

LazyLog: A Shared Log Abstraction for Low-Latency Applications

Xuhao Luo, Shreesha G. Bhat, Jiyu Hu, Ram Alagappan, Aishwarya Ganesan
University of Illinois Urbana-Champaign

ANKITA SINGH
2024MCS2891

ABHIYA M JOSE
2024JCS2046

VADDI JHANCY
2024JCS2040

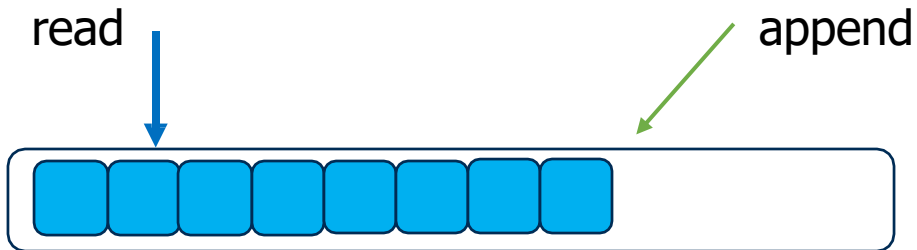
Shared Log: Abstraction and Interface

Abstraction:

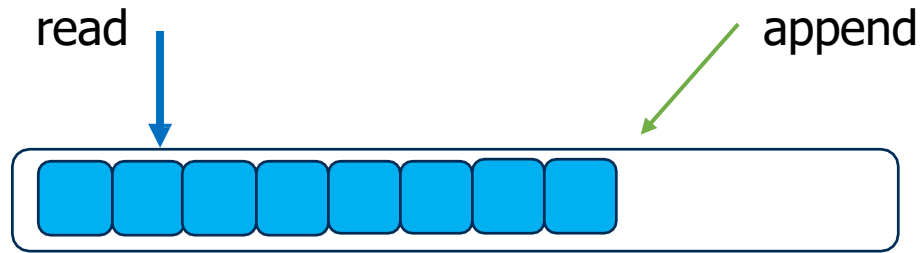
Fault-tolerant, durable,
linearizably ordered sequence of records

Interface:

append
read



Shared Log: Abstraction and Interface



```
// append to log; return log position  
uint64_t append(record r);  
  
// read 'len' records starting at 'from'  
list read(logpos_t from, uint64_t len);
```

The Problem with Current Shared Logs

- high append latencies
 - Append takes multiple RTTs
- Low append latency is critical to applications
- Eager ordering nature of shared logs:

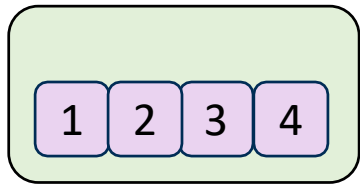
The Problem with Current Shared Logs

- **high ingestion latencies**
 - Append takes multiple RTTs
- Low ingestion latency is critical to applications
- **Eager ordering** nature of shared logs:
 - Order is established eagerly upon appends
 - Position of record is decided by the time append completes

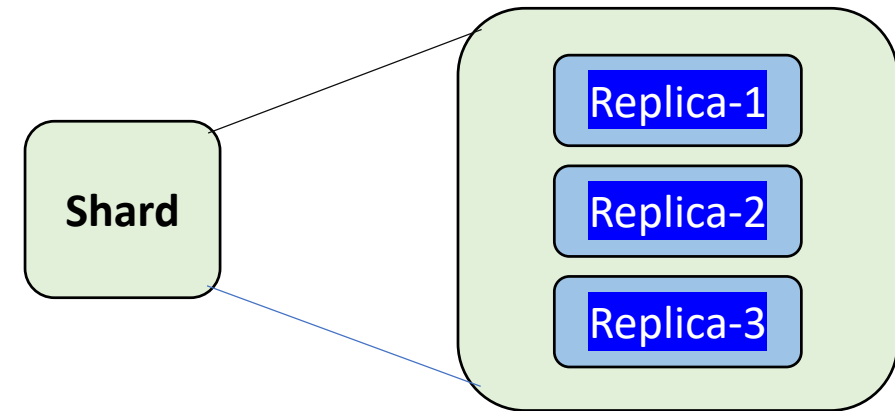
Total Ordering in Shared Logs

Linearizable order:

if `append(B)` starts after `append(A)` completes, then B appears after A in the shared log



Shard-1

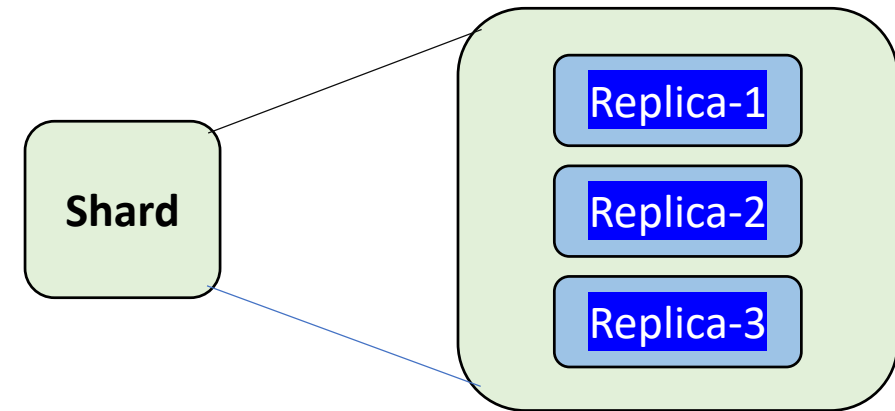
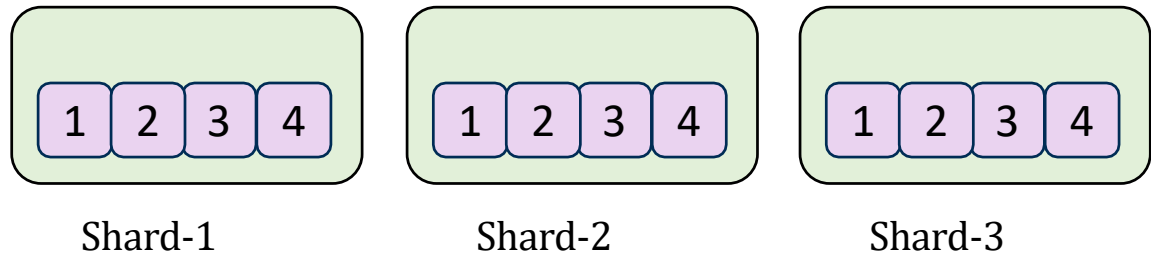


Fault Tolerance via Replication

Total Ordering in Shared Logs

Linearizable order:

if `append(B)` starts after `append(A)` completes, then B appears after A in the shared log



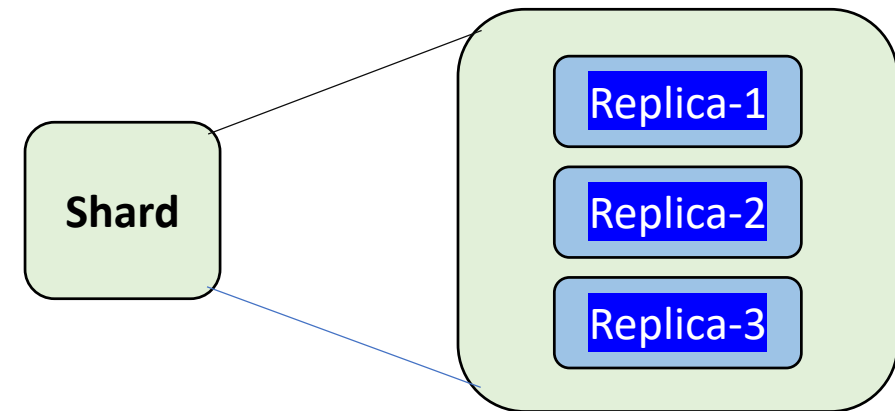
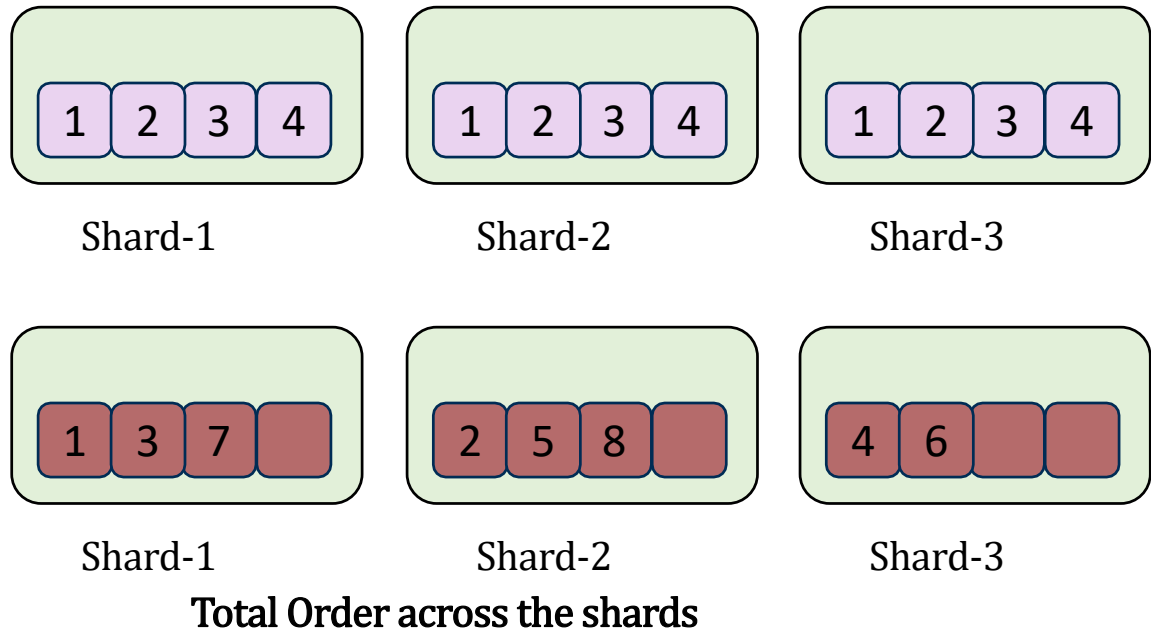
Fault Tolerance via Replication

Requires Total Ordering across the shards

Total Ordering in Shared Logs

Linearizable order:

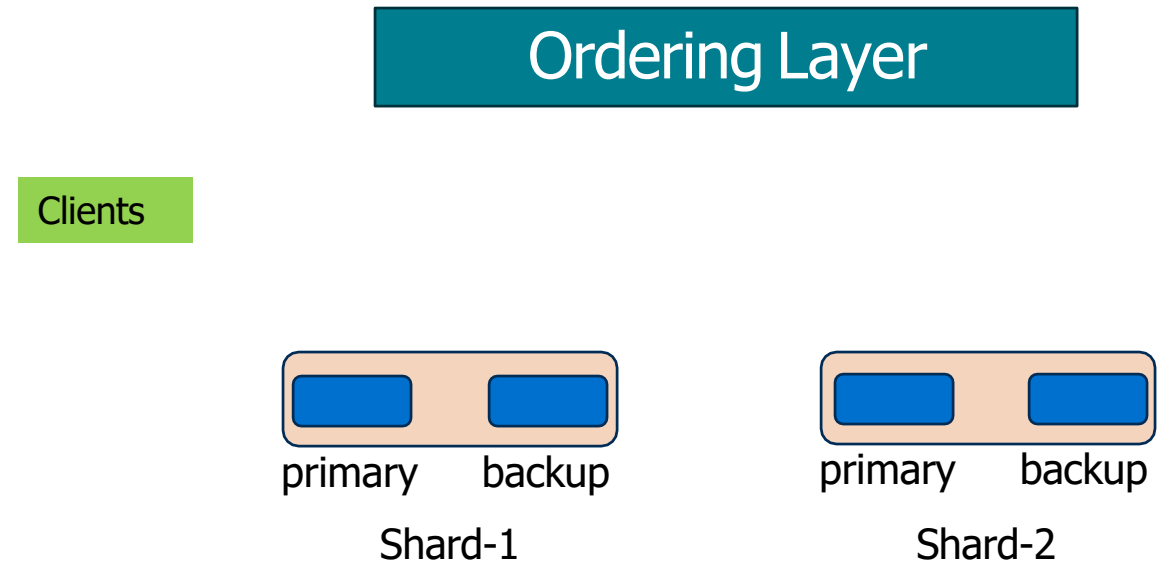
if `append(B)` starts after `append(A)` completes, then B appears after A in the shared log



Fault Tolerance via Replication

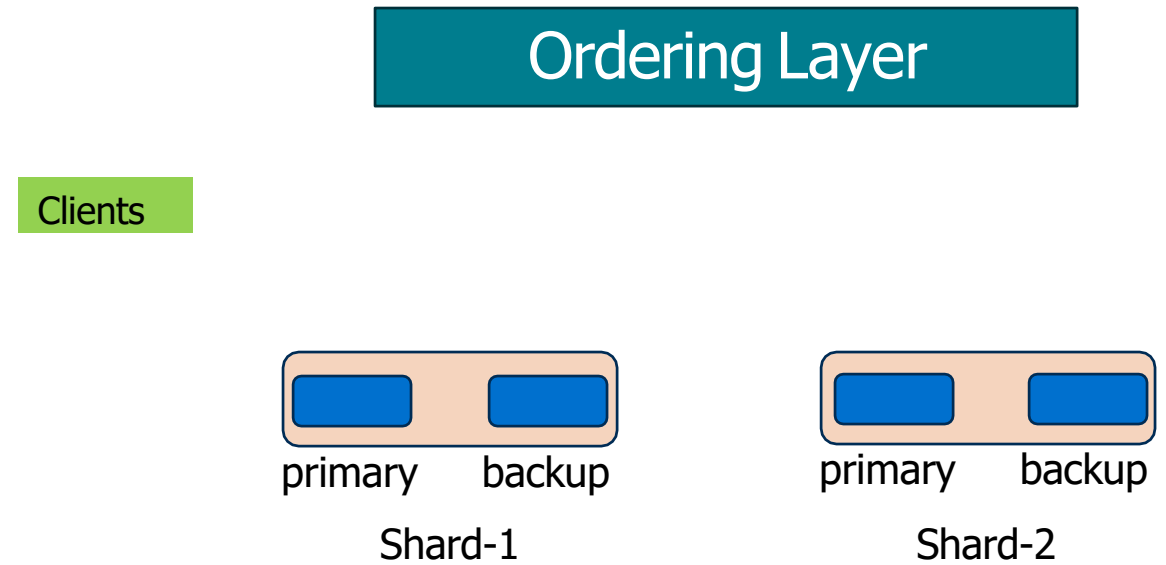
Scalability with total order: Scalog & Corfu

Eager Ordering → High Latency



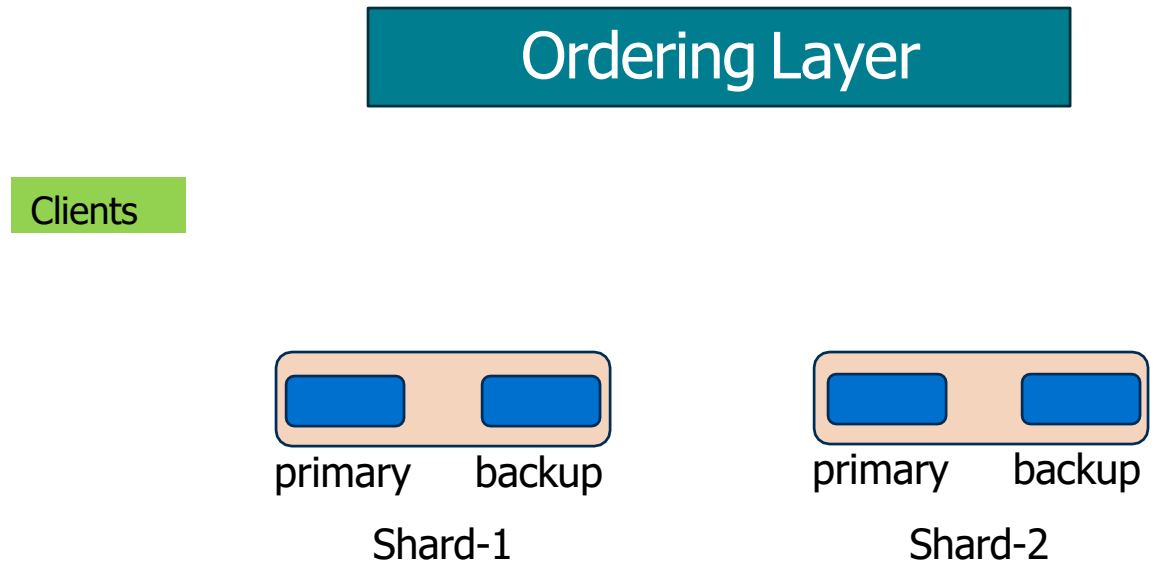
Eager Ordering → High Latency

- Incur high append latency



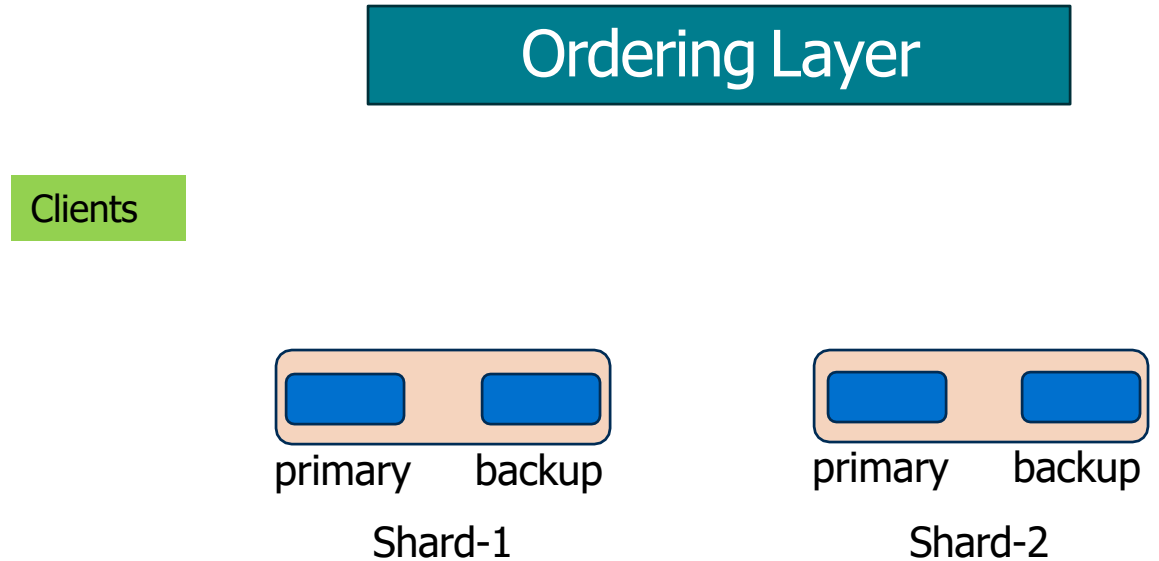
Eager Ordering → High Latency

- Incur high append latency
- Rooted in **eager ordering**



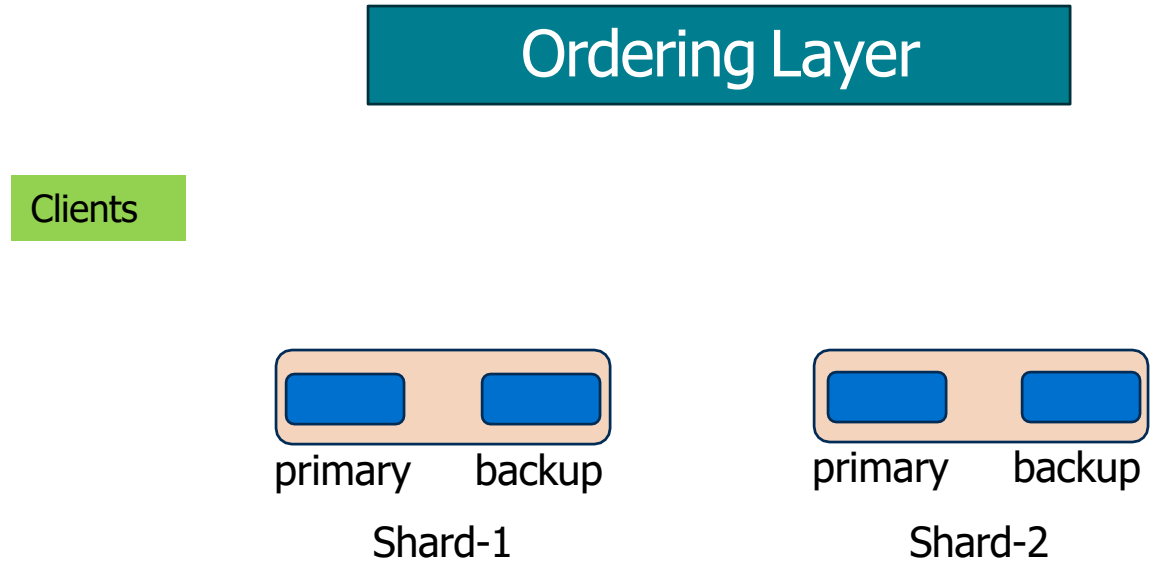
Eager Ordering → High Latency

- Incur high append latency
- Rooted in **eager ordering**
- Both durability and global ordering are completed before getting back to clients



