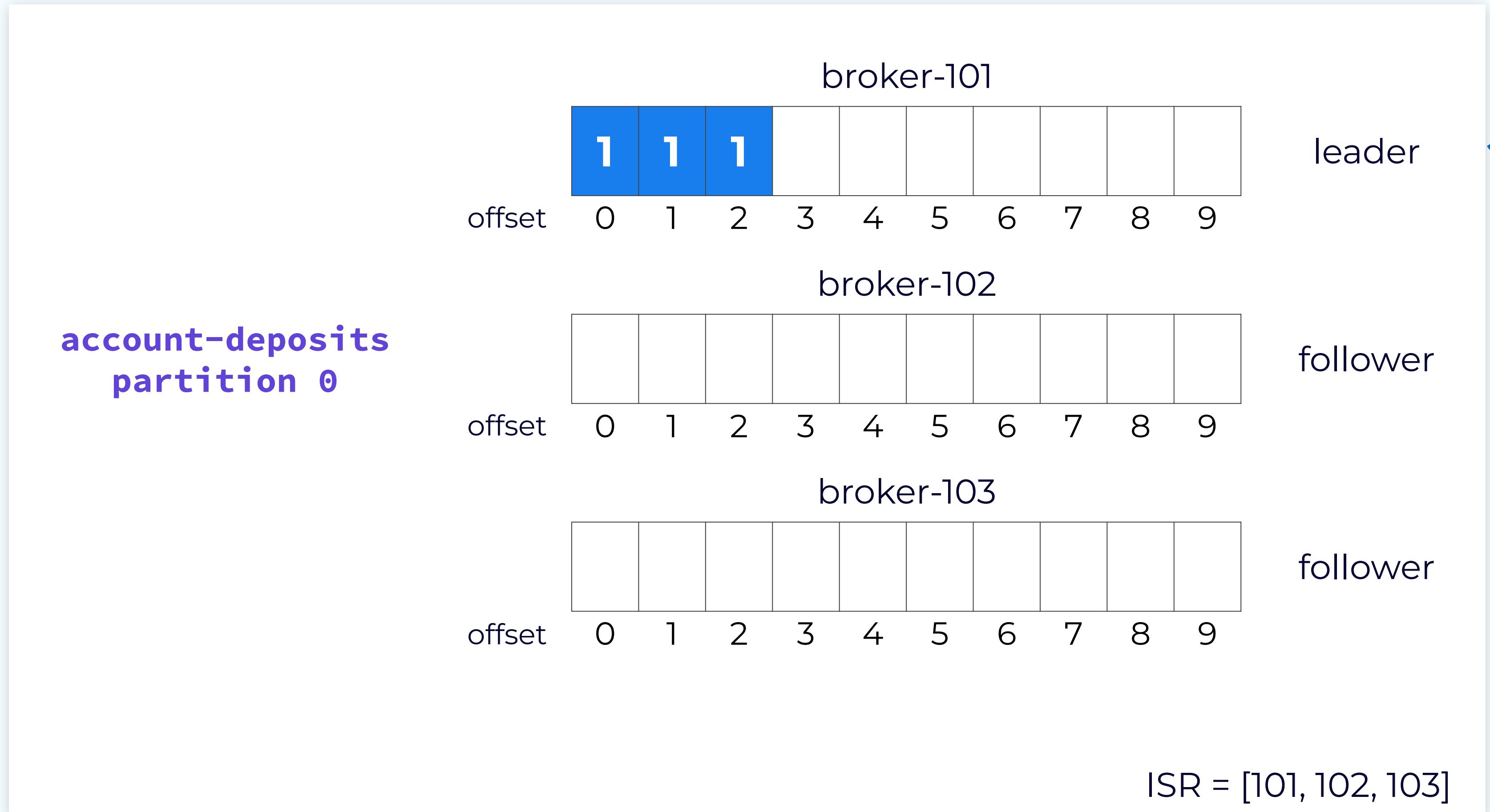
Leader, Follower, and In-Sync Replica (ISR) List



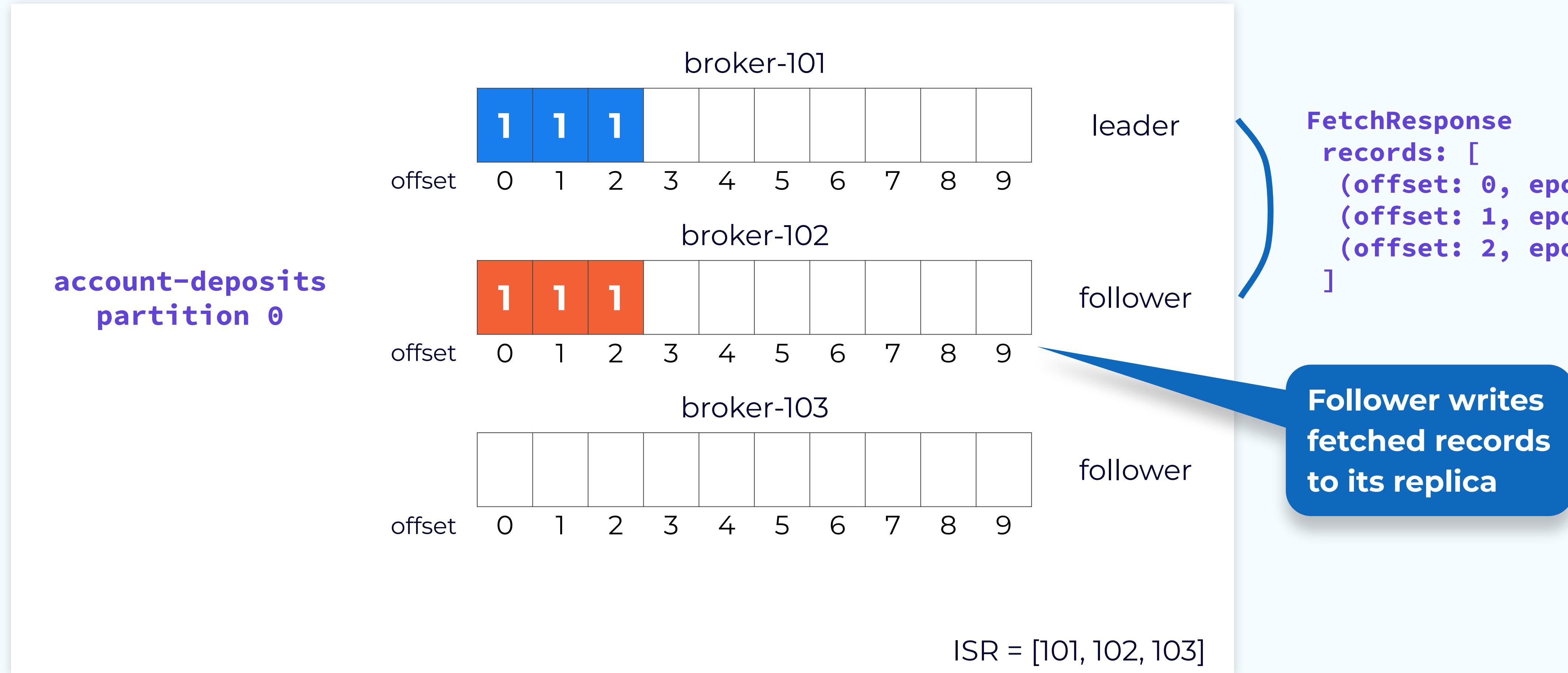
Leader Epoch



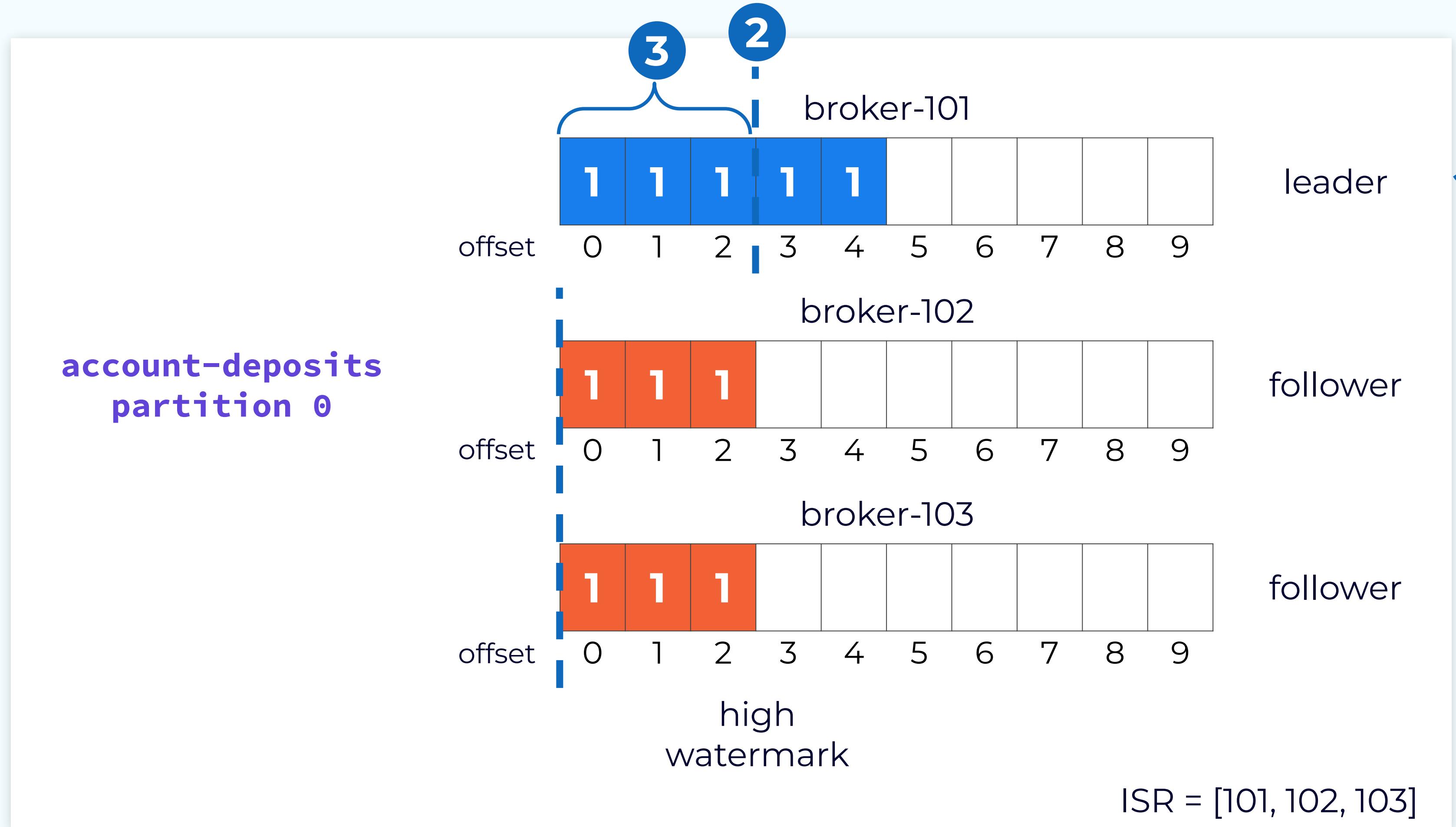
Follower Fetch Request



Follower Fetch Response

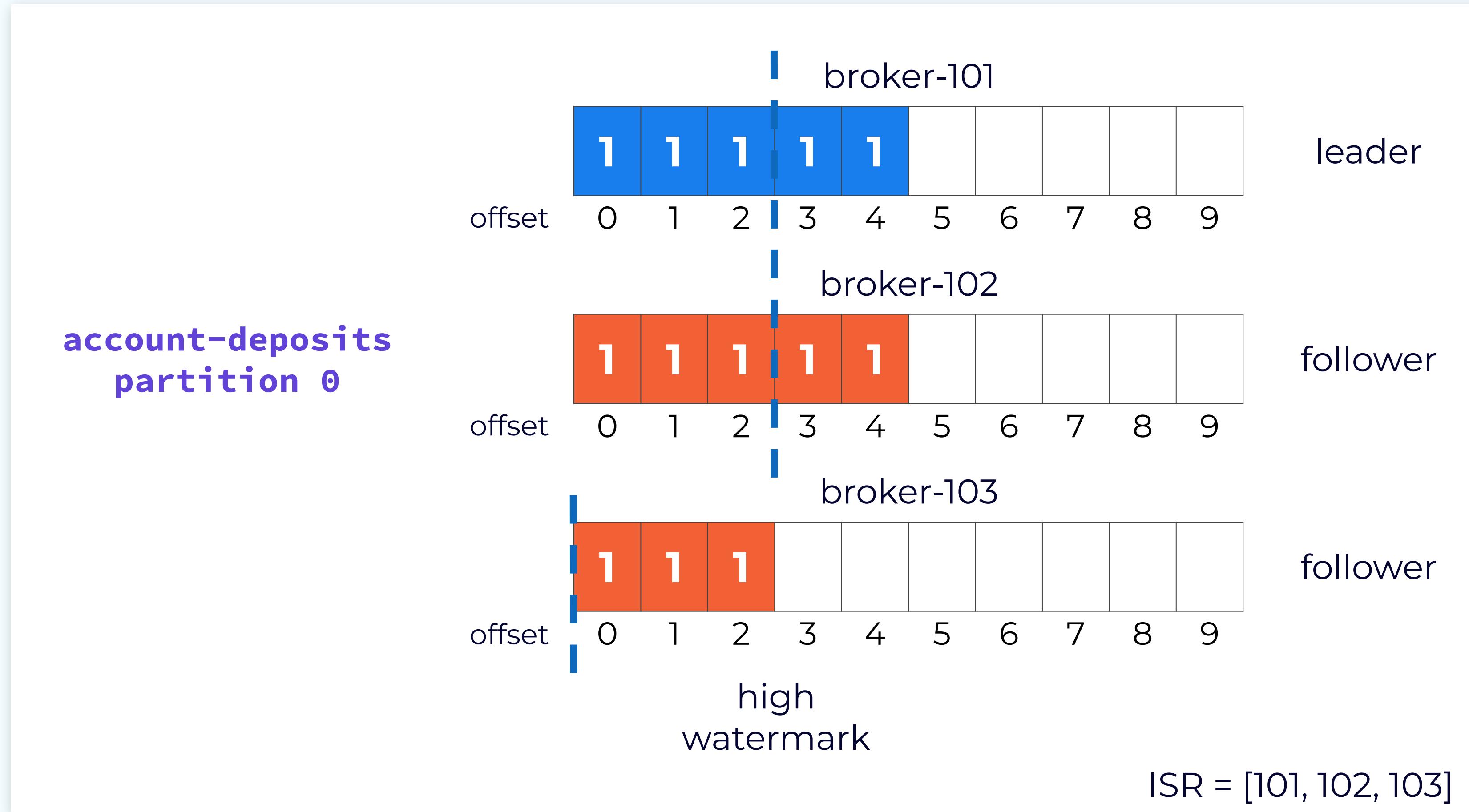


Committing Partition Offsets



- 1
 - FetchRequest
offset: 3
 - 1) Subsequent fetch request implicitly confirms receipt of previously fetched records
 - 2) Leader commits records once all followers in ISR have confirmed receipt
 - 3) Only committed records are exposed to consumers

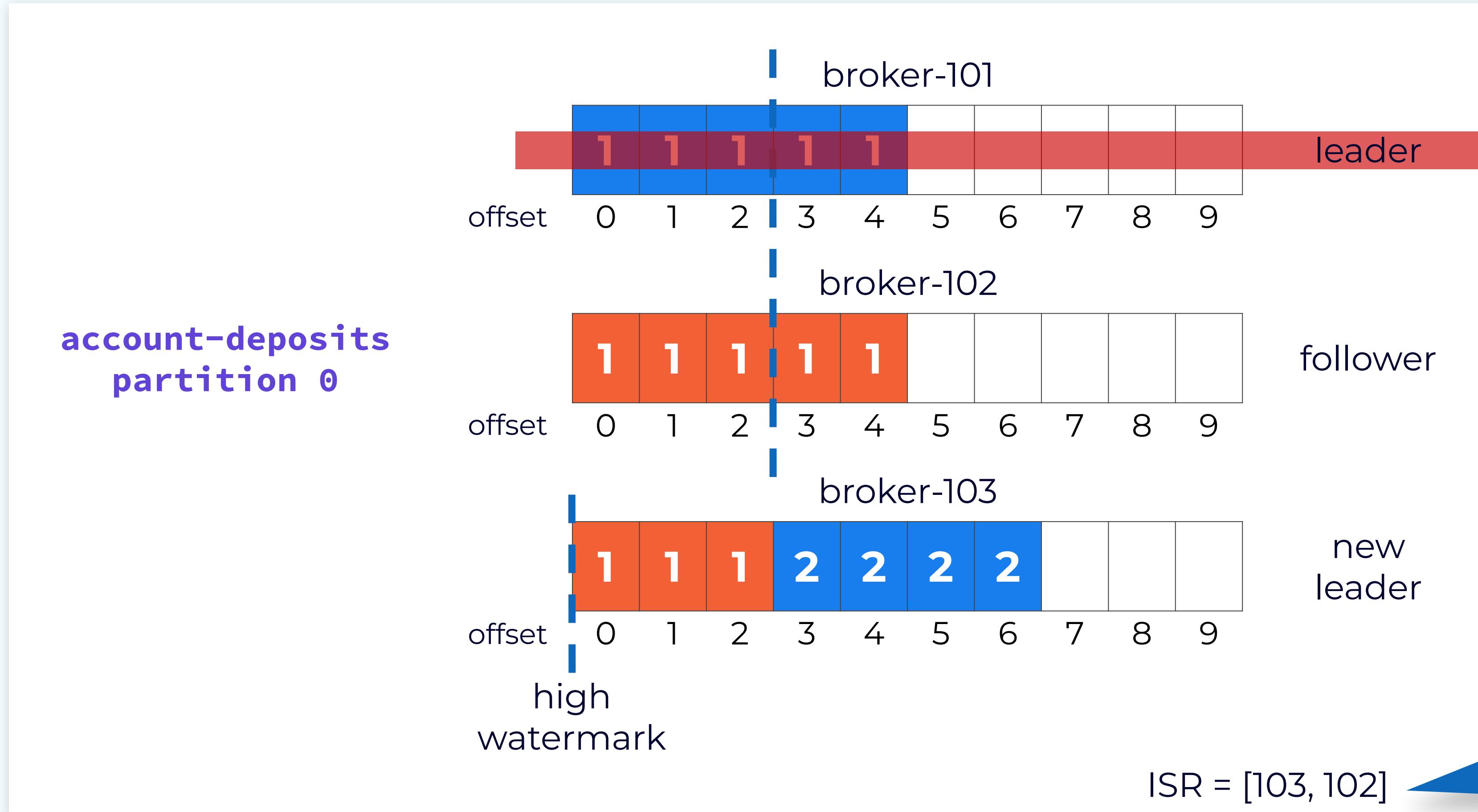
Advancing the Follower High Watermark



```
FetchResponse  
HighwaterMarkOffset: 3  
records: [  
  (offset: 3, epoch: 1),  
  (offset: 4, epoch: 1)  
]
```

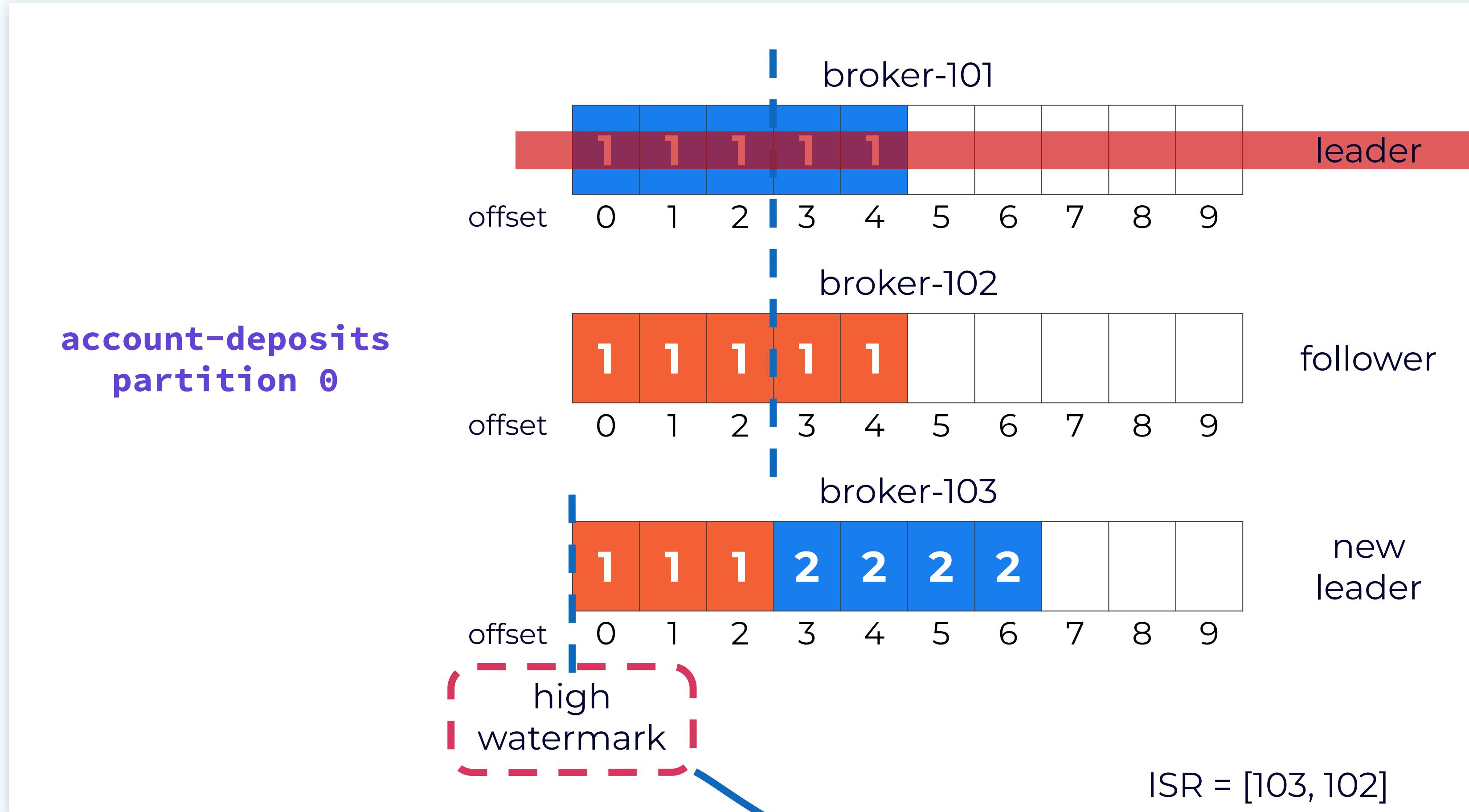
Subsequent fetch responses include new high watermark

Handling Leader Failure



New leader elected from ISR and propagated through control plane

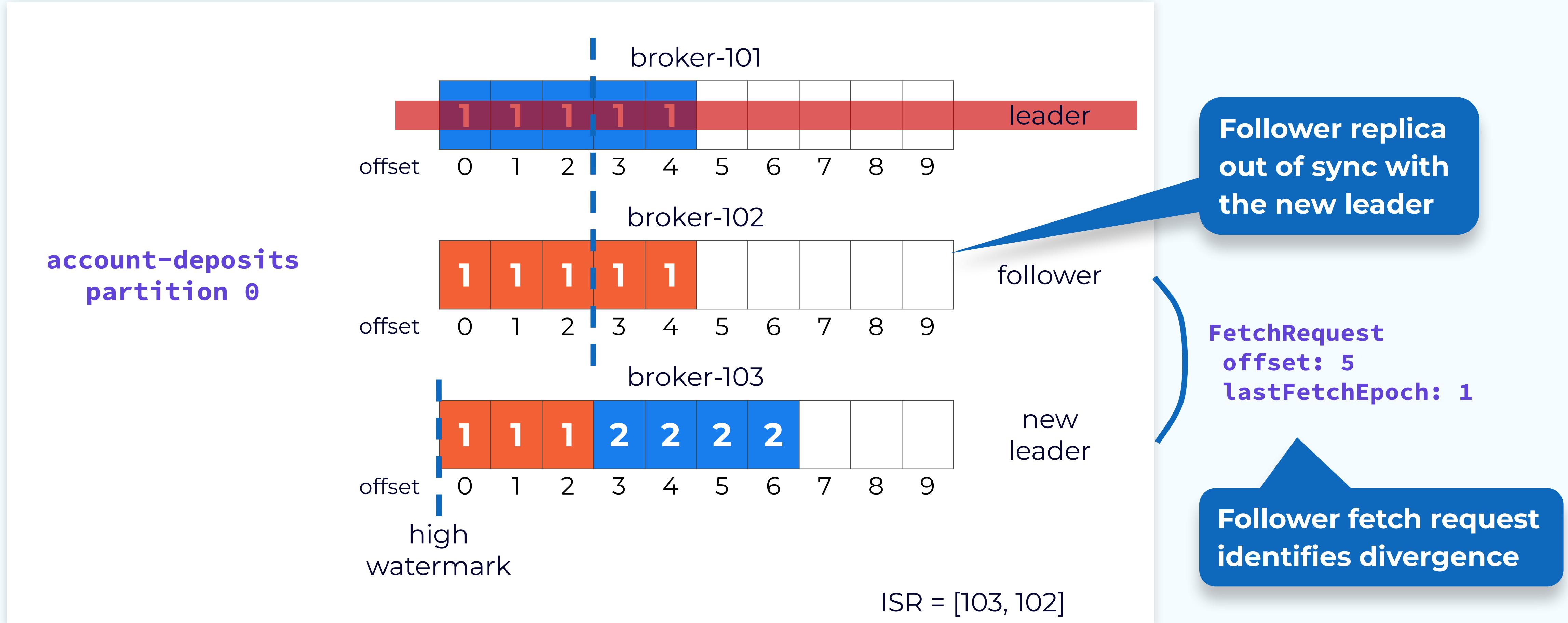
Temporary Decreased High Watermark



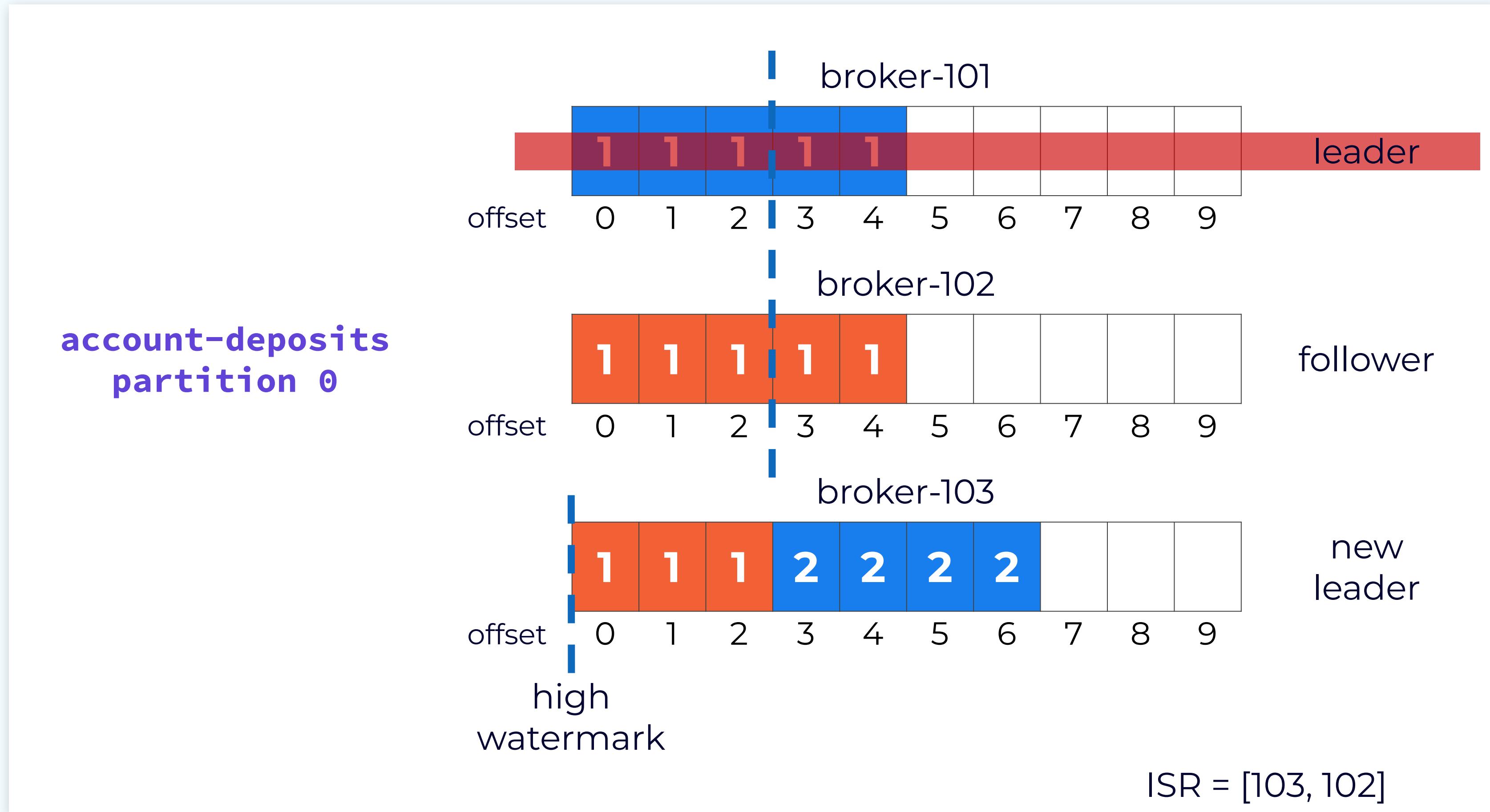
New leader high watermark could be less than true high watermark

- If consumer tries to read data in between, it gets a retrievable `OFFSET_NOT_AVAILABLE` error
- Once high watermark catches up, normal consumption continues

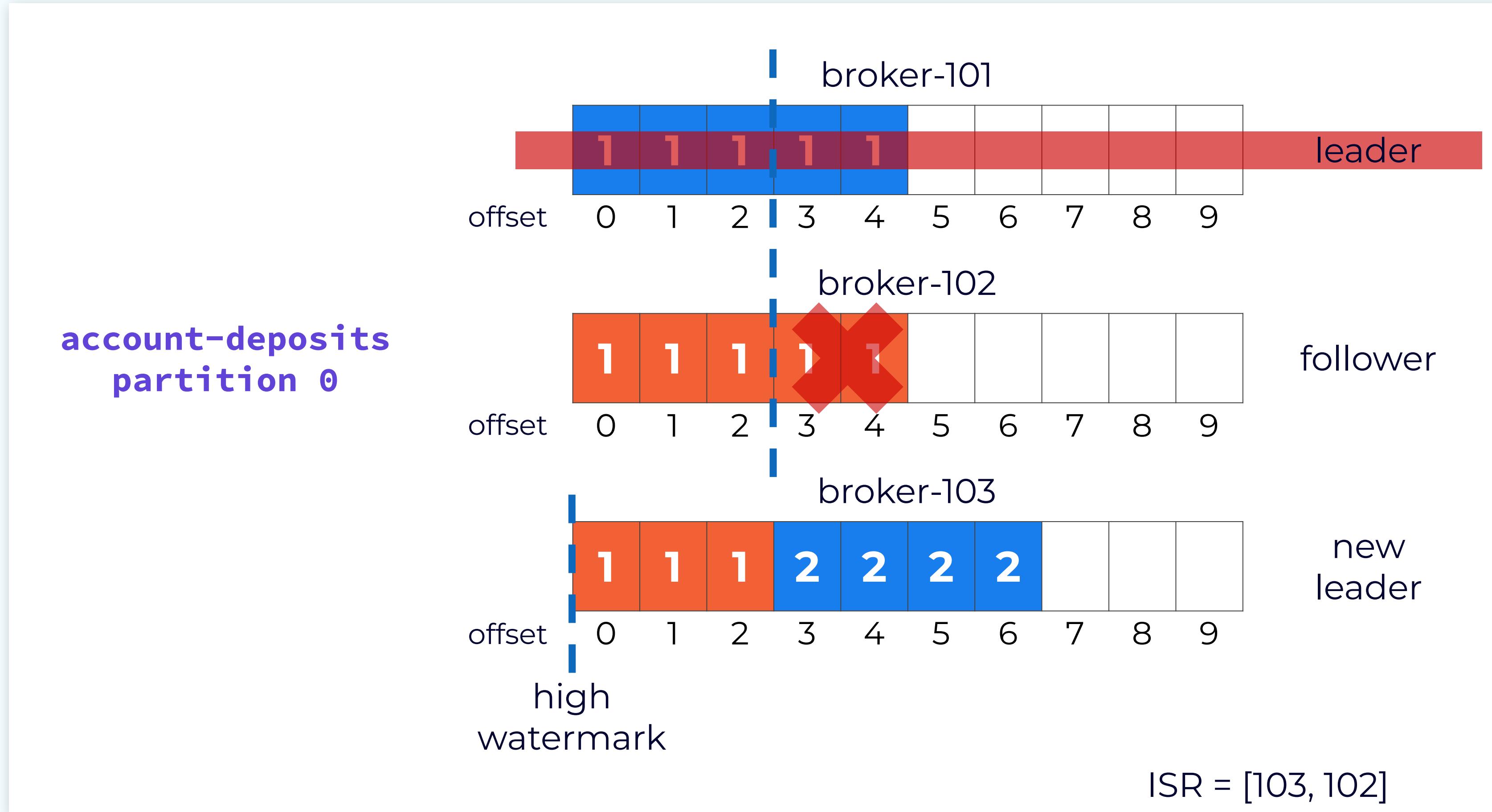
Partition Replica Reconciliation



Fetch Response Informs Follower of Divergence

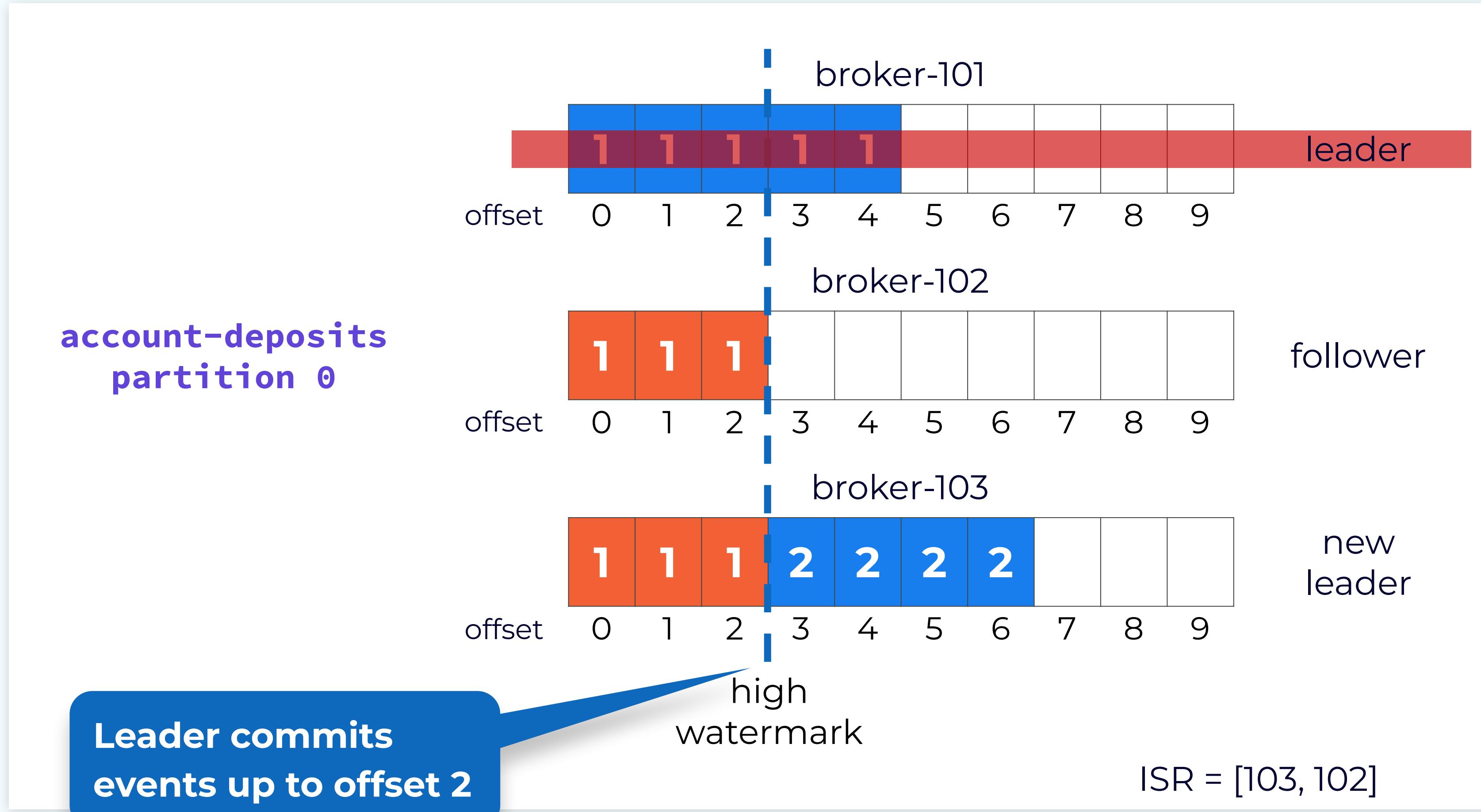


Follower Truncates Log To Match Leader Log

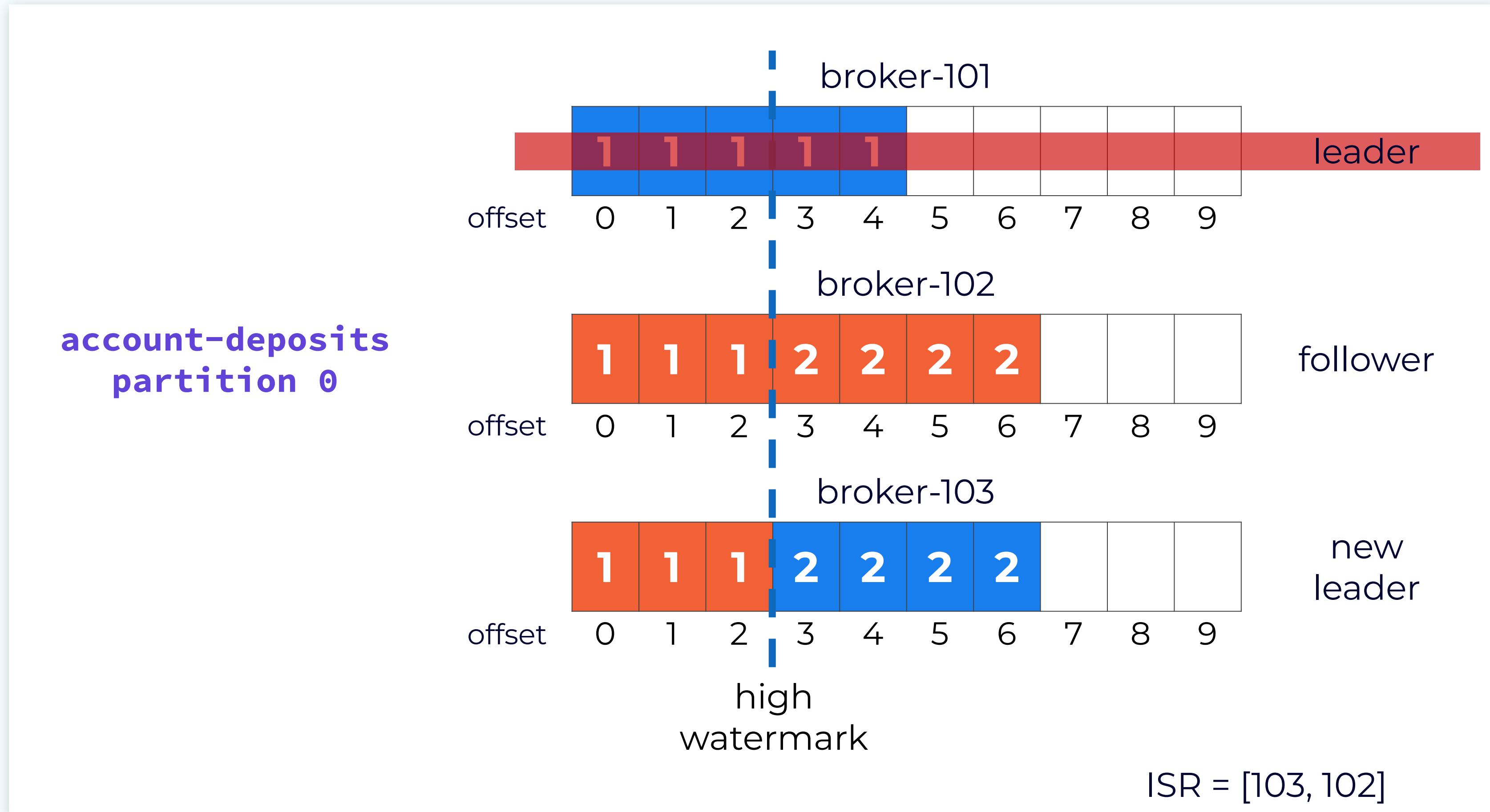


FetchResponse
HighwaterMarkOffset: 0
divergingEpoch:
epoch: 1
endOffset: 3

Subsequent Fetch with Updated Offset & Epoch



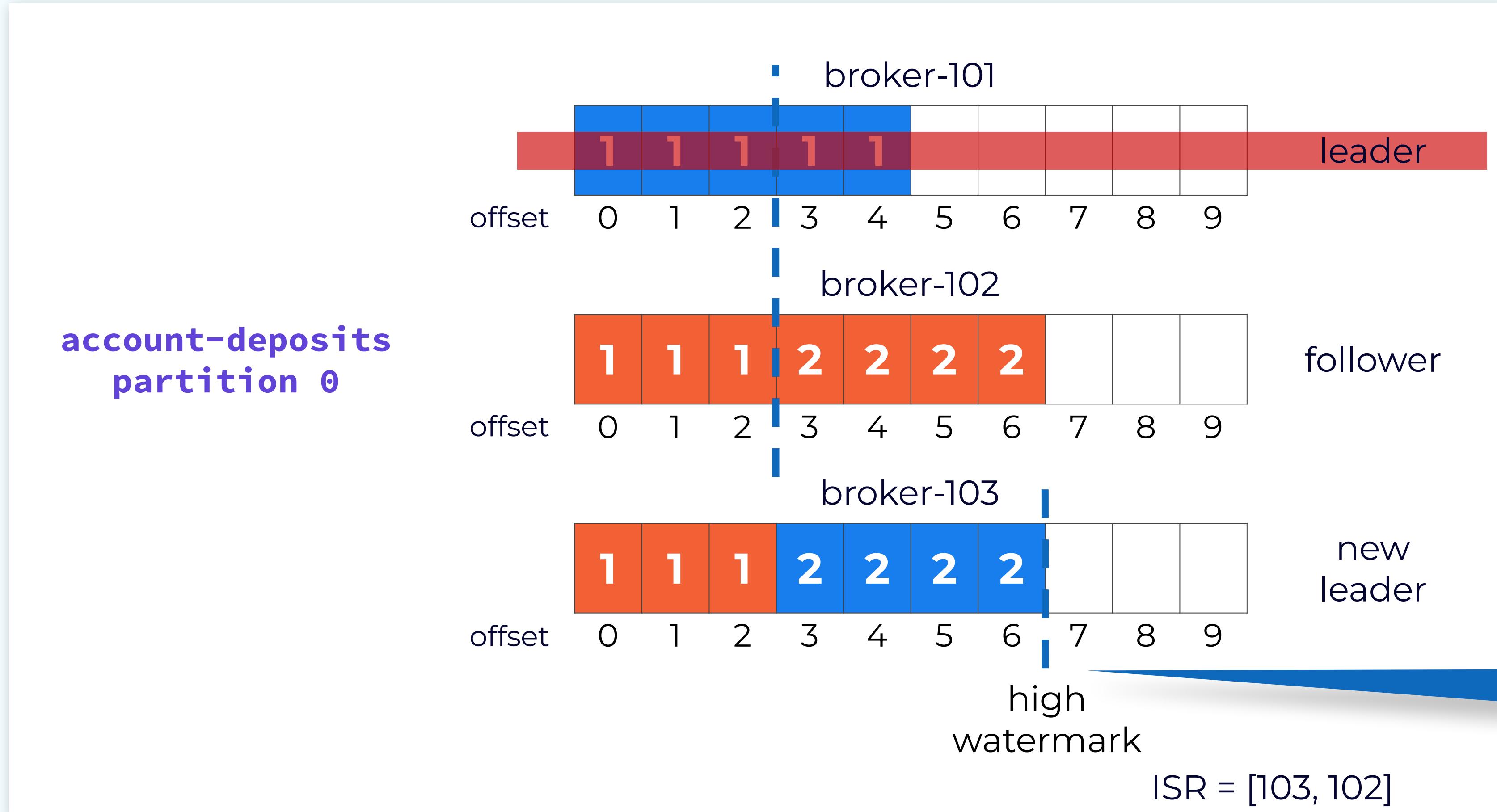
Follower 102 Reconciled



Leader sends requested offsets

FetchResponse
HighwaterMarkOffset: 3
records: [
 (offset: 3, epoch: 2),
 (offset: 4, epoch: 2),
 (offset: 5, epoch: 2),
 (offset: 6, epoch: 2)
]