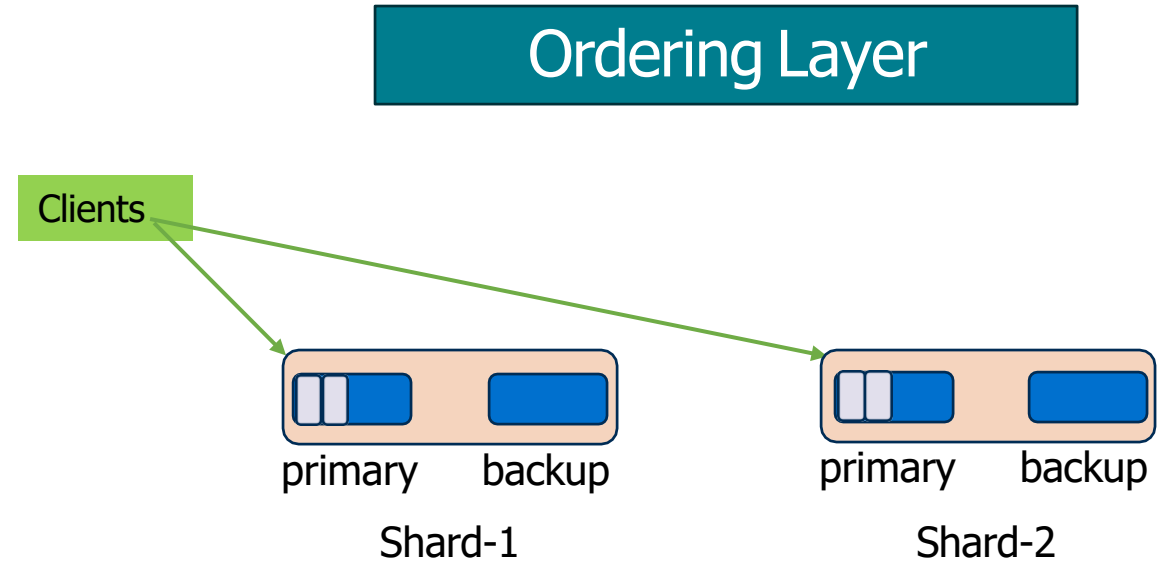
Eager Ordering → High Latency

- Scalog



Eager Ordering → High Latency

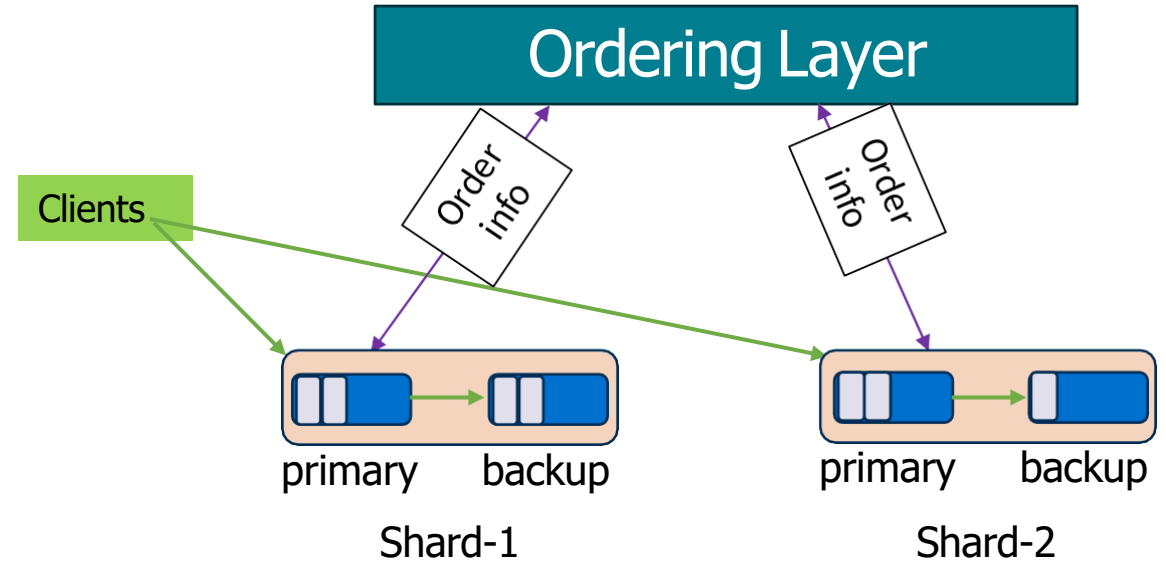
- Scalog
 - durability first



Eager Ordering → High Latency

■ Scalog

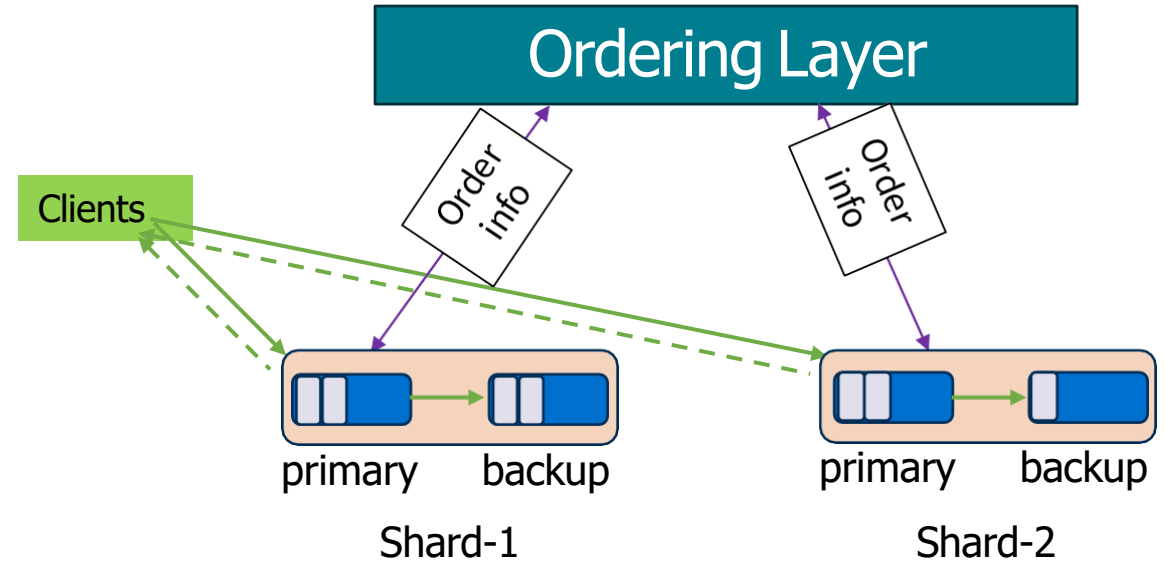
- durability first
- then global ordering



Eager Ordering → High Latency

■ Scalog

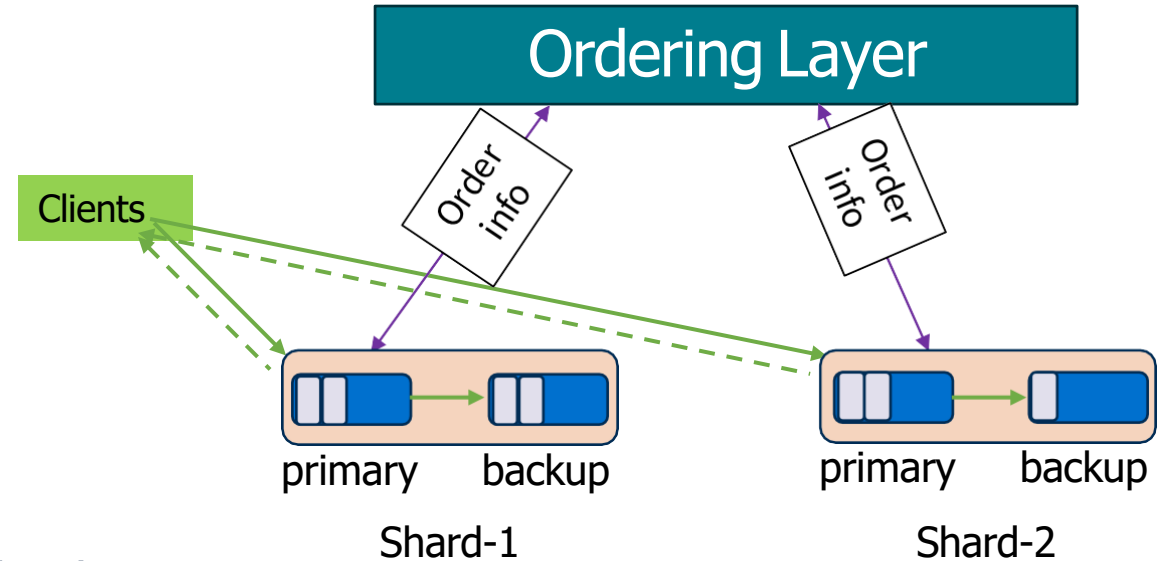
- durability first
- then global ordering



Eager Ordering → High Latency

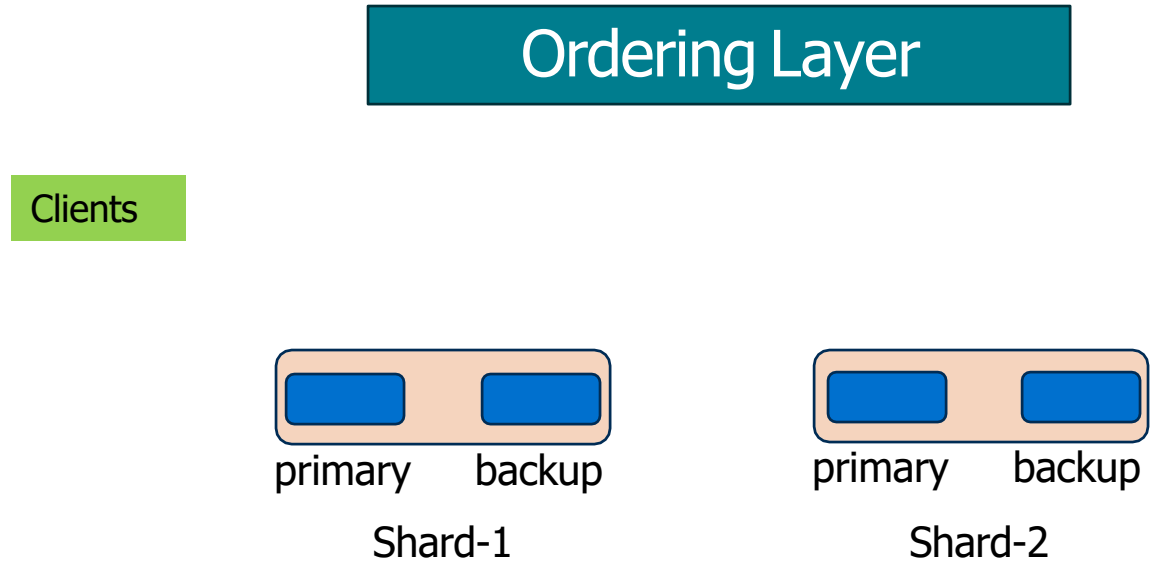
■ Scalog

- durability first
- then global ordering
- multiple RTT + batch interval
- Results in high ingestion latency for applications