Follower 102 Acknowledges New Records

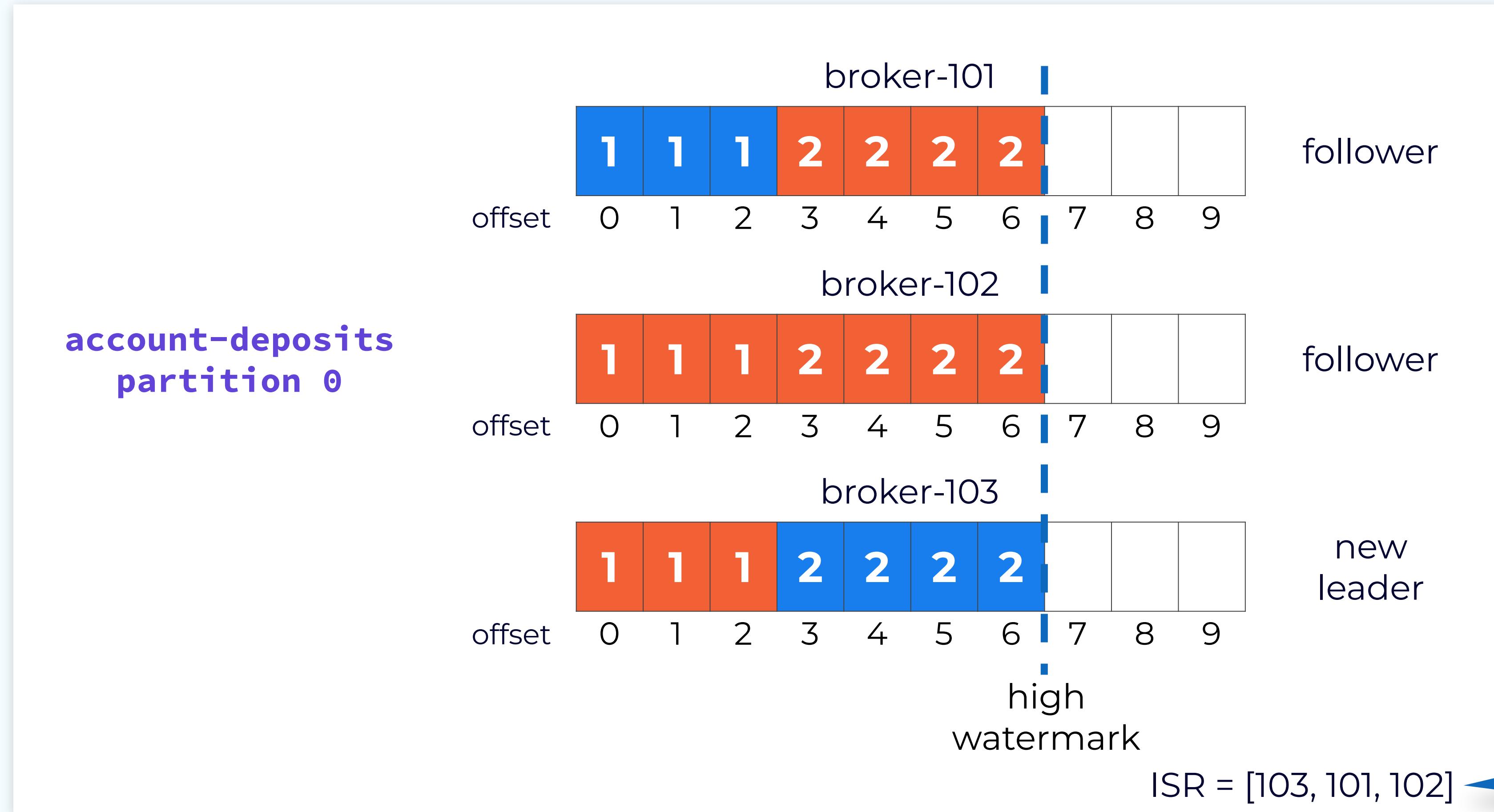


New request with latest offset and epoch values

FetchRequest
offset: 7
lastFetchEpoch: 2

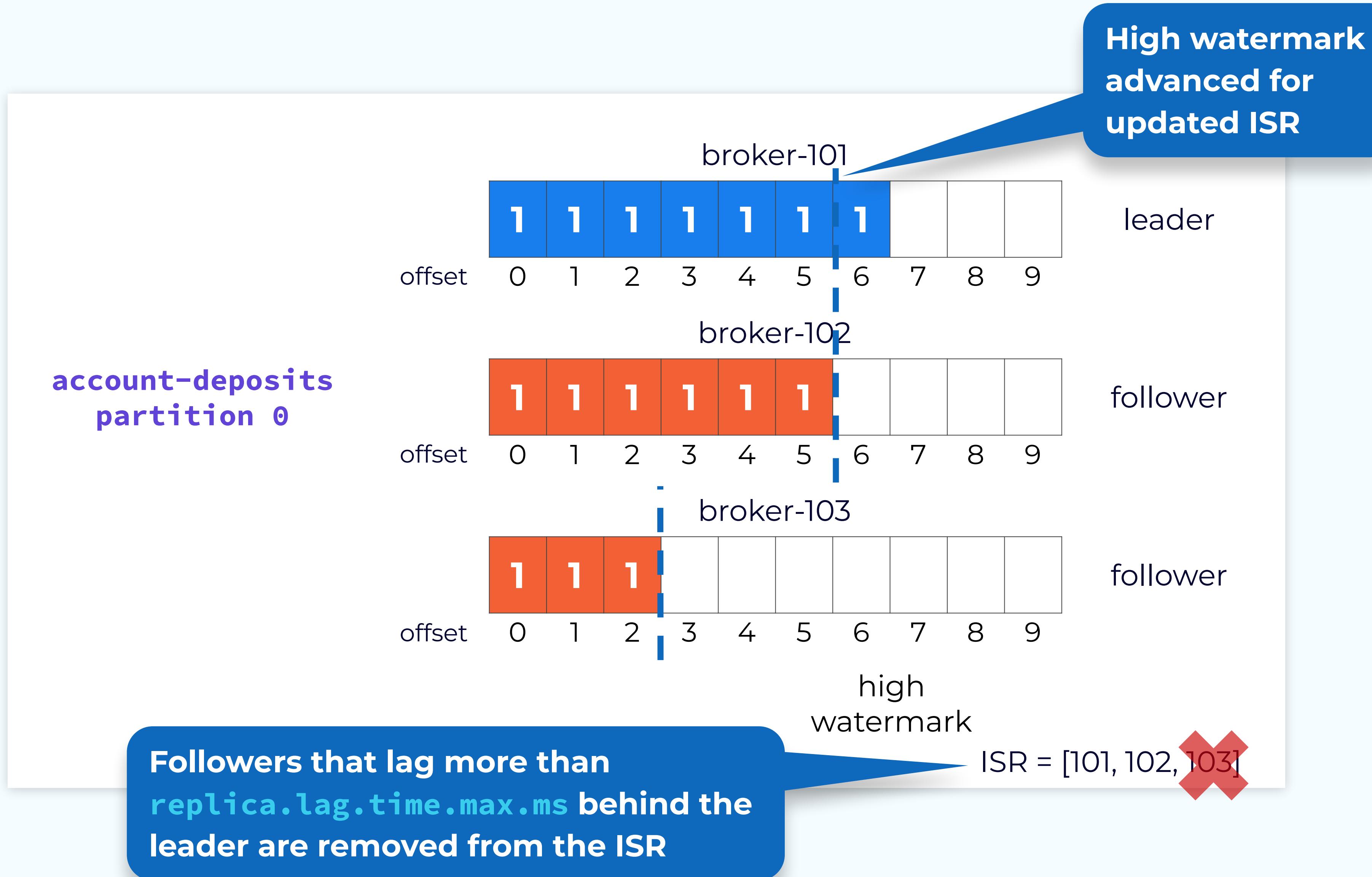
Leader advances its high watermark

Follower 101 Rejoins the Cluster



Broker-101 added
back to ISR

Handling Failed or Slow Followers

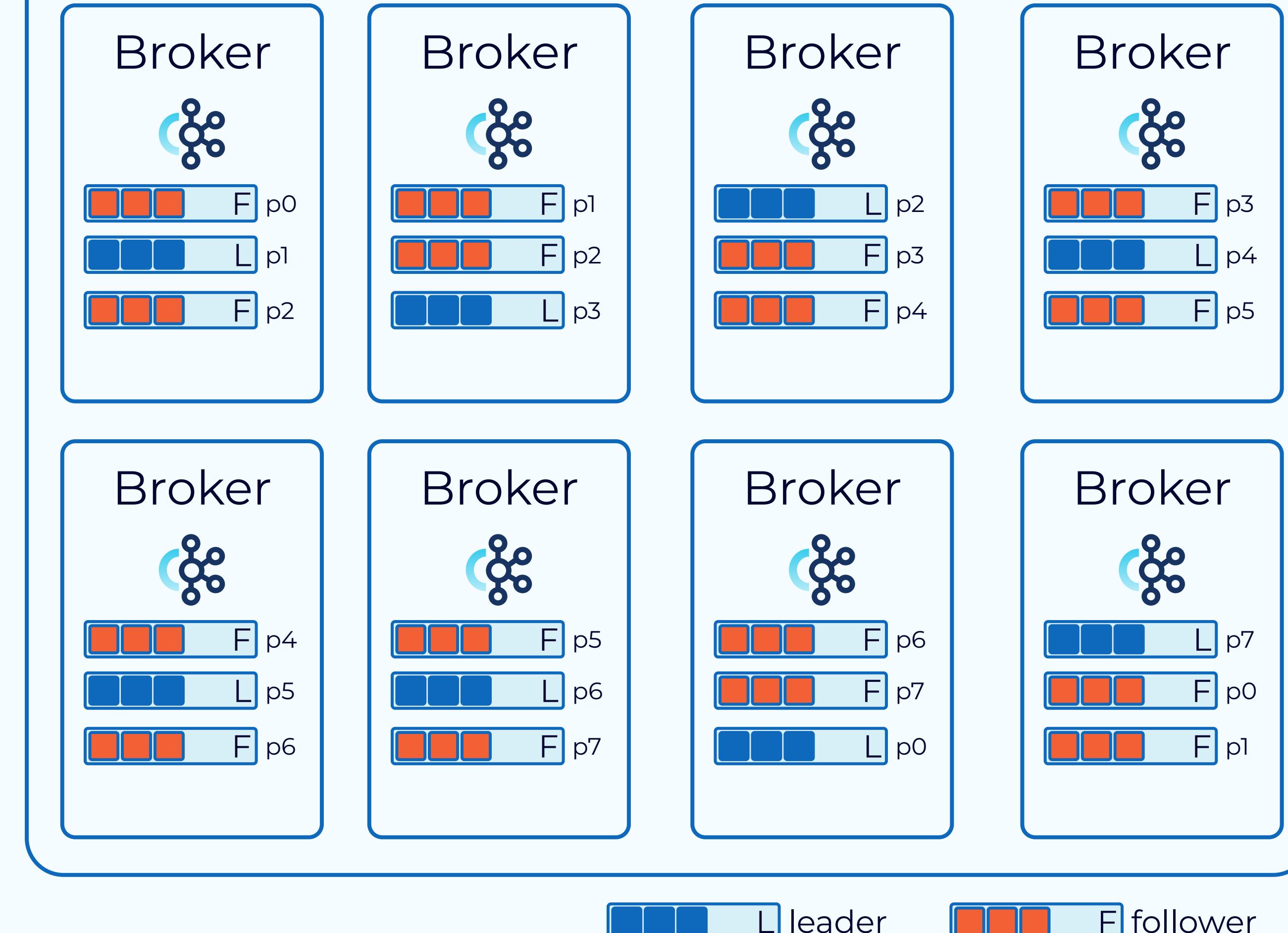


Partition Leader Balancing



- First replica considered preferred
- Preferred replica distributed evenly during assignment
- Background thread moves leader to preferred replica when it's in-sync

Kafka Cluster A



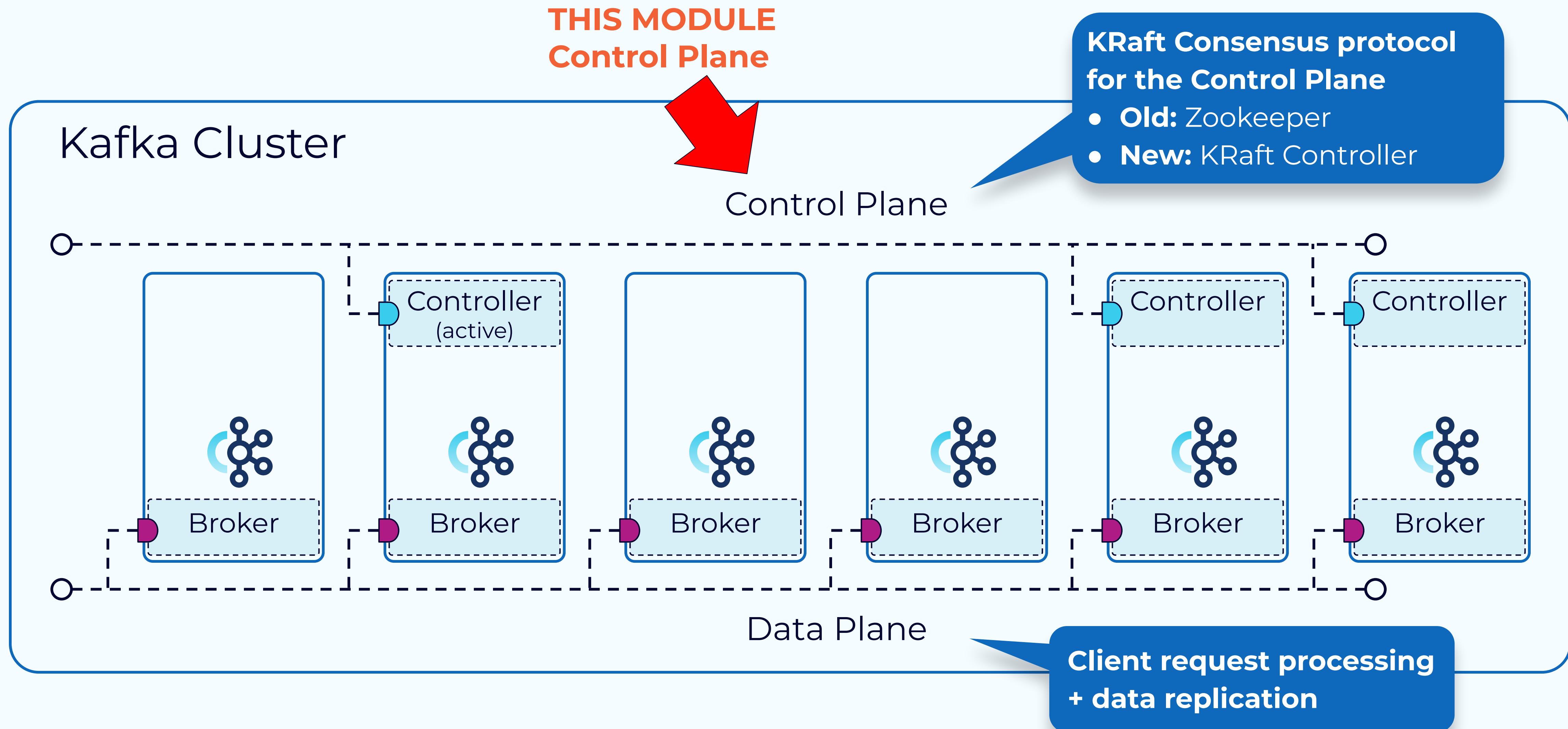
L leader

F follower

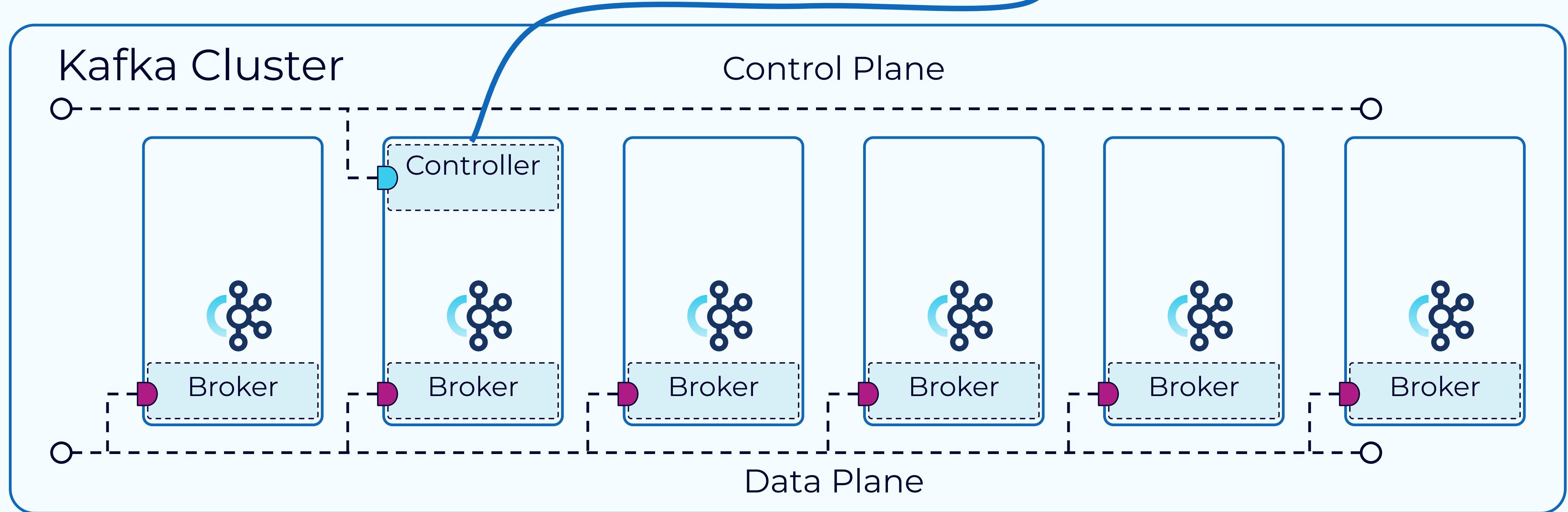
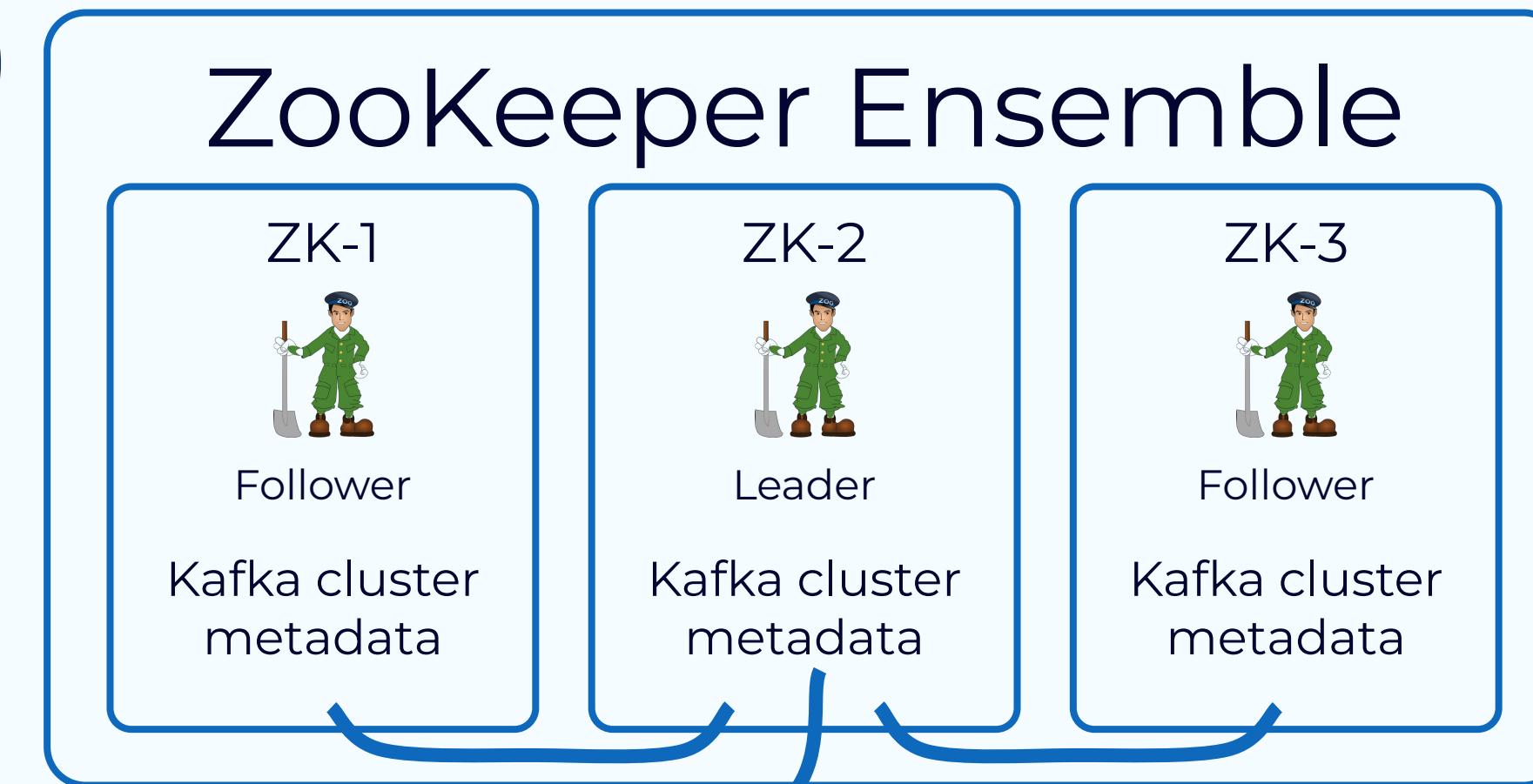


The Apache Kafka® Control Plane

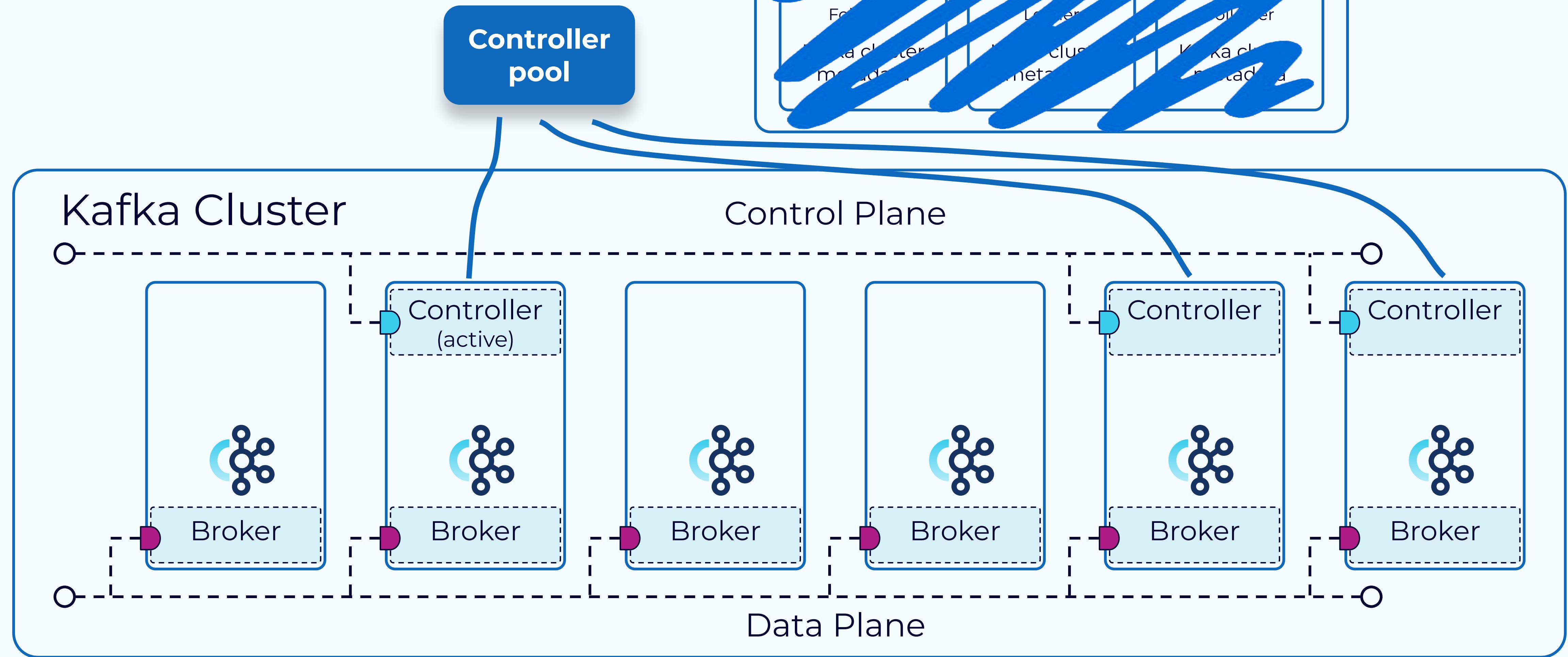
Kafka Manages Data & Metadata Separately



ZooKeeper Mode (Legacy)

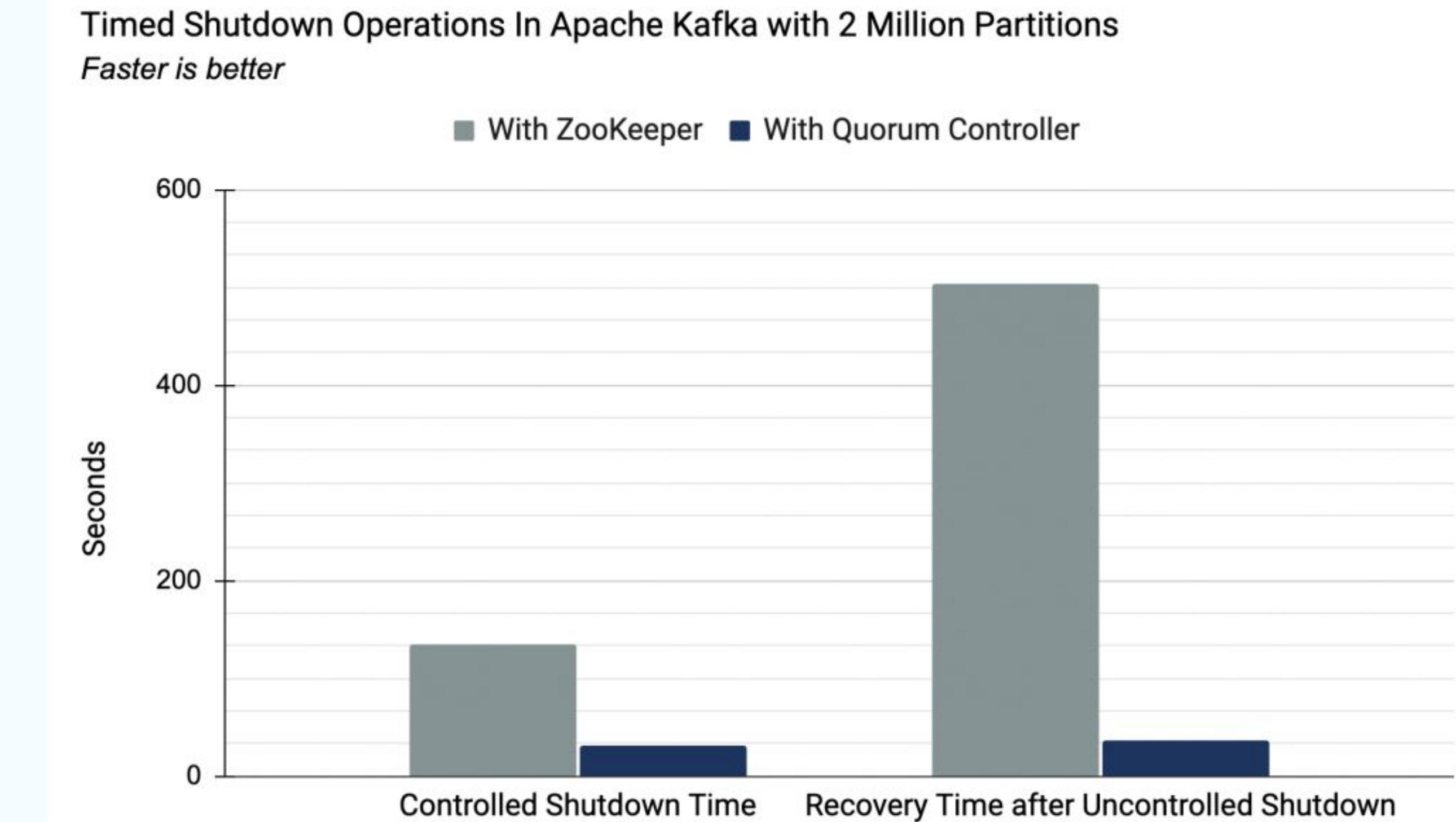


KRaft Mode (Nascent)



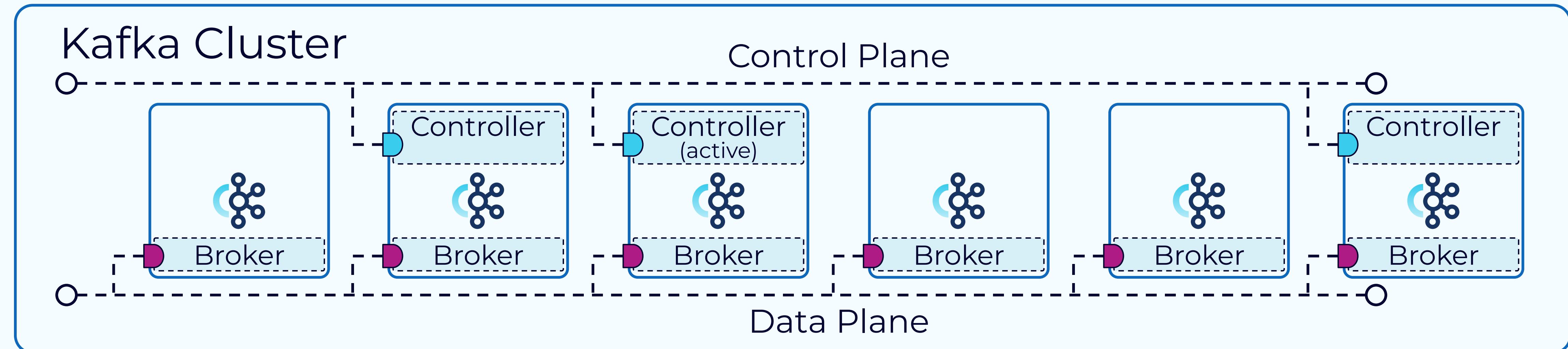
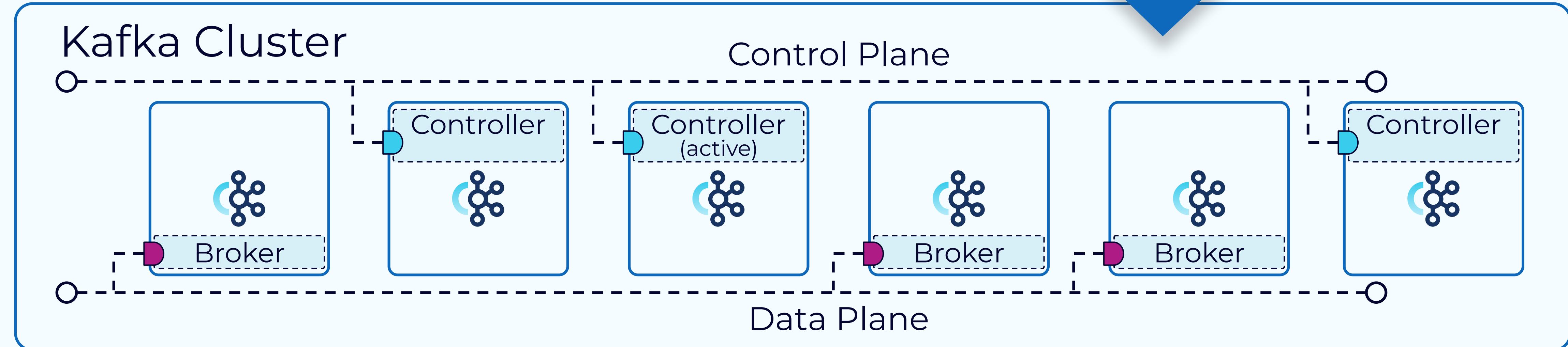
KRaft Mode Advantages

- Fewer moving parts make simpler operation
- Improved cluster scalability due to faster recovery (see chart)
- More efficient, incremental propagation of metadata to brokers

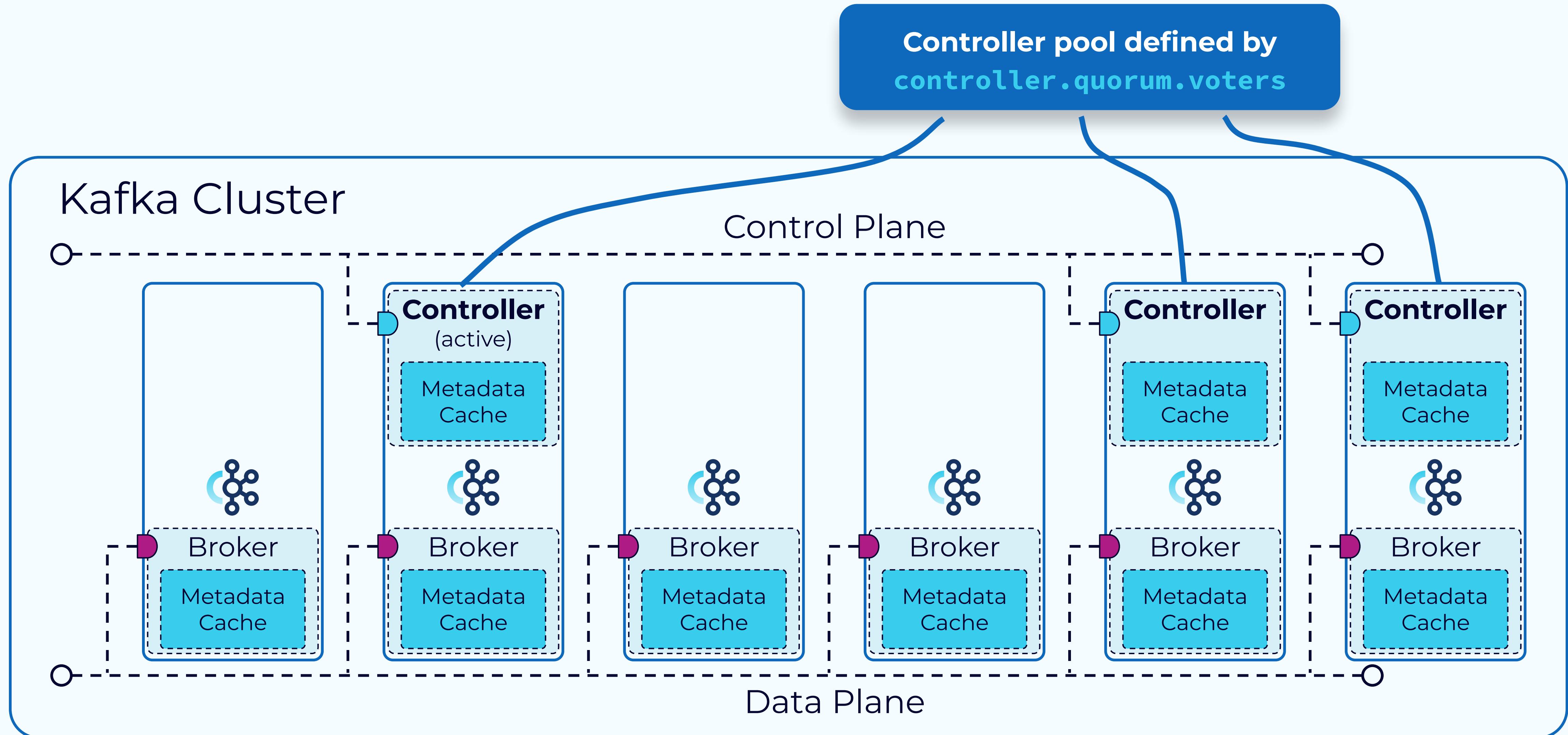


KRaft Cluster Node Roles

Cluster node role can be controller or broker or controller, broker



KRaft Mode Controller

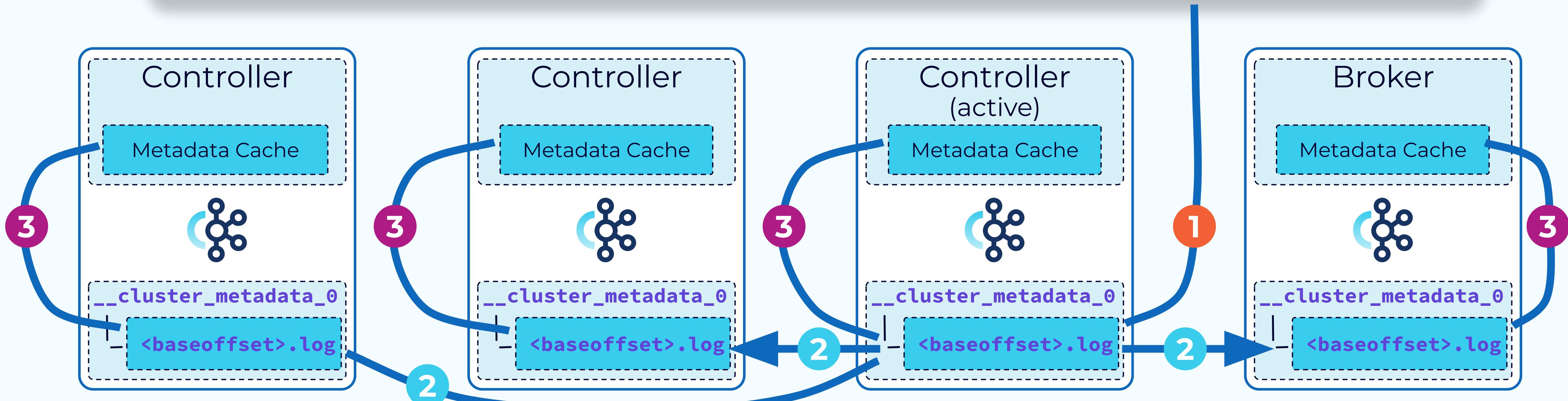


KRaft Cluster Metadata



- 1 Active controller writes the metadata change records to the `__cluster_metadata` topic, e.g.

```
{"topicId":"JJmWDIT9Q9eXVDvCaBIJhg", "partitionId":0, "leader":1, "isr": [1,2,3]}
```



- 2 Controllers (followers) and brokers (observers) replicate `__cluster_metadata` topic events.

- 3 Controllers and brokers update metadata cache from committed log records

KRaft Metadata Replication

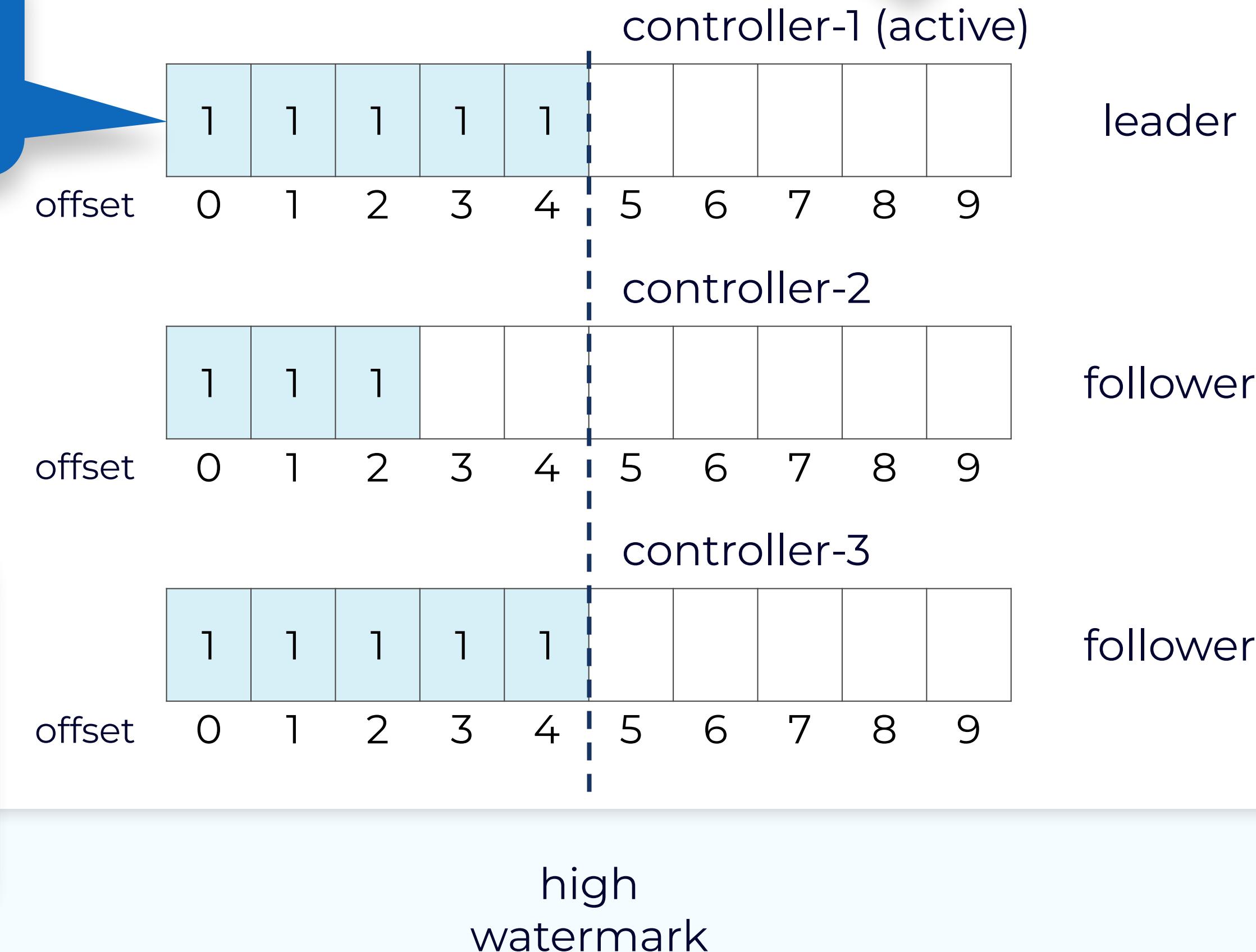
The active controller is the
`--cluster_metadata-0`
partition leader

Leader epoch is
included for a
batch of records

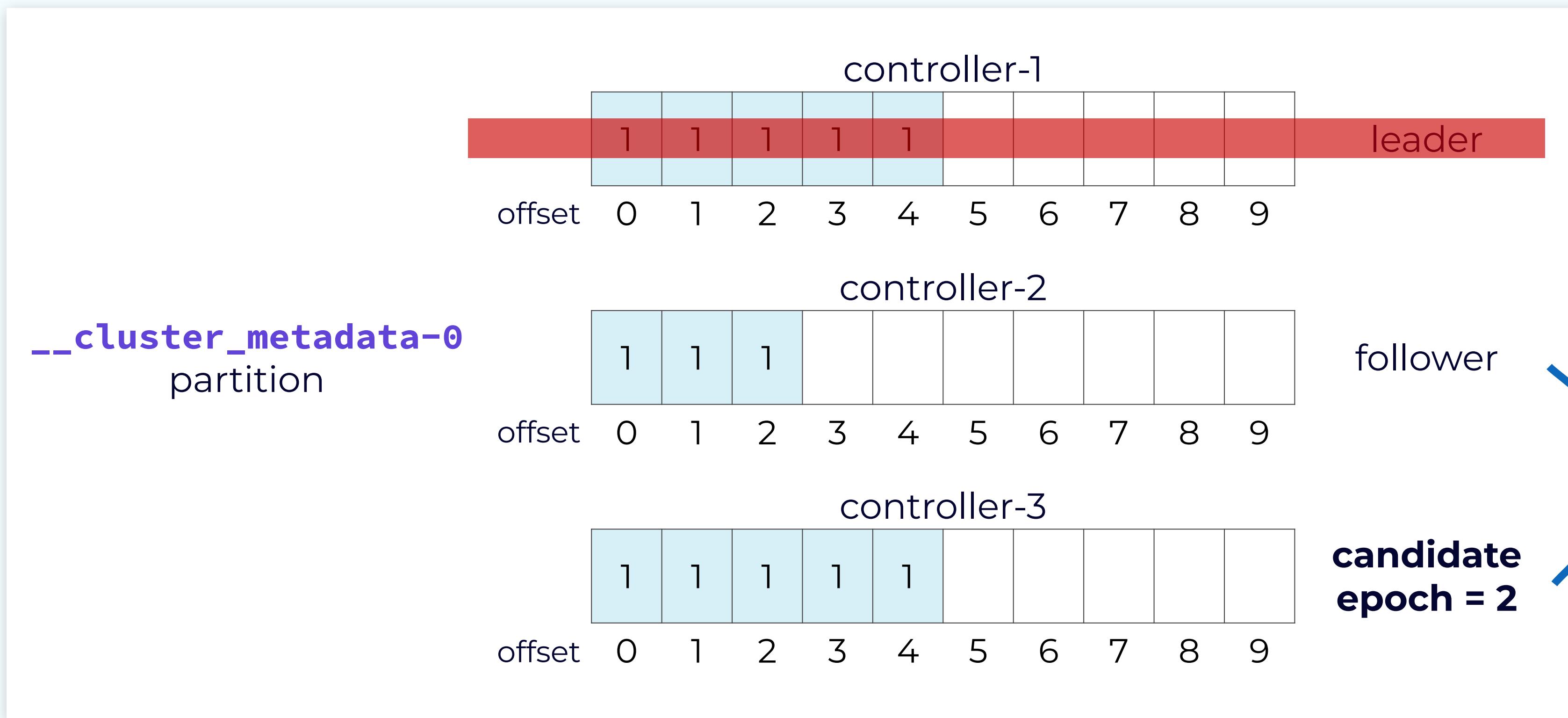
`--cluster_metadata-0`
partition

Difference to Kafka data replication:

- Leader election and committing record: quorum instead of ISR
- Records are fsync to disk prior to commit



Leader Election Step 1: Vote Request

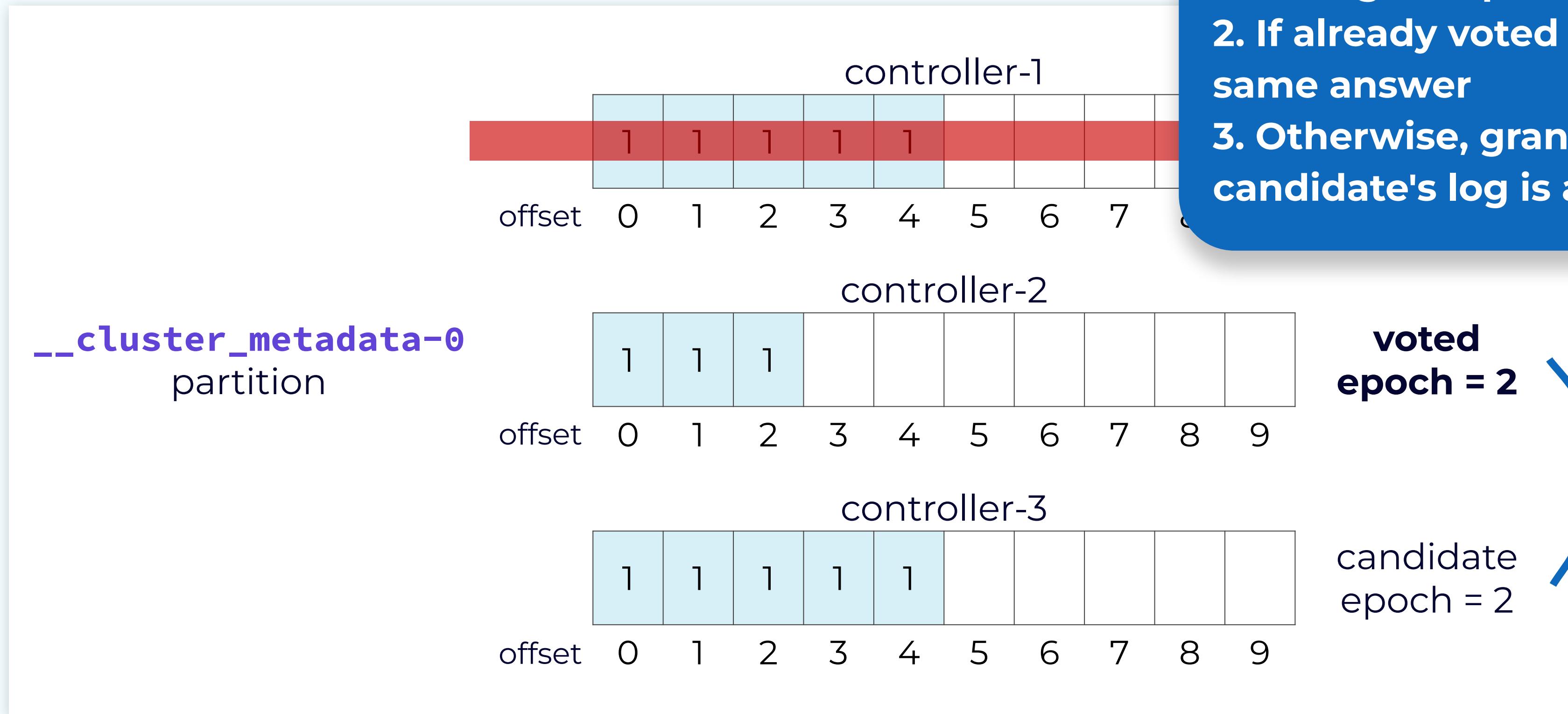


In response to an absent leader, a controller:

- Increases its epoch
- Transitions to candidate
- Votes for itself
- Requests votes from other controllers

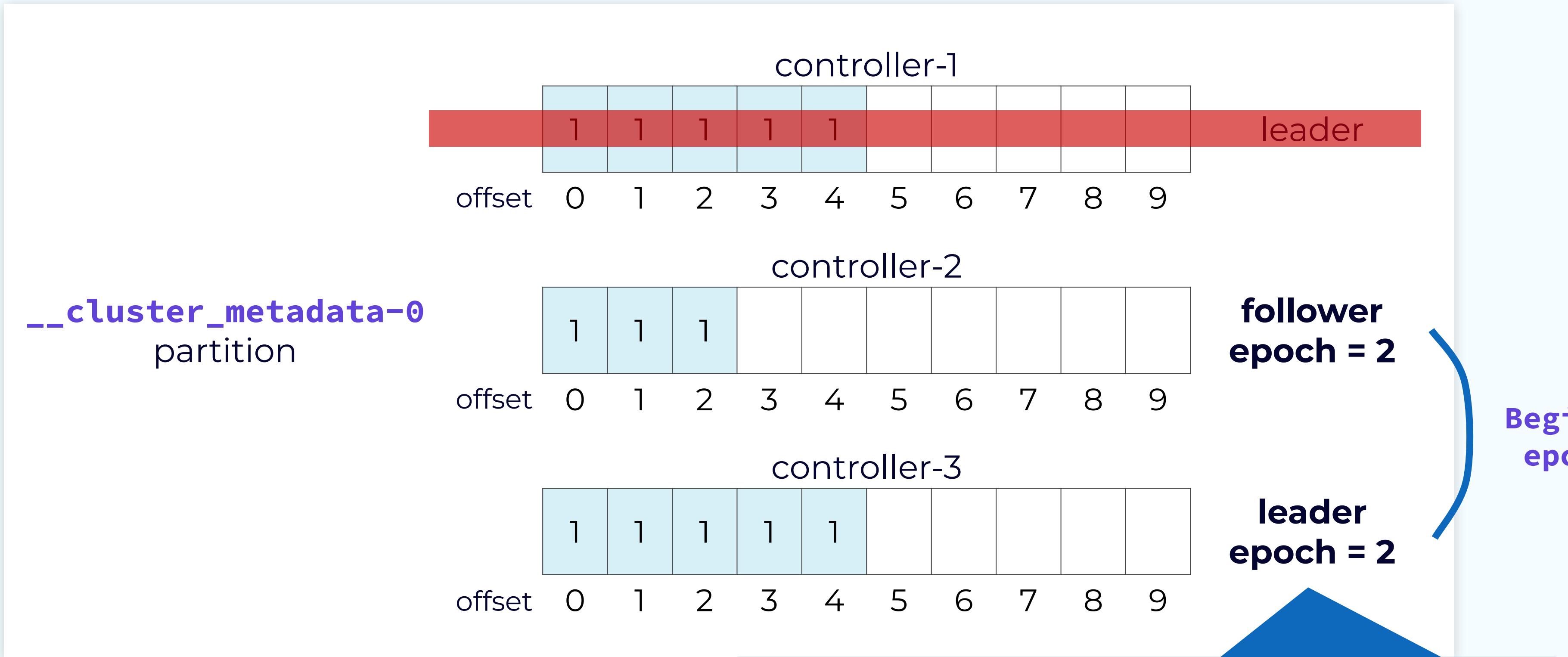
VoteRequest
candidateEpoch: 2
lastOffset: 4
lastOffsetEpoch: 1

Leader Election Step 2: Vote Response



③ controller-3 lastOffset: 4 lastOffsetEpoch: 1 \geq controller-2 lastOffset: 2 lastOffsetEpoch: 1

Leader Election Step 3: Completion

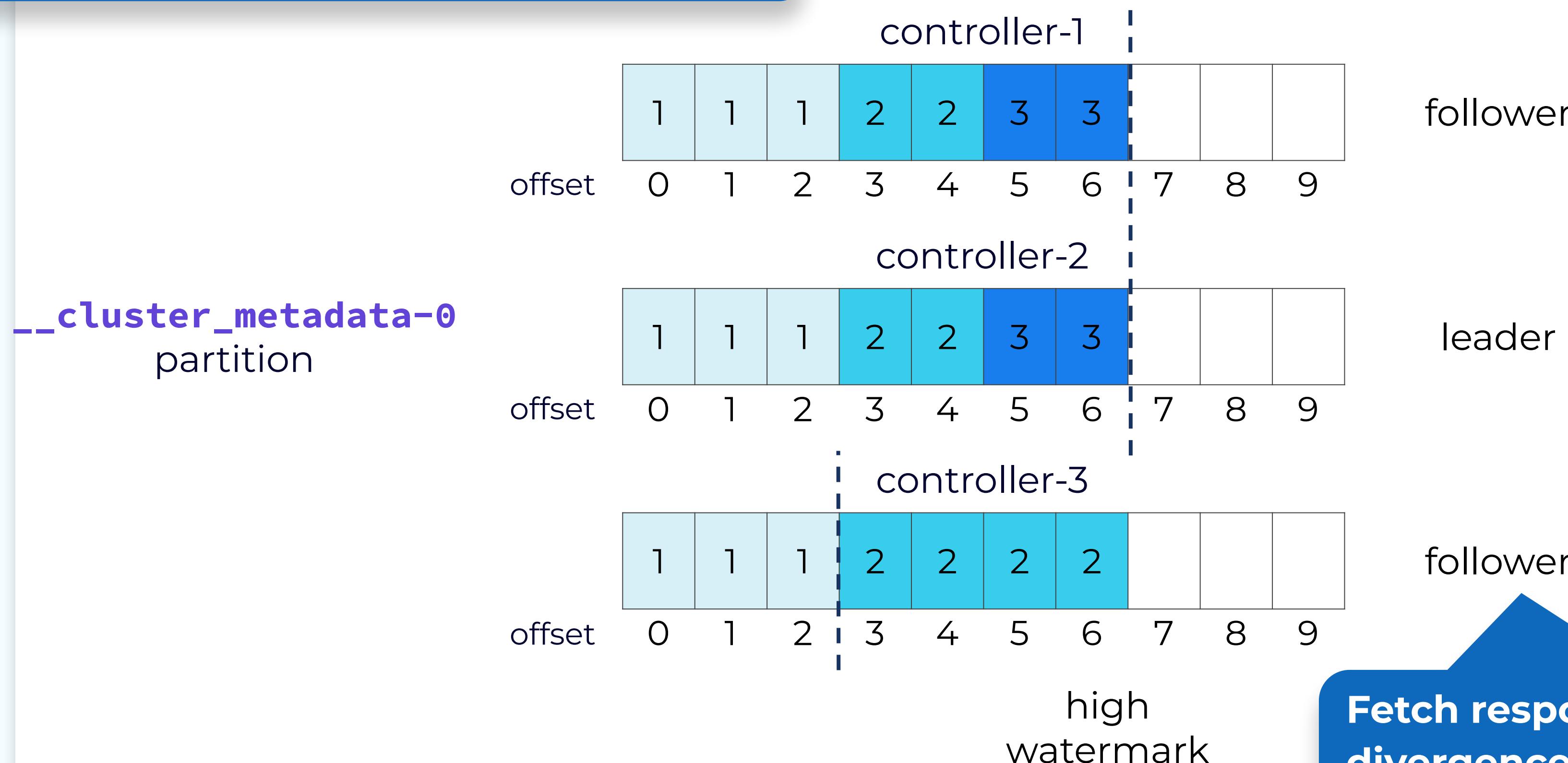


When a candidate receives the majority votes, it becomes the new leader and notifies other controllers the election is complete

Metadata Replica Reconciliation

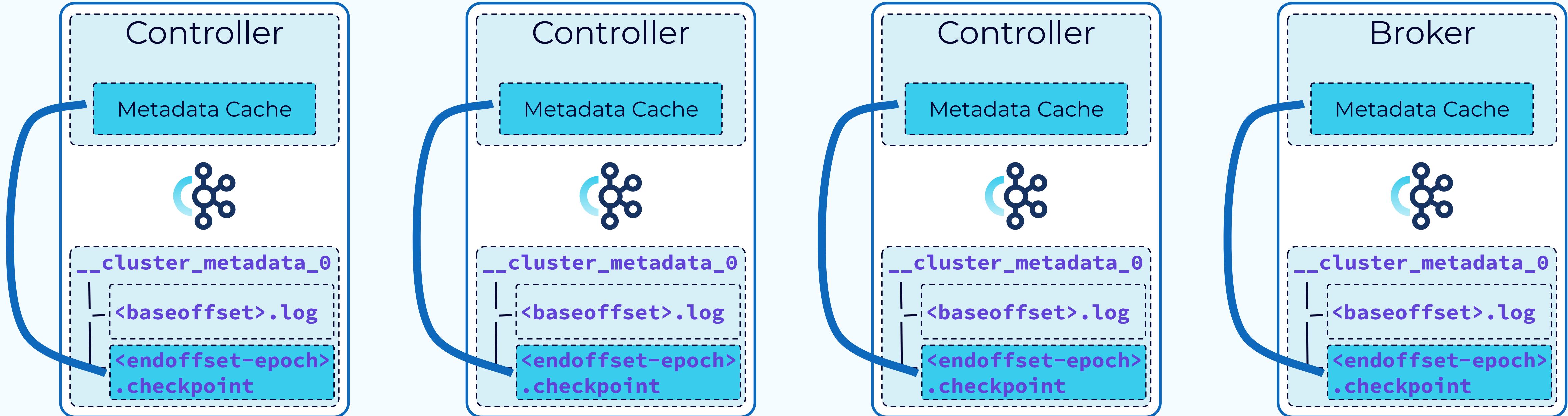


The same reconciliation process is used for both KRaft metadata and Kafka topic data



Fetch response will tell controller-3 divergence exists and it will truncate its replica to offset 5

KRaft Cluster Metadata Snapshot



0000000000000000128-0000000001.checkpoint

Periodic point-in-time snapshot of in-memory metadata cache

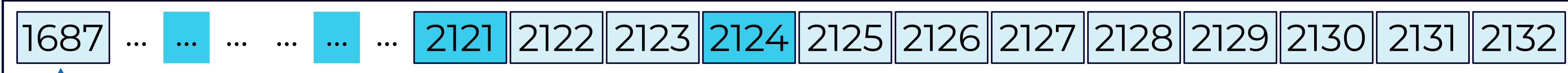
- Uniquely identified by its **EndOffset** and **Epoch**
- Stores only one entry per key
- A bit like a “compacted topic”

Metadata Log Truncation with Snapshot



Log of events, i.e. multiple events,
potentially for the same key, in order

`__cluster_metadata-0/00000000000000001687.log`

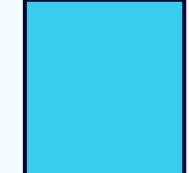


LogStartOffset

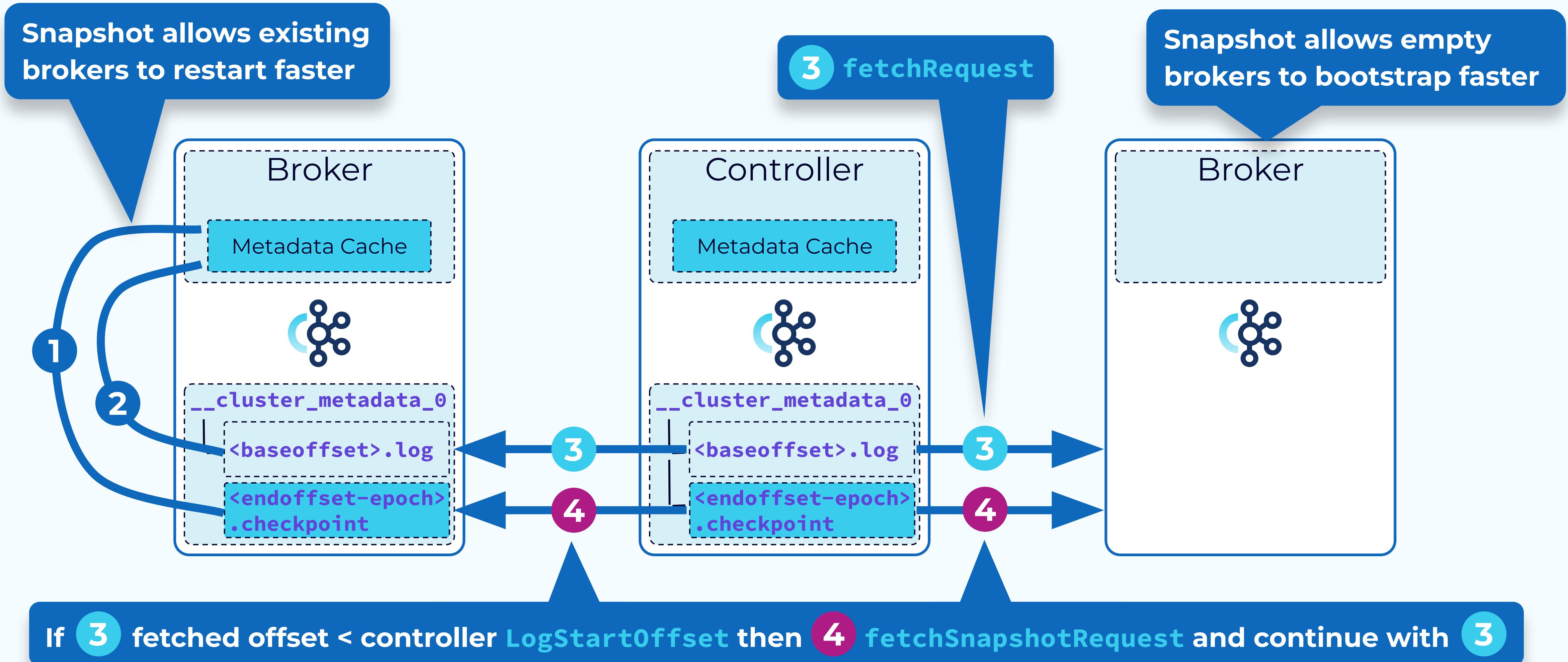
Snapshot → 2129

Log growth is limited since it is safe
to truncate the beginning of the log
up to the end of the latest snapshot

`00000000000000002130-0000000003.checkpoint`

 - redundant record for its key
- not included in snapshot

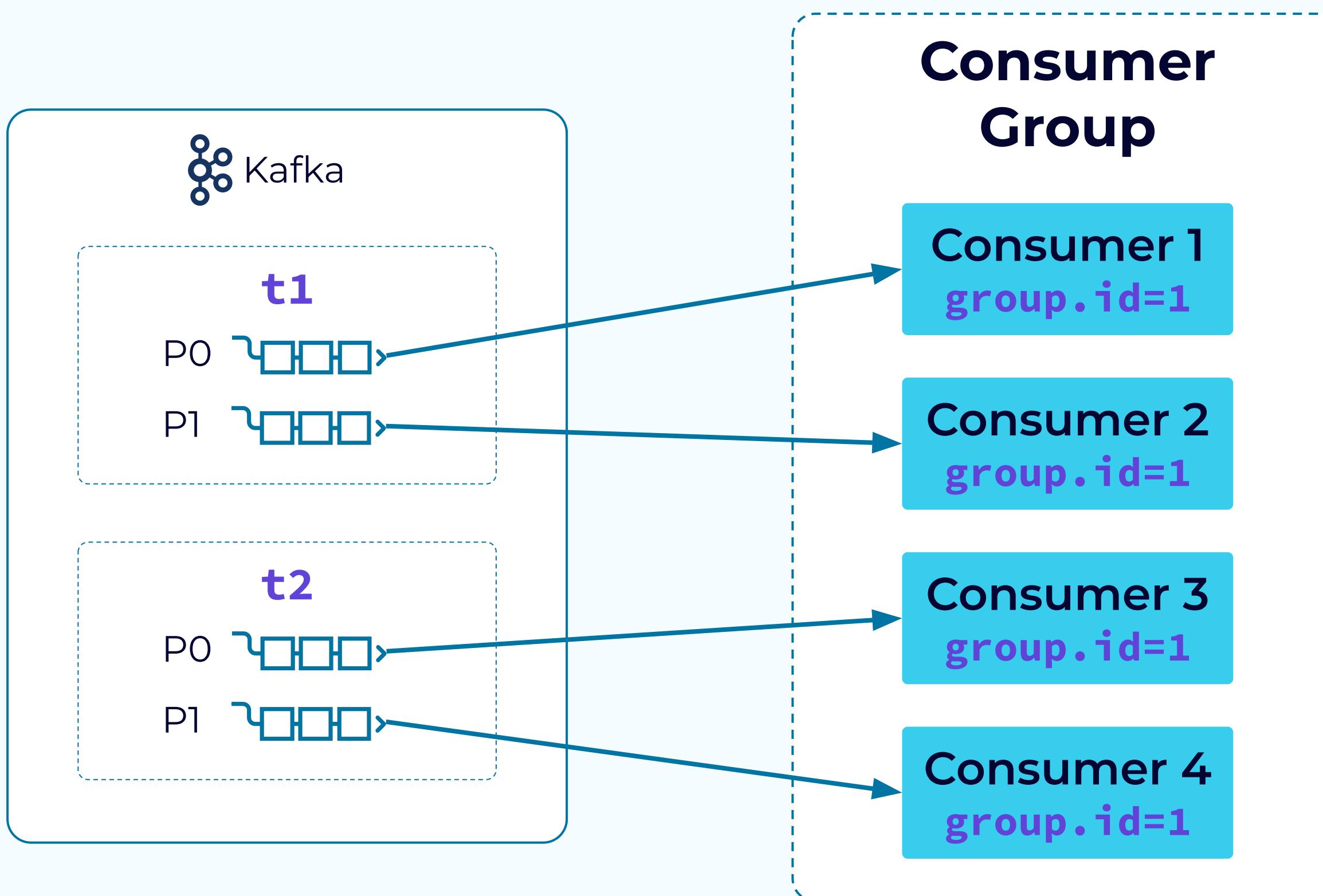
When a Snapshot Is Read





Consumer Group Protocol

Kafka Consumer Group



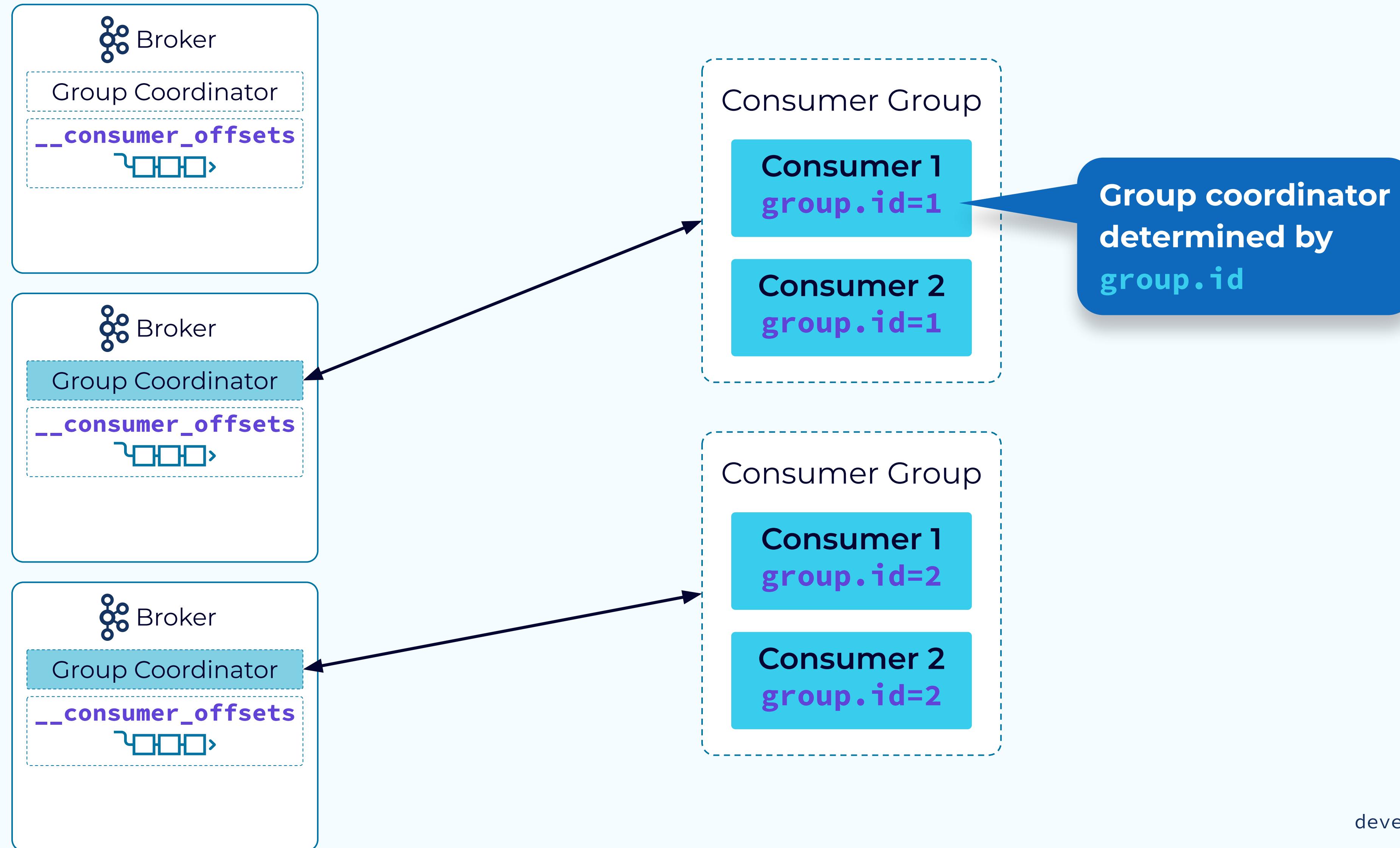
To enable group consumption:

- 1) **Configure group.id**
- 2) **topics={"t1", "t2"};**
consumer.subscribe(topics)

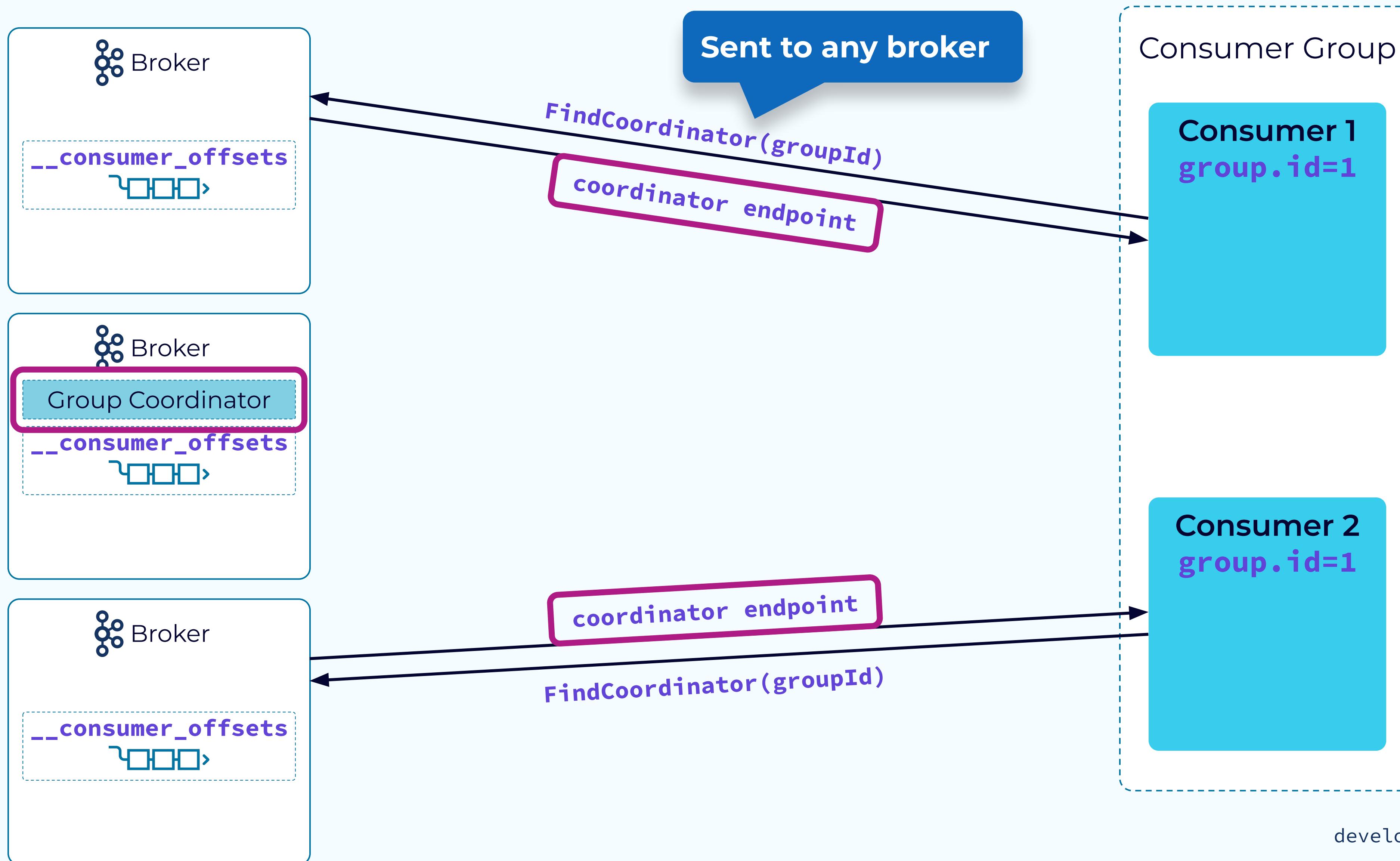
Benefits:

- 1) **Scalability**
- 2) **Elasticity**
- 3) **Fault tolerance**

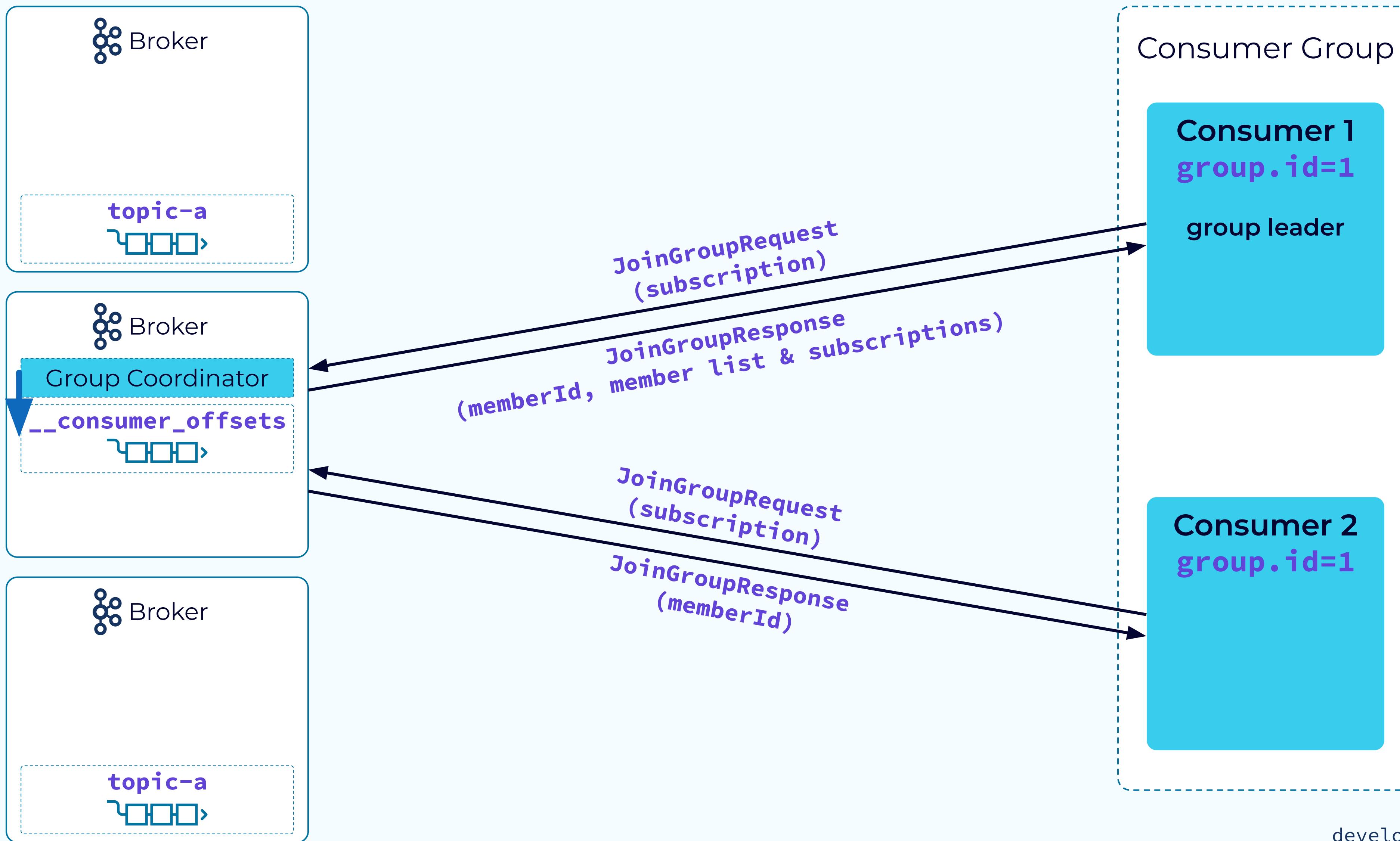
Group Coordinator



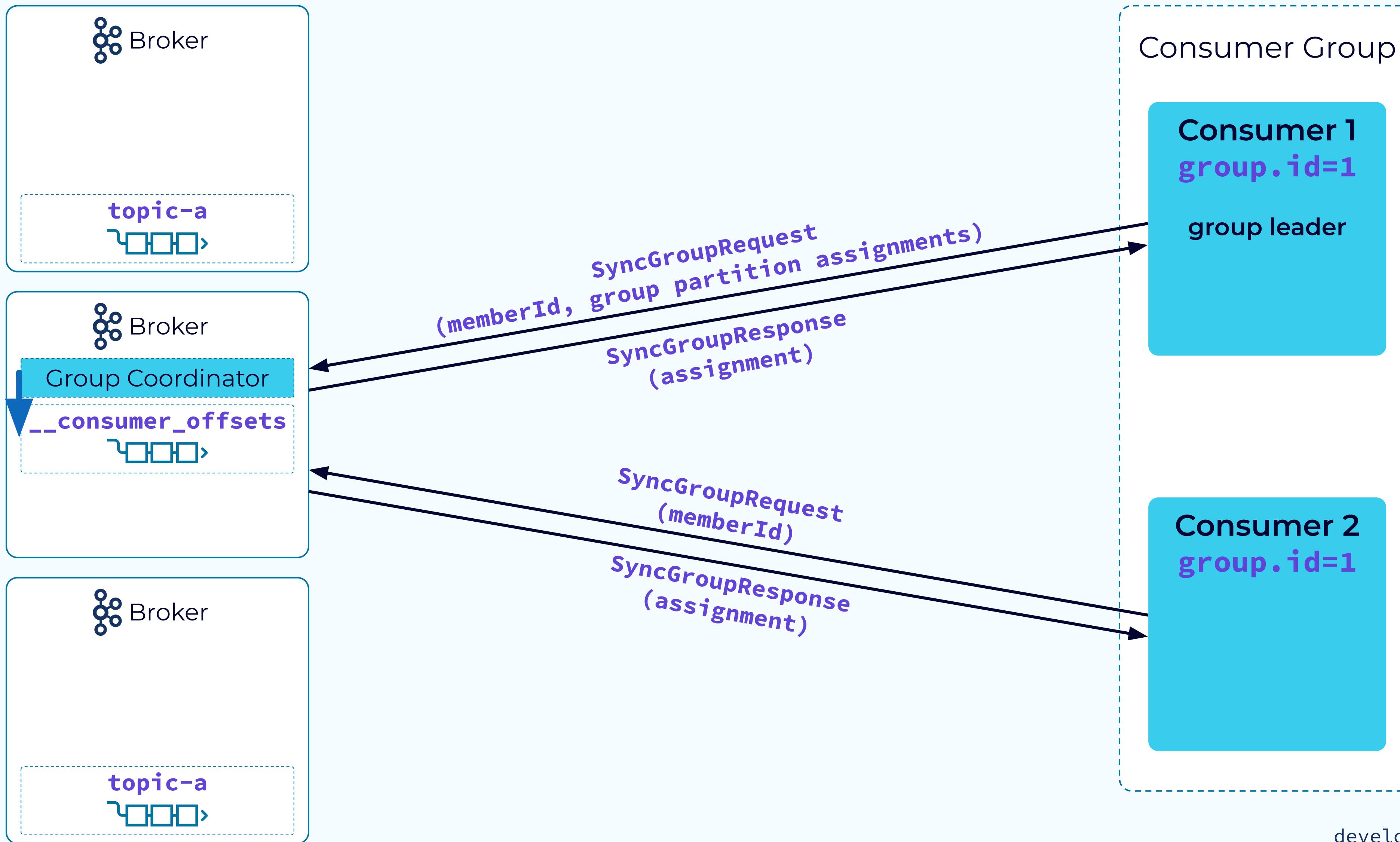
Group Startup: Step 1 - Find Group Coordinator



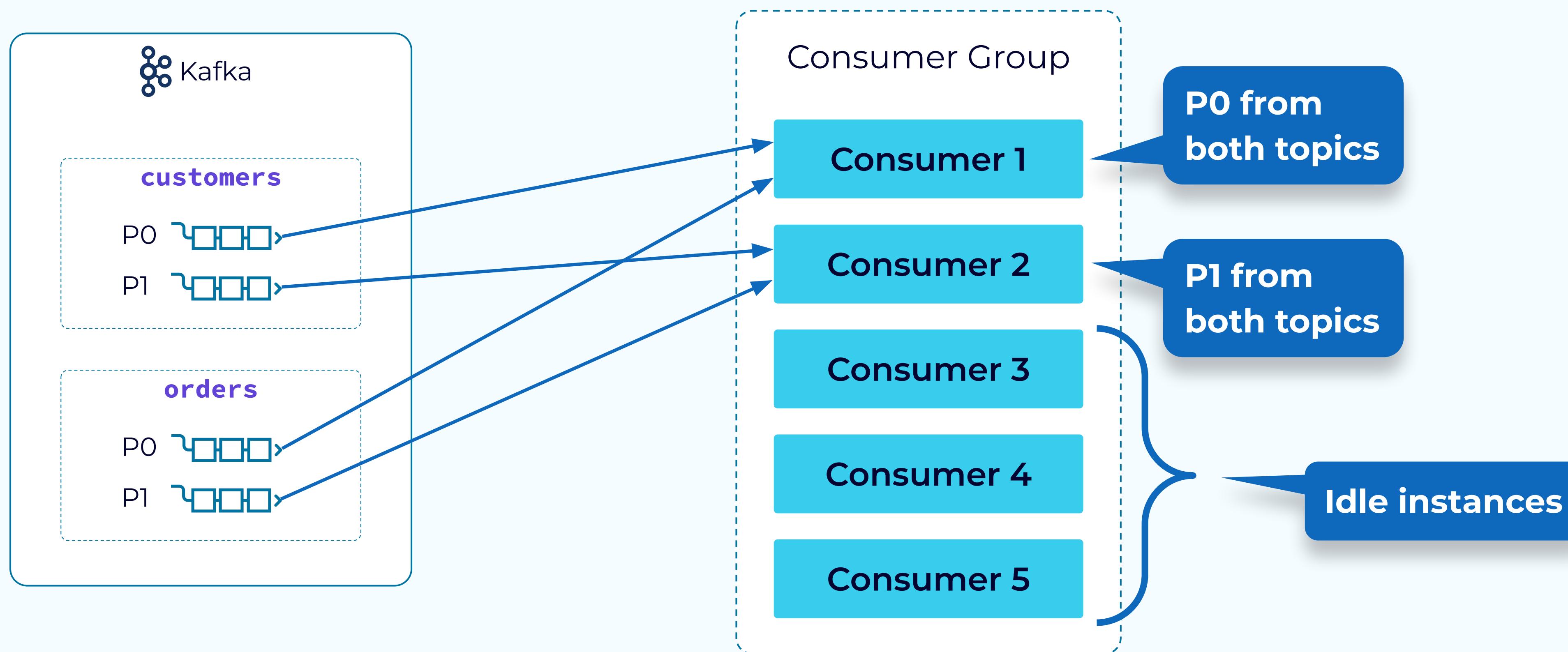
Group Startup: Step 2 - Members Join



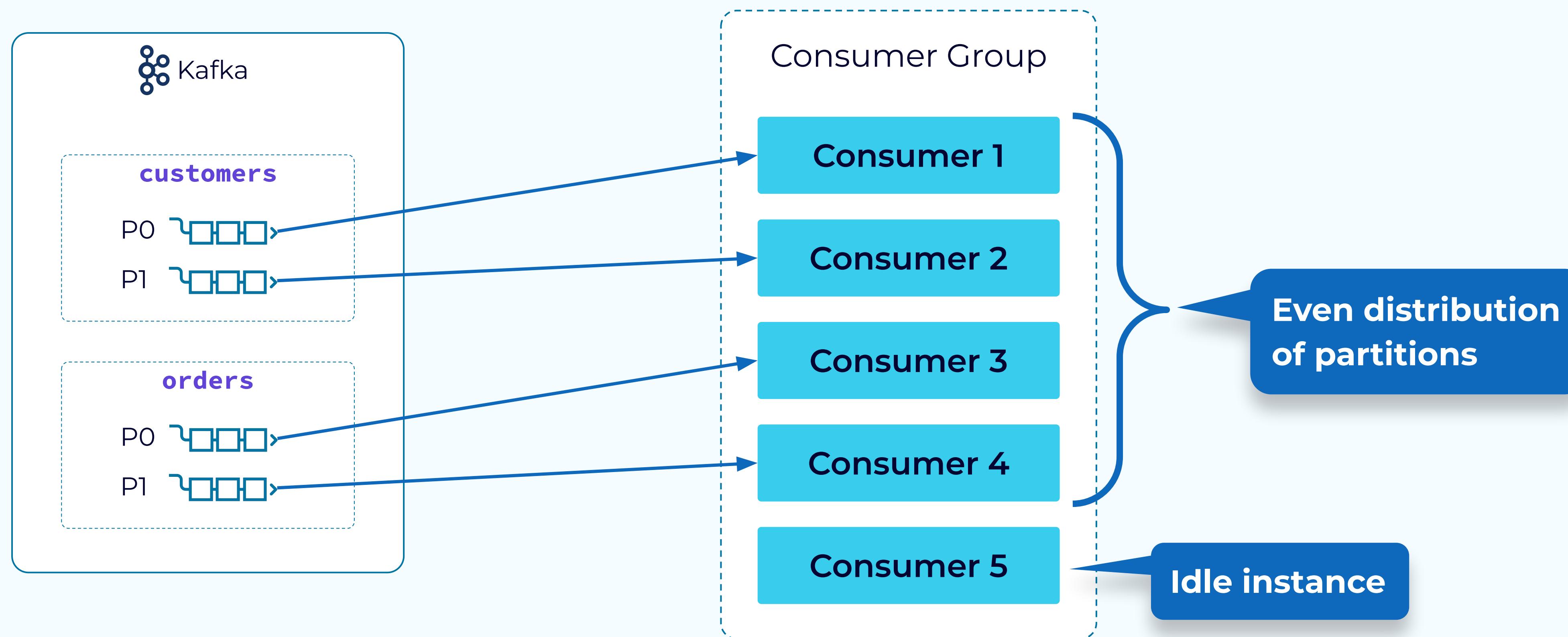
Group Startup: Step 3 - Partitions Assigned



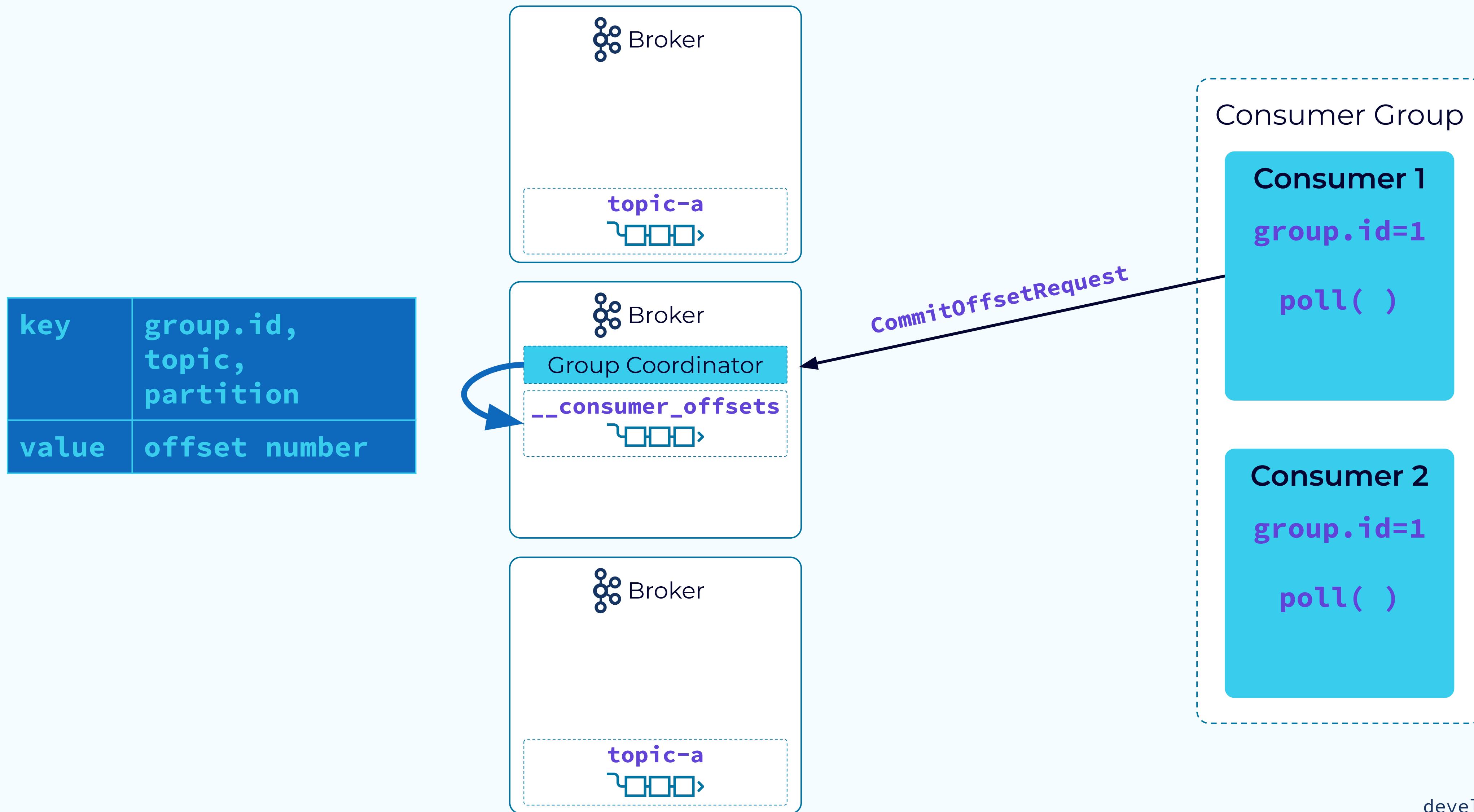
Range Partition Assignment Strategies



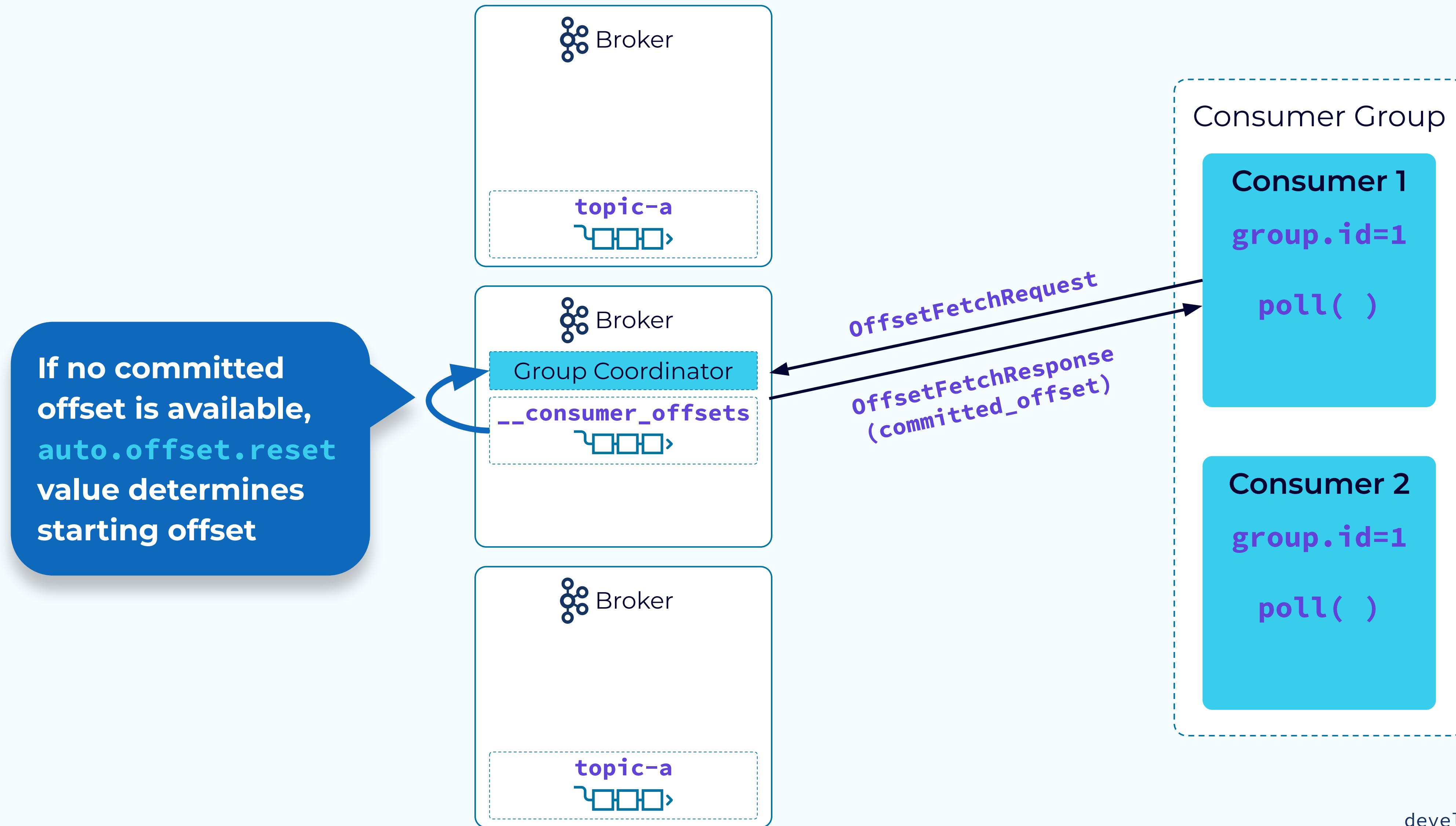
Round Robin and Sticky Partition Assignment Strategies