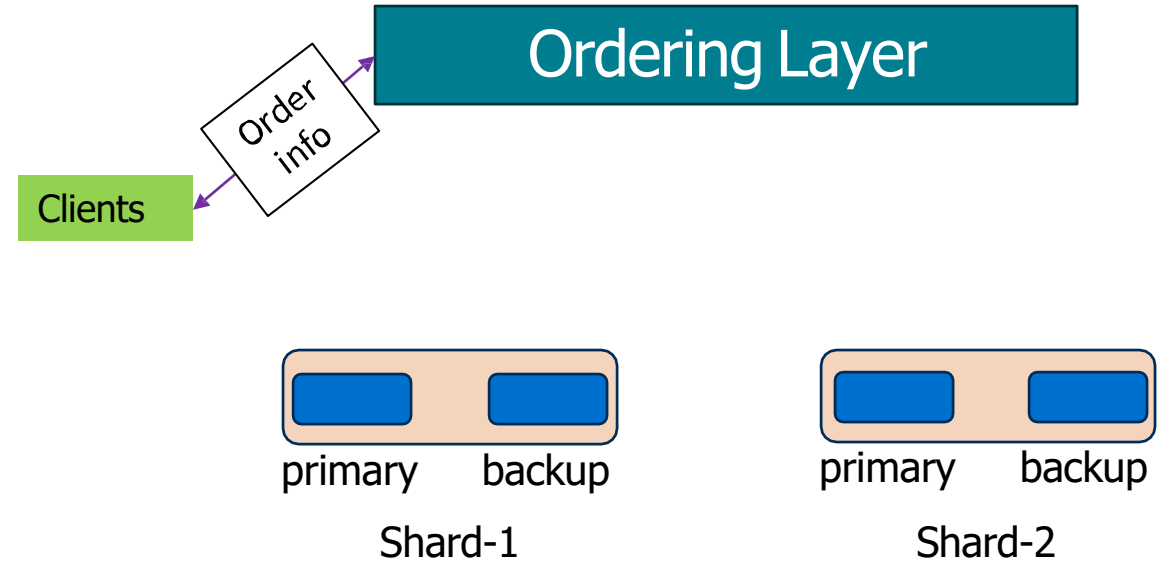
Eager Ordering → High Latency

- Corfu



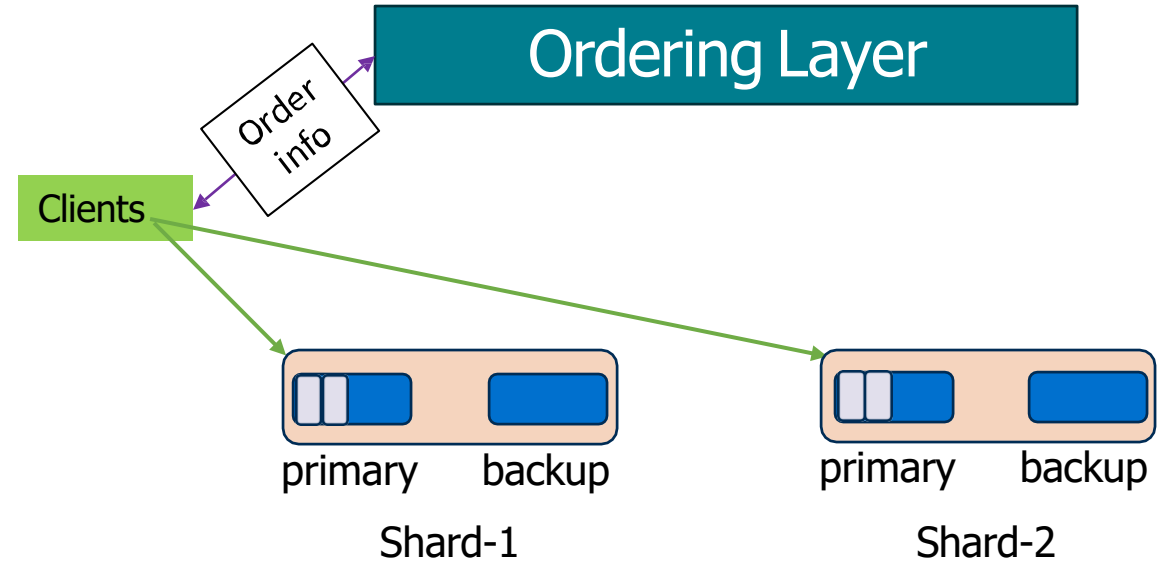
Eager Ordering → High Latency

- **Corfu**
 - global ordering first



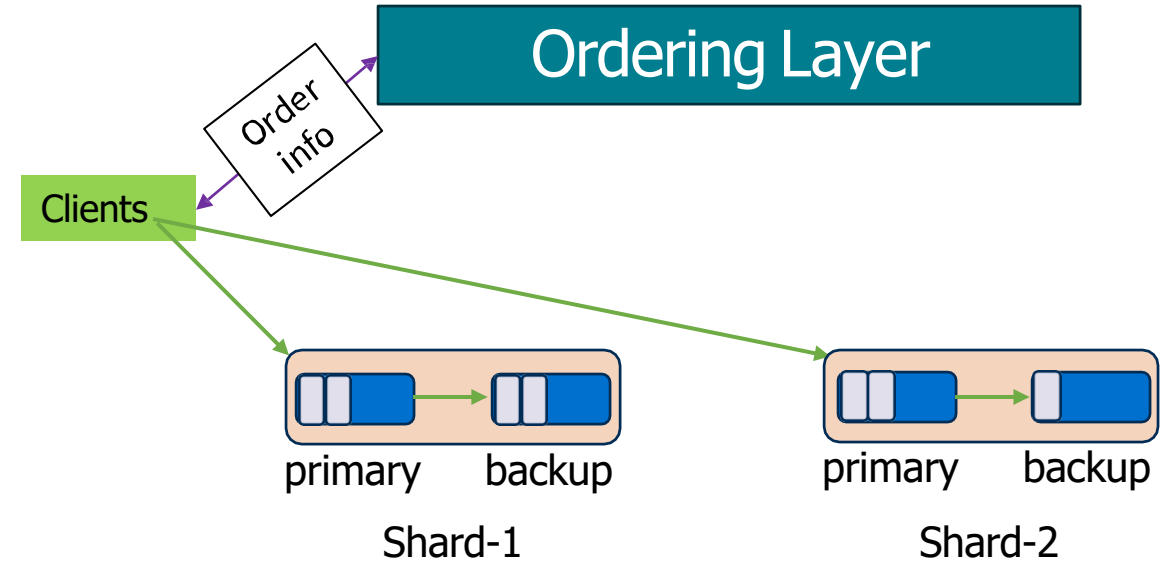
Eager Ordering → High Latency

- **Corfu**
 - global ordering first
 - then durability



Eager Ordering → High Latency

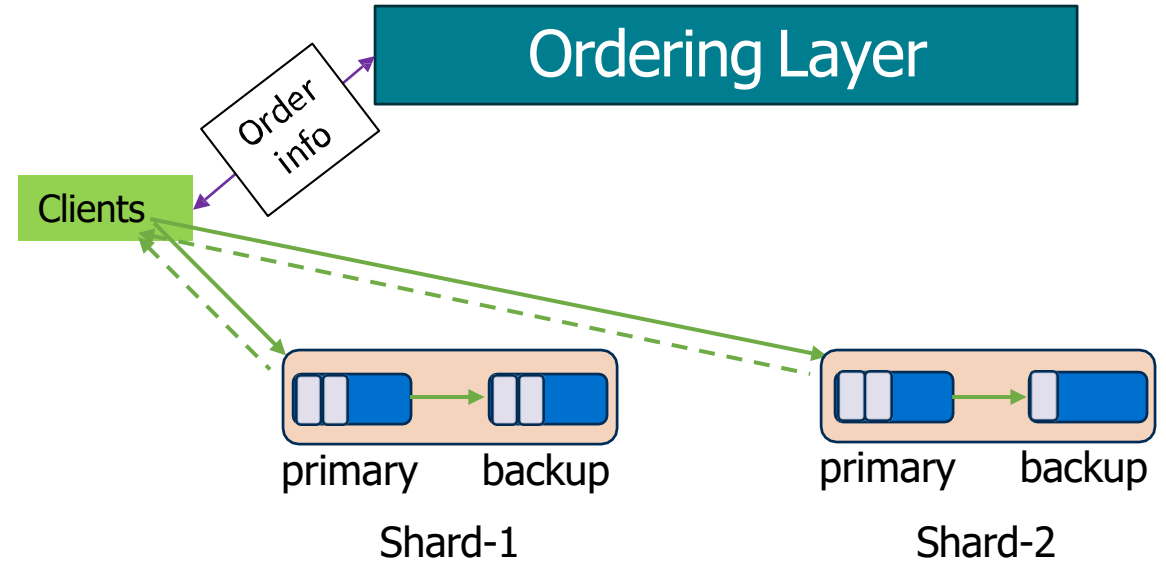
- Corfu
 - global ordering first
 - then durability



Eager Ordering → High Latency

■ Corfu

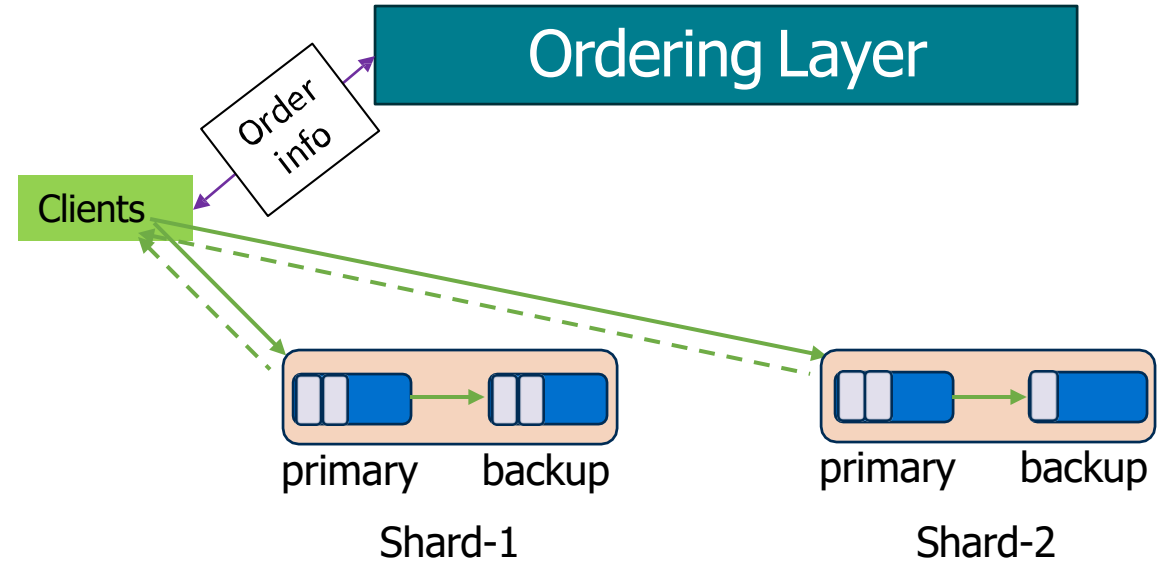
- global ordering first
- then durability



Eager Ordering → High Latency

■ Corfu

- global ordering first
- then durability
- multiple RTT
- Results in high ingestion latency for applications



Can a shared log **avoid eager ordering**, yet also **preserve the ordering guarantees** of conventional shared logs?