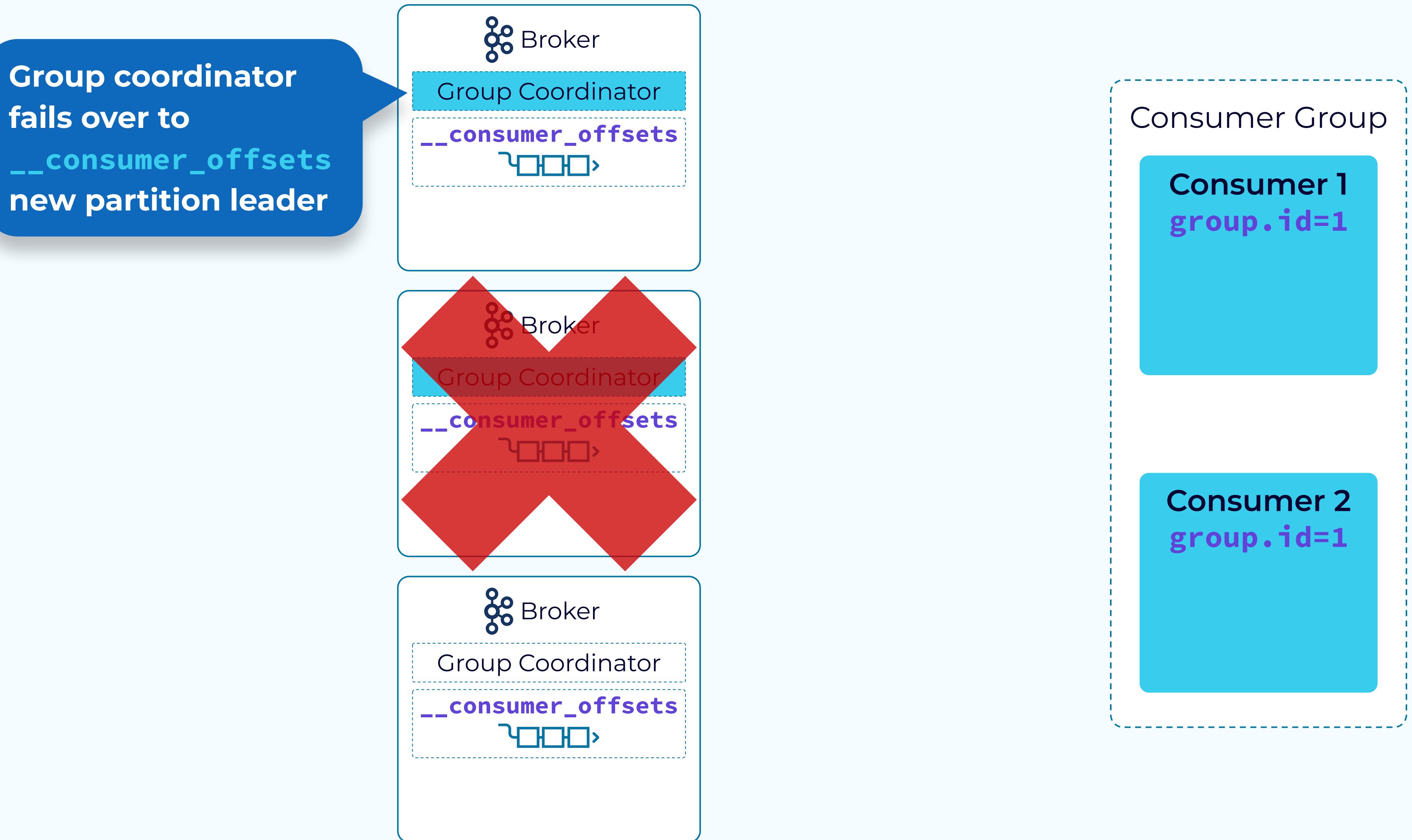
Tracking Partition Consumption



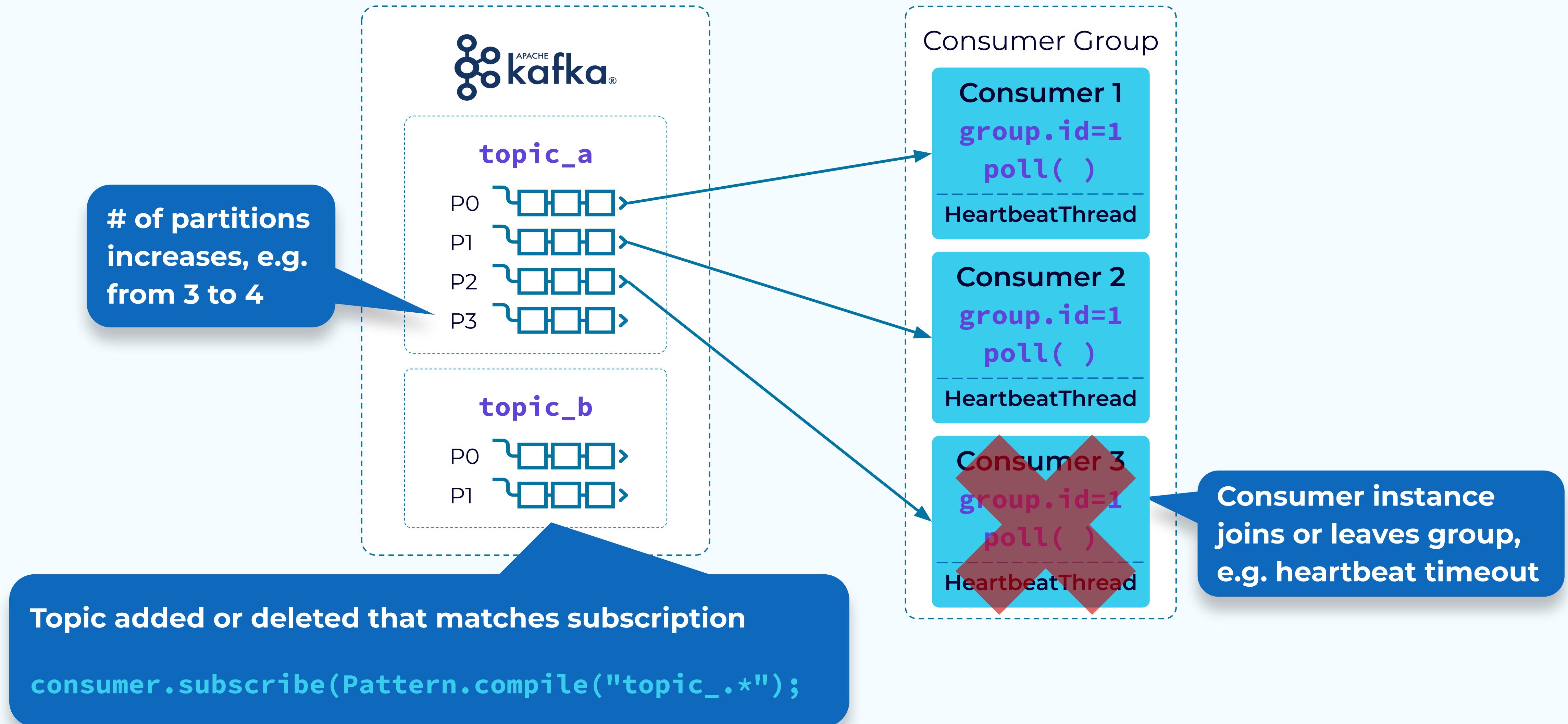
Determining Starting Offset to Consume



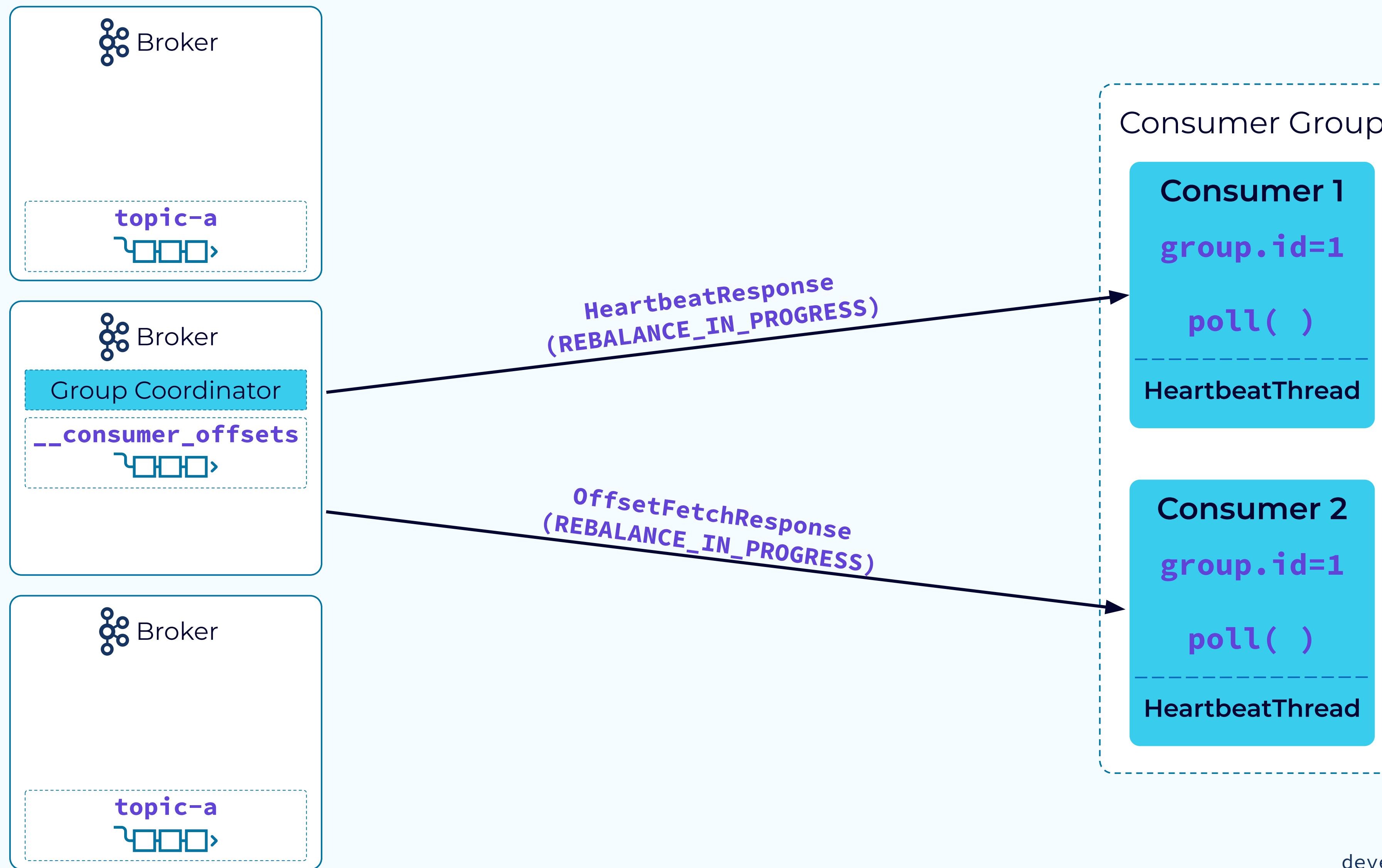
Group Coordinator Failover



Consumer Group Rebalance Triggers



Consumer Group Rebalance Notification

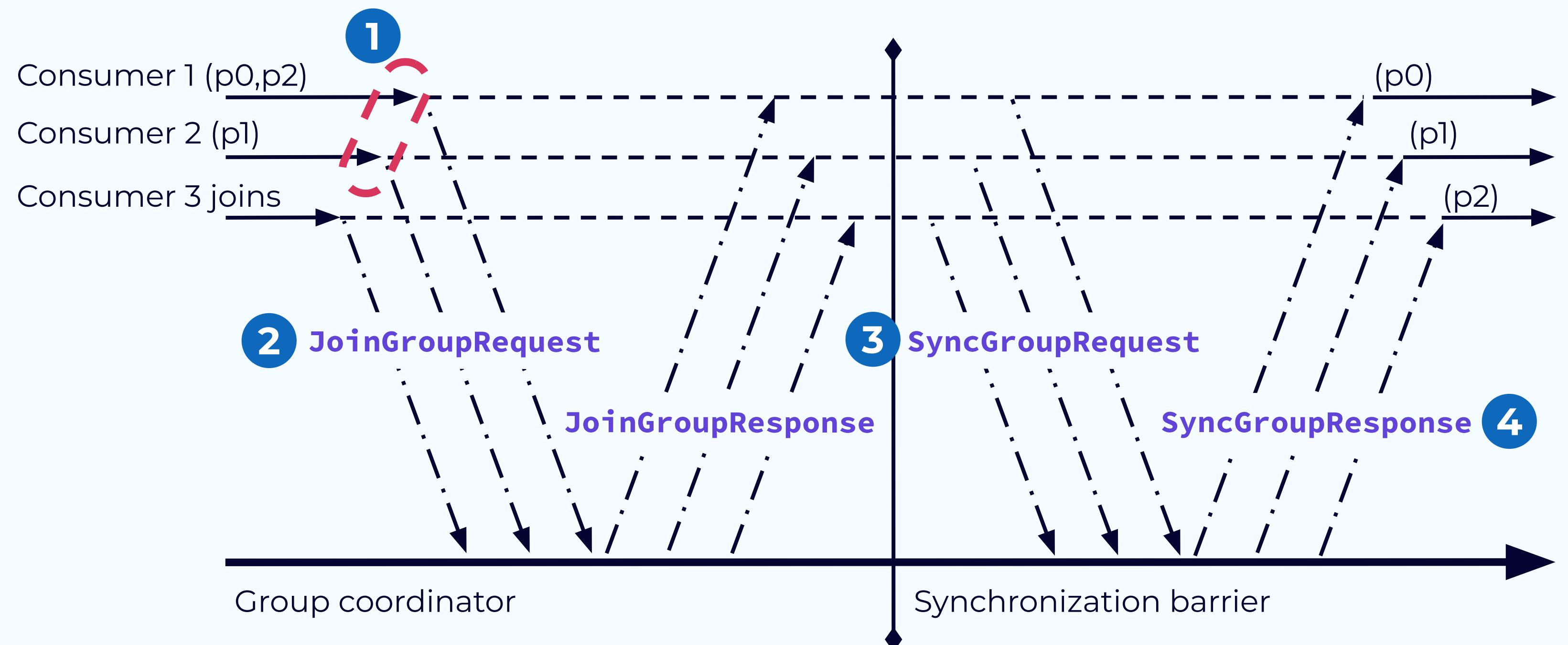


Stop-the-world Rebalance



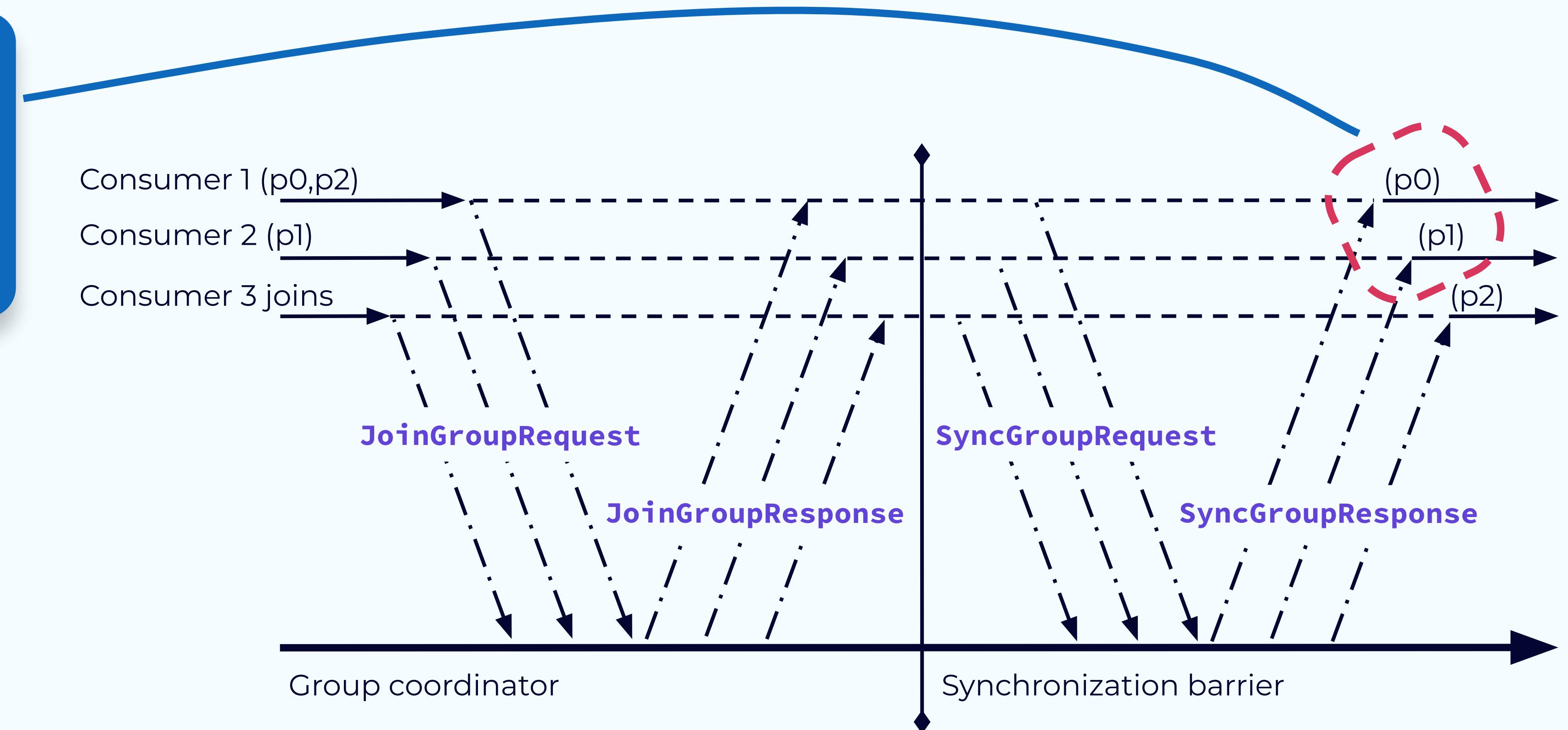
Consumers:

- 1) Revoke current partition assignment and clean up the partition states
- 2) Join the group
- 3) Sync with the group
- 4) Receive new partition assignments
 - a) Build the partition state
 - b) Resume consumption



Stop-the-world Problem 1 - Rebuilding the State

Since partitions p0 and p1 are assigned to the same consumer instance, rebuilding the state is unnecessary

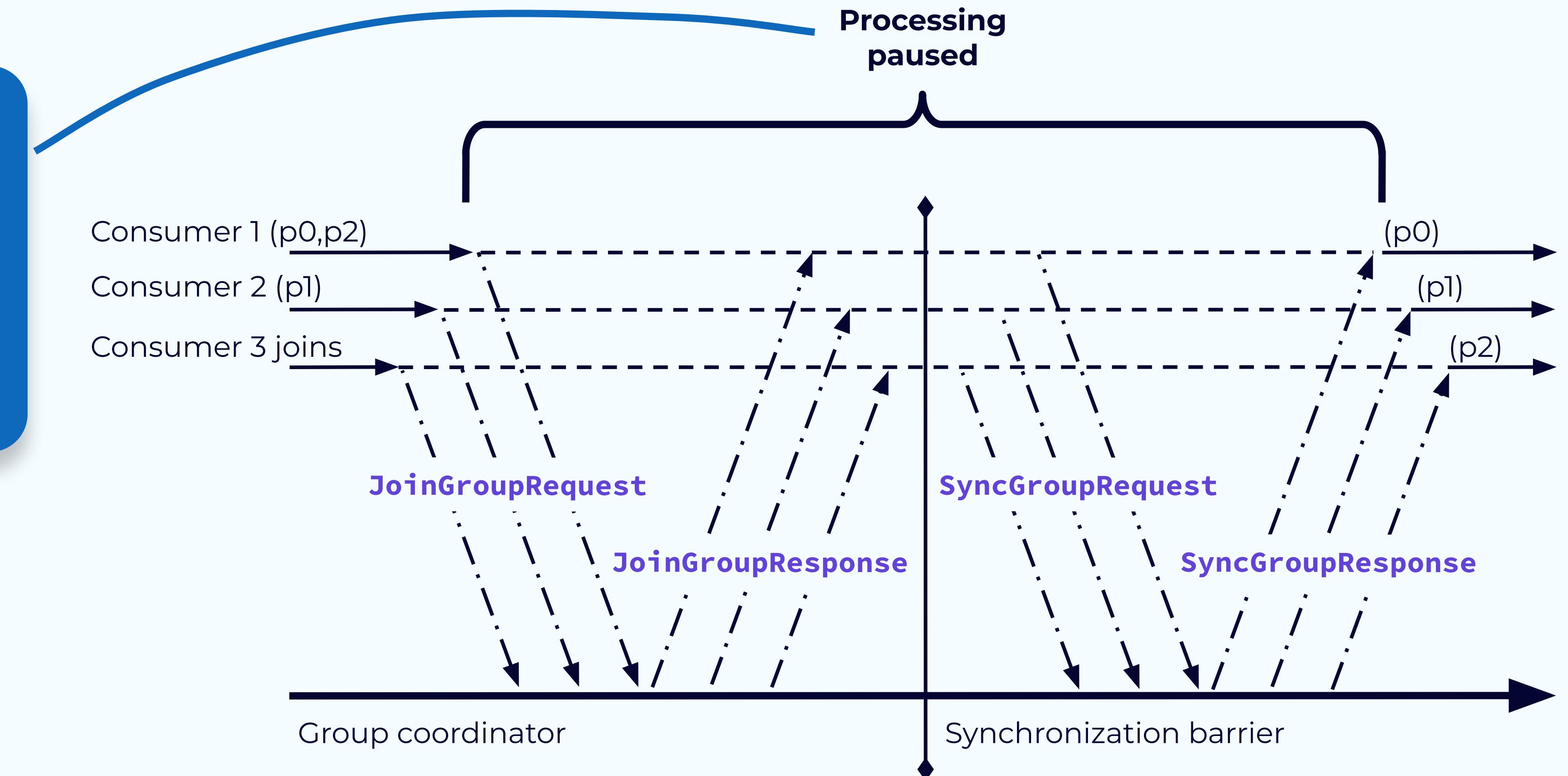


Stop-the-world Problem 2 - Paused Processing

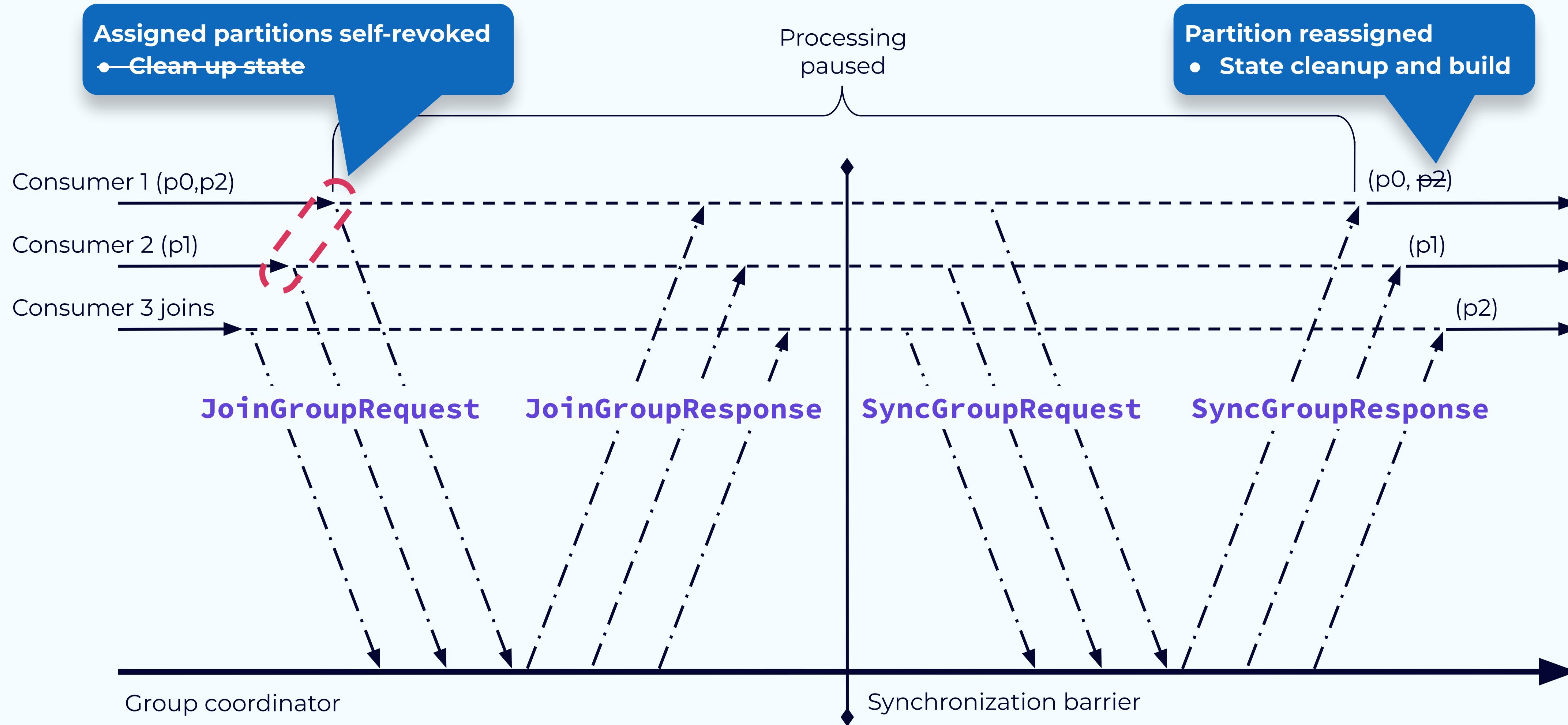


Processing pauses for all subscribed partitions for the duration of the rebalance

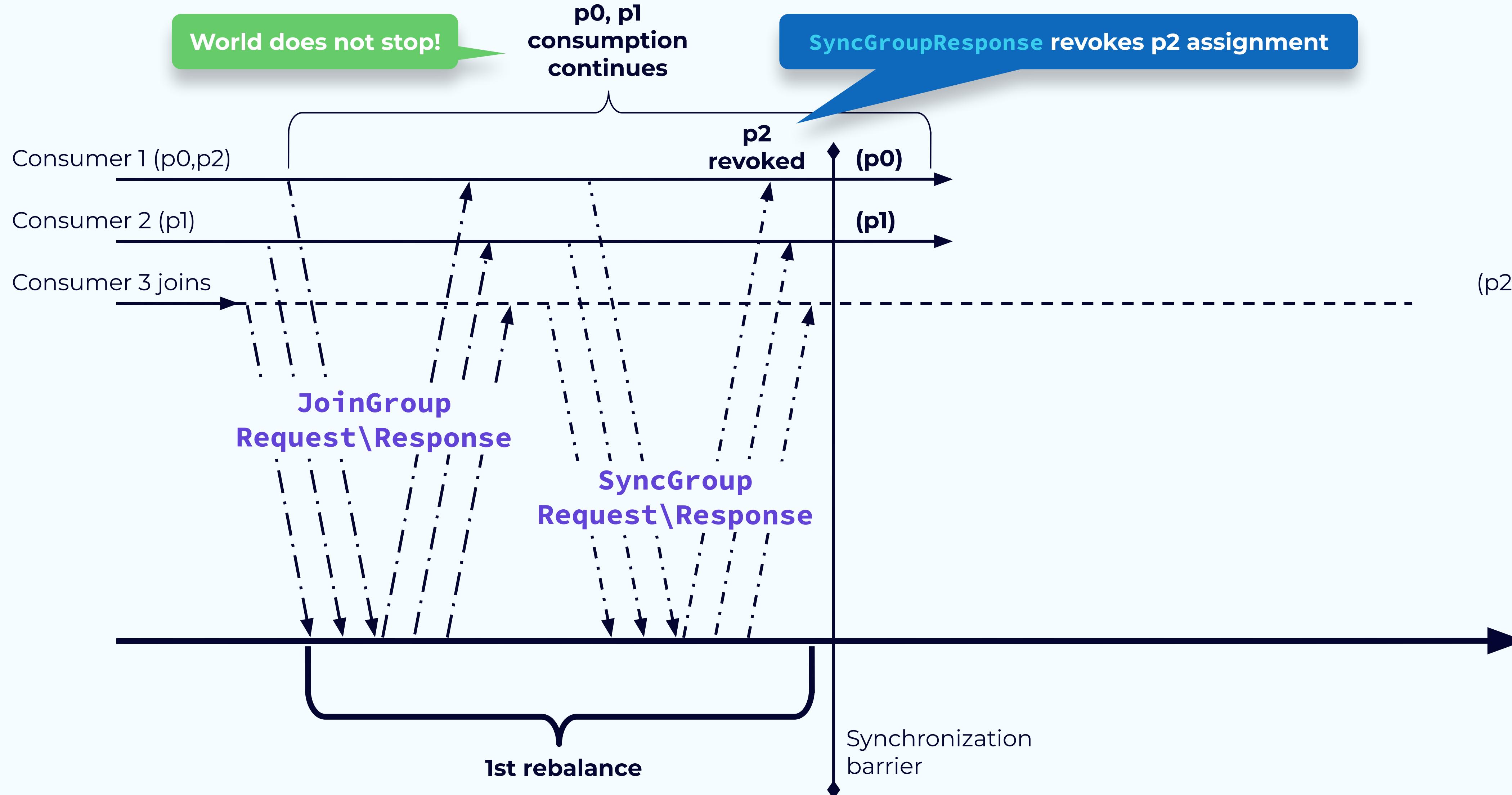
- The pausing for p0 and p1 is unnecessary



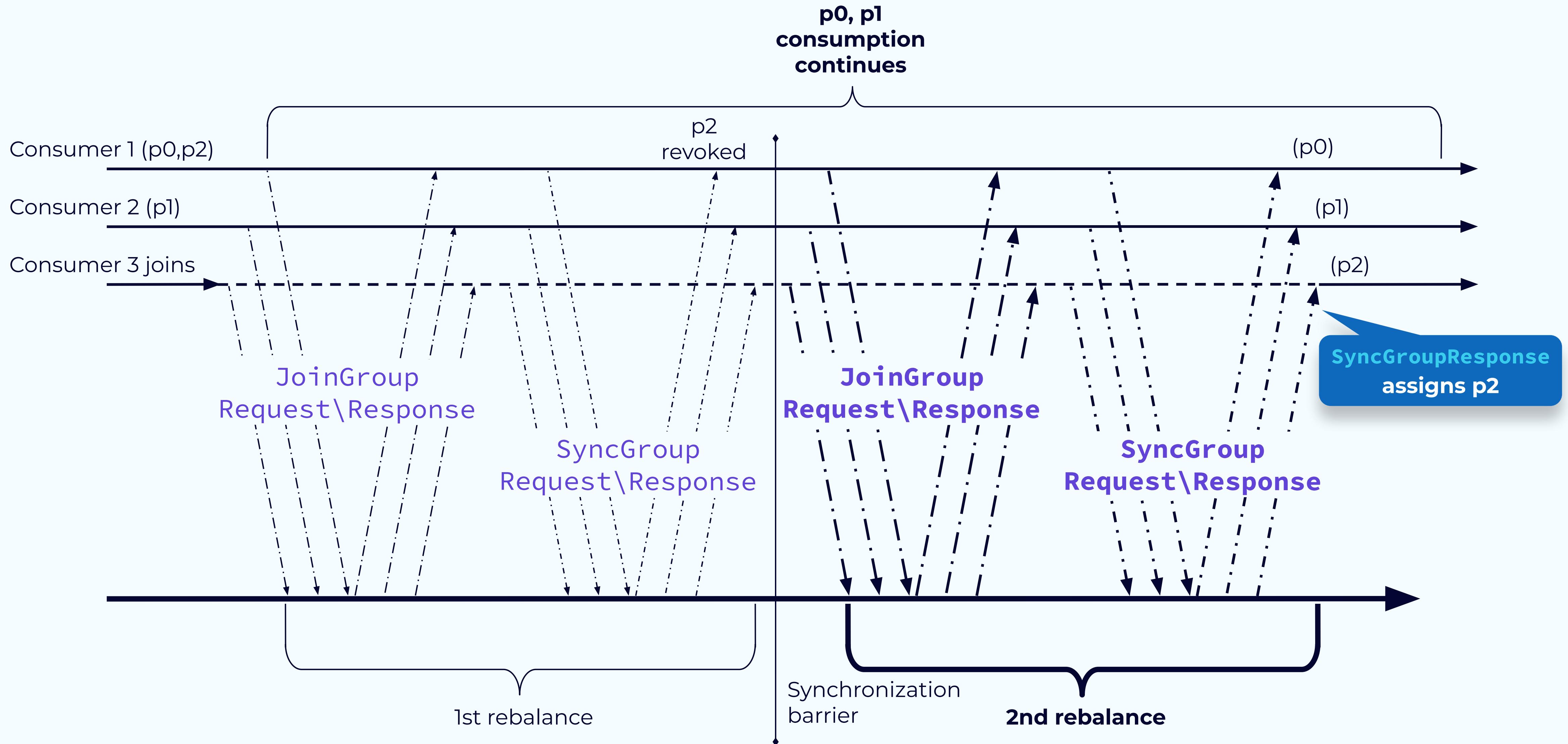
Avoid Needless State Rebuild with StickyAssignor



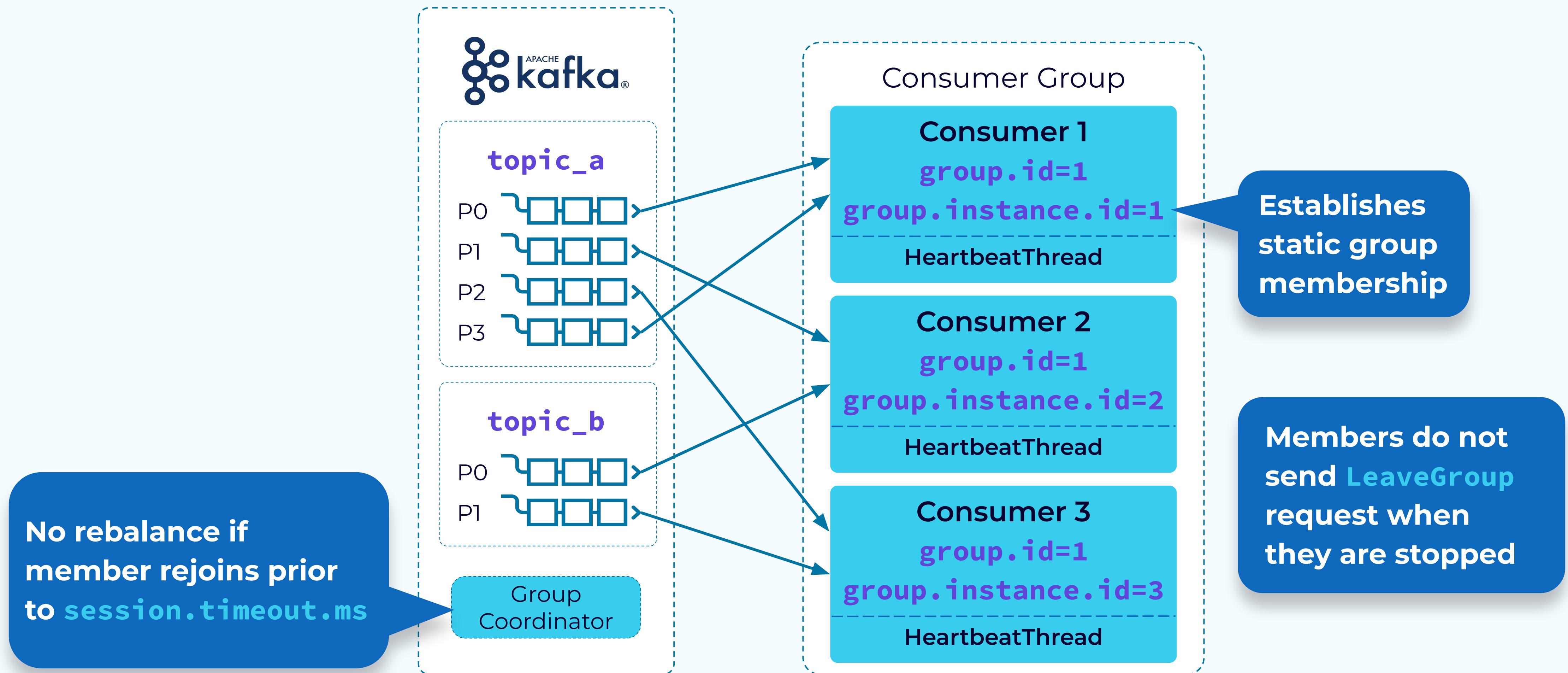
Avoid Processing Pause with CooperativeStickyAssignor



Avoid Processing Pause with CooperativeStickyAssignor



Avoid Rebalance with Static Group Membership

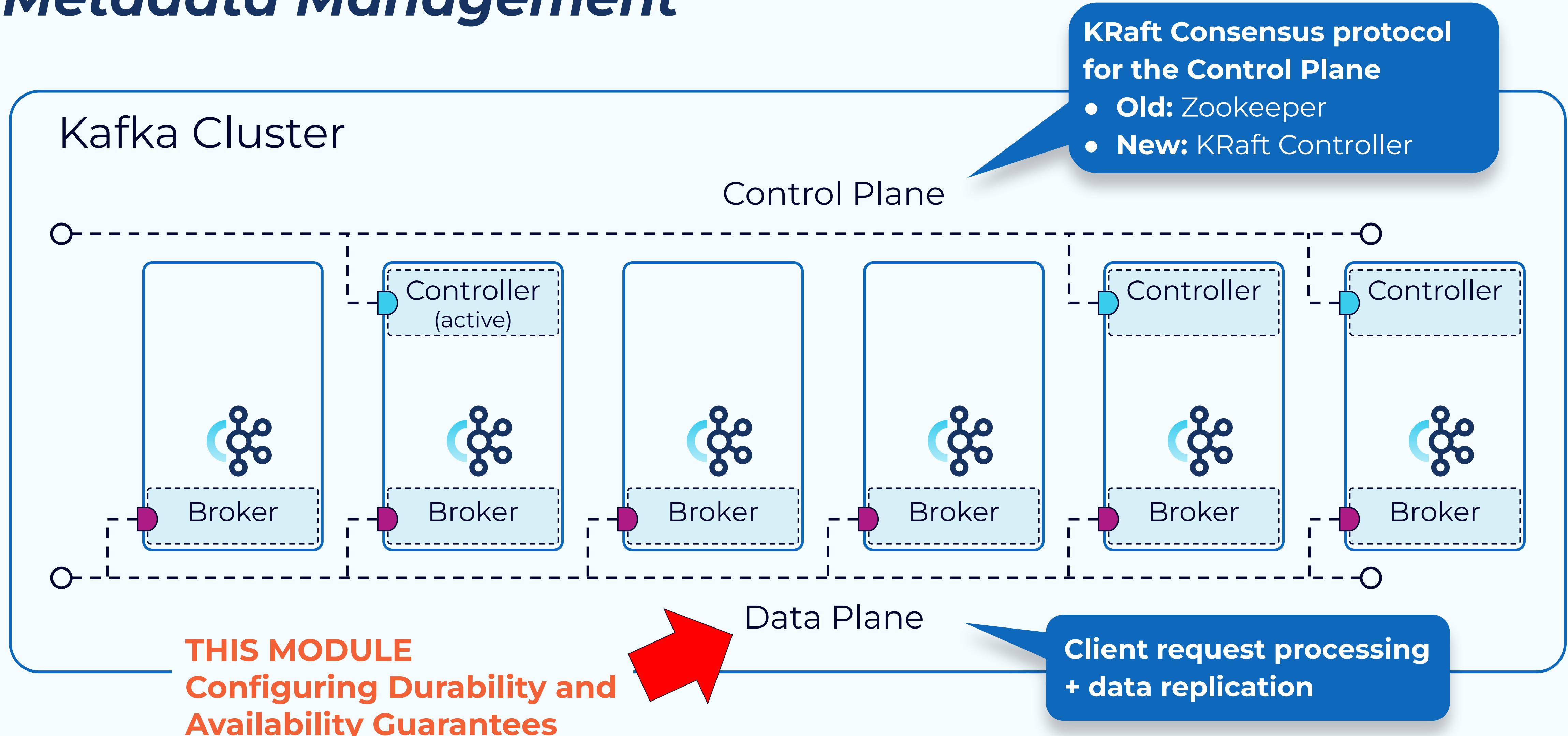




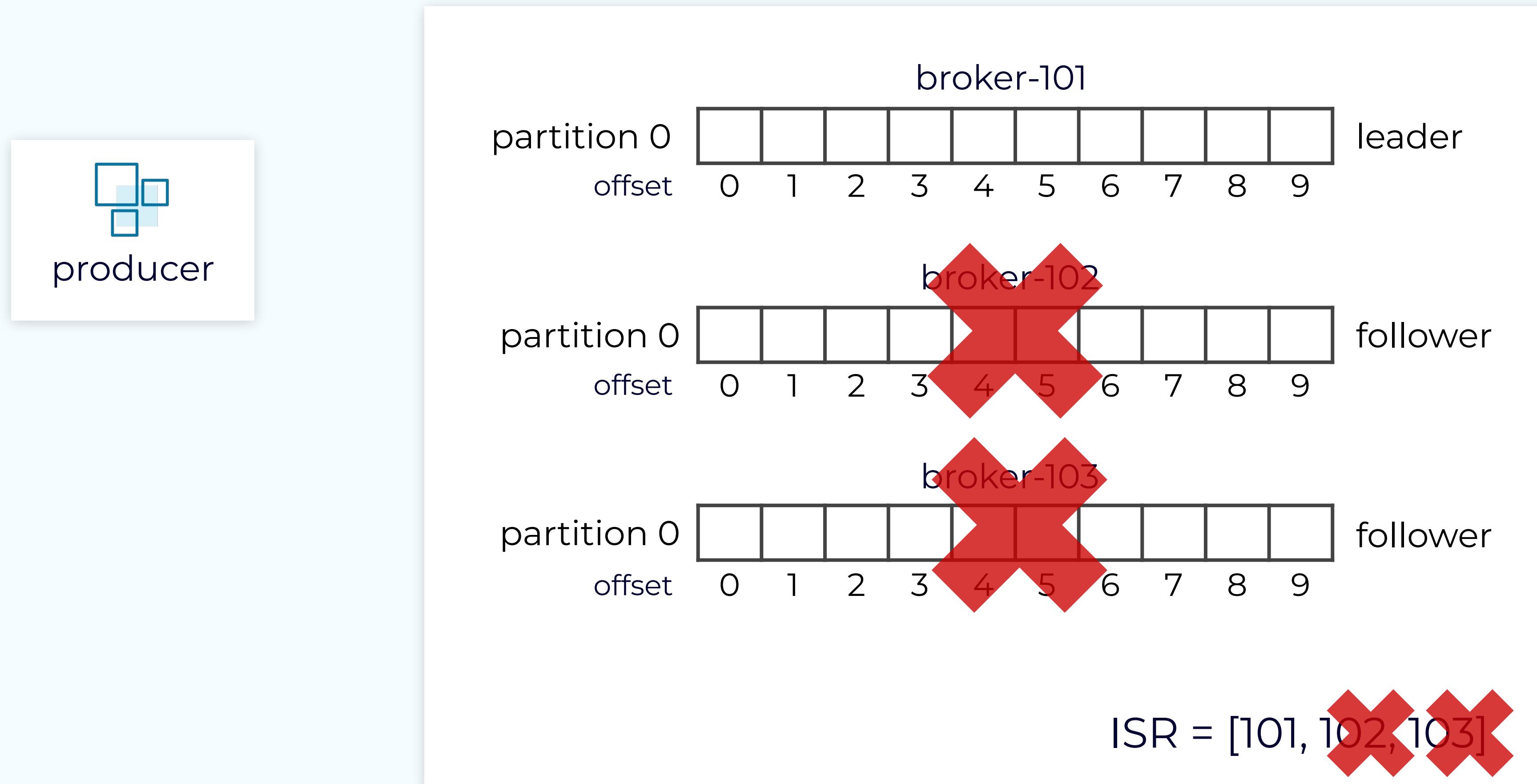
Configuring Durability, Availability and Ordering Guarantees



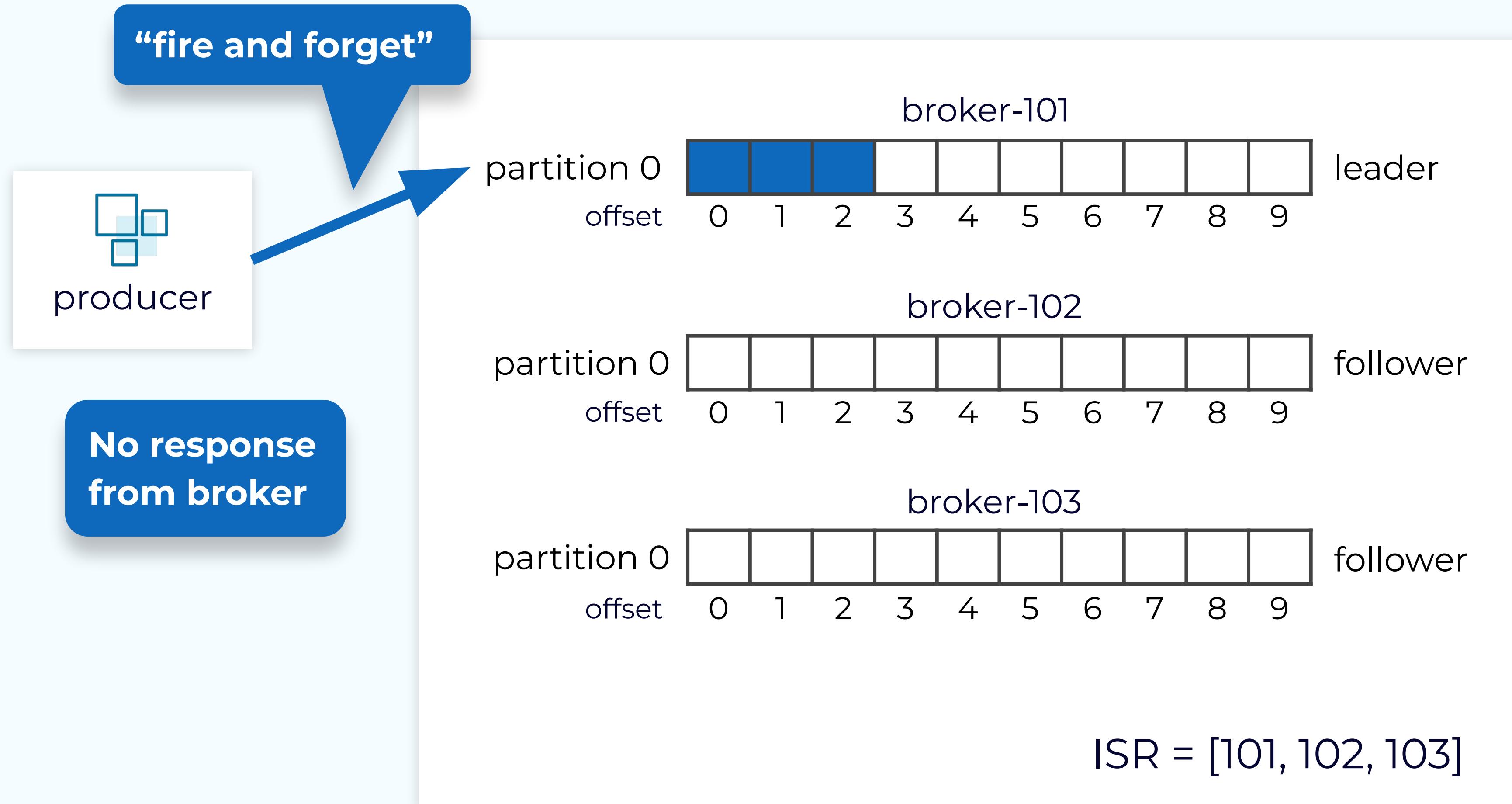
Kafka Separates Data & Metadata Management