LazyLog: Idea and Abstraction

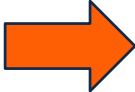
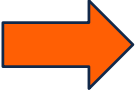
Idea



- No need for **eagerly bind** upon an append
 - But only make it durable

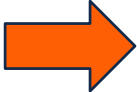


LazyLog: Idea and Abstraction

Idea

-  ○ No need for **eagerly bind** upon an append
 - But only make it durable
-  ○ Bind the records lazily, it must **enforce** ordering before positions can be read

LazyLog: Idea and Abstraction

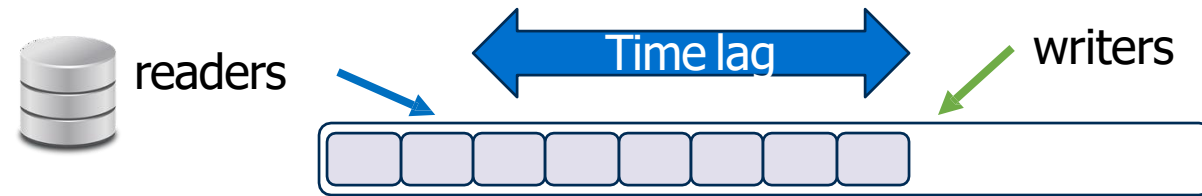
Idea

- 
 - No need for **eagerly bind** upon an append
 - But only make it durable
- 
 - Bind the records lazily, it must **enforce** ordering before positions can be read
- 
 - Shared log can do the ordering comfortably **in the background**

Insight

A shared log can **defer ordering** upon appends But
establish it before reads arrive

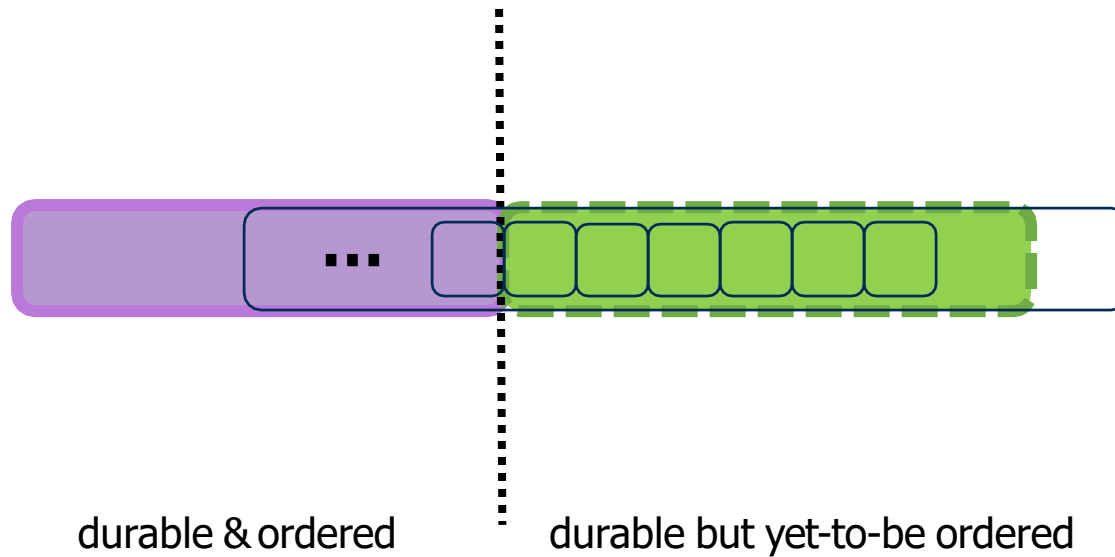
Holds for Many Apps



- Readers and writers are **time-decoupled**: readers typically **lag behind** writers