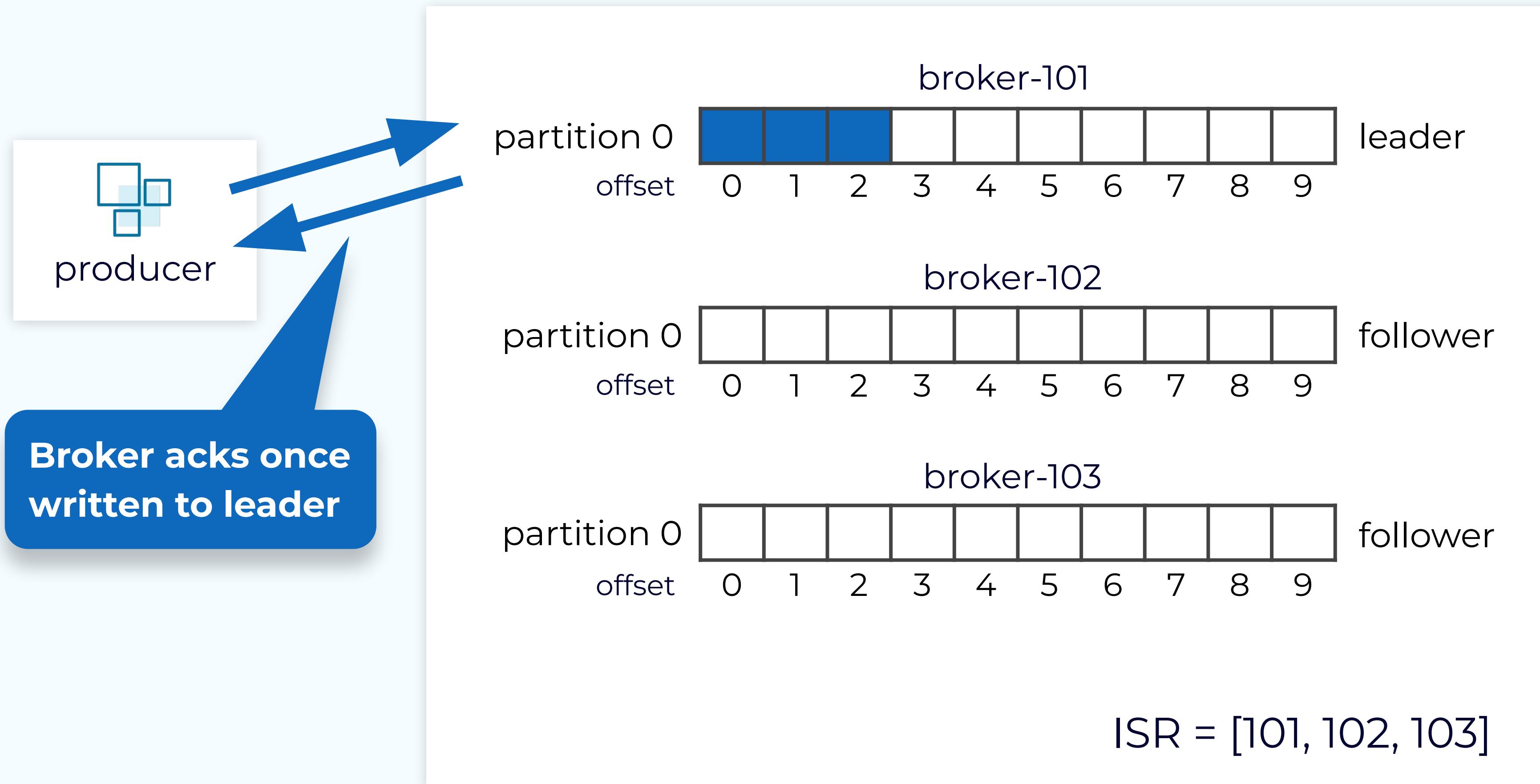
Data Durability and Availability Guarantees



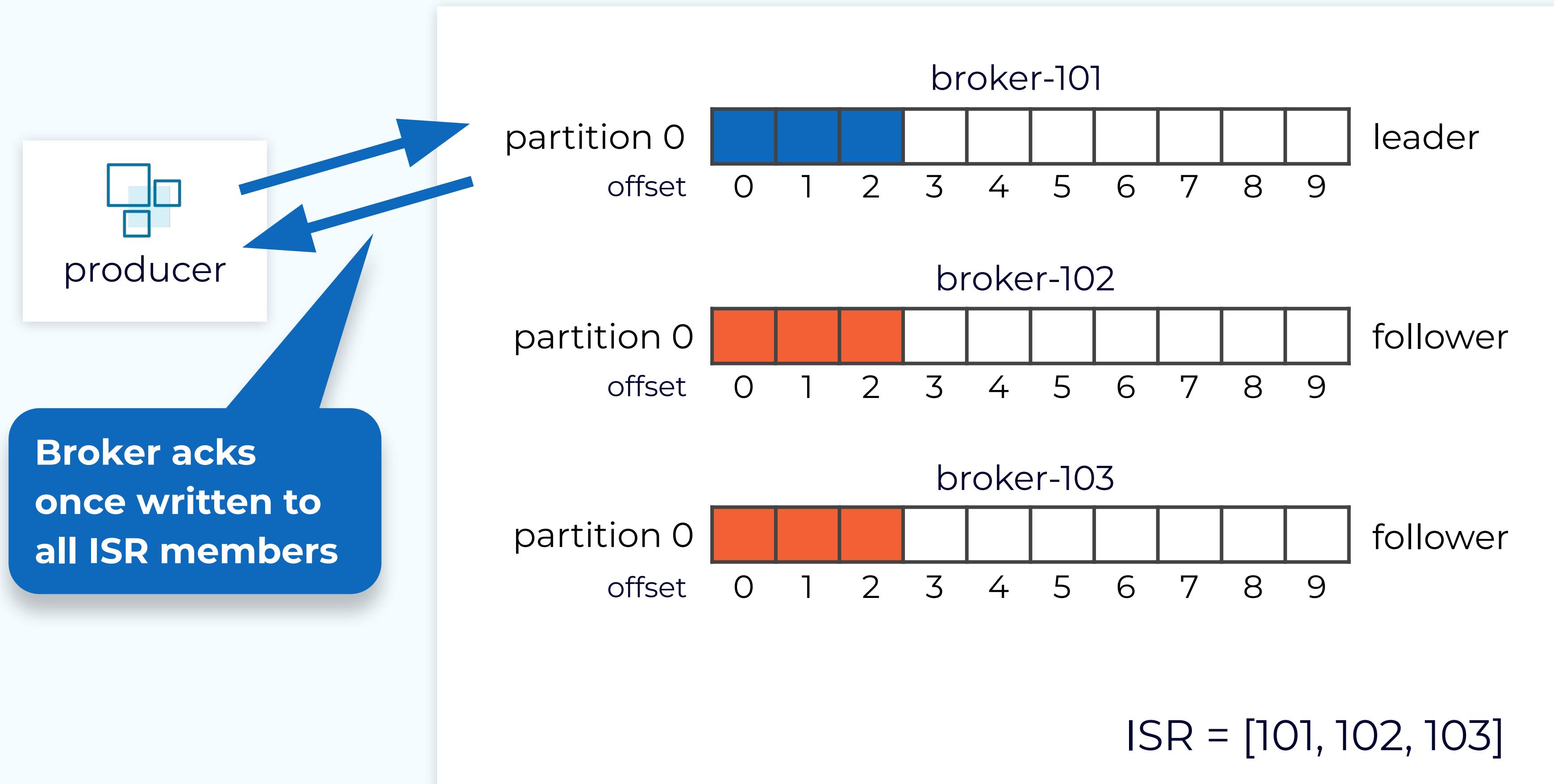
Producer acks=0



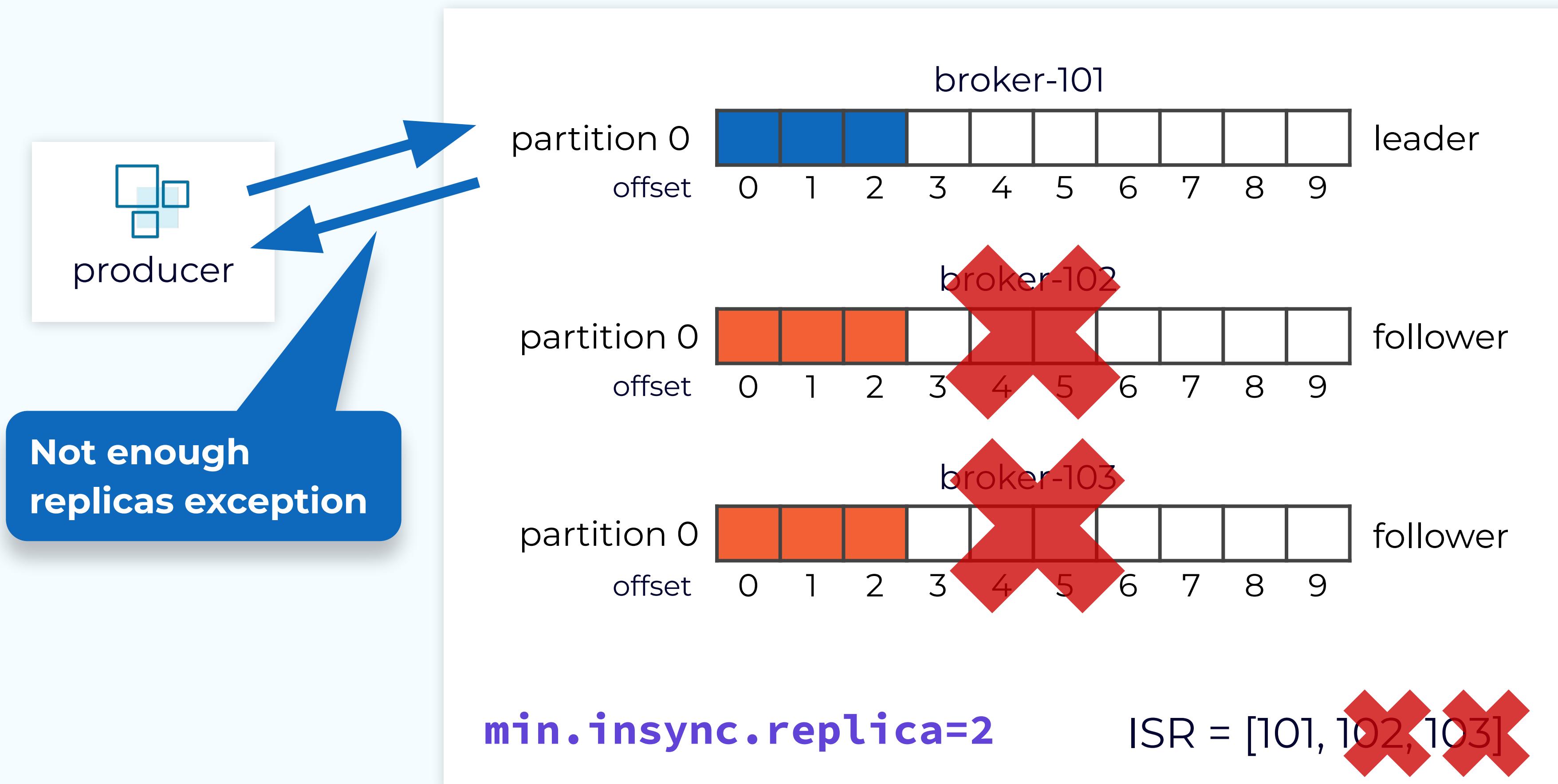
Producer acks=1



Producer acks=all



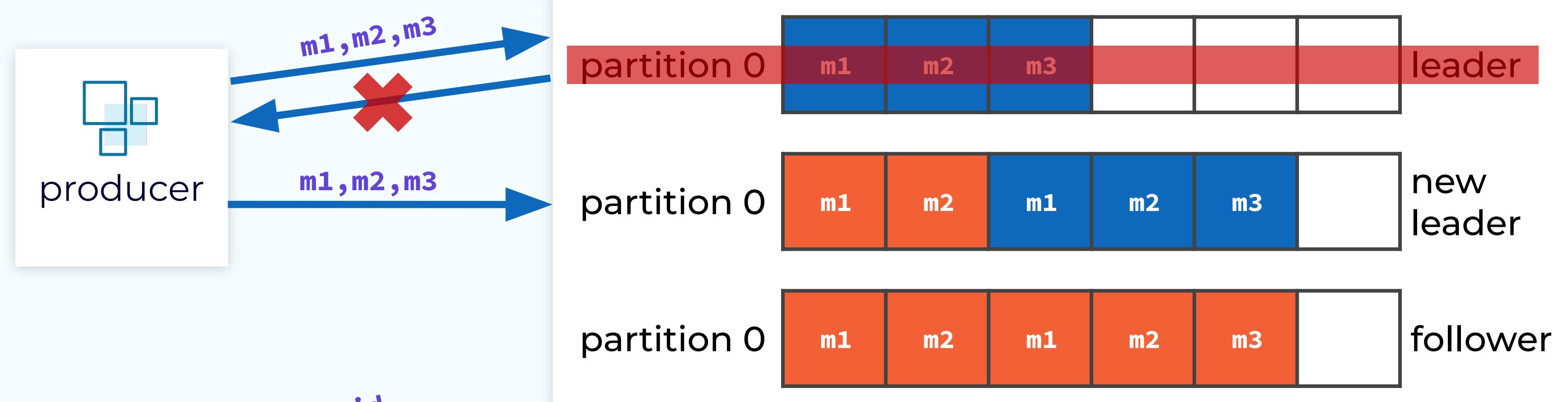
Topic min.insync.replicas



Producer Idempotency

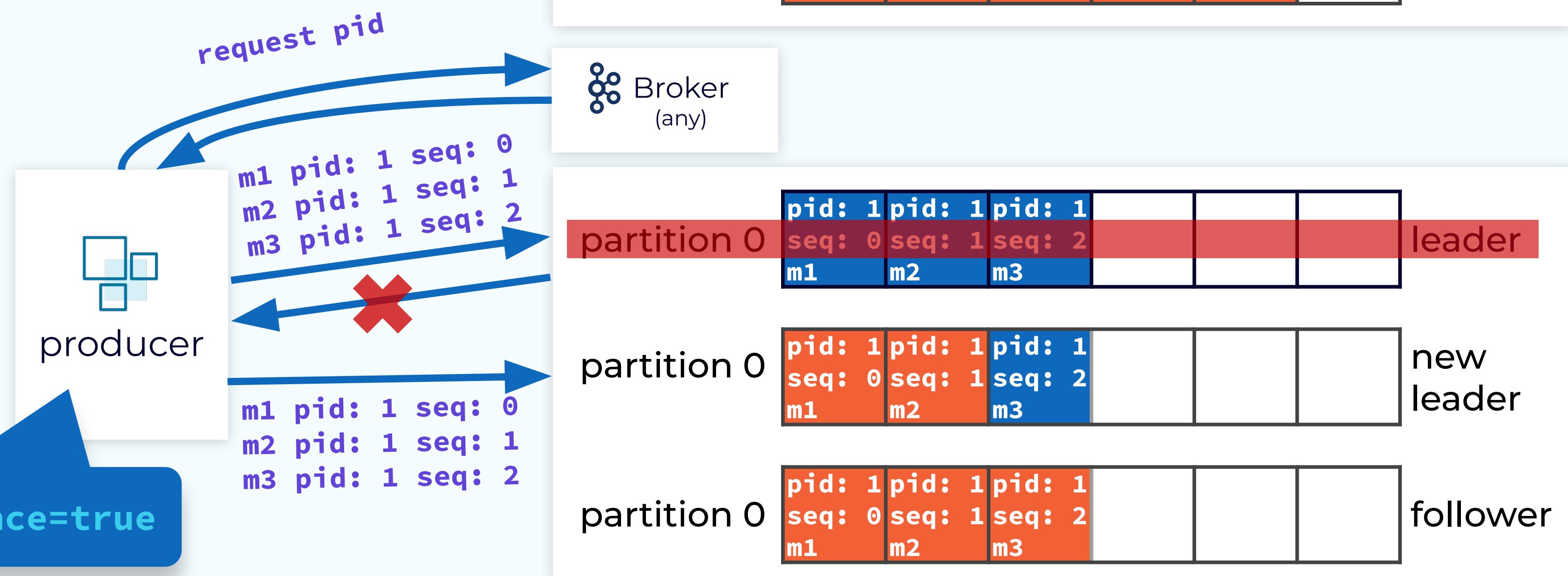


**Without idempotency,
duplicate and out of
order records can occur**

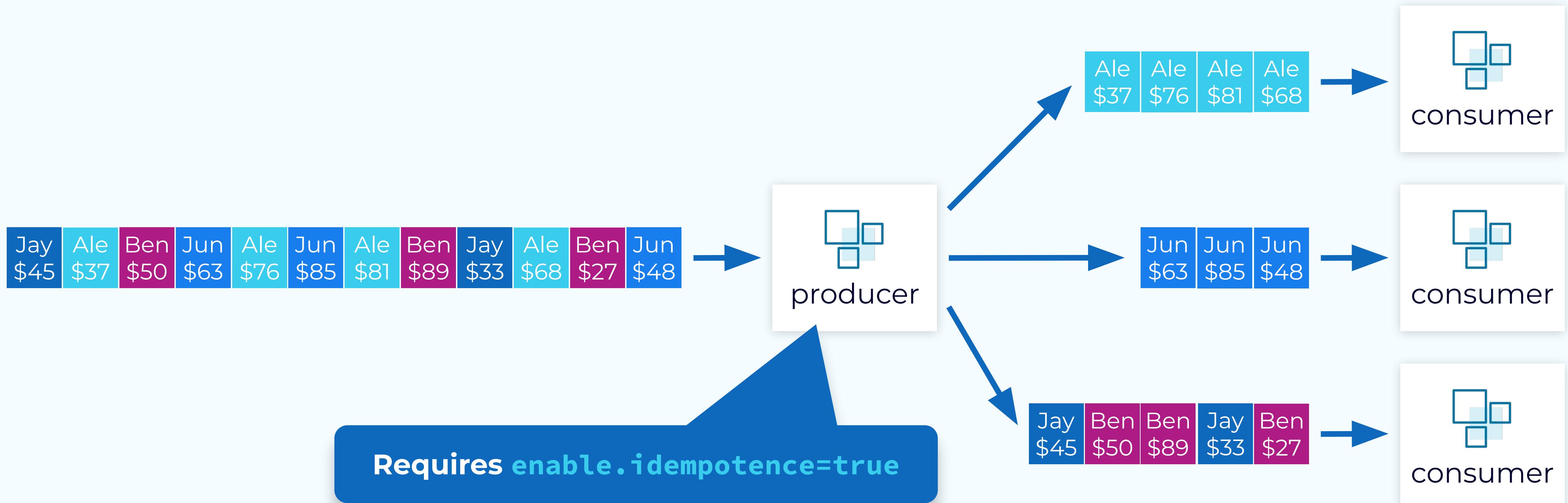


**Idempotency's seq #'s
prevent both of these
issues**

`enable.idempotence=true`



End-to-End Ordering Guarantee



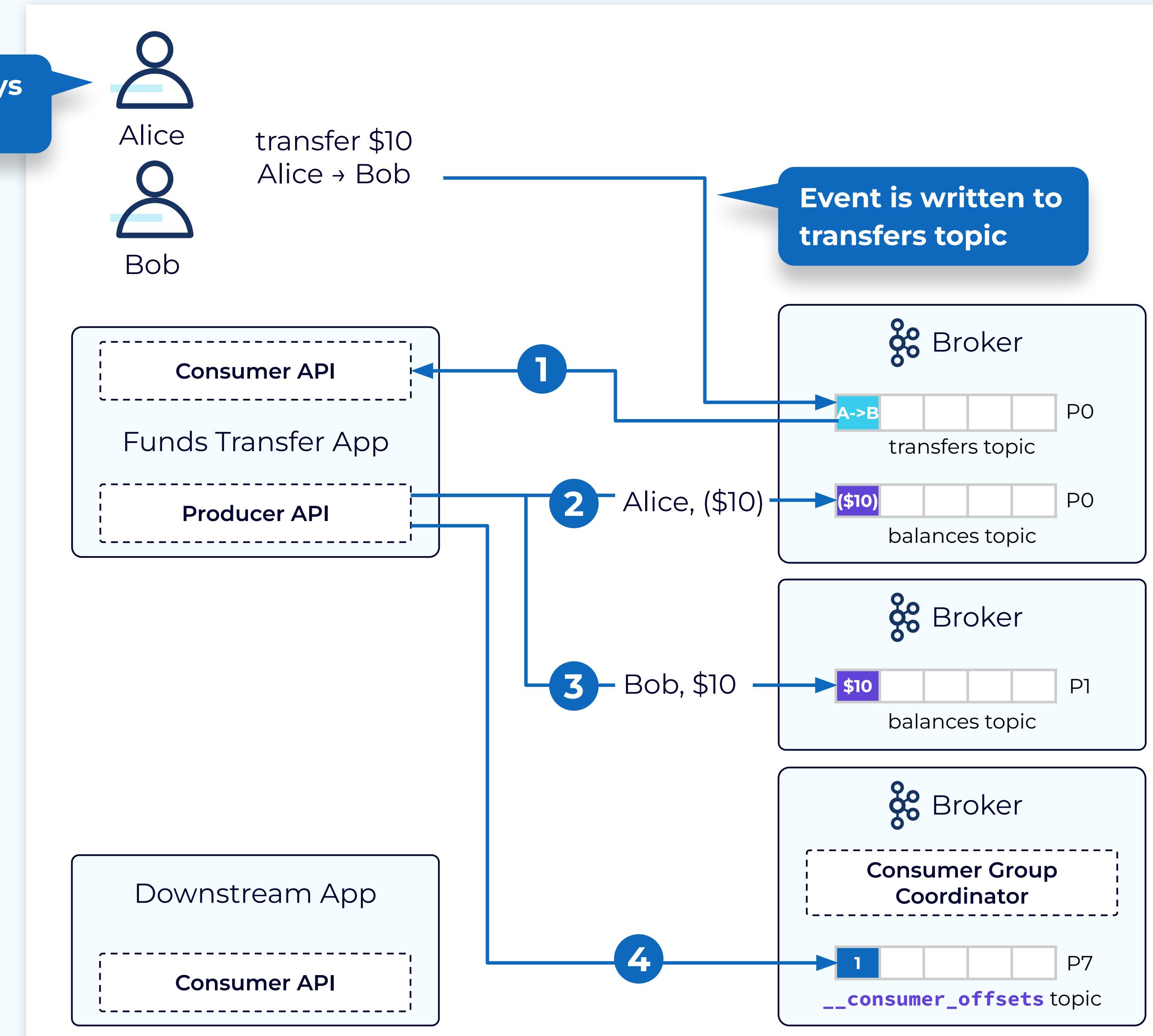


CONFLUENT
Developer

Transactions

Why Are Transactions Needed?

1. Event is fetched by the consumer
2. Debit event is written
3. Credit event is written
4. Transfer event offset is committed



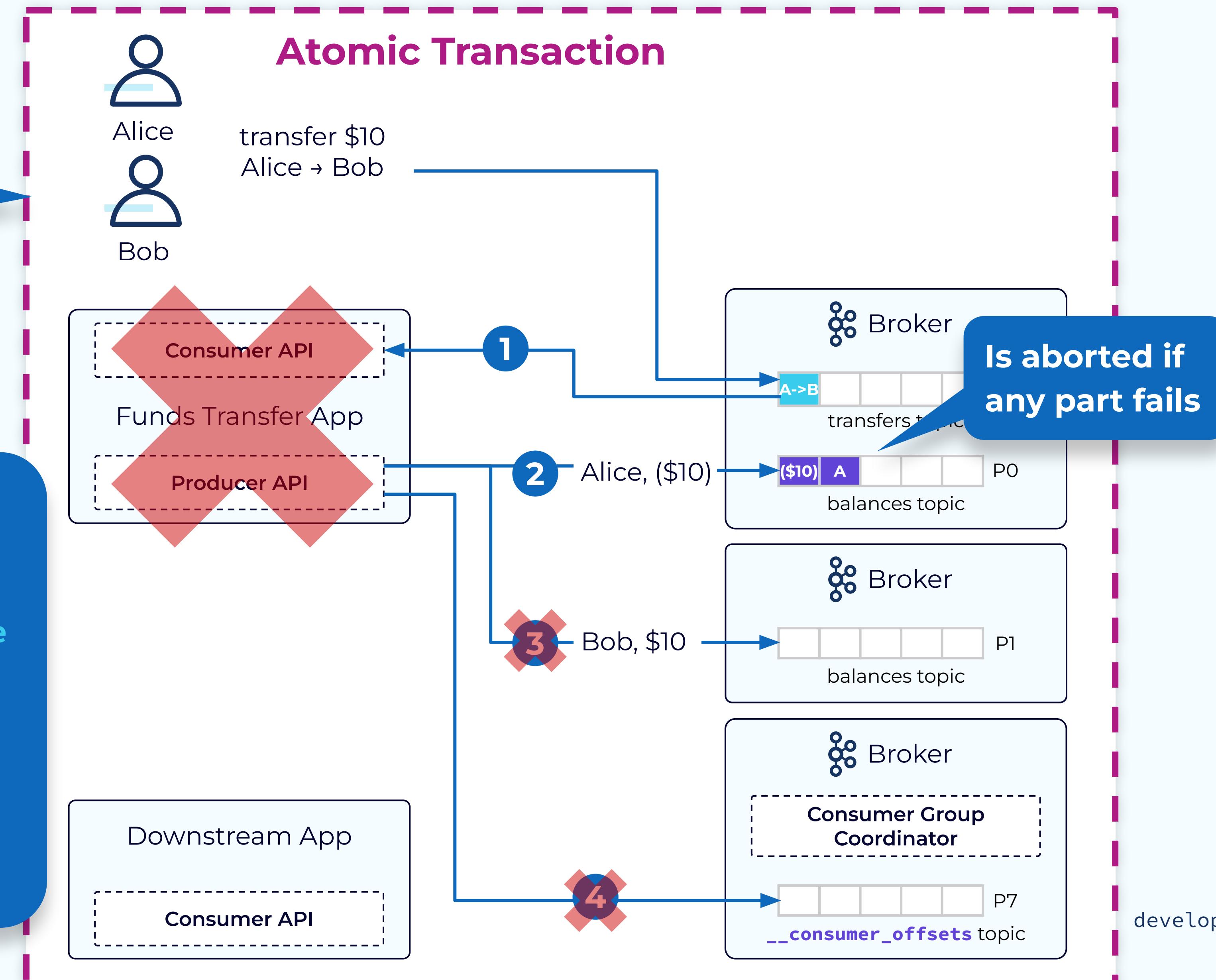
Kafka Transactions Deliver Exactly Once



Transaction is only committed if all parts succeed

Using transactions with Kafka Streams is quite simple:

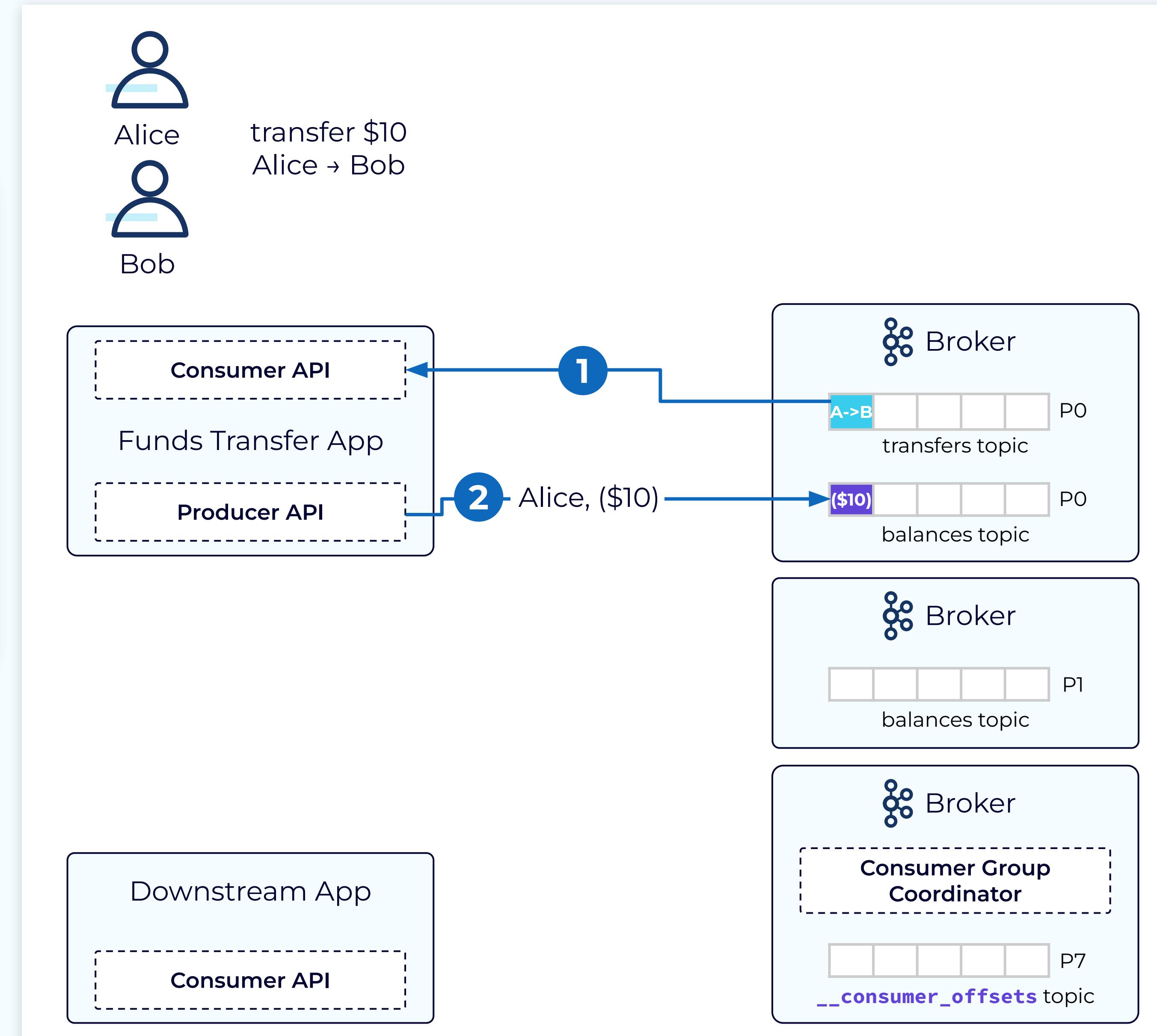
- 1) Set `processing.guarantee` to `exactly_once_v2` in `StreamsConfig`
- 2) Set `isolation.level` to `read_committed` in the Consumer configuration



System Failure Without Transactions



1. Event fetched by consumer
2. Alice's account debited



System Failure Without Transactions

1. Event fetched by consumer

2. Alice's account debited

Application instance fails without committing offset and new application instance starts

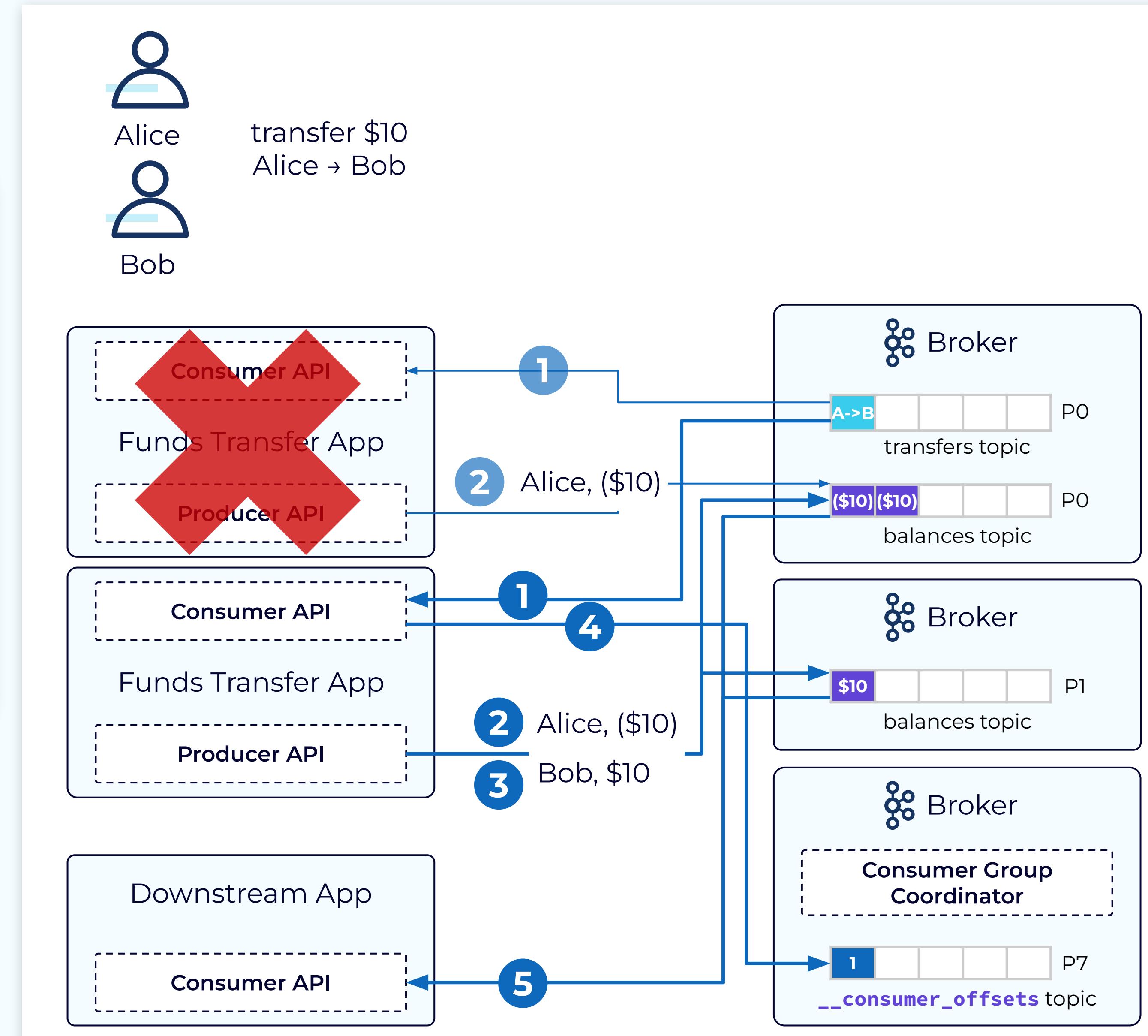
1. Event fetched by consumer

2. Alice's account is debited a second time

3. Bob's account is credited

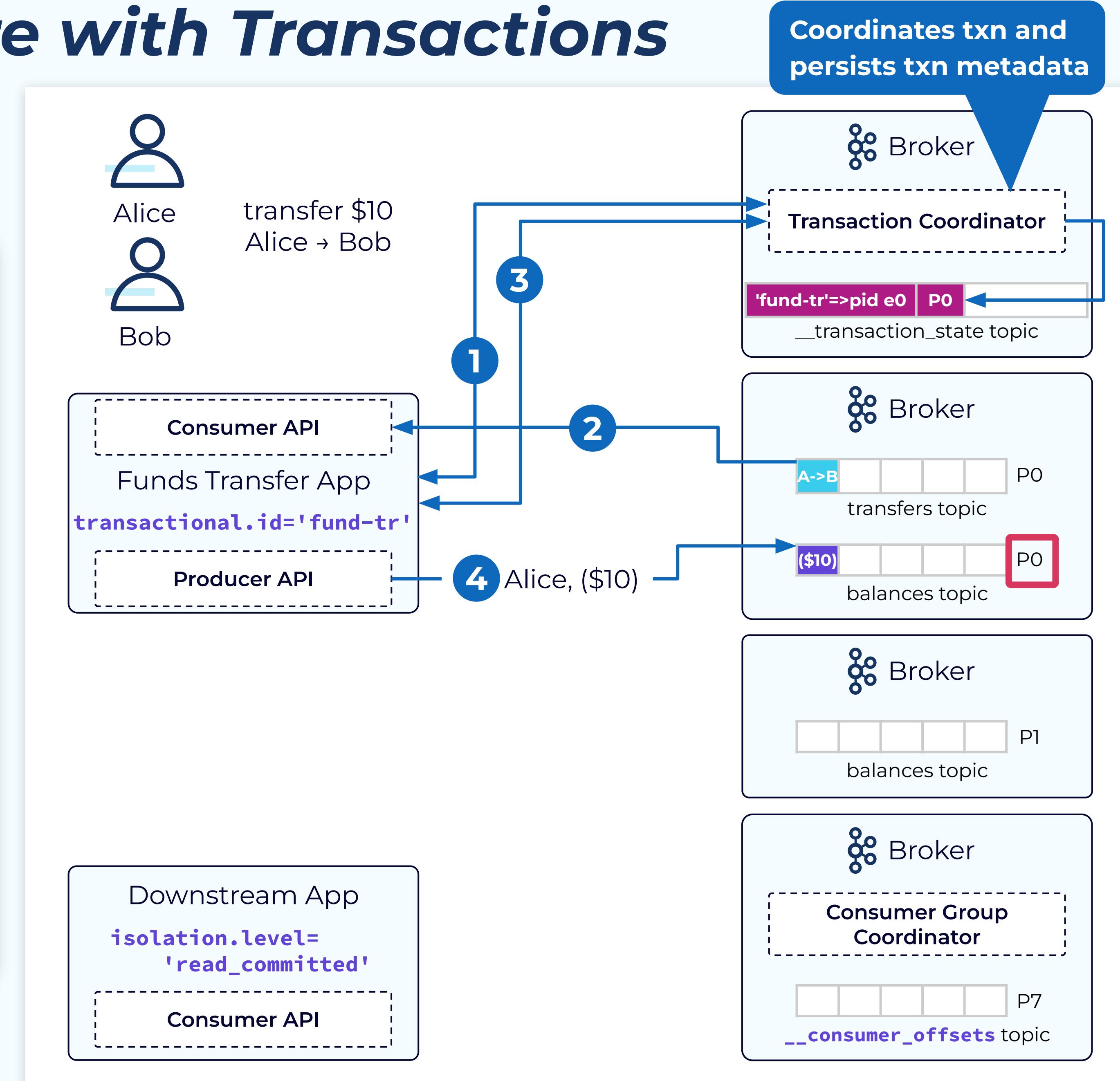
4. Consumer offset committed

5. Two debit events processed by downstream consumer



System Failure with Transactions

1. Requests txn ID, is returned PID and txn epoch
2. Event fetched by consumer
3. Notifies coordinator of partition being written to
4. Alice's account debited

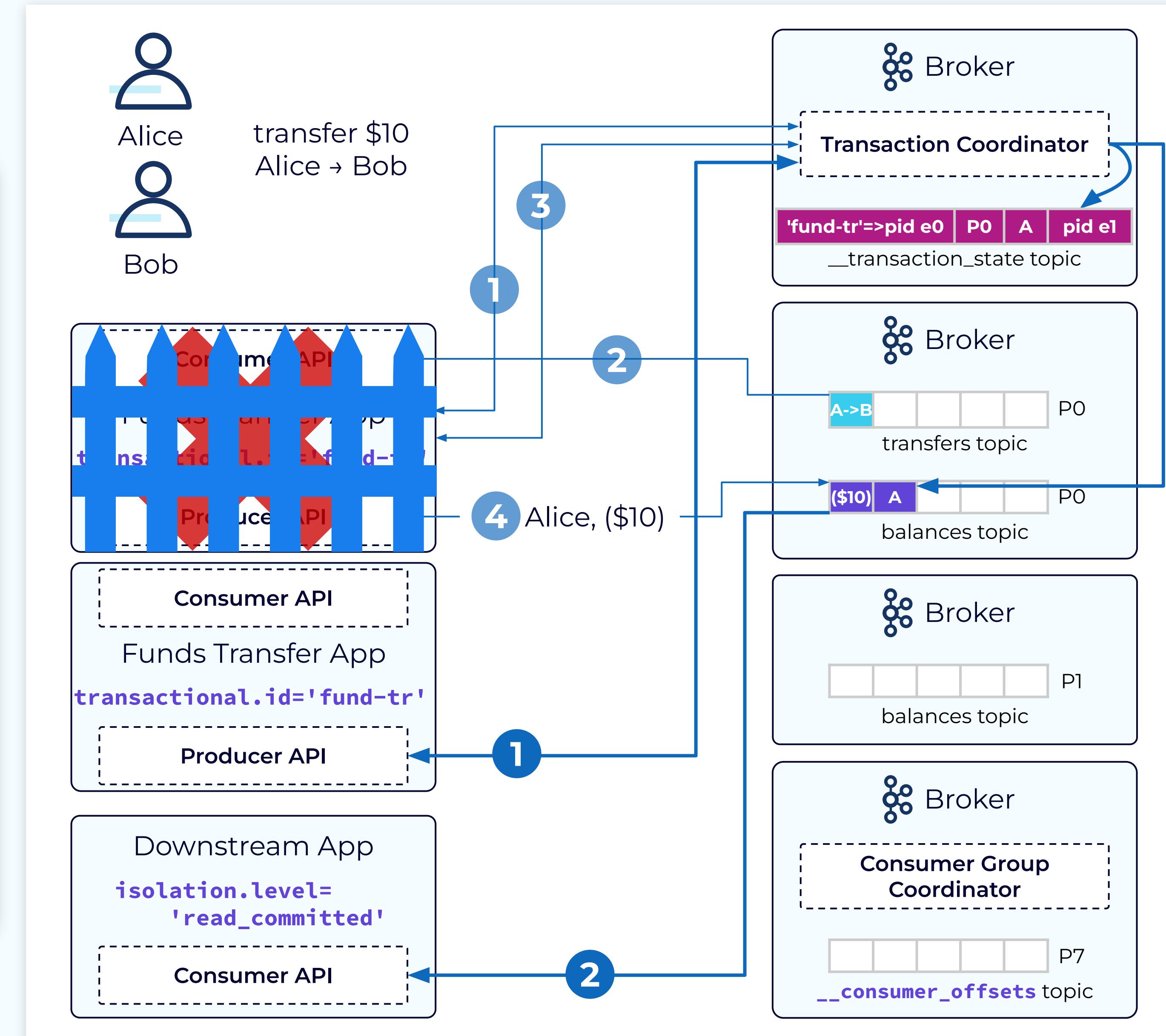


System Failure with Transactions

1. Requests txn ID, is returned PID and txn epoch
2. Event fetched by consumer
3. Notifies coordinator of partition being written to
4. Alice's account debited

Application instance fails without committing offset and new application instance starts

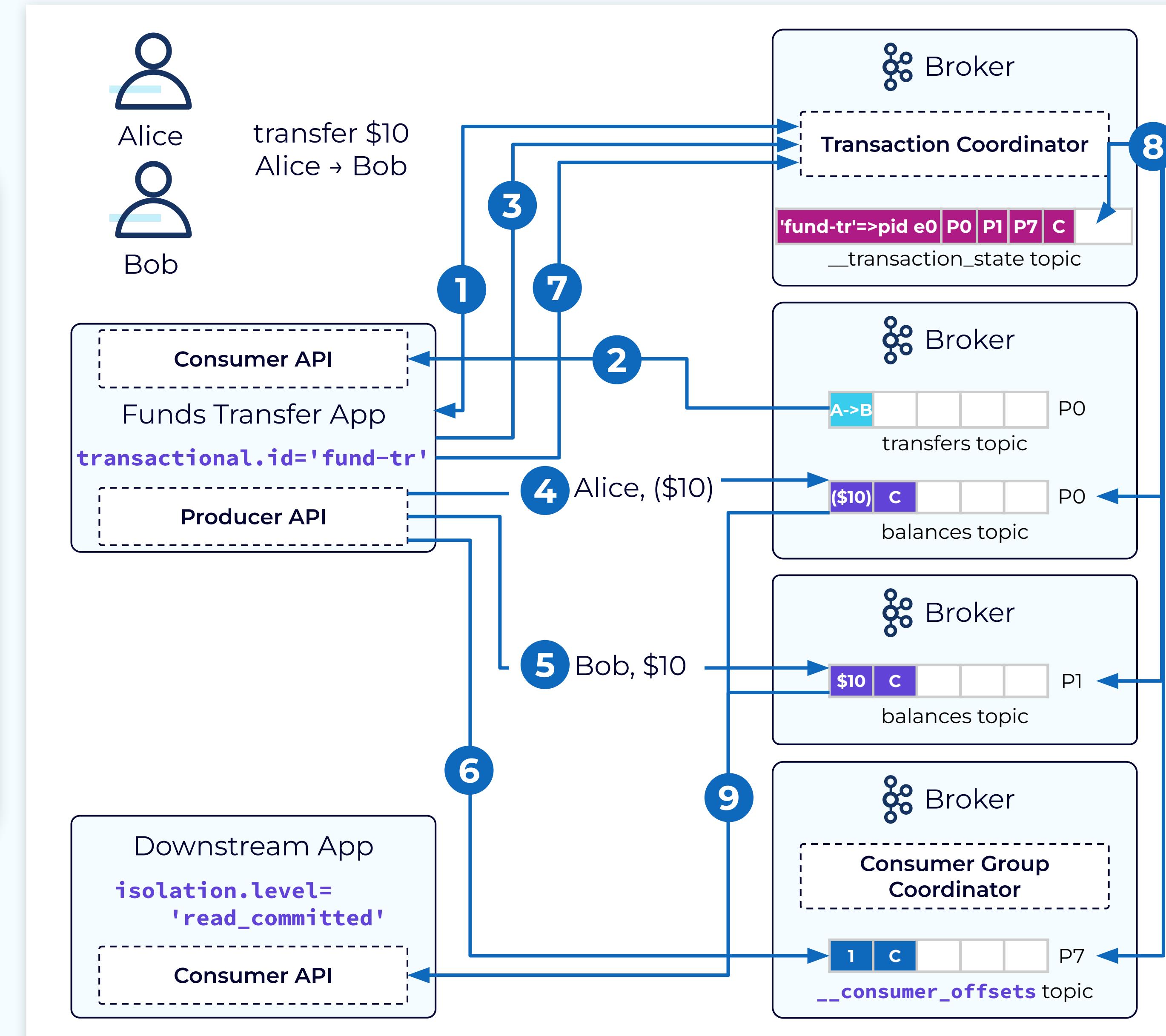
1. New instance requests txn ID
 - a. Coordinator fences previous instance by aborting pending txn and bumping up epoch
2. Downstream consumer with `read_committed` discards aborted events