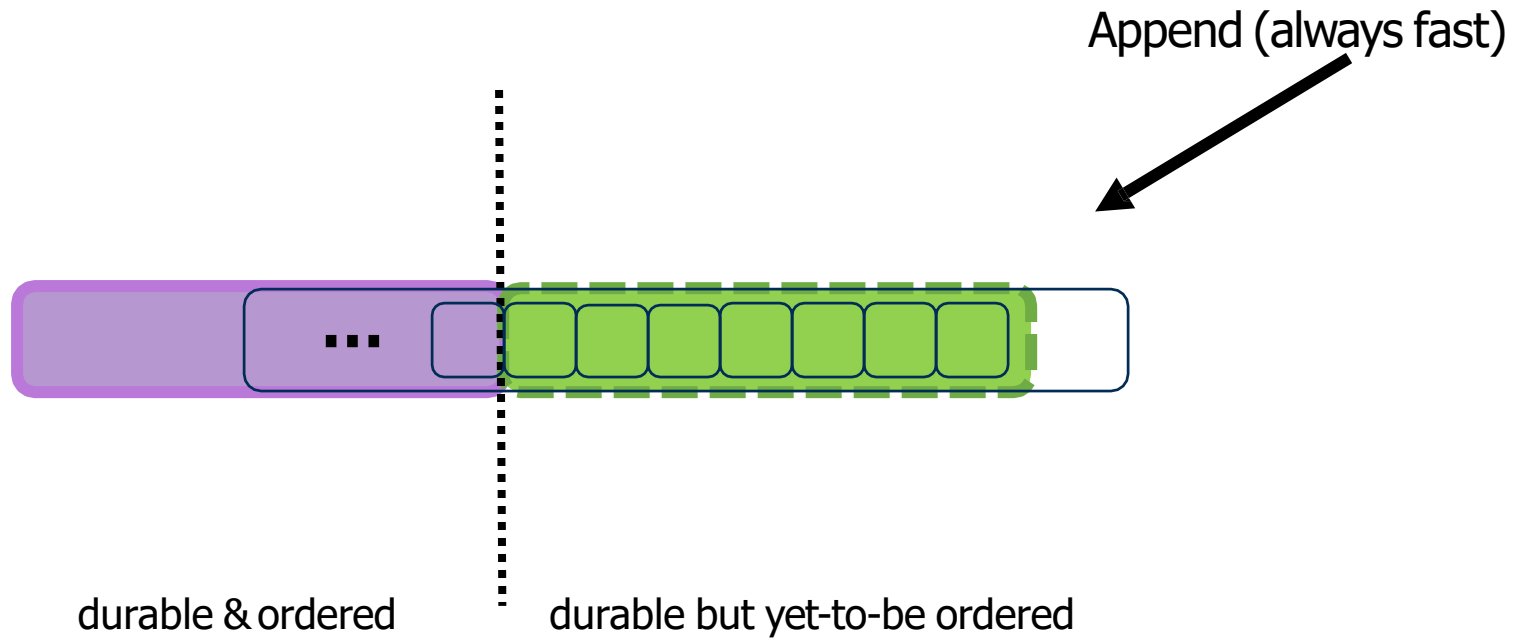
Performance Property

- Cannot be too lazy – keep ordering in the background



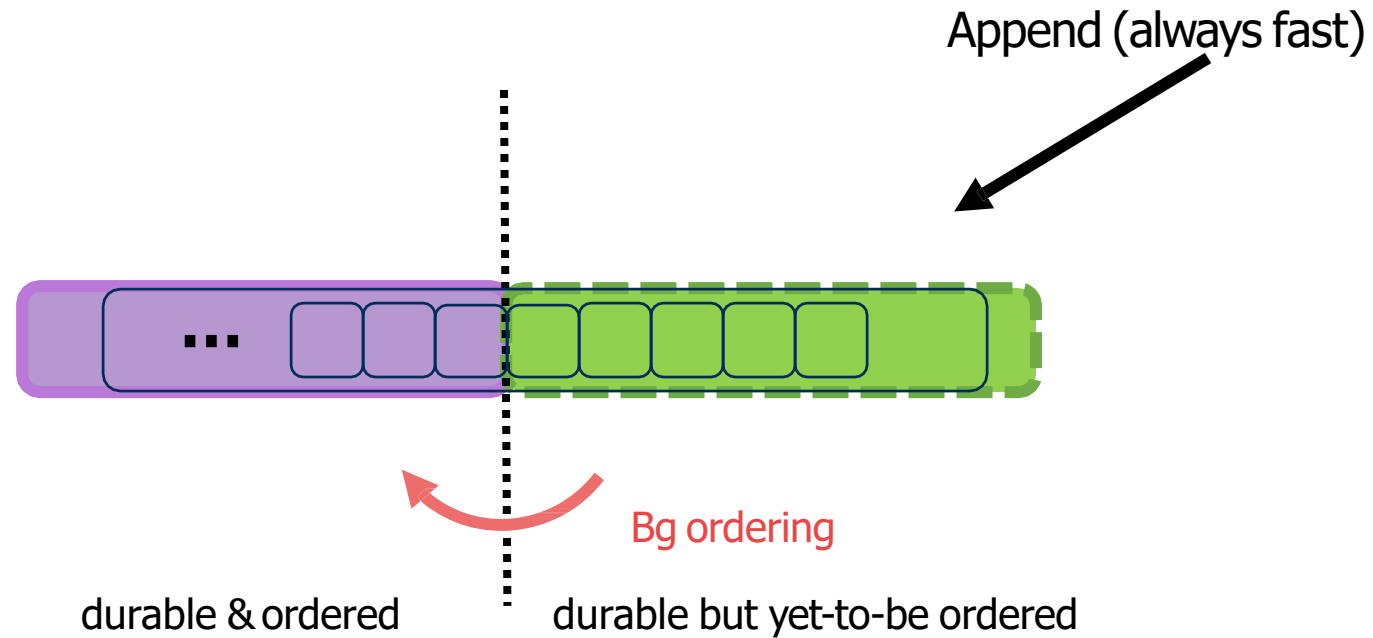
Performance Property

- Cannot be too lazy – keep ordering in the background



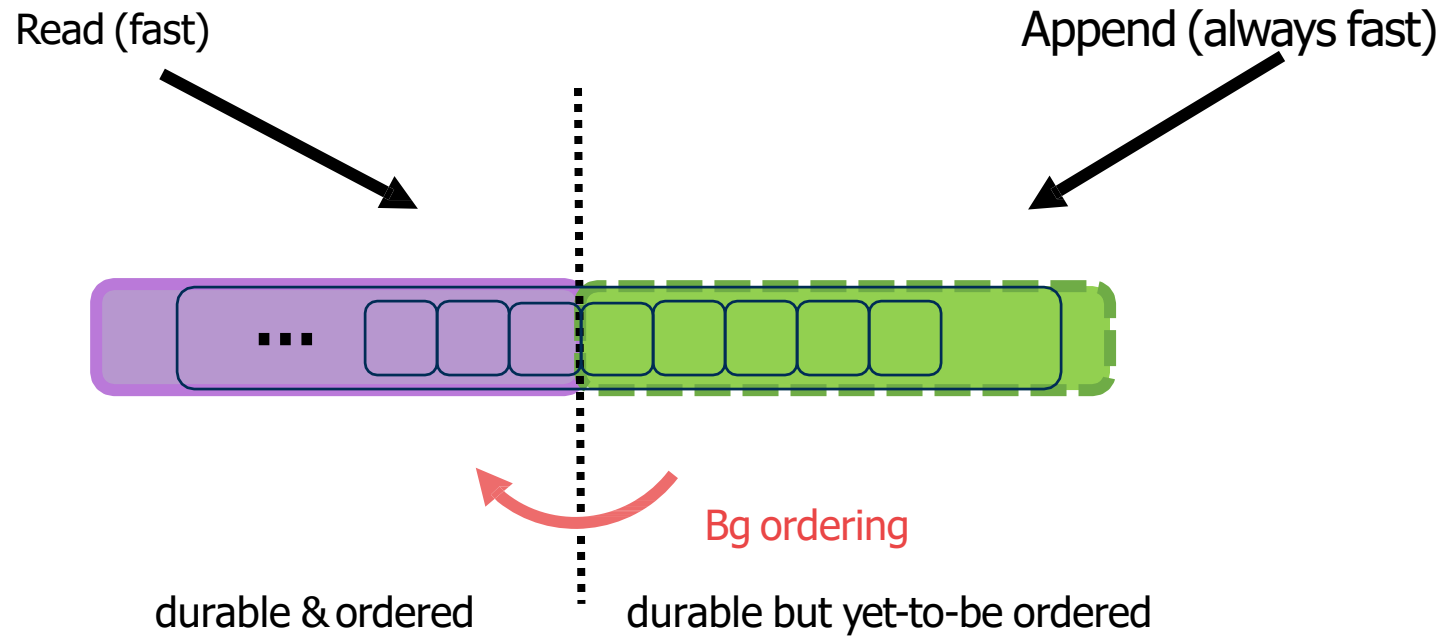
Performance Property

- Cannot be too lazy – keep ordering in the background



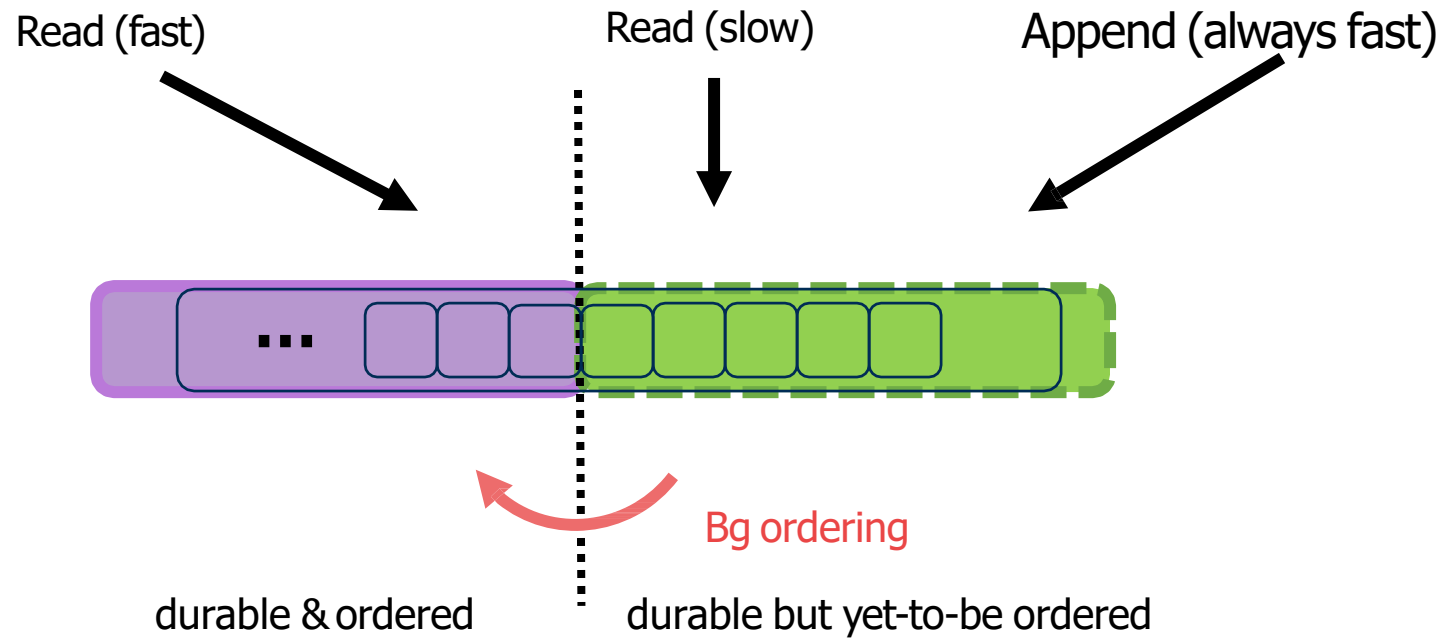
Performance Property

- Cannot be too lazy – keep ordering in the background



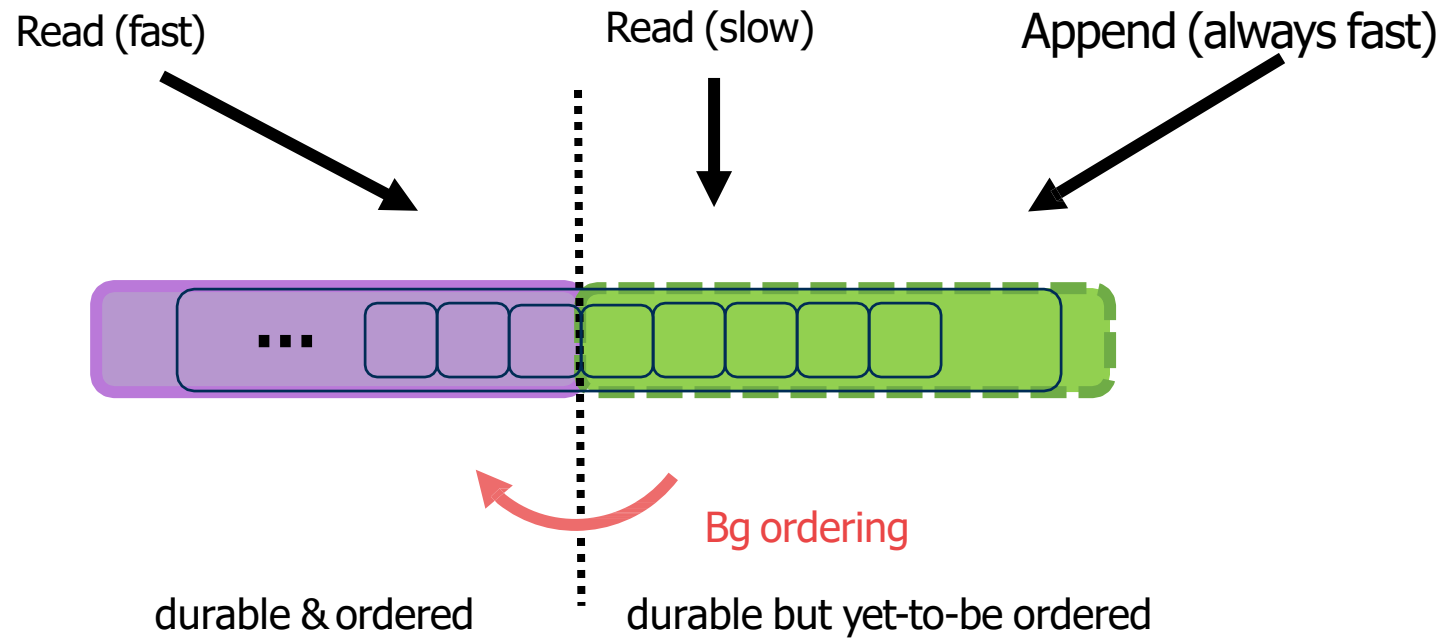
Performance Property

- Cannot be too lazy – keep ordering in the background



Performance Property

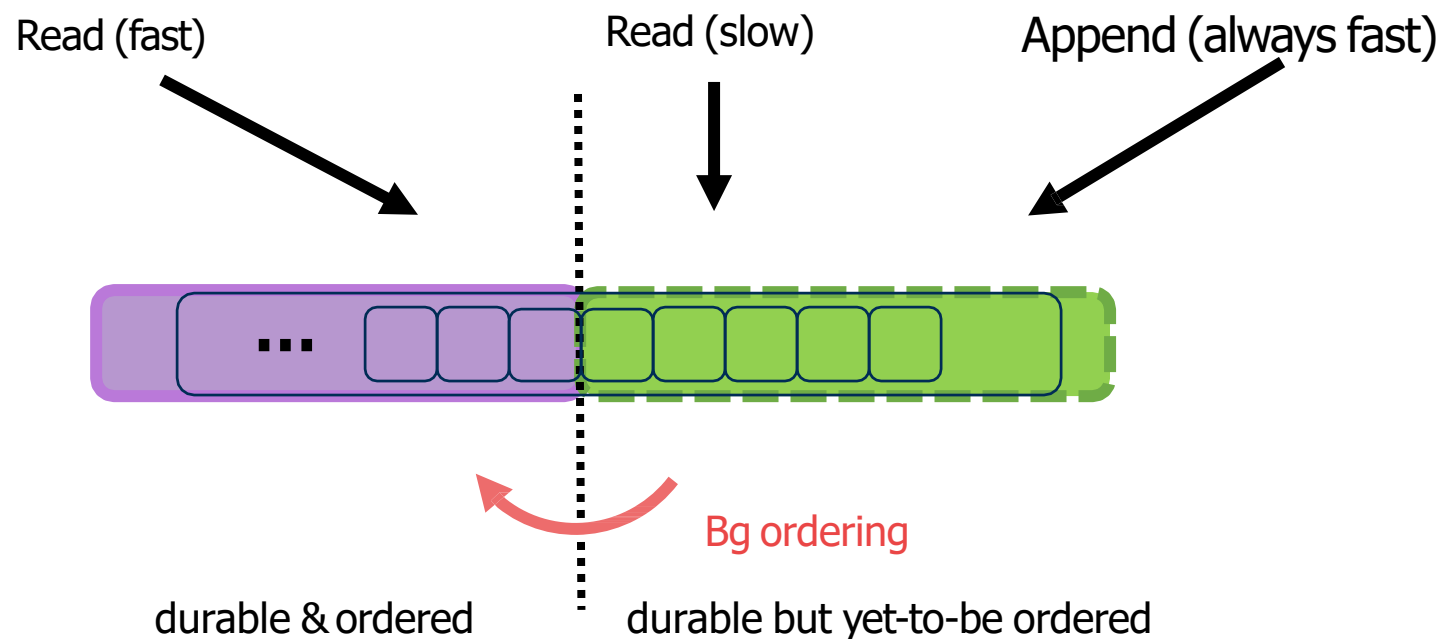
- Cannot be too lazy – keep ordering in the background



- For many apps – reads are always fast

Performance Property

- Cannot be too lazy – keep ordering in the background



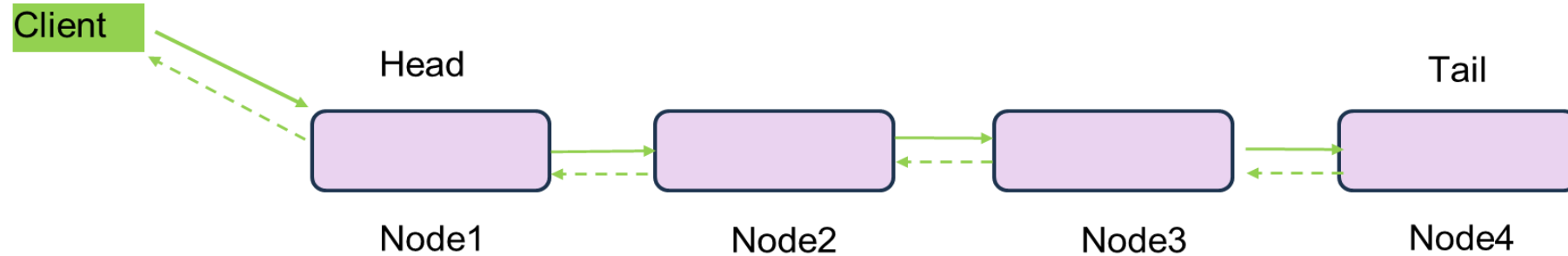
- For many apps – reads are always fast
- For immediate read, LazyLog preserves the performance of eager shared logs
 - never worse than an eager-ordering shared log!

Insight

Why **1-RTT Appends** are hard to achieve in Systems like Chain Paxos or Raft?

Chain Paxos

Requires multi-step coordination between replicas:



Client → Node1

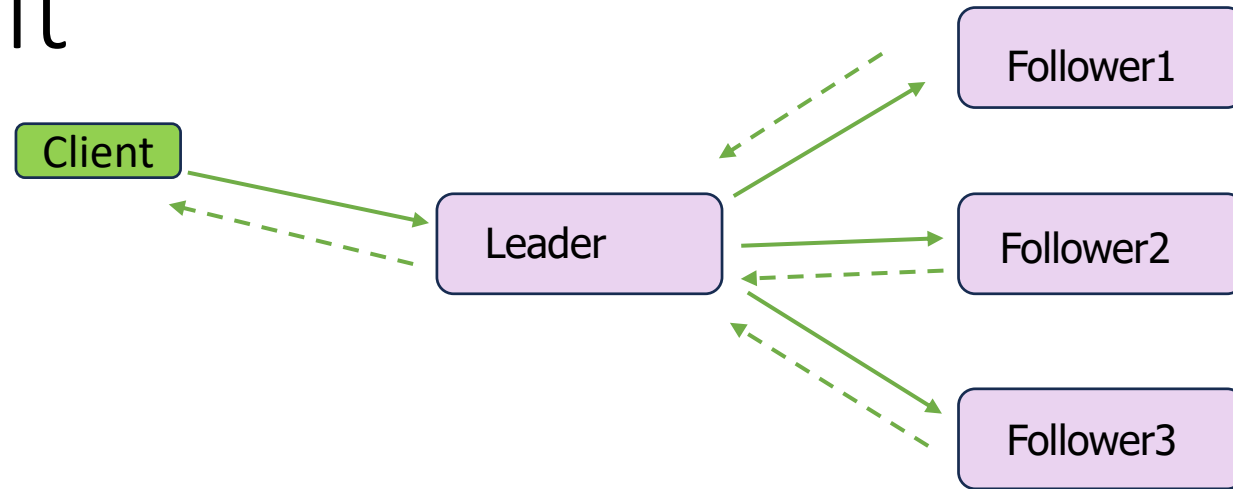
Node1 → Node2

Node2 → Node3

Node3 → Node4

Node4 → Client

Raft



Client → Leader (Send Request)

Leader → Followers (Replicate Entry)

Followers → Leader (Acks)

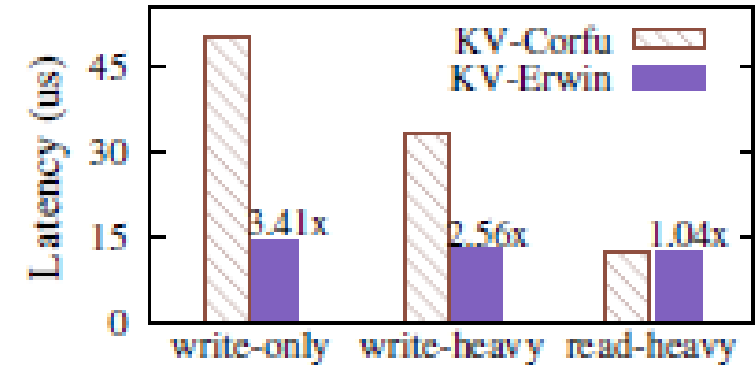
Leader → Client (Reply)

Insight

Application: Lazy-Log

Application: Key-Value Store

- Shared-log based key-value store
- **Readers – Writers are decoupled**
- Write-processing server- **handles append request**
- Read server – **consume the log**
- **One writer server, one reader server.**
- Shared log with **1 shard (3 replicas)**.
- Compared **Corfu** vs **Erwin** (LazyLog)



(a) KV Store

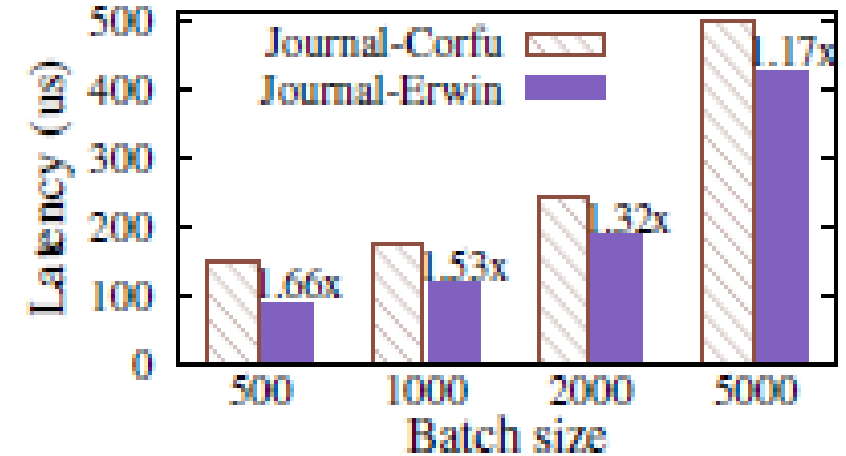
3 Tests

- | | | | |
|----------|------------------------|-------------|--------------------|
| • YCSB | 100% writes | write only | 3.4× lower latency |
| • YCSB-A | 50% writes & 50% reads | write heavy | 2.5× lower latency |
| • YCSB-B | 5% writes & 95% reads | read heavy | No change |

Application: Journaling for Stream Processing

- Task Worker-
 - reads the incoming records
 - use shared log to store their states
 - output the result
 - processes in batch

To **recover** during failure



(c) Journaling

Latency = Record Reading + Processing + checkpointing + emitting
Compared **Corfu** vs **Erwin**

Tests: stream-processing word count app with 5 workers For Different batches

- 5,000 records more computation **1.17x** faster
- 500 records logging dominates **1.66x** faster

System Design

Designed an implementation of the LazyLog interface: **Erwin**

Goal- Offers linearizable ordering across shards with 1-RTT appends

System Design

Designed an implementation of the LazyLog interface: **Erwin**

Offers linearizable ordering across shards with 1-RTT appends

Treats shards as black boxes and requires them to support:

- append an entry and
- read the entry at a specified index / position

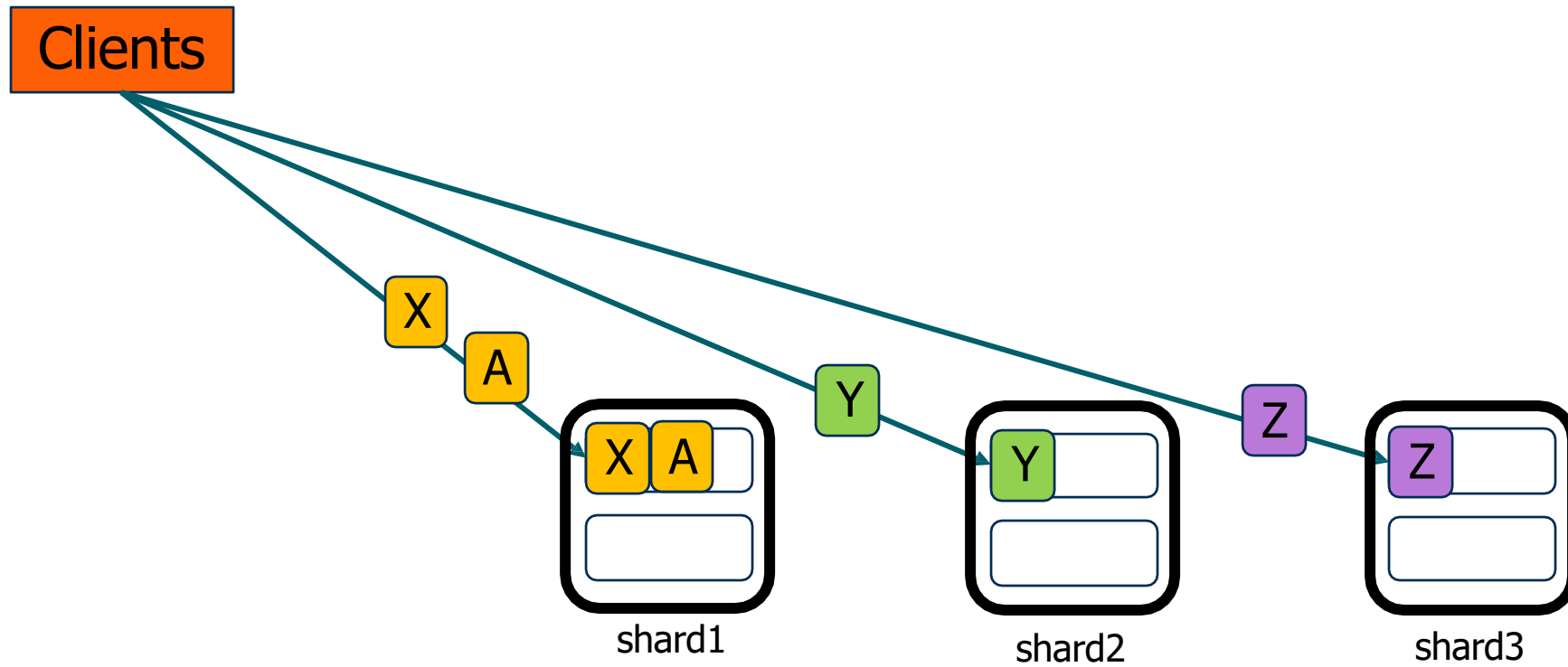
Why Erwin?