System with Successful Committed Transaction



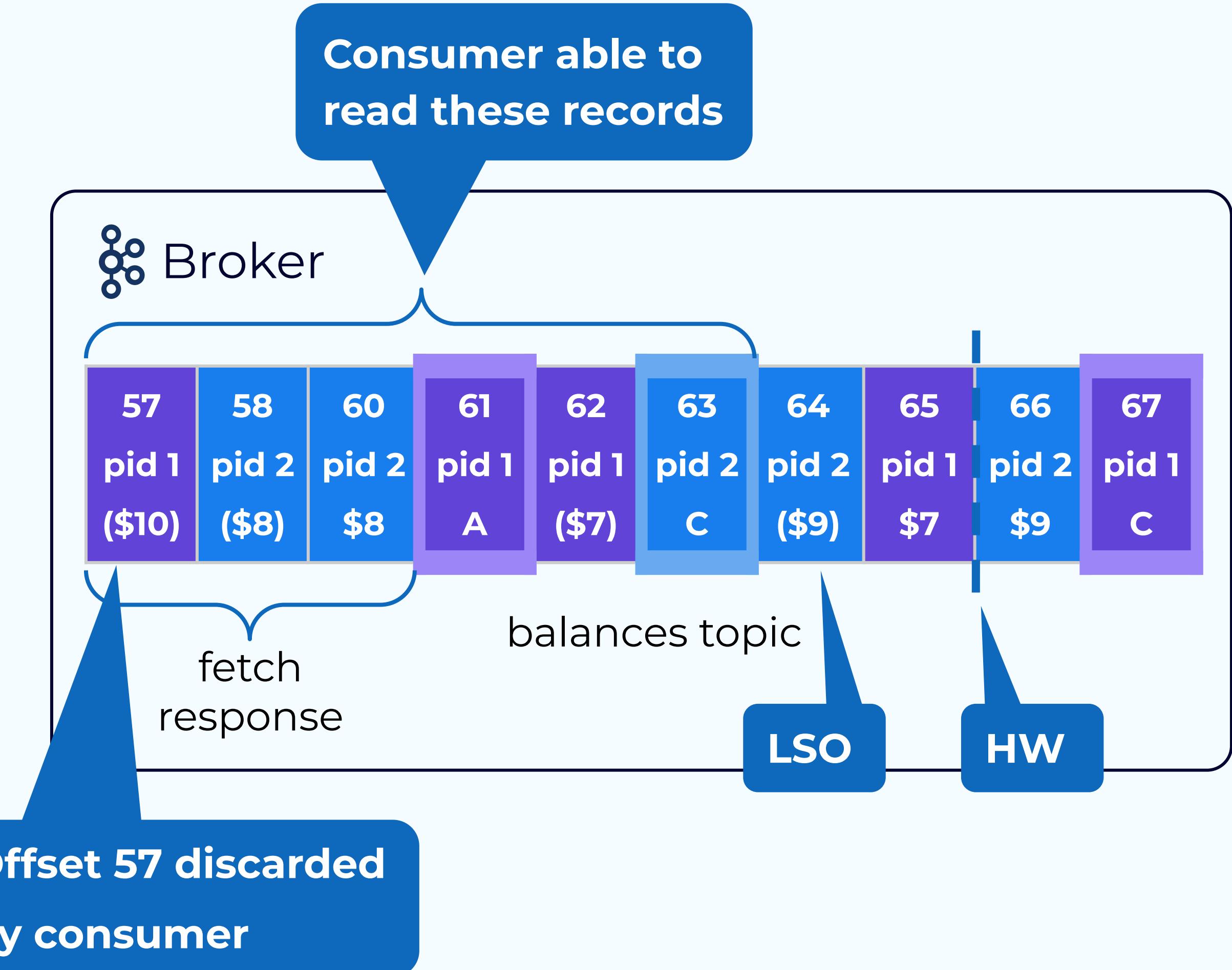
1. Requests txn ID and assigned PID and epoch
2. Event fetched by consumer
3. Notifies coordinator of partition being written to
4. Alice's account debited
5. Bob's account credited
6. Consumer offset committed
7. Notify coordinator that transaction is complete
8. Coordinator writes commit markers to p0, p1, p7
9. Downstream consumer with `read_committed` processes committed events



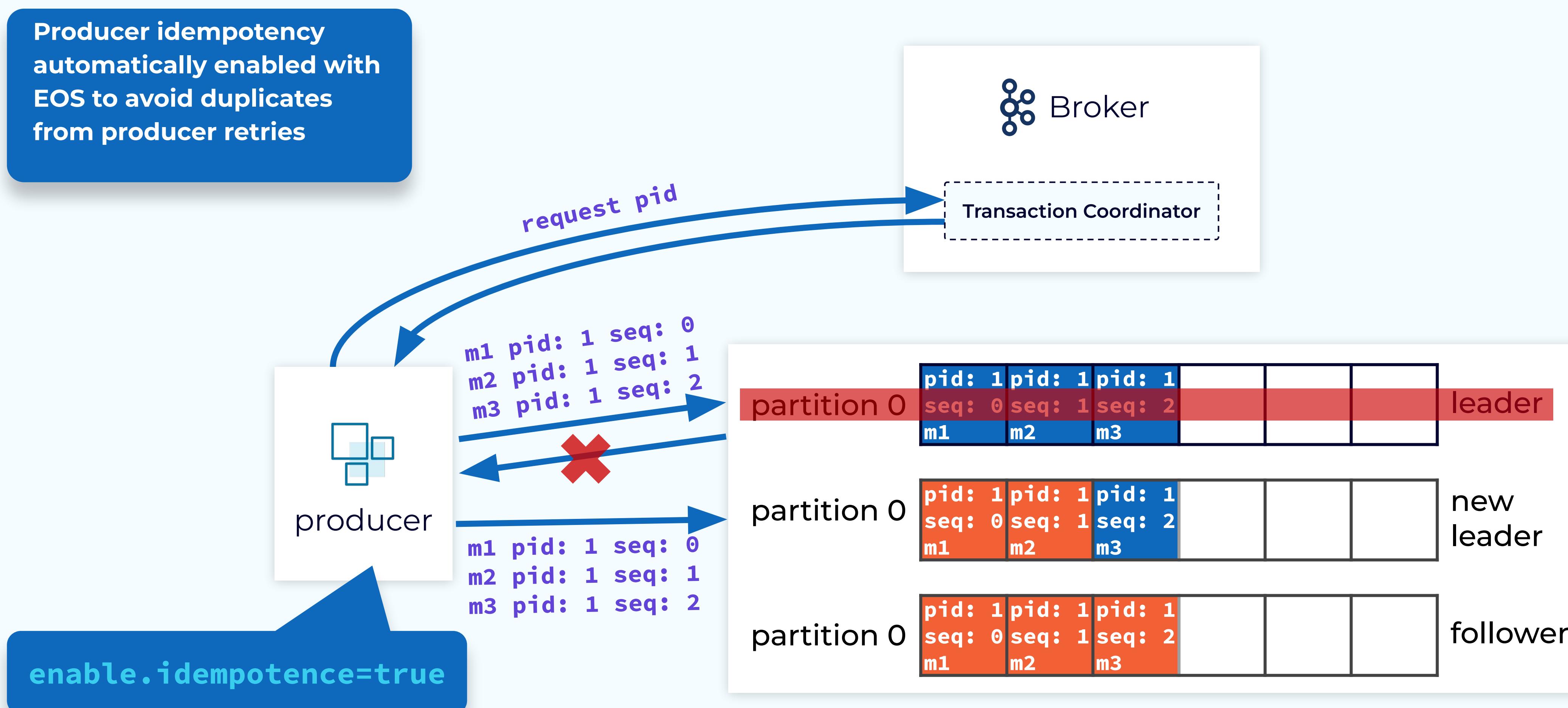
Consuming Transactions with `read_committed`



- Leader maintains last stable offset (LSO), the smallest offset of any open transaction
- Fetch response includes
 - only records up to LSO
 - metadata for skipping aborted records



Transactions: Producer Idempotency



Transactions: Balance Overhead with Latency

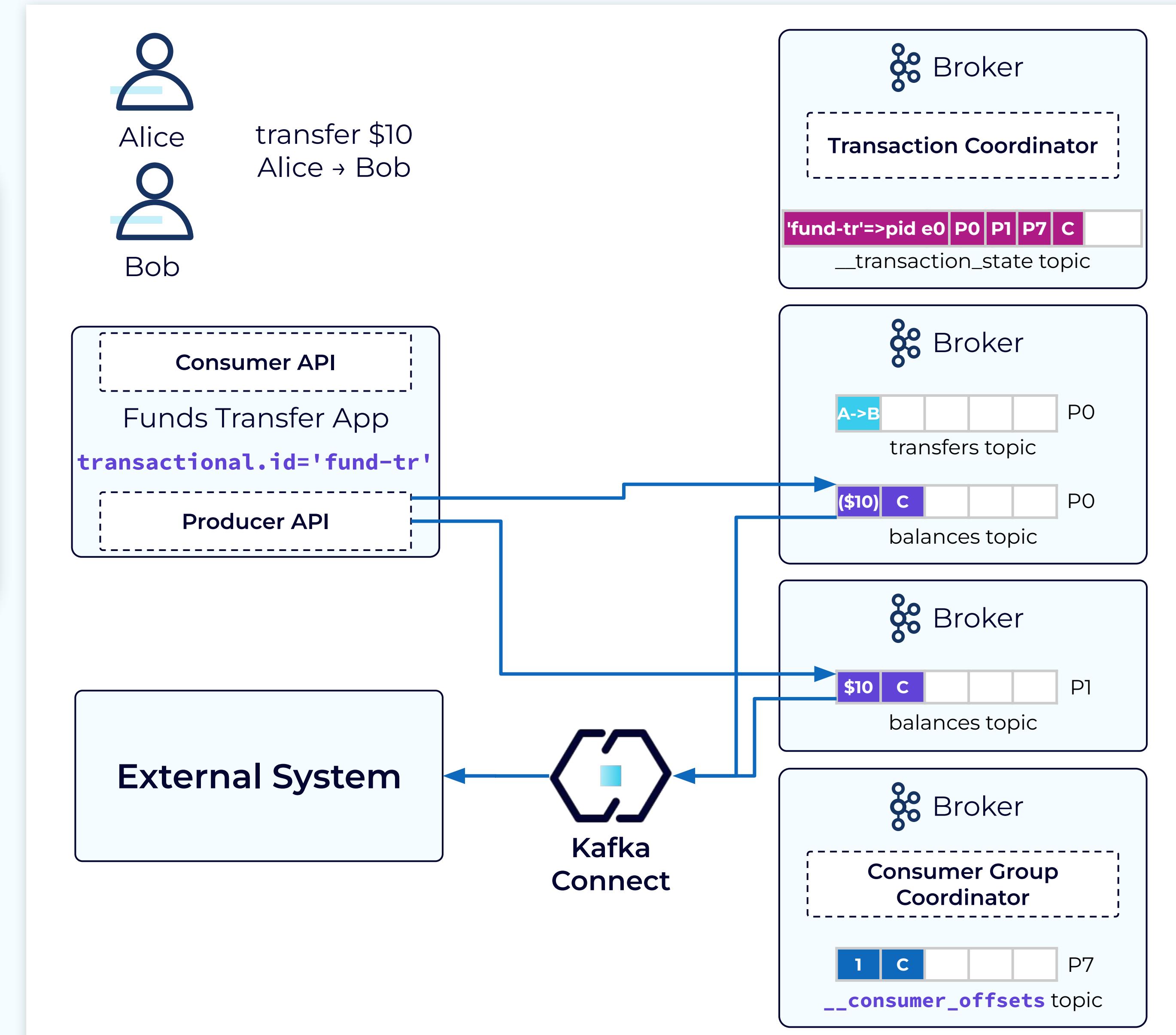
- Too few records per txn adds noticeable overhead
- Too many records per txn delays the exposure of output
- Controllable through `commit.interval.ms` in KStreams

Interacting with External Systems



Atomic writes to Kafka and external systems are not supported

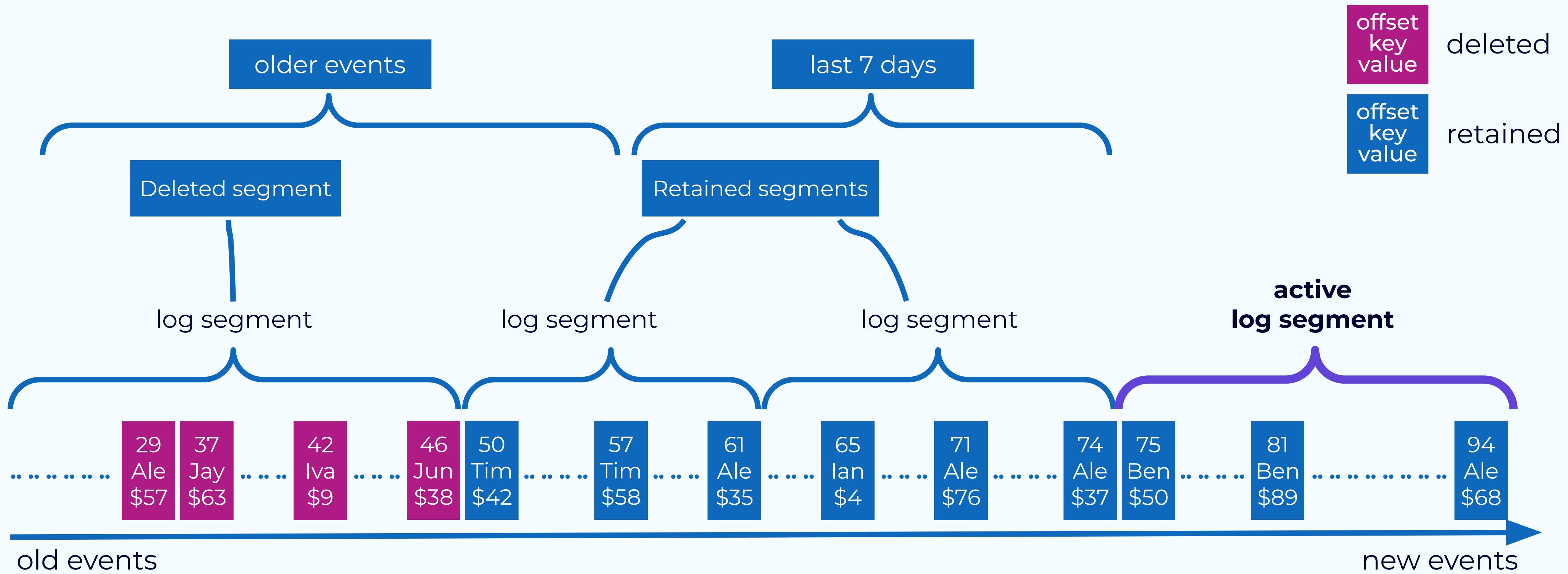
- Instead, write the transactional output to a Kafka topic first
- Rely on idempotency to propagate the data from the output topic to the external system





Topic Compaction

Time-Based Retention



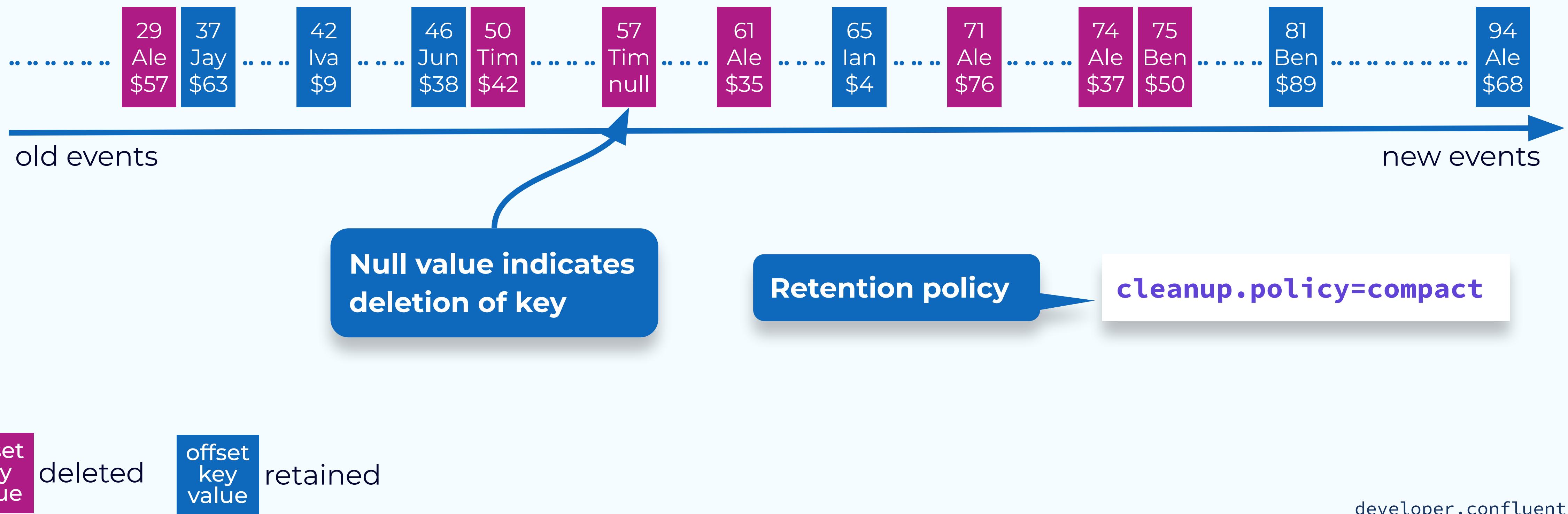
Time based retention
applied at the log segment
level of granularity

`cleanup.policy=delete`
`retention.ms=60480000`
(7 days default)

Topic Compaction: Key-Based Retention



- Most recent occurrence per key is kept
- All other occurrences marked for deletion over time



Usage and Benefits of Topic Compaction

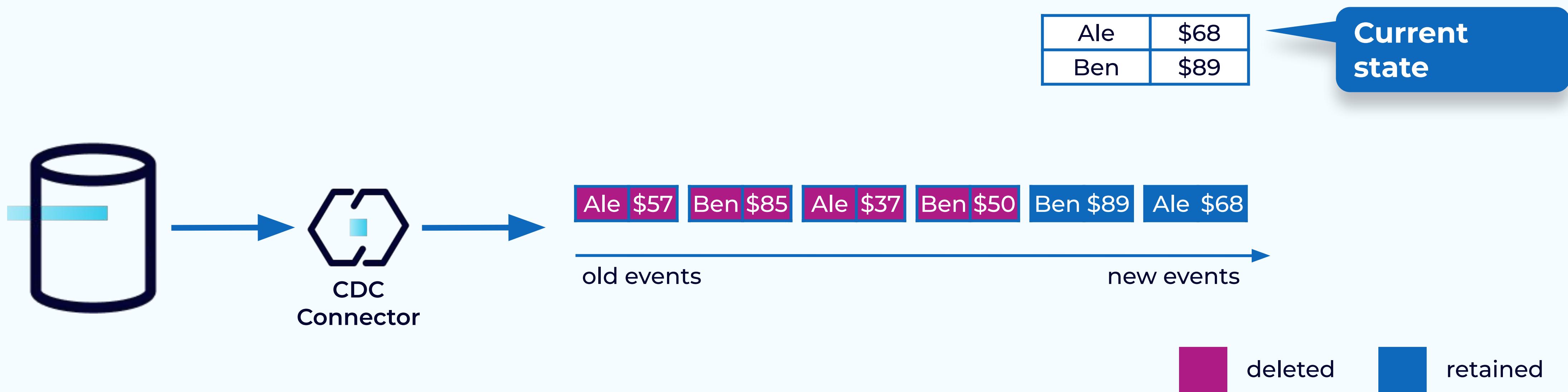


Usage:

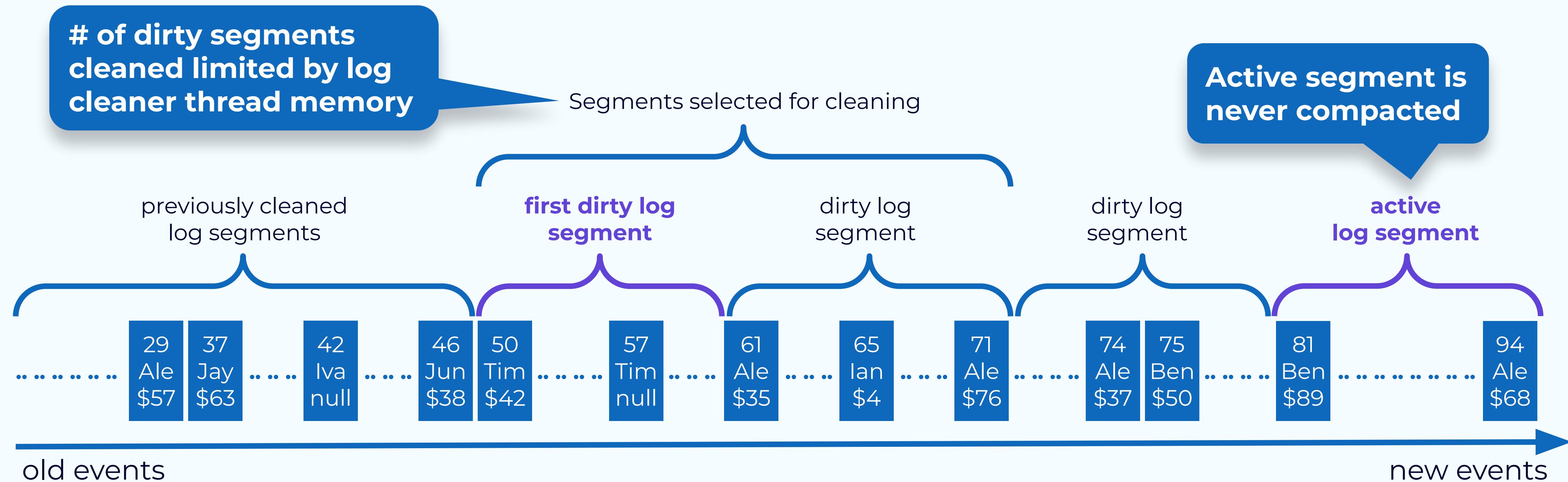
- Streaming updatable data sets (e.g. profiles, catalogue)
- ksqlDB state and tables

Benefits:

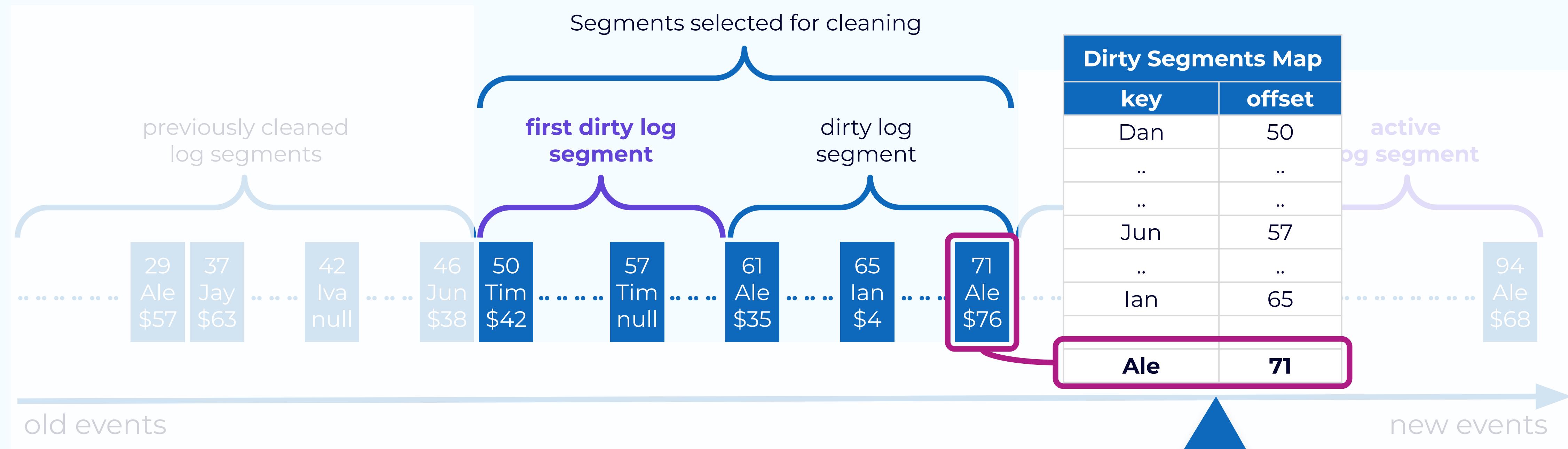
- Keeping latest info with bounded storage
- Allows incremental consumption and bootstrap to be done the same way



Compaction Process - Segments to Clean

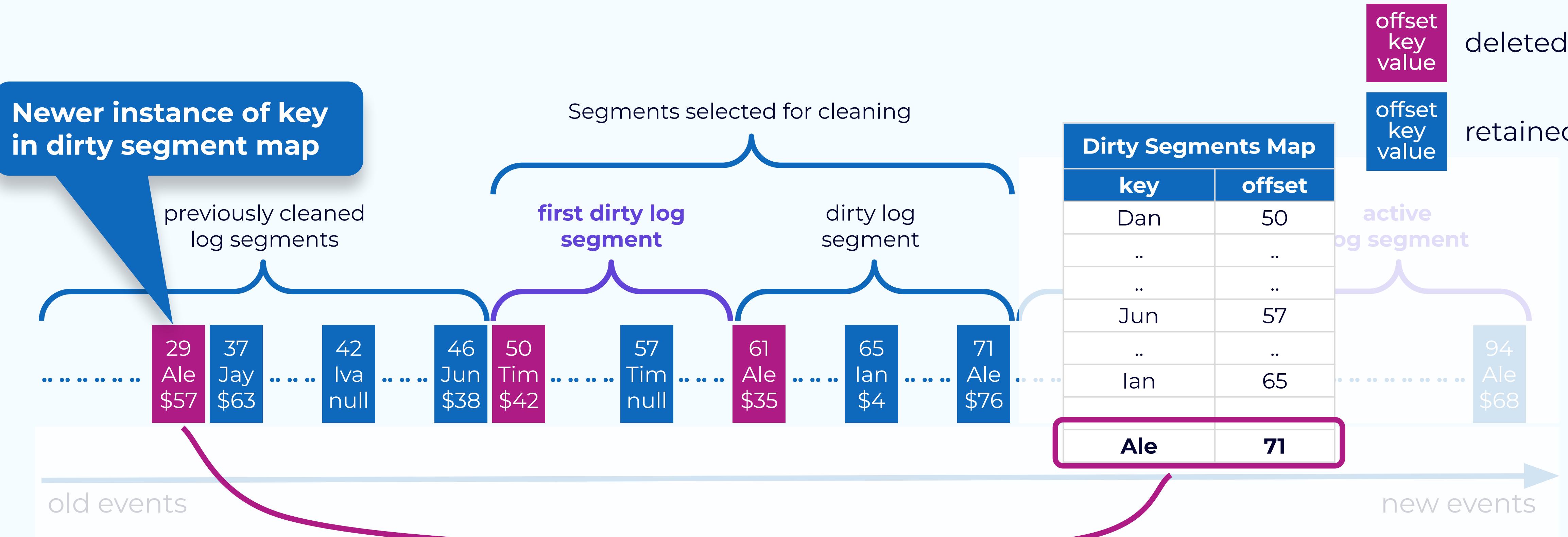


Compaction Process - Build Dirty Segment Map

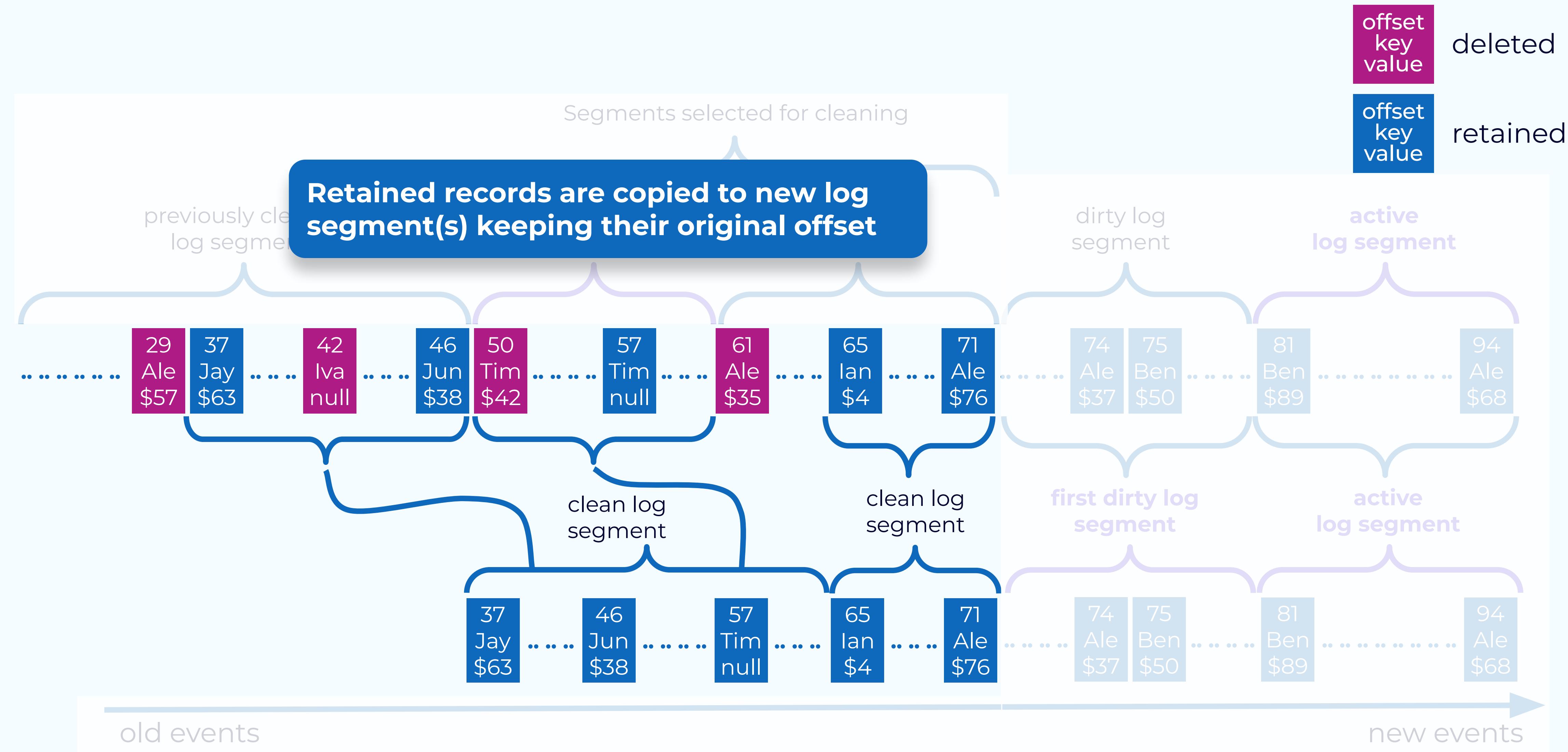


In-memory map of latest offset for each key in selected dirty segments

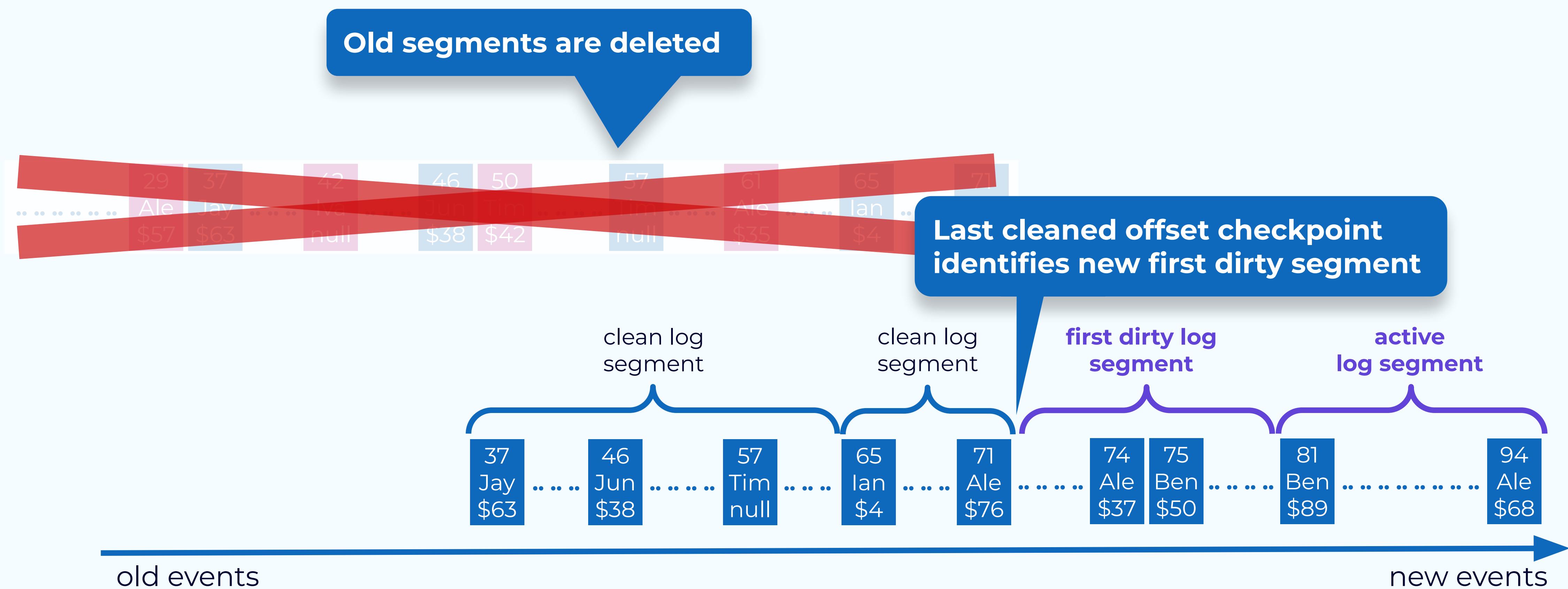
Compaction Process - Deleting Events



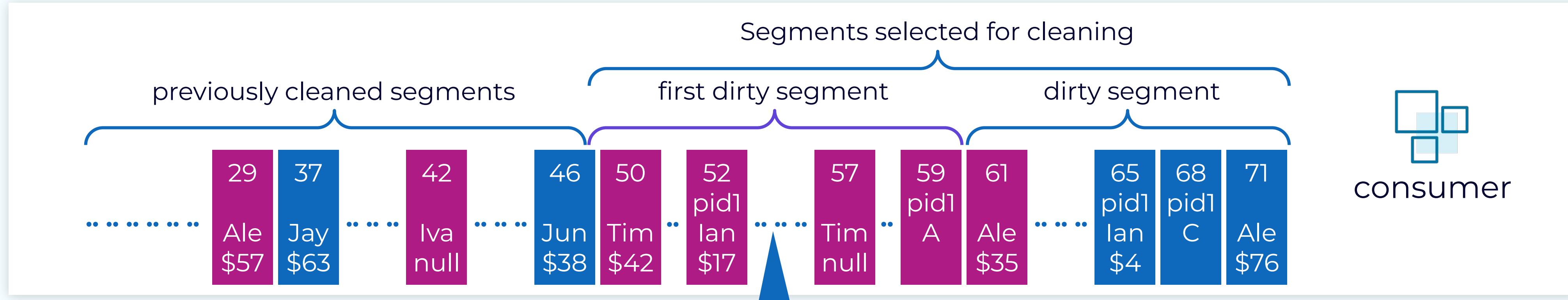
Compaction Process - Retaining Events



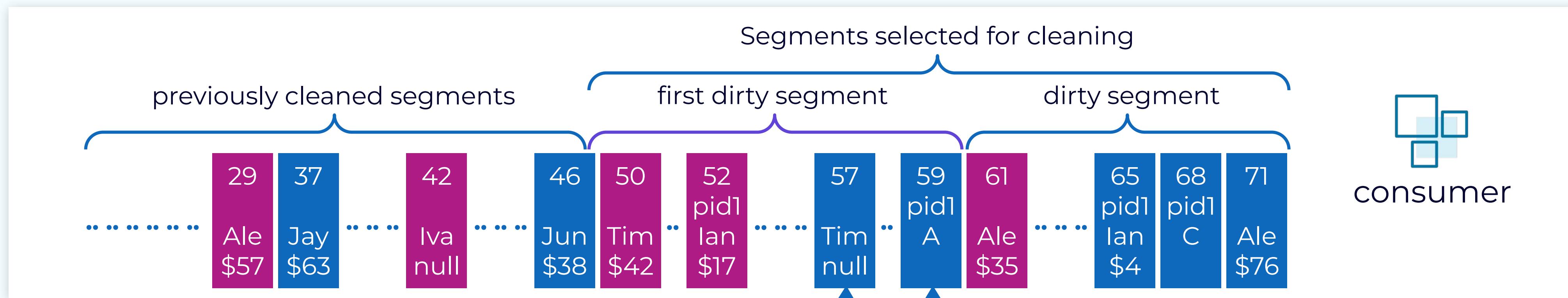
Compaction Process - Replace Old Segments



Cleaning Tombstone and Transaction Markers



Paused consumer with committed offset = 54 misses offset 57 null and 59 txn commit records



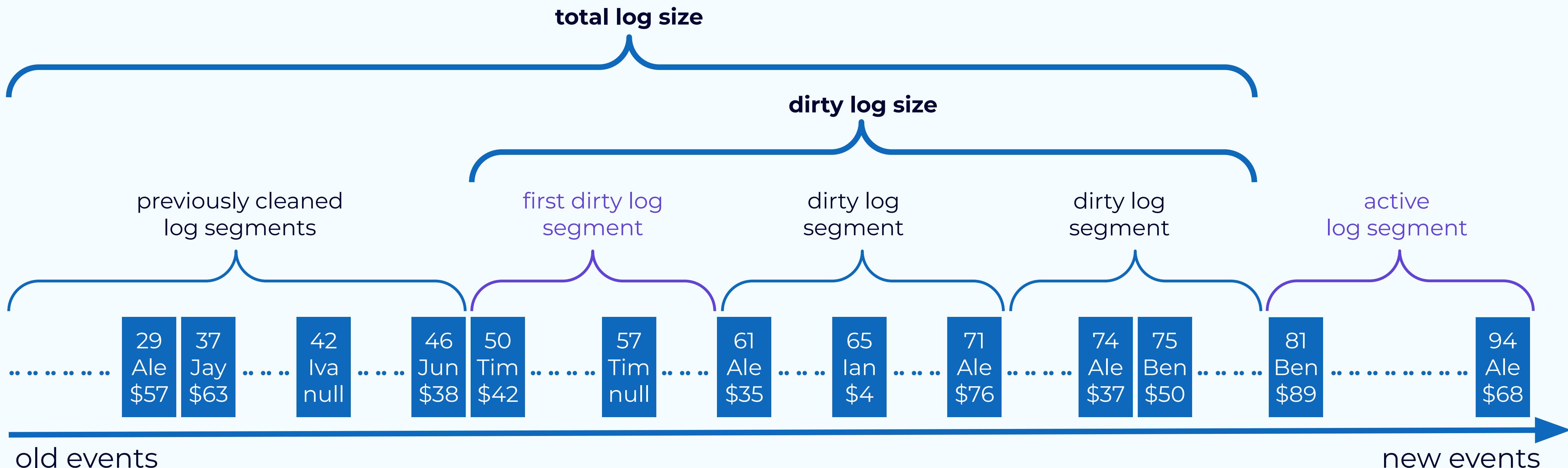
Solution: delay removal of tombstone and txn markers by `first-clean-time + delete.retention.ms`

When Compaction Is Triggered



A topic partition is compacted if one of the following is true:

1. `dirty / total > min.cleanable.dirty.ratio` and
`message timestamp >= current time - min.compaction.lag.ms`
2. `message timestamp < current time - max.compaction.lag.ms`



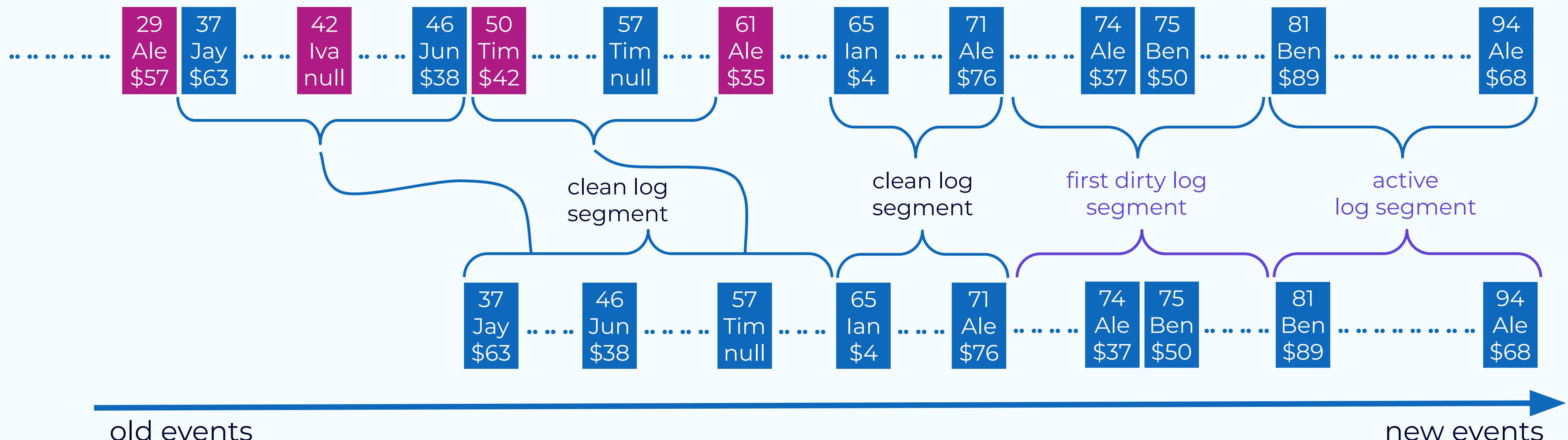
Topic Compaction Guarantees



A consumer:

- 1) is not guaranteed to see every record
- 2) is guaranteed to see the latest record for a key

offset key value	deleted
offset key value	retained



Topic Compaction



developer.confluent.io

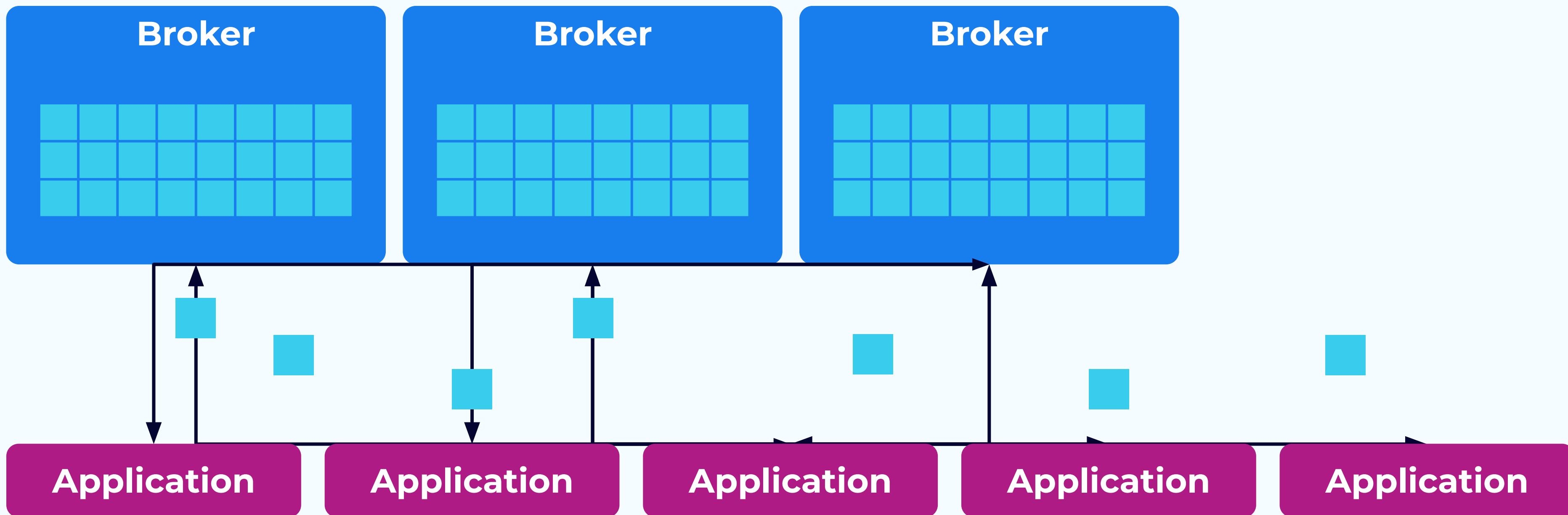


CONFLUENT

Developer

Tiered Storage

Current Kafka Storage and Its Shortcomings

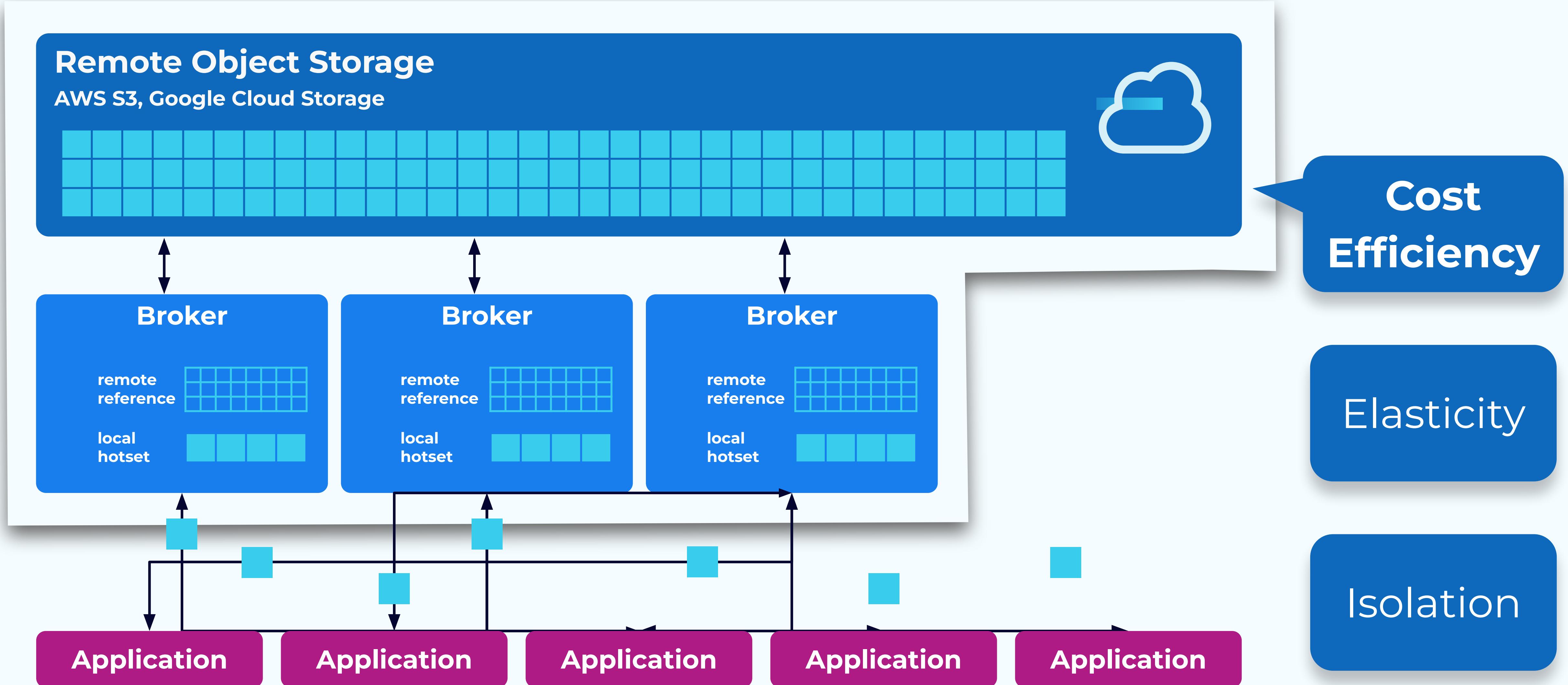


Cost

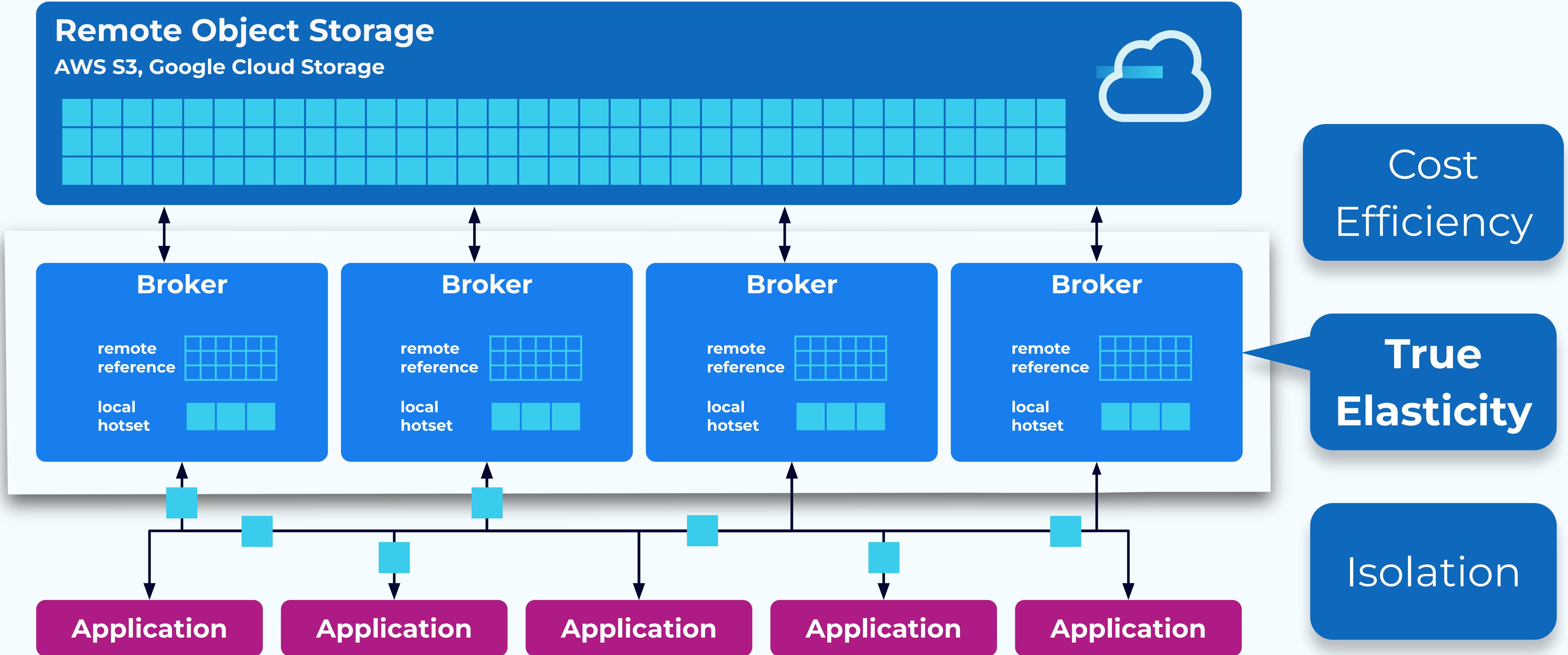
Elasticity

Isolation

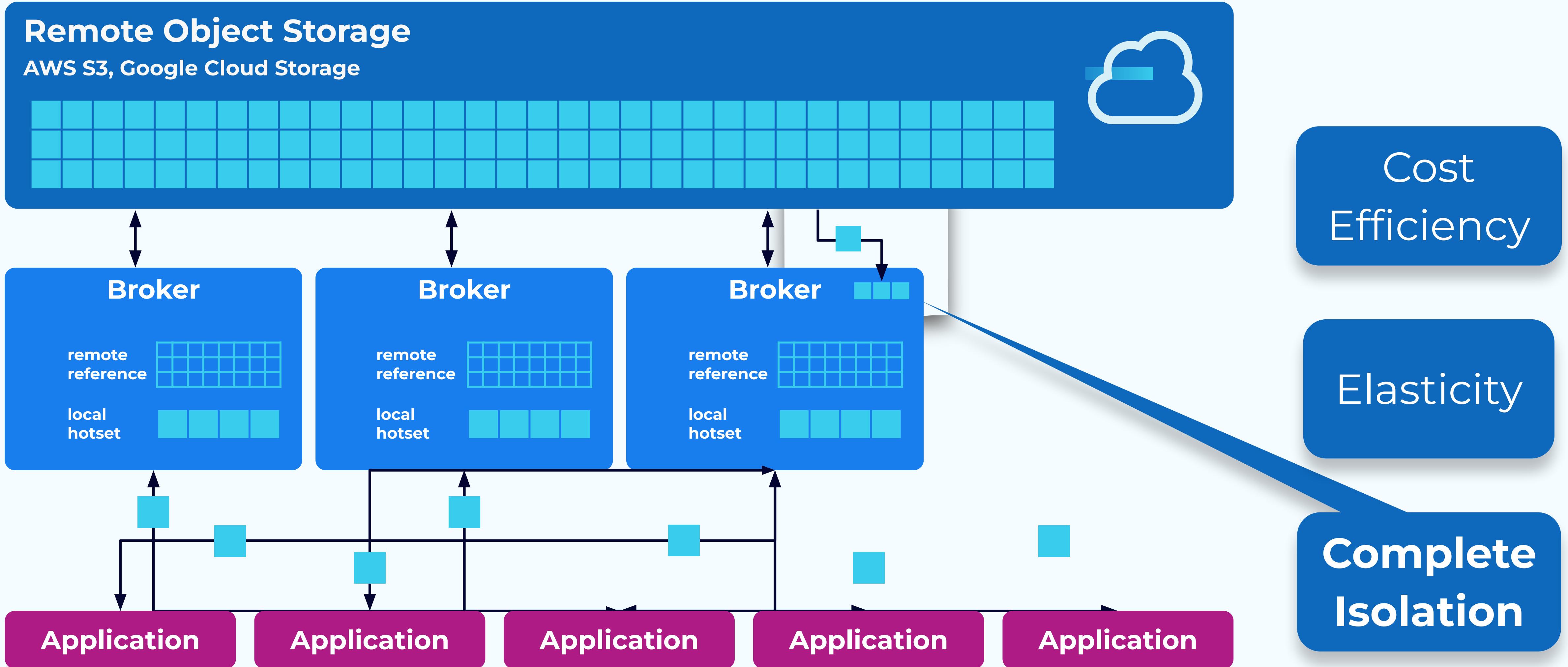
Tiered Storage



Tiered Storage



Tiered Storage

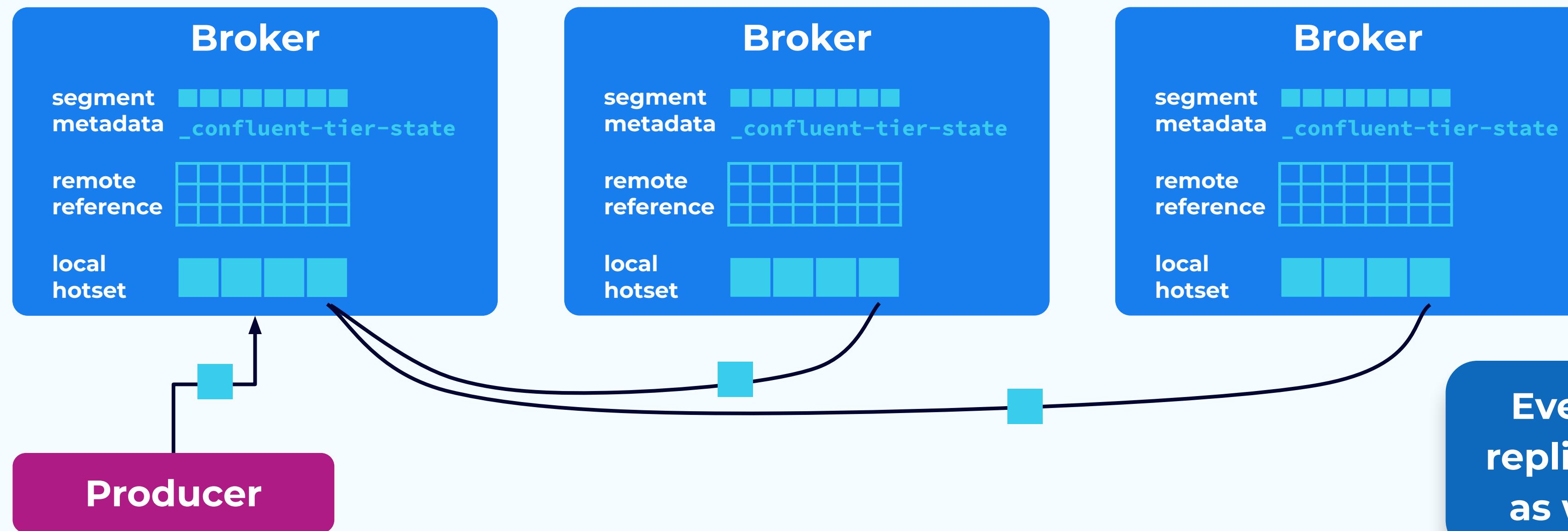


Writing Events to a Tiered Topic

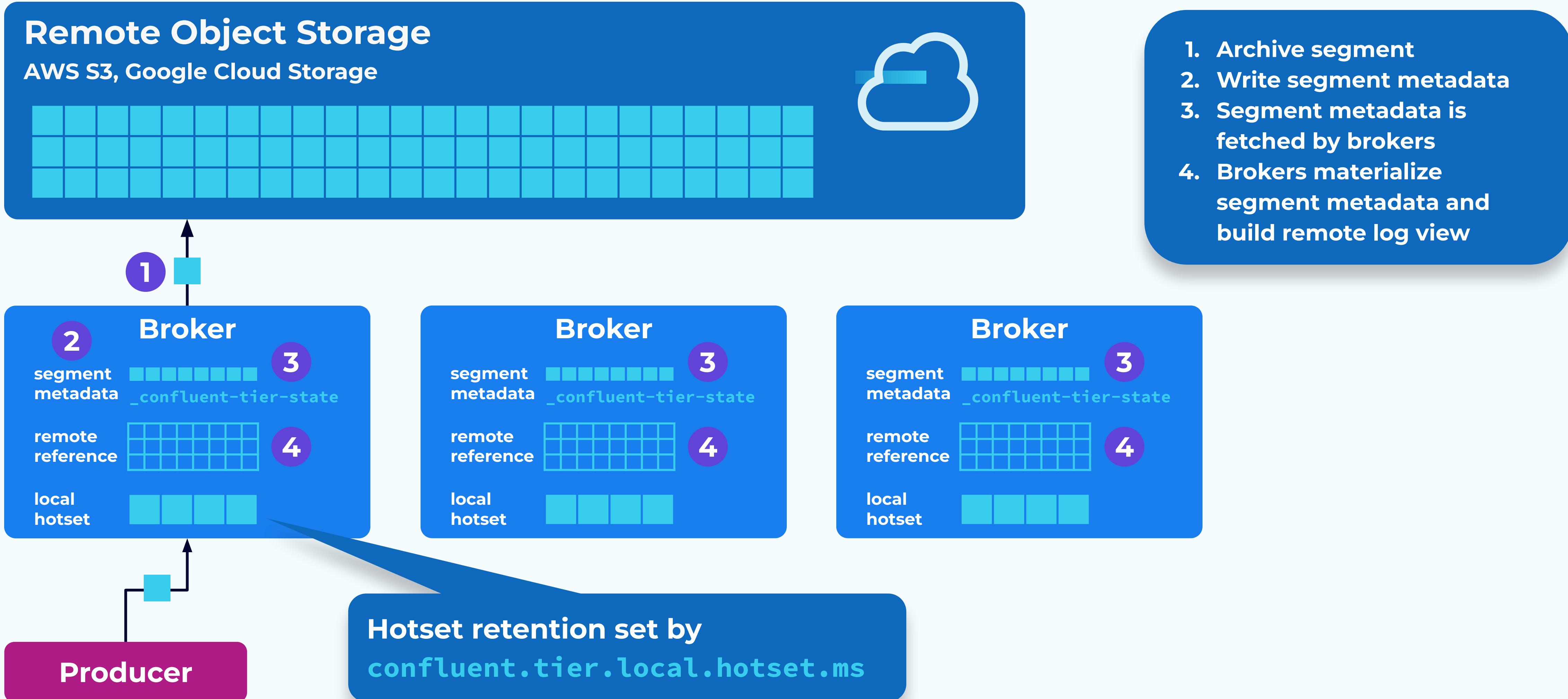


Remote Object Storage

AWS S3, Google Cloud Storage



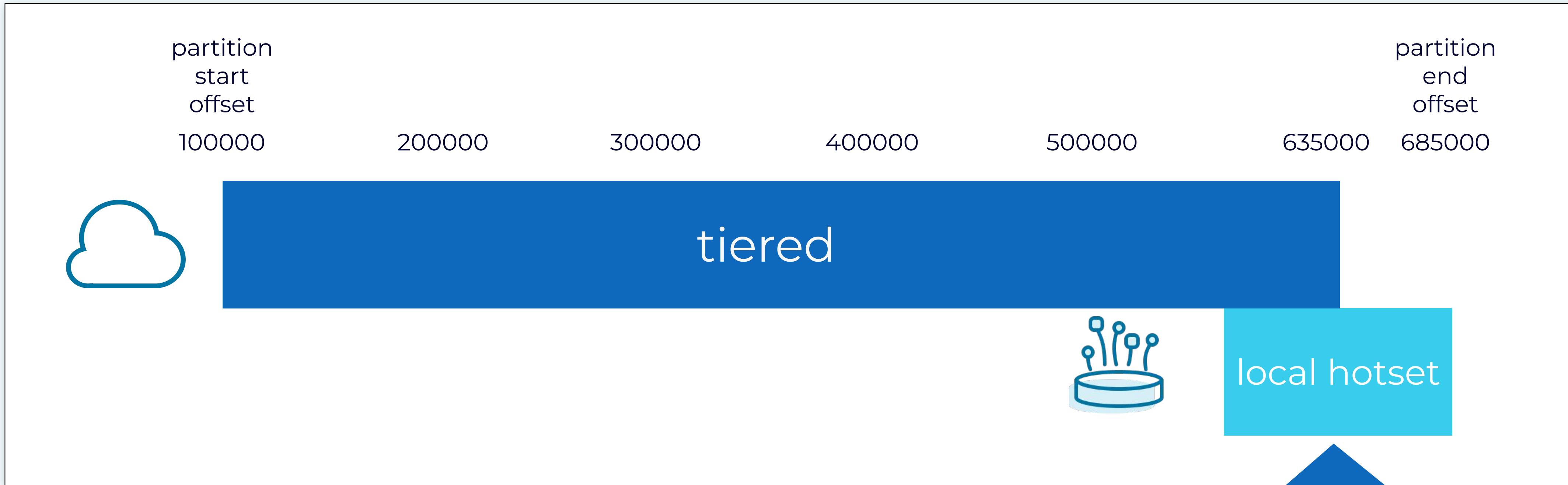
Tiering Events to the Remote Object Store



Broker Logical View of Tiered Partition



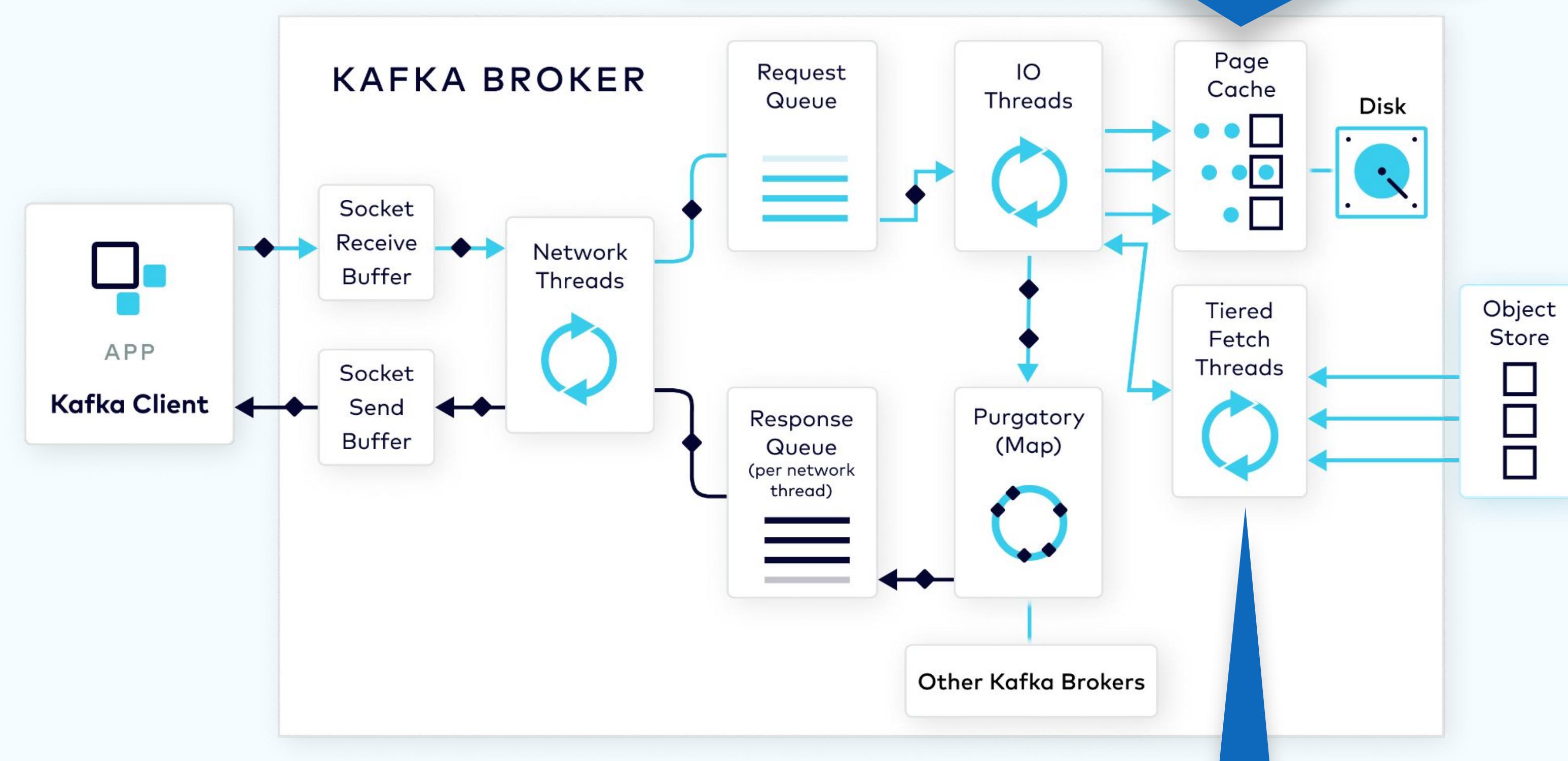
Logical view of log based upon metadata contained in `_confluent-tier-state`



Recent events can exist in both the hotset and object storage

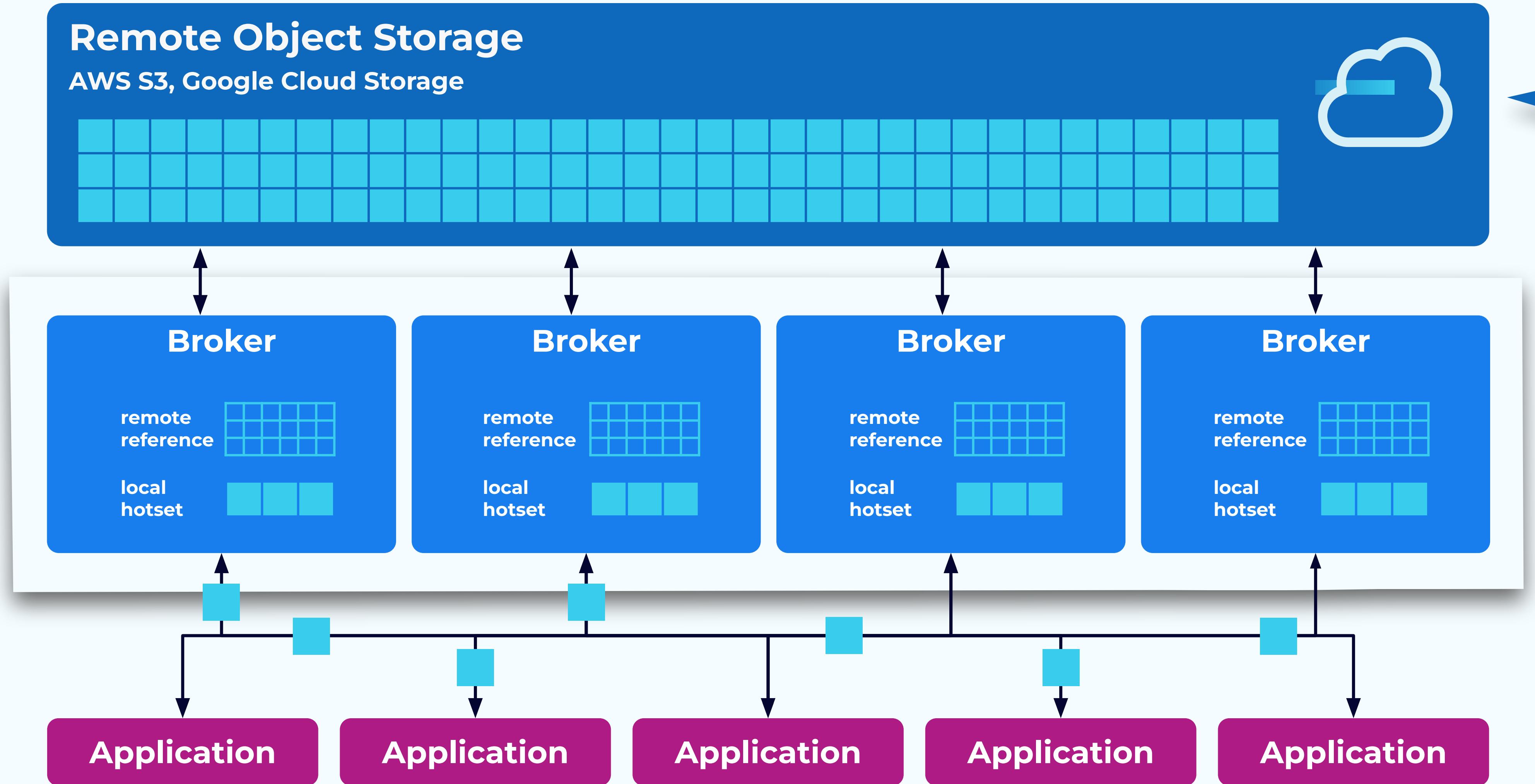
Fetching Tiered Data

If data is in the hotset, it is returned using standard Kafka methods



If data is not in the hotset, broker retrieves it async from remote storage and serves it to the client

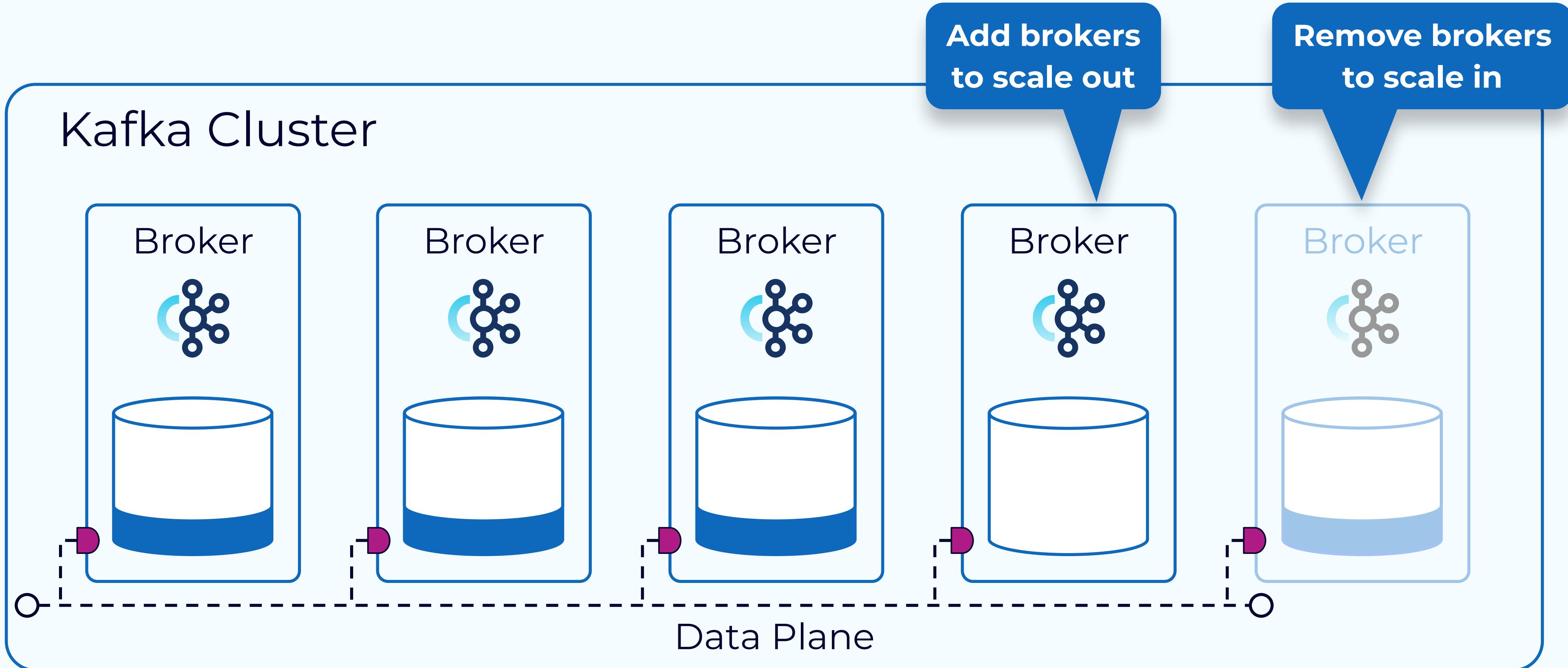
Tiered Storage Portability



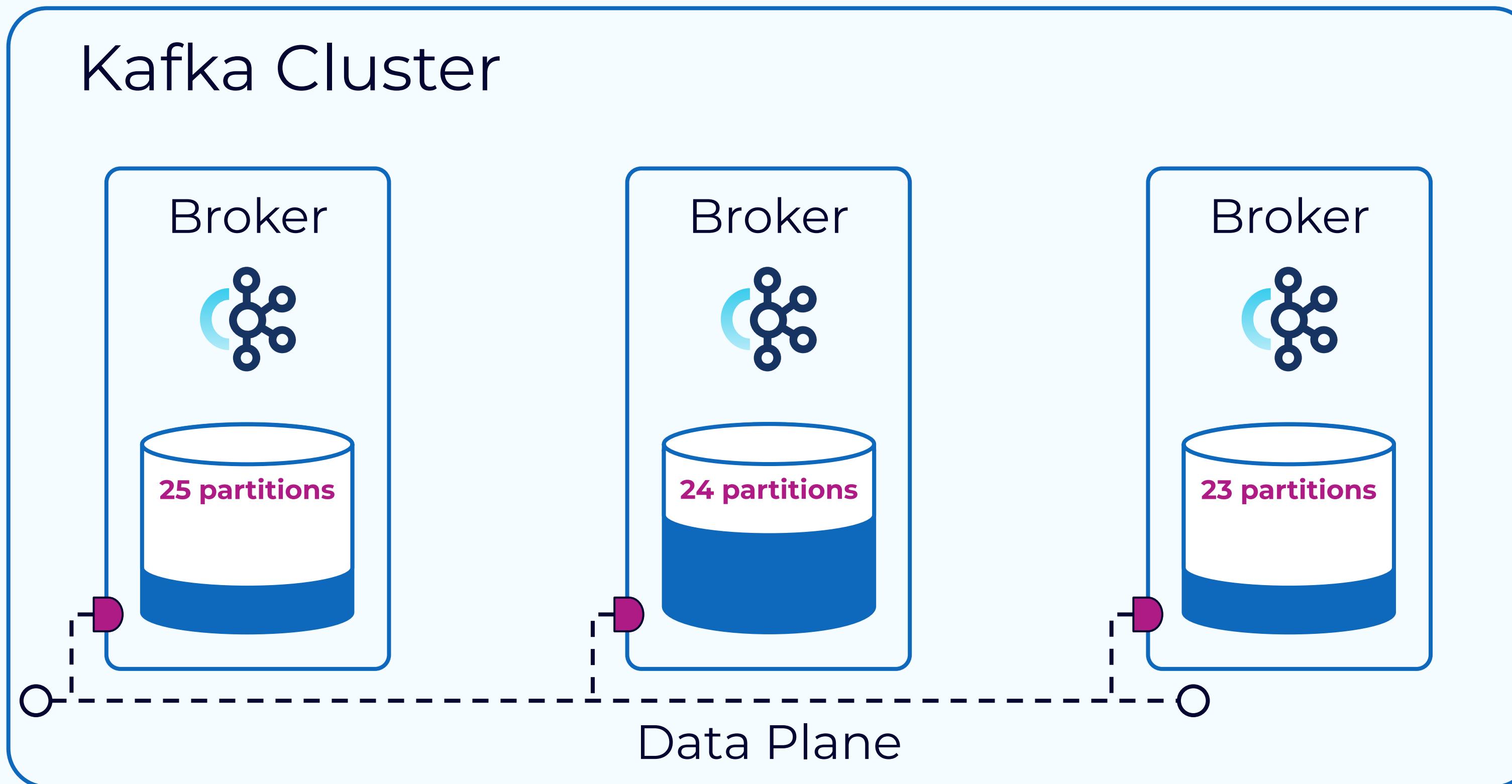


Cluster Elasticity

Cluster Scaling



Unbalanced Data Distribution



Automatic Cluster Scaling with Confluent Cloud



Confluent Cloud

kafka-reassign-partitions.sh

Included with Kafka core

- 1) Relocates partitions based upon JSON configuration file
- 2) Provides granular control
- 3) Manual process to create JSON configuration file

```
kafka-reassign-partitions \
--bootstrap-server server-1:39094 \
--reassignment-json-file reassign-topic-1.json \
--execute
```

```
{"version":1,
"partitions":[
    {"topic":"topic-1","partition":0,"replicas":[1,2,3]},
    {"topic":"topic-1","partition":1,"replicas":[2,3,1]},
    {"topic":"topic-1","partition":2,"replicas":[3,1,2]},
    {"topic":"topic-1","partition":3,"replicas":[1,2,3]},
    {"topic":"topic-1","partition":4,"replicas":[2,3,1]},
    {"topic":"topic-1","partition":5,"replicas":[3,1,2]}
]}
```

Confluent Auto Data Balancer



Included with Confluent Platform

- 1) Automatically generates redistribution plan based on stats
- 2) When the plan is executed, partitions are rebalanced across brokers in the cluster
- 3) These are both manual steps done with the `confluent-rebalancer` command

```
confluent-rebalancer execute \
--bootstrap-server kafka-1:9092 \
--metrics-bootstrap-server kafka-1:9092 \
--throttle 1000000 \
--verbose
```

```
Computing the rebalance plan (this may take a while) ...
You are about to move 512 replica(s) for 512 partitions to 1 broker(s) with total size 233.9 MB.
The preferred leader for 192 partition(s) will be changed.
In total, the assignment for 515 partitions will be changed.
The minimum free volume space is set to 20.0%.
```

```
The following brokers will have less than 40% of free volume space during the rebalance:
```

Broker	Current Size (MB)	Size During Rebalance (MB)	Free % During Rebalance	Size After Rebalance (MB)	Free % After Rebalance
102	311.7	311.7	37.6	233.8	37.6
104	0	233.9	36.8	233.9	36.8
103	311.7	311.7	37.6	233.7	37.6
101	311.7	311.7	37.6	233.8	37.6

```
Min/max stats for brokers (before -> after):
```

Type	Leader Count	Replica Count	Size (MB)
Min	0 (id: 104) -> 171 (id: 101)	0 (id: 104) -> 510 (id: 102)	0 (id: 104) -> 233.7 (id: 103)
Max	229 (id: 102) -> 171 (id: 101)	684 (id: 102) -> 515 (id: 101)	311.7 (id: 101) -> 233.9 (id: 104)

```
No racks are defined.
```

```
Broker stats (before -> after):
```

Broker	Leader Count	Replica Count	Size (MB)	Free Space (%)
101	226 -> 171	683 -> 515	311.7 -> 233.8	37.6 -> 37.8
102	229 -> 171	684 -> 510	311.7 -> 233.8	37.6 -> 37.8
103	229 -> 171	683 -> 513	311.7 -> 233.7	37.6 -> 37.8
104	0 -> 171	0 -> 512	0 -> 233.9	37.6 -> 36.8

```
Would you like to continue? (y/n):
```

```
The rebalance has been started, run `status` to check progress.
```

```
Warning: You must run the `status` or `finish` command periodically, until the rebalance completes, to ensure the throttle is removed. You can also alter the throttle by re-running the execute command passing a new value.
```

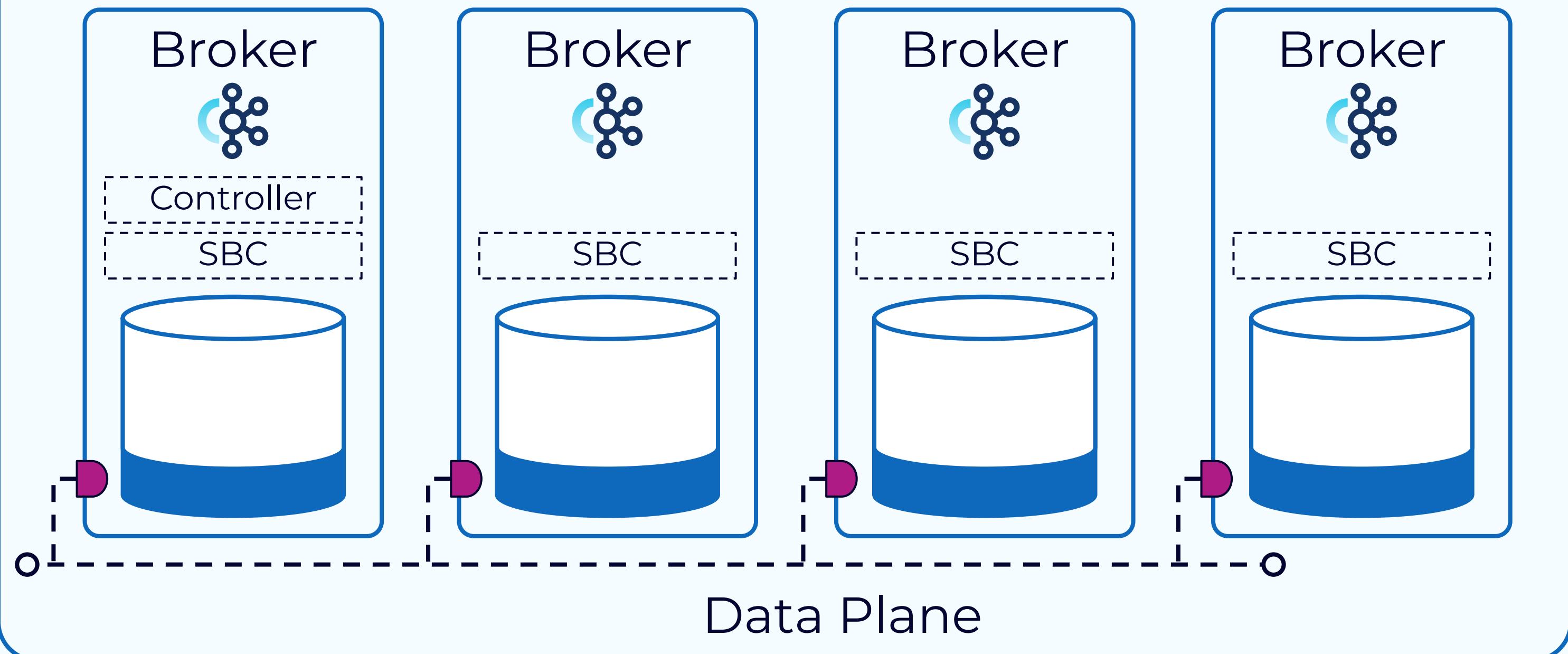
Confluent Self-Balancing Clusters (SBC)



Advantages over Auto Data Balancer

- 1) No additional tools to run (built into the brokers)
- 2) Cluster balance is continuously monitored so that rebalances run whenever they are needed
- 3) Much faster rebalancing (several orders of magnitude)

Kafka Cluster



SBC Metrics Collection and Processing



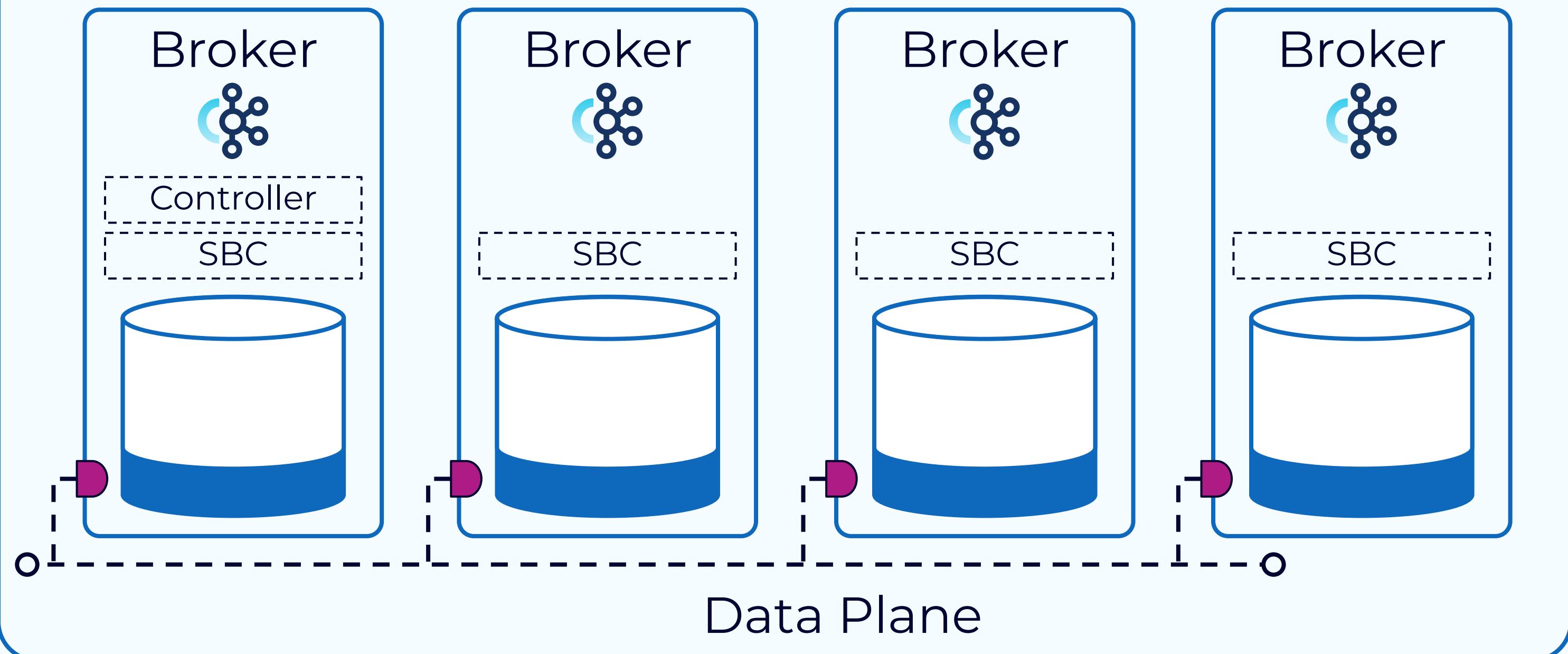
SBC runs on each broker:

- 1) Collects metrics for that broker
- 2) Writes the metrics to an internal metrics collection topic

SBC on the controller node:

- 1) Aggregates the metrics
- 2) Generates load balance plans based on goals
- 3) Executes the plan and exposes monitoring data to Control Center

Kafka Cluster



SBC Rebalance Triggers



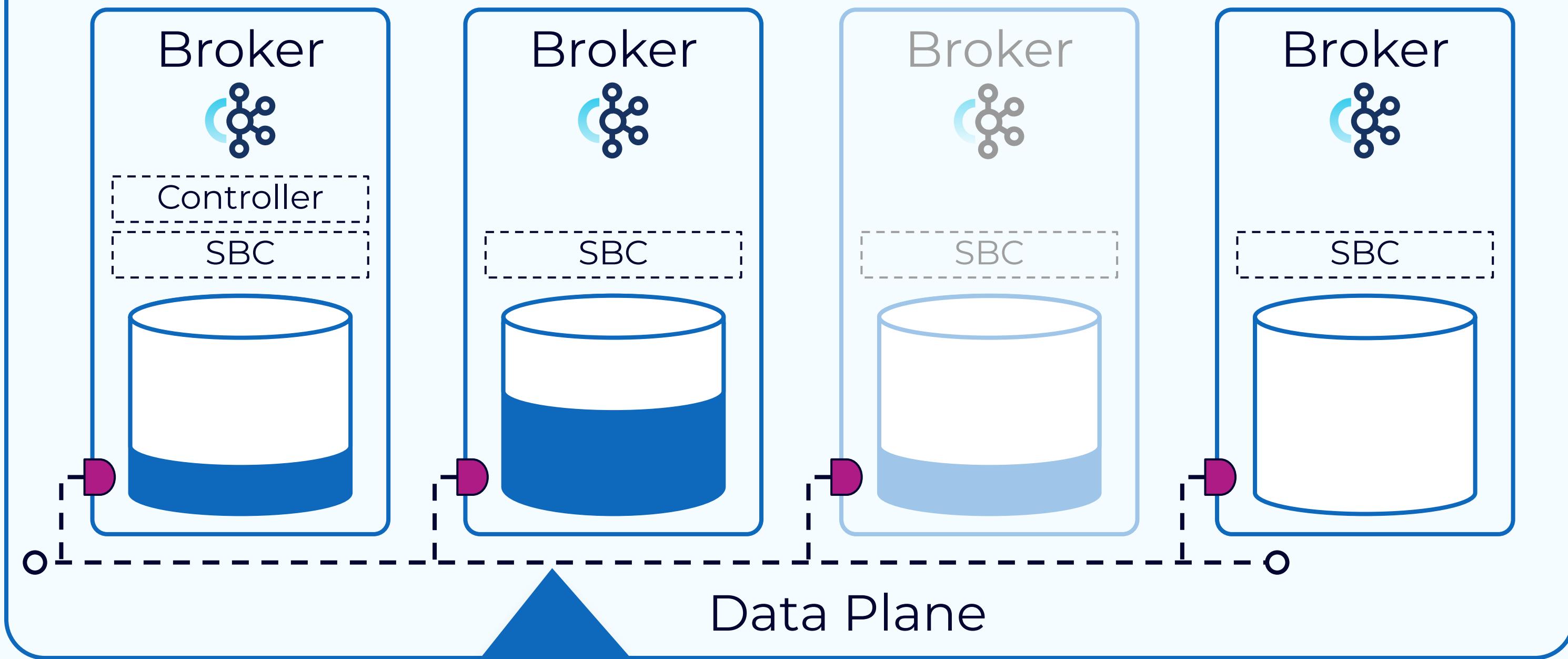
Auto rebalance trigger options:

- 1) Added and removed brokers
- 2) Any uneven load

Uneven load condition is based upon:

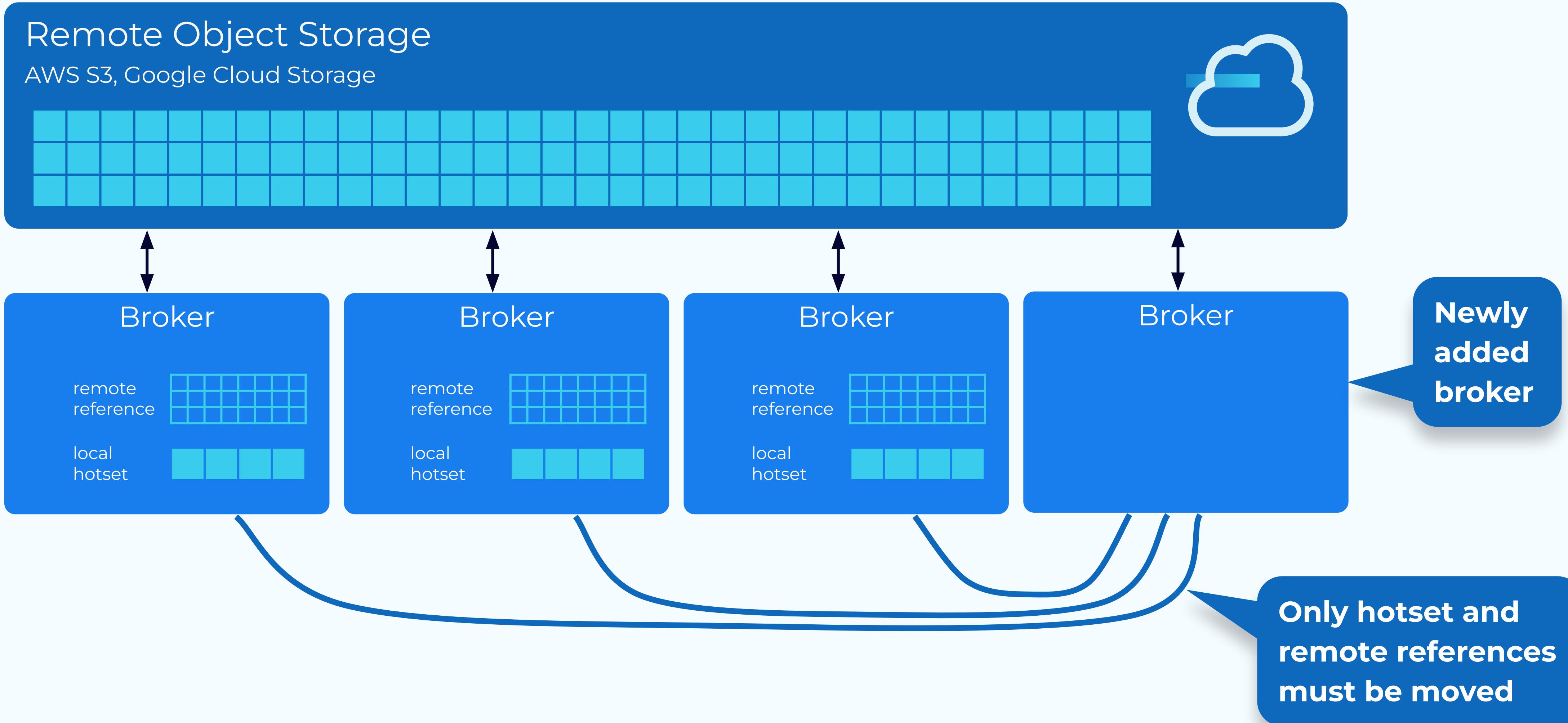
- 1) Disk usage and network usage
- 2) Number of partitions(replicas per broker)
- 3) Leadership and rack awareness