Motivation

Motivation



Record



Motivation



Record

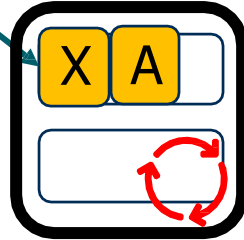
Clients

Problem:

Require coordination within
a shard

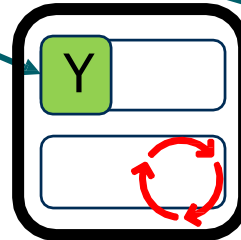
X

A



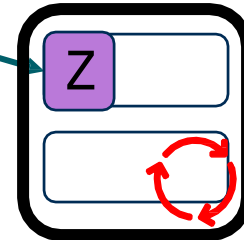
shard1

Y



shard2

Z

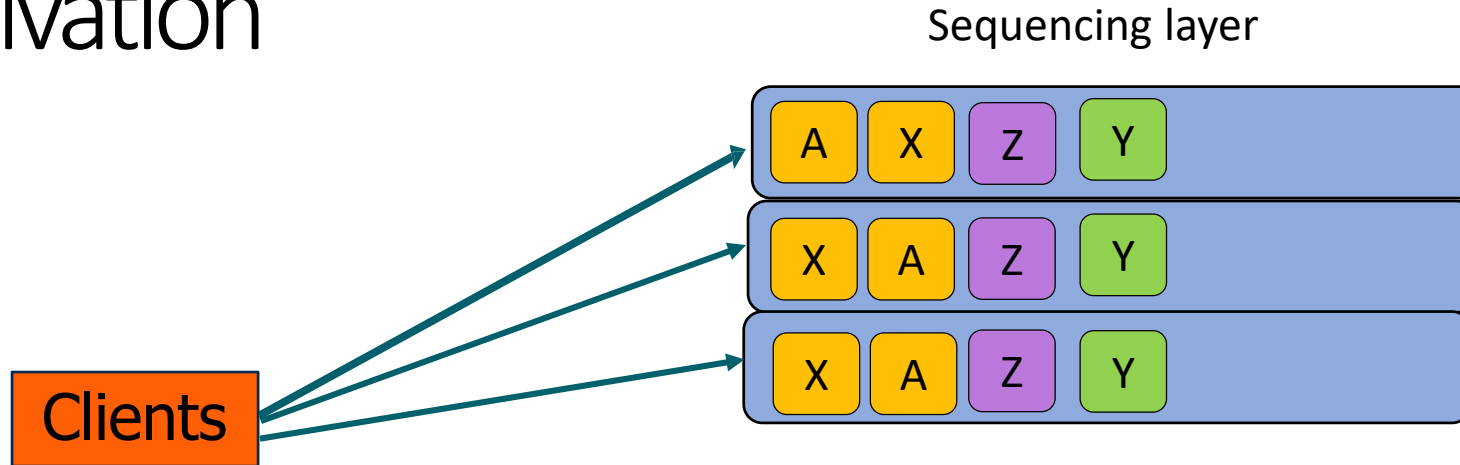


shard3

Motivation



Record



Problem:

Require coordination within a shard

Solution:

Use a buffer and send requests to this buffer



shard1

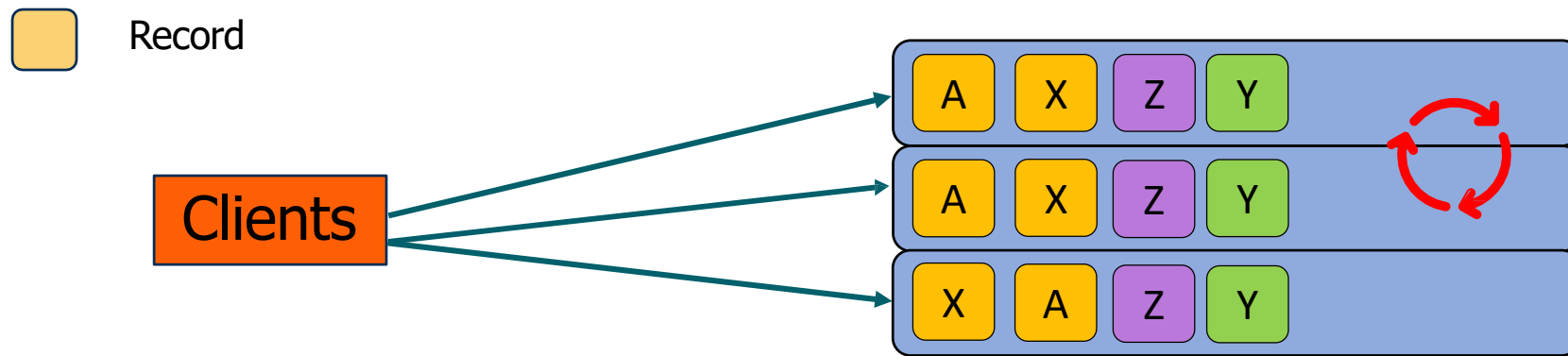


shard2



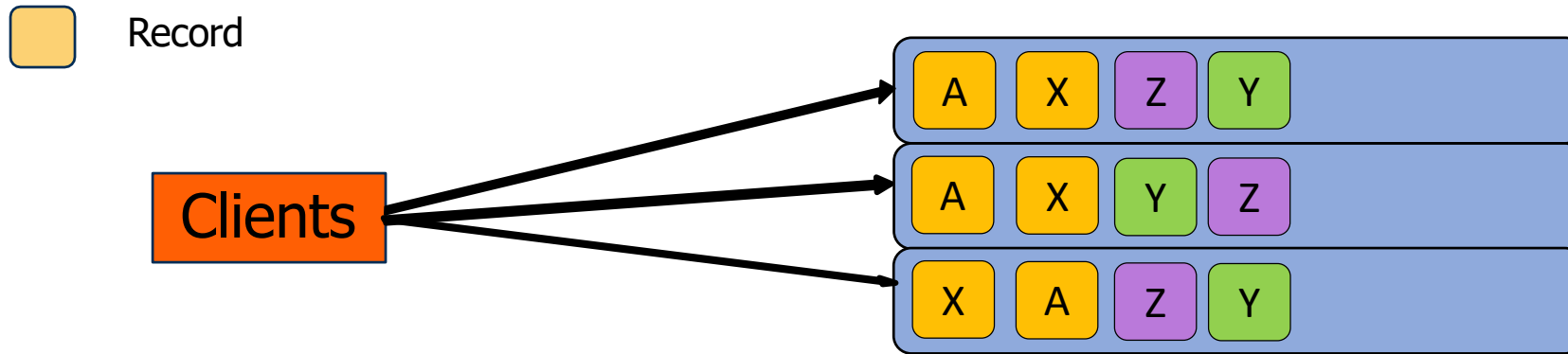
shard3

Erwin: 1-RTT Append



Problem:
Sequencing layer must run consensus to make ordering fault-tolerant → Incurs coordination within replicas

Erwin: 1-RTT Append



Problem:

Sequencing layer must run consensus to make ordering fault-tolerant → Incurs coordination within replicas

Solution:

Coordination-free sequencing

- Clients write to shard replicas in 1RTT; in same RTT, write metadata to **all** seq replicas
 - Appends complete in 1 RTT
- Erwin allows different orders across sequencing replicas
 - but without violating the linearizability

Erwin's Architecture

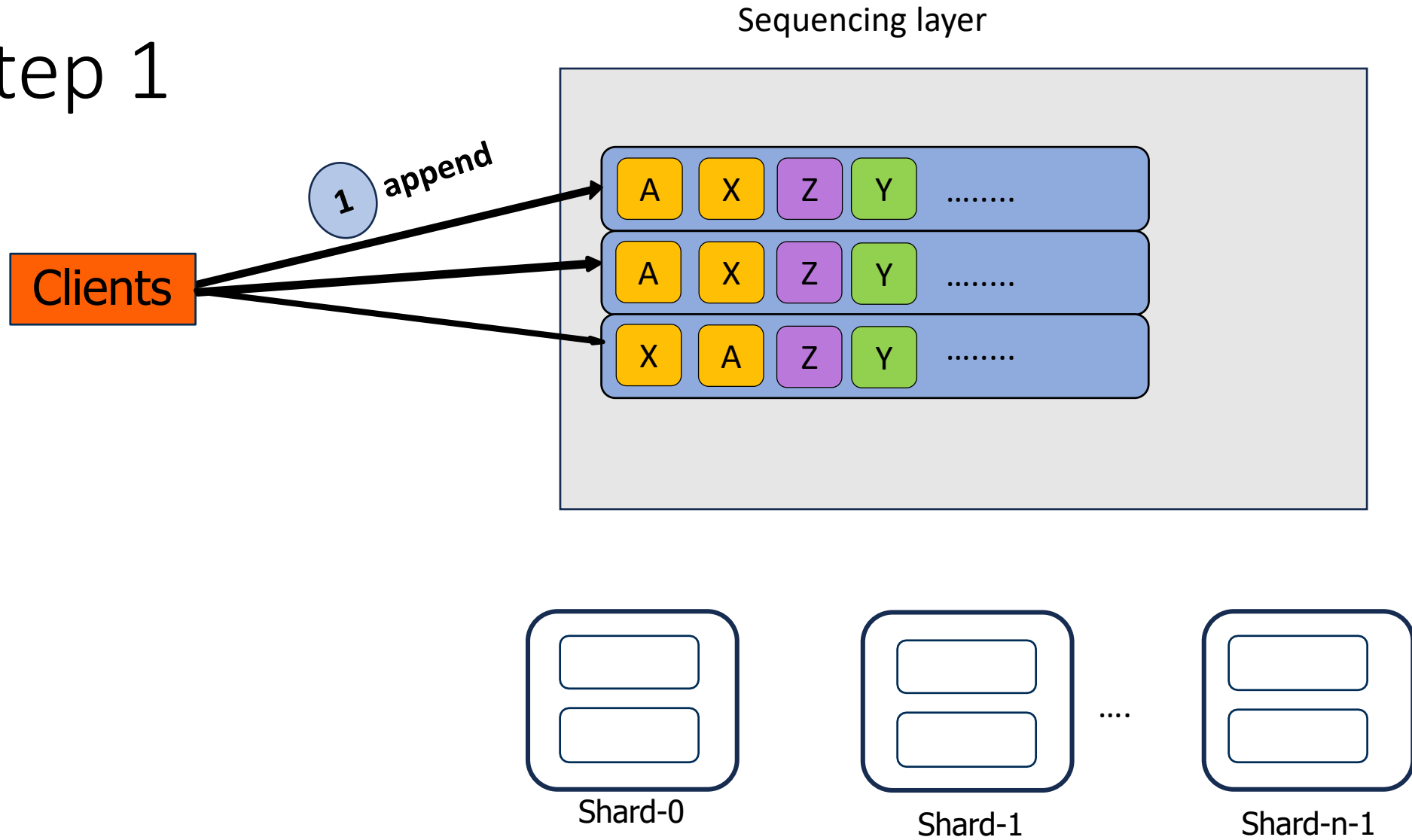
- Has set of unmodified primary-backup shards and **Sequencing layer**
- **Sequencing layer:**
 - Fault-tolerant : $f+1$ replicas to tolerate f failures
 - Coordination-free
 - Provide only short-term durability

How does Erwin Append work?

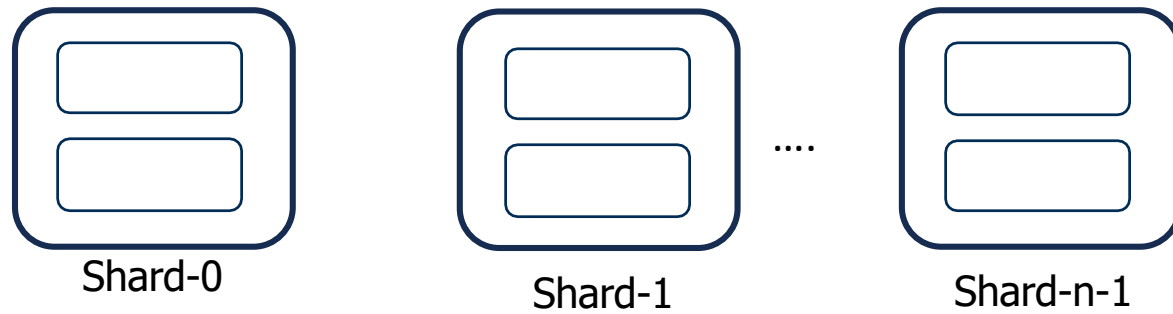