Kafka Cluster



**SBC throttles replication during rebalance
to prevent negative impact to clients**

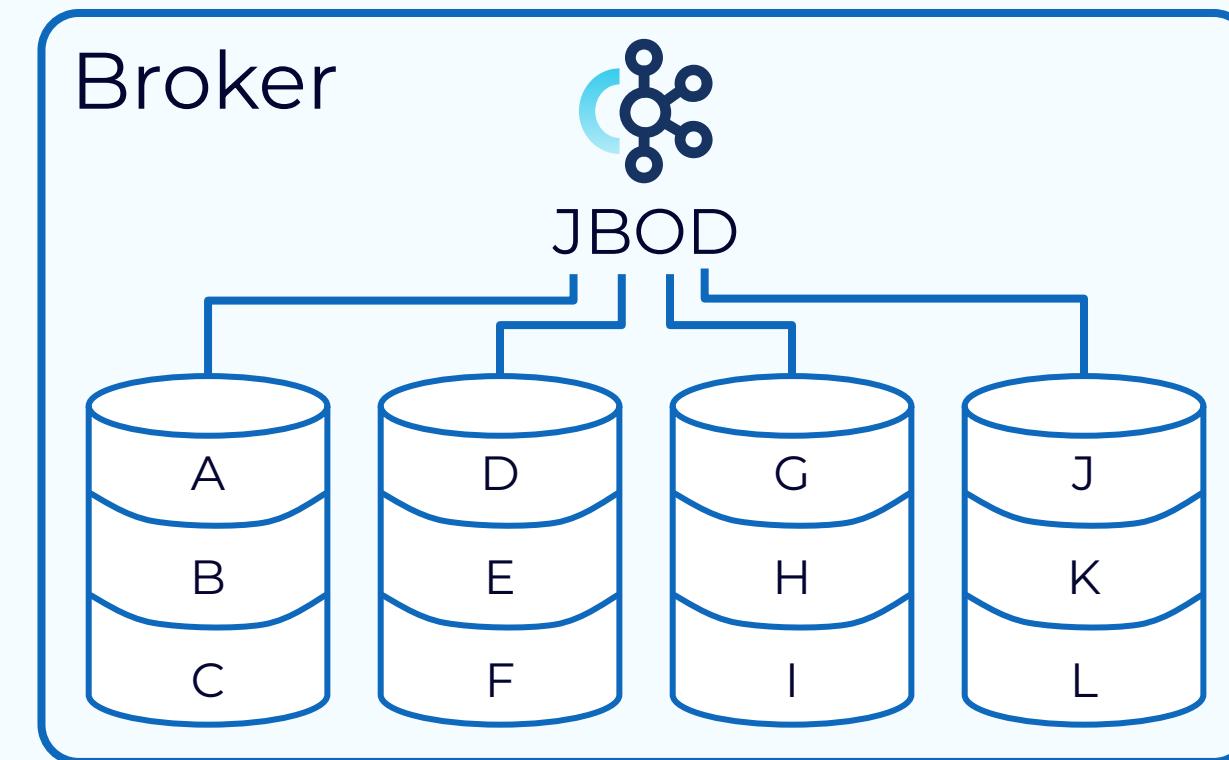
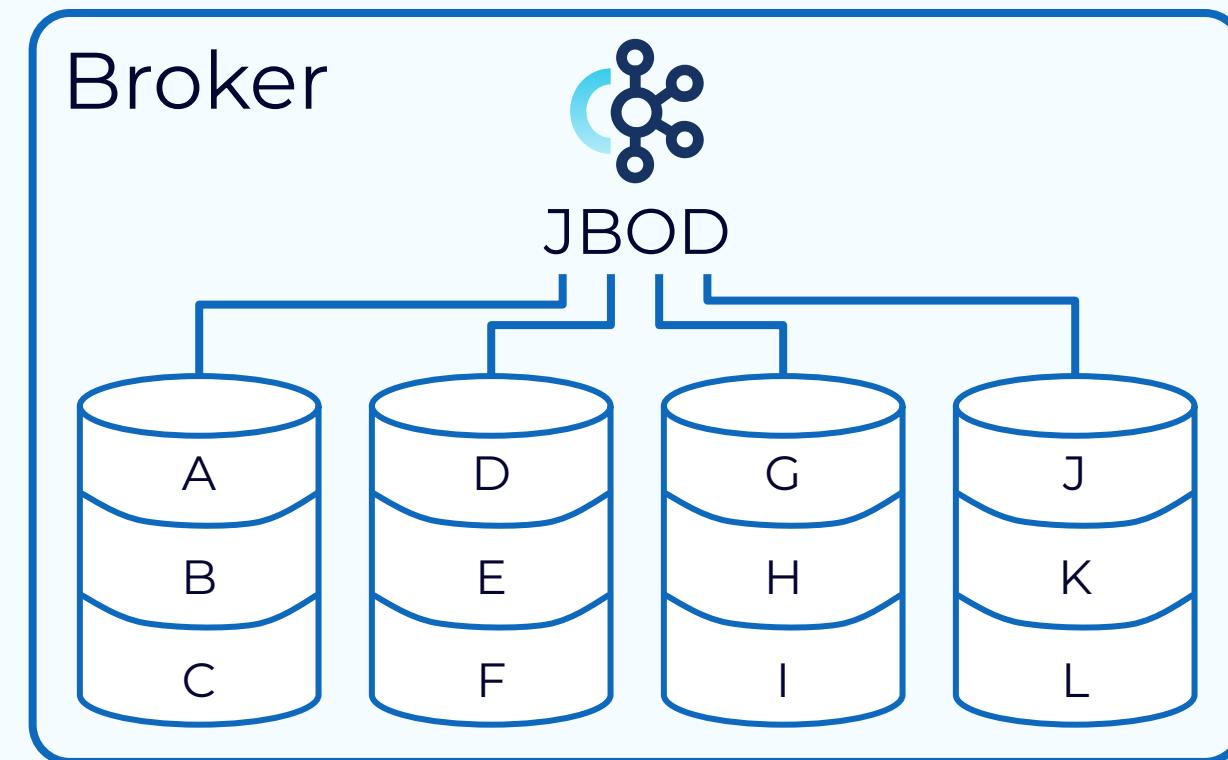
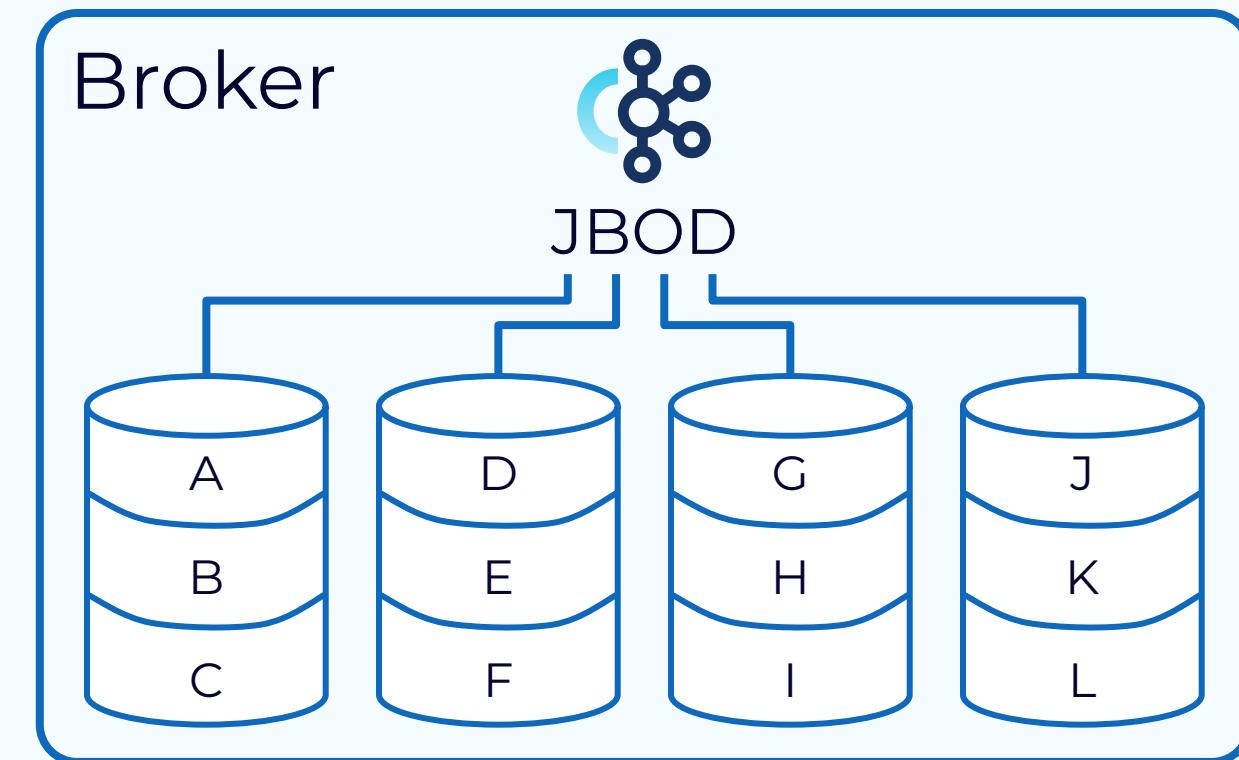
Fast Rebalancing with Tiered Storage



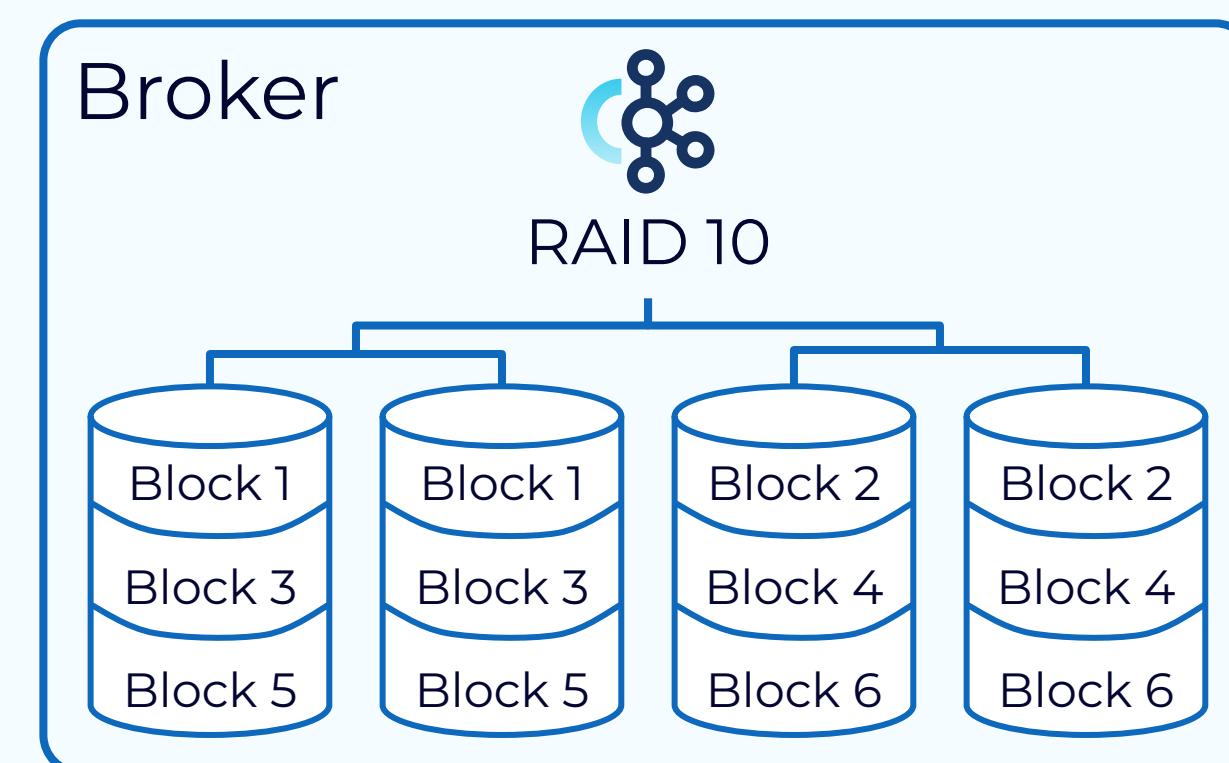
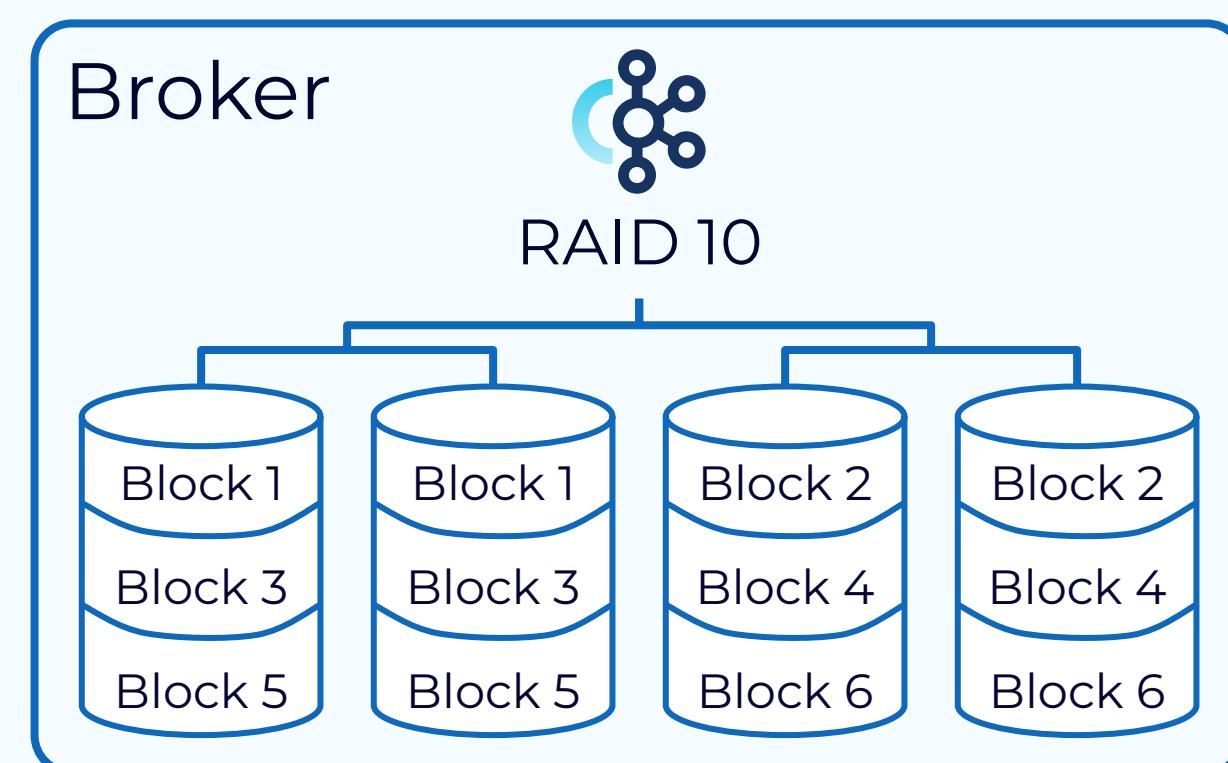
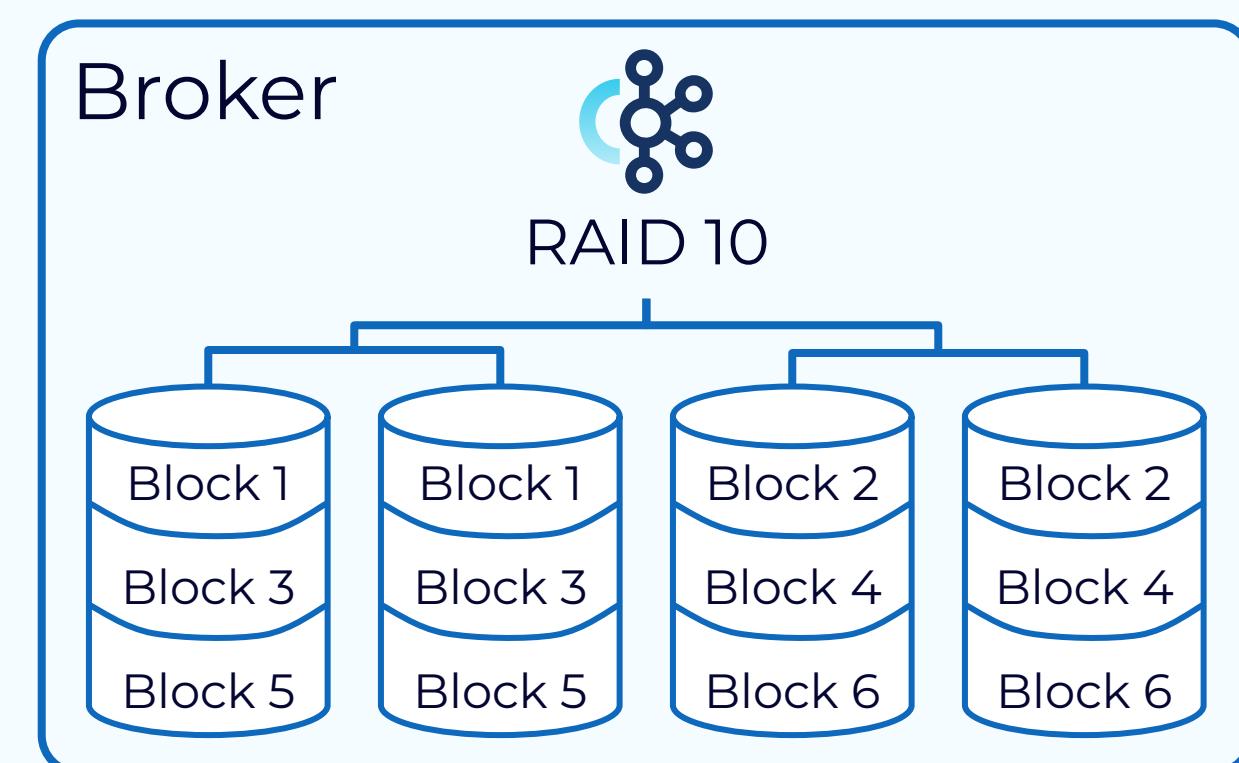
JBOD vs RAID



Kafka Cluster



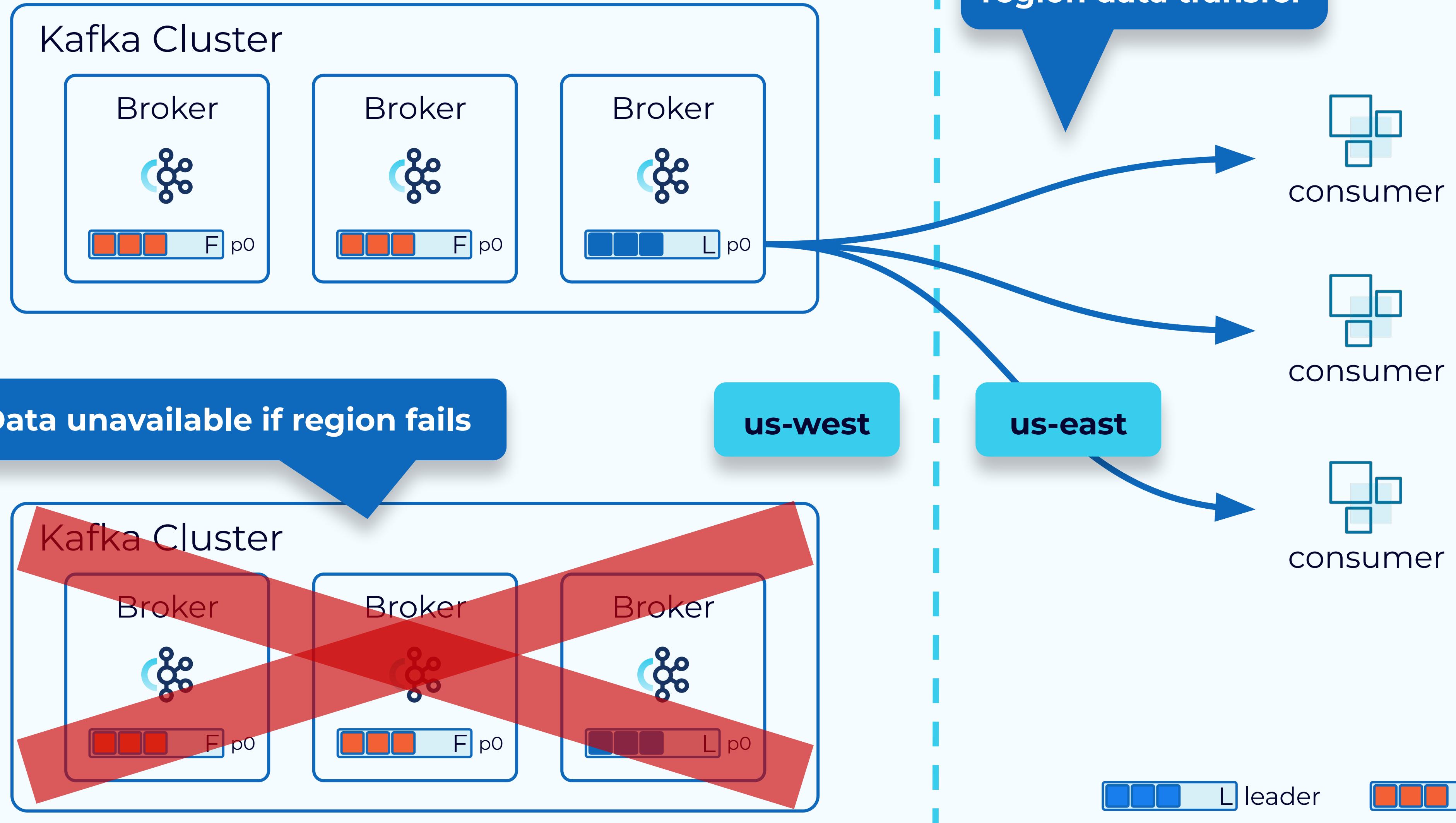
Kafka Cluster



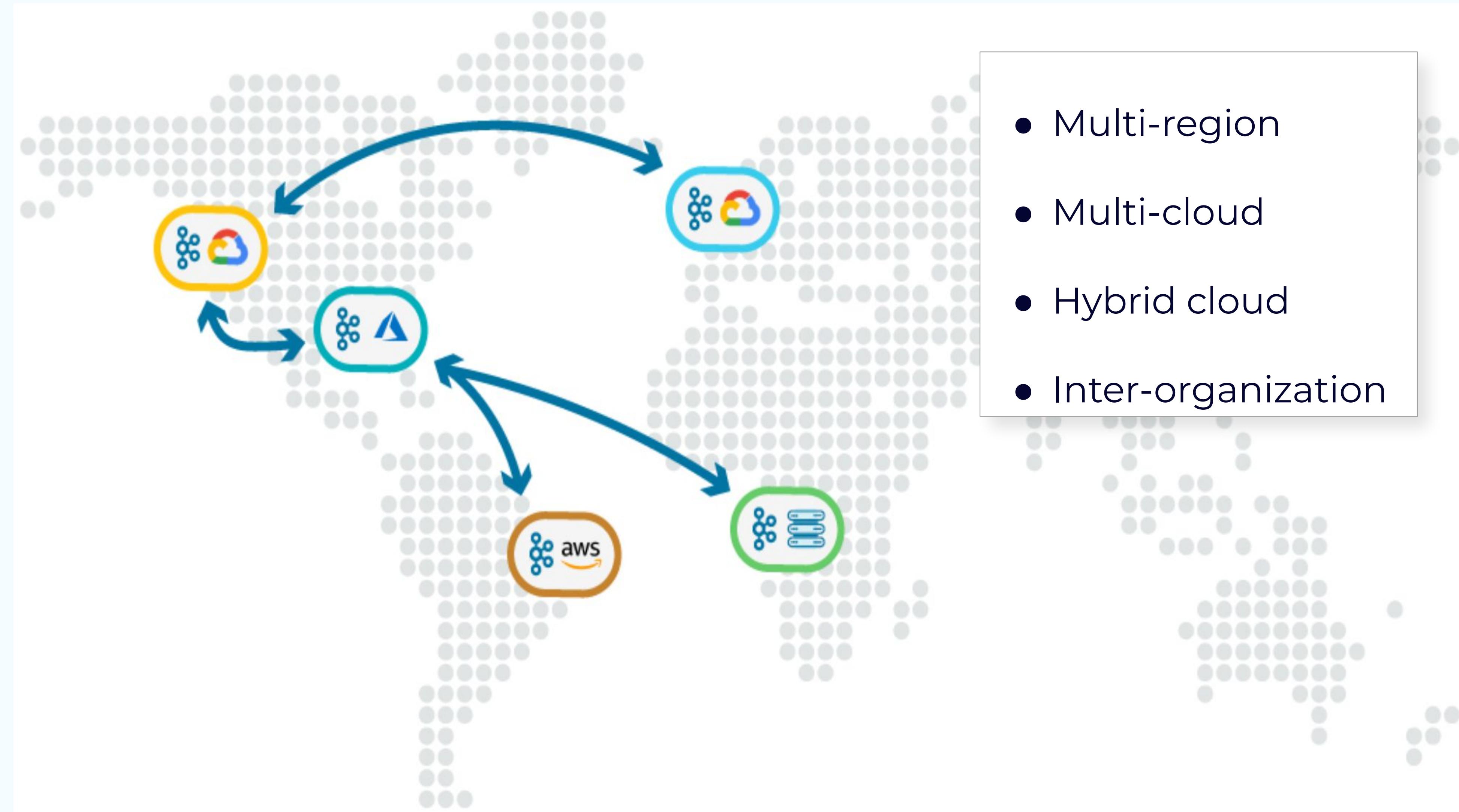


Geo-Replication

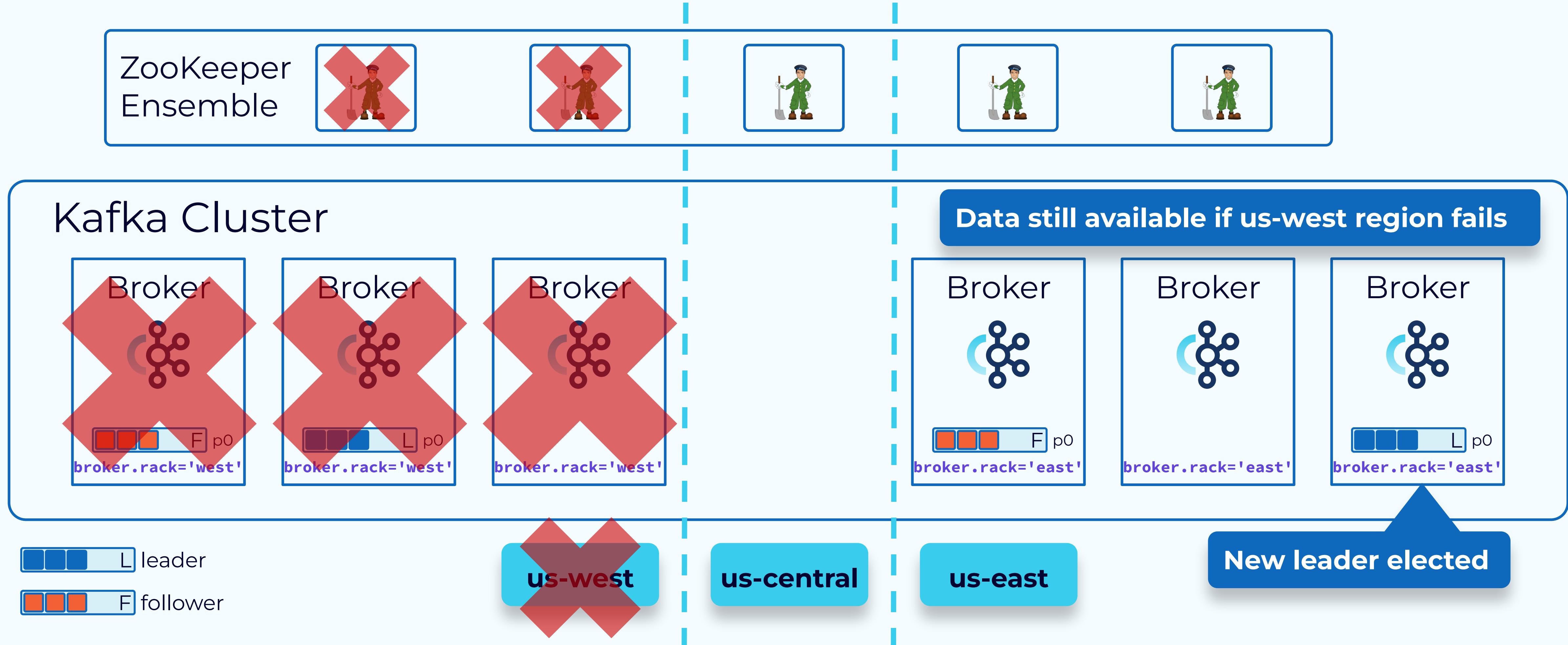
Single-Region Cluster Concerns



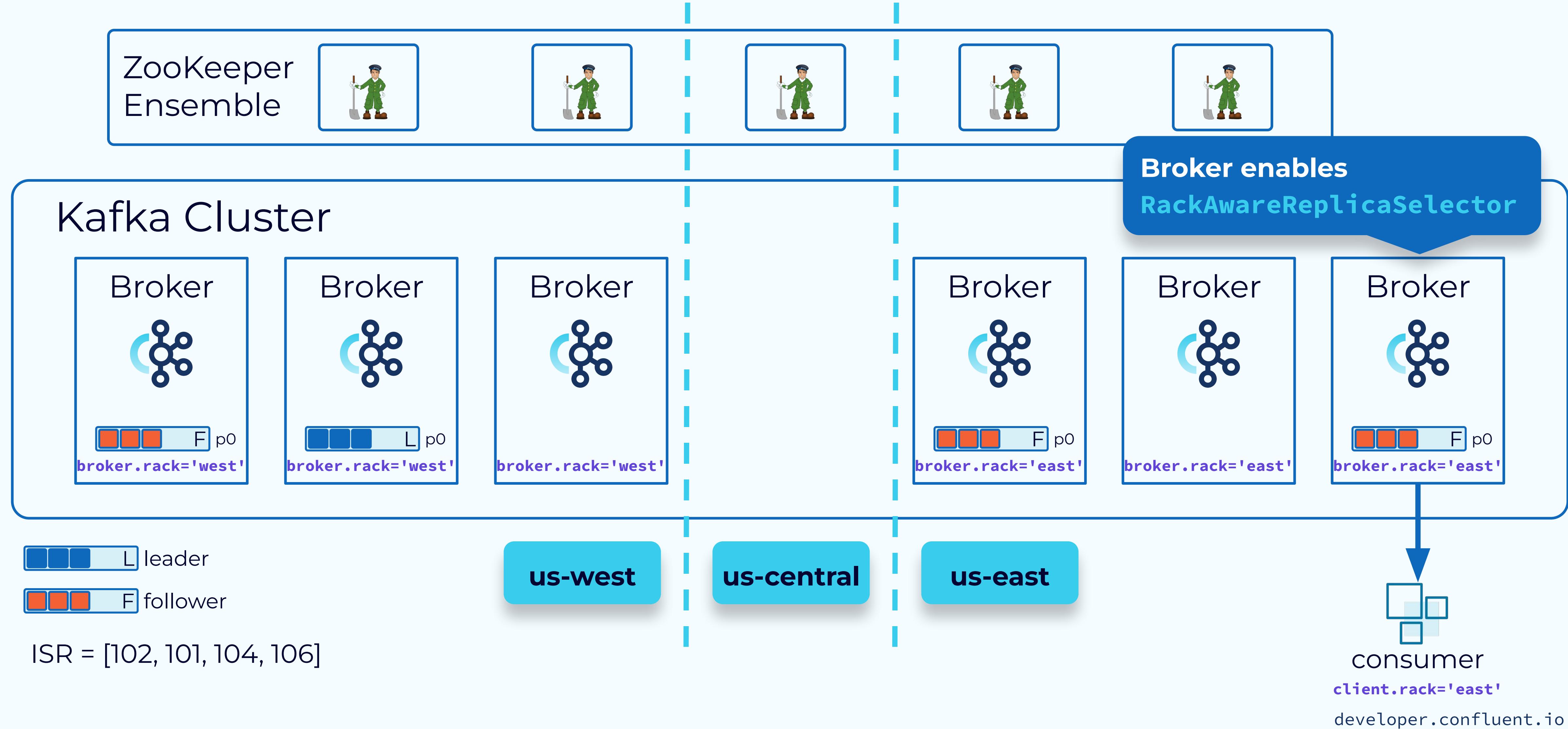
Geo-Replication



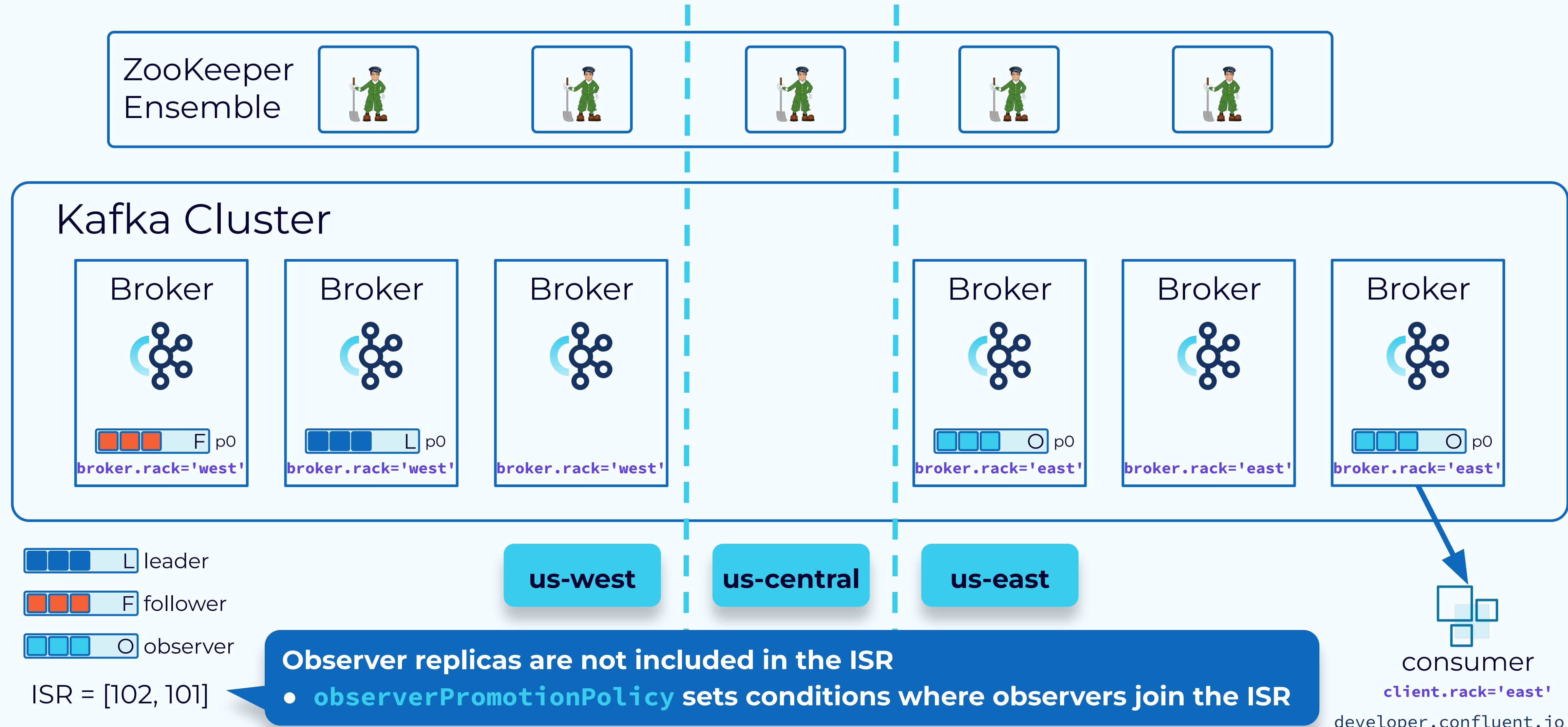
Confluent Multi-Region Cluster



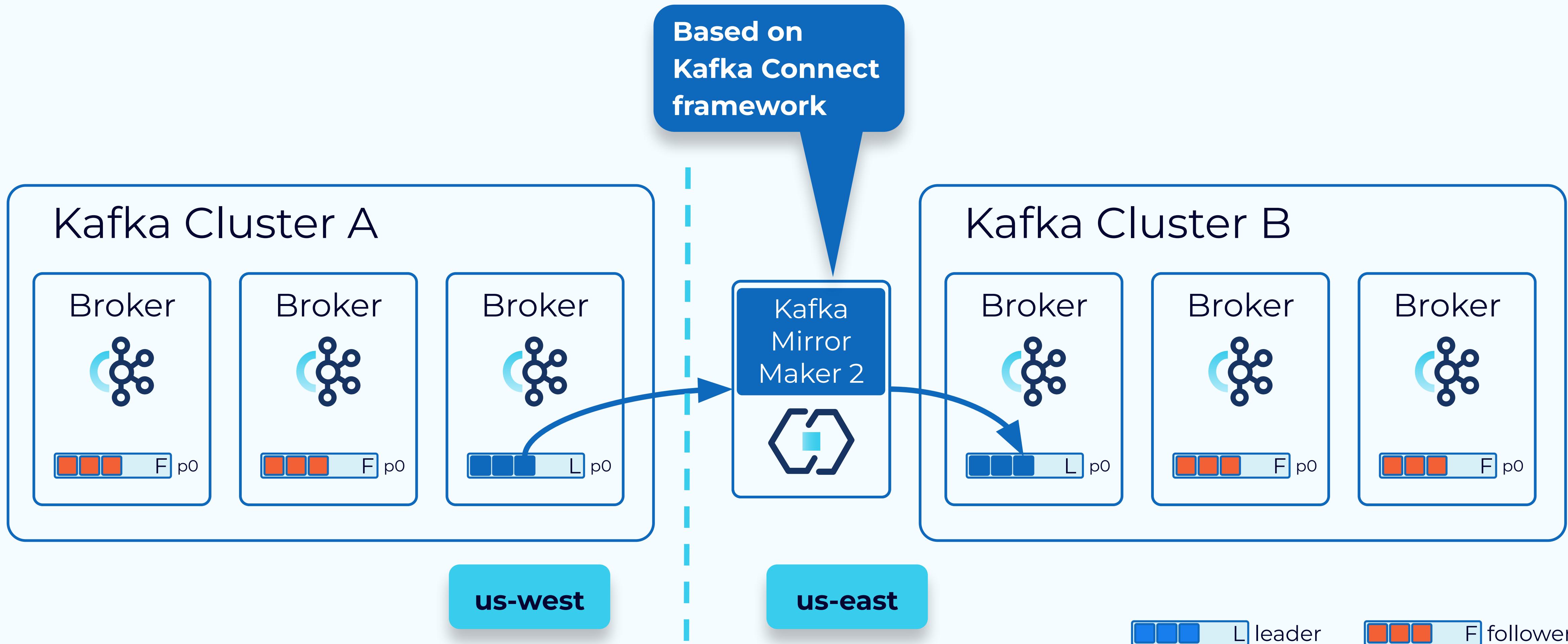
Better Locality with Fetch From Follower



Async Replication with Observers



Kafka MirrorMaker 2

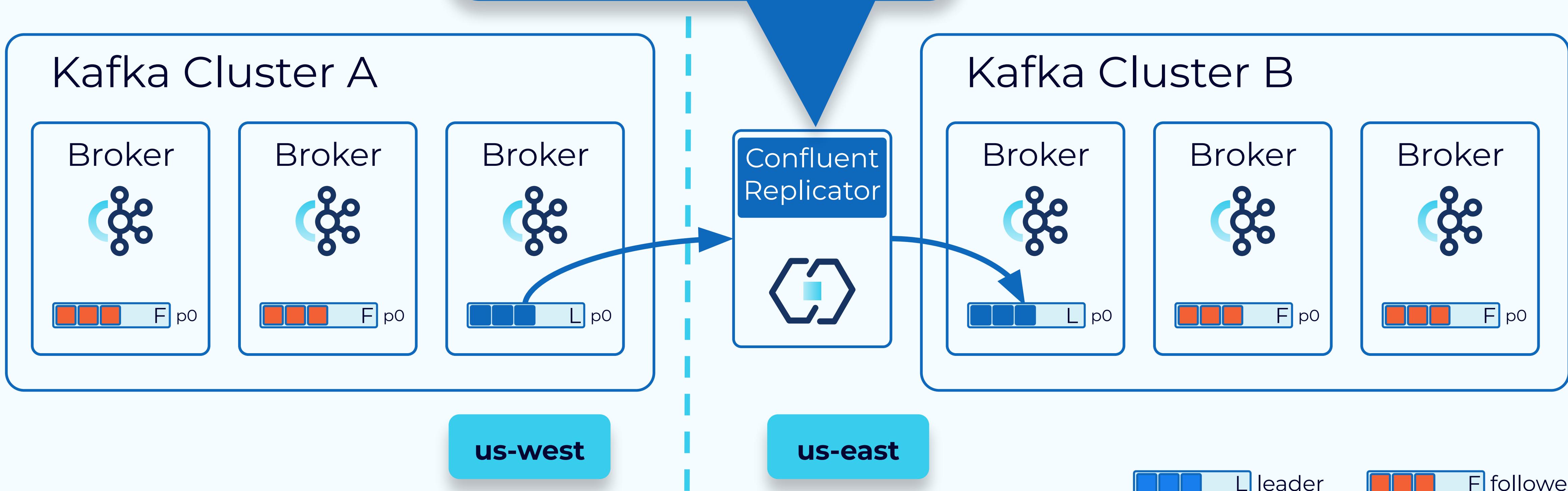


Confluent Replicator



Advantages over MirrorMaker 2

- 1) Metadata replication
- 2) Auto offset translation for Java consumers

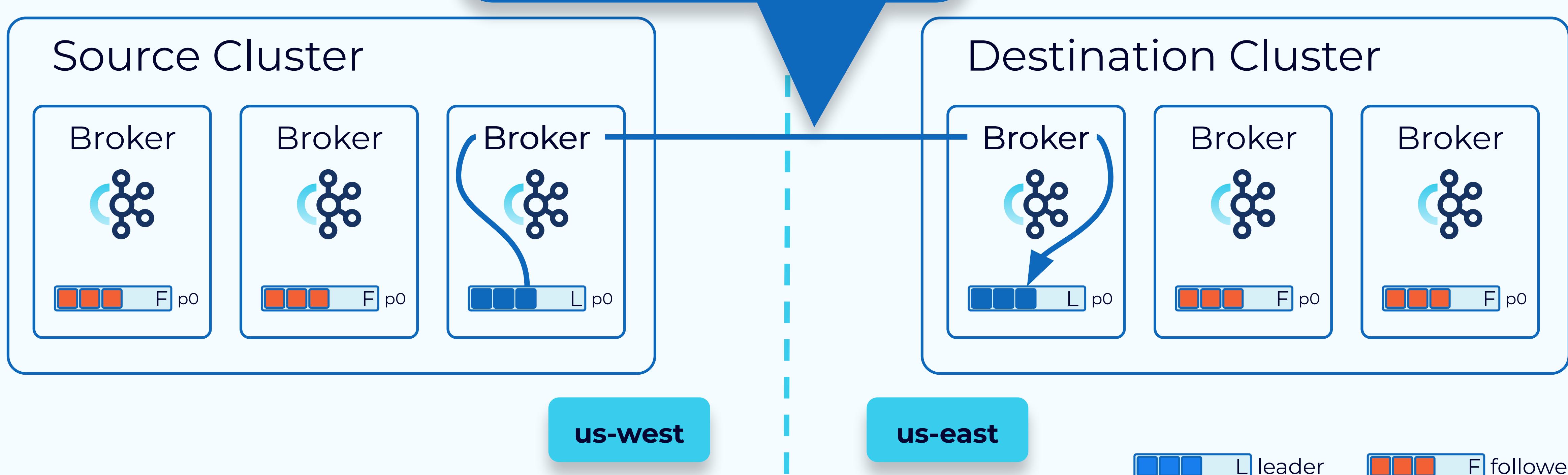


Confluent Cluster Linking

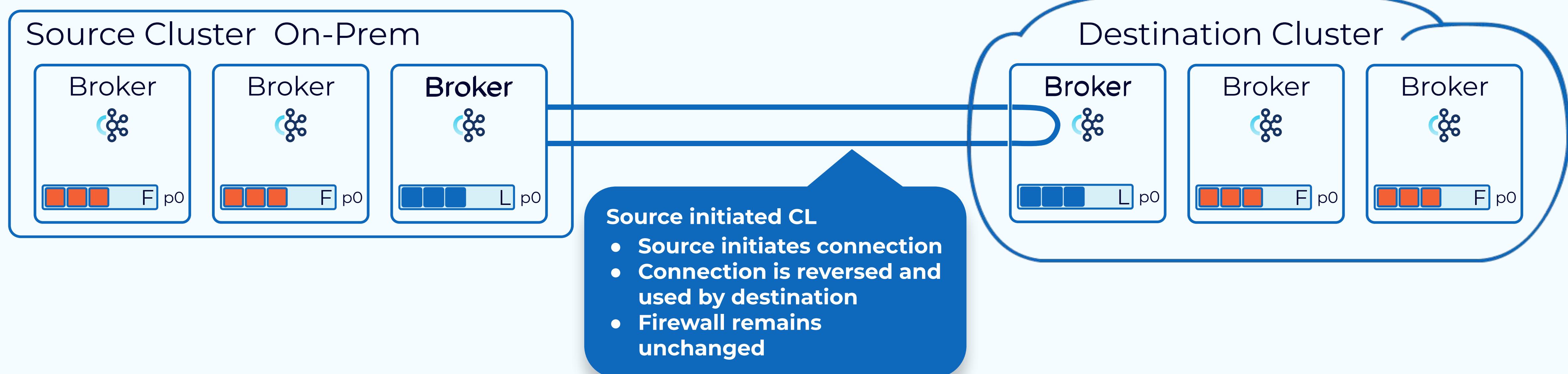


Benefits:

- 1) No separate process
- 2) Offset preserving
- 3) More efficient
(no recompression overhead)



Destination vs Source Initiated Cluster Linking



Geo-Replication Recap



Feature	Sync vs. Async	Offset Preserving	Built-In to the Cluster	Latency Tolerant
Confluent Multi-Region Cluster	Sync or Async	Yes	Yes	No (dark fiber recommended)
Kafka MirrorMaker 2	Async	No (manual offset translation)	No	Yes
Confluent Replicator	Async	No (auto offset translation for Java consumers)	No	Yes
Confluent Cluster Linking	Async	Yes	Yes	Yes



Hands On: Cluster Linking



*Your Apache Kafka
journey begins here*

developer.confluent.io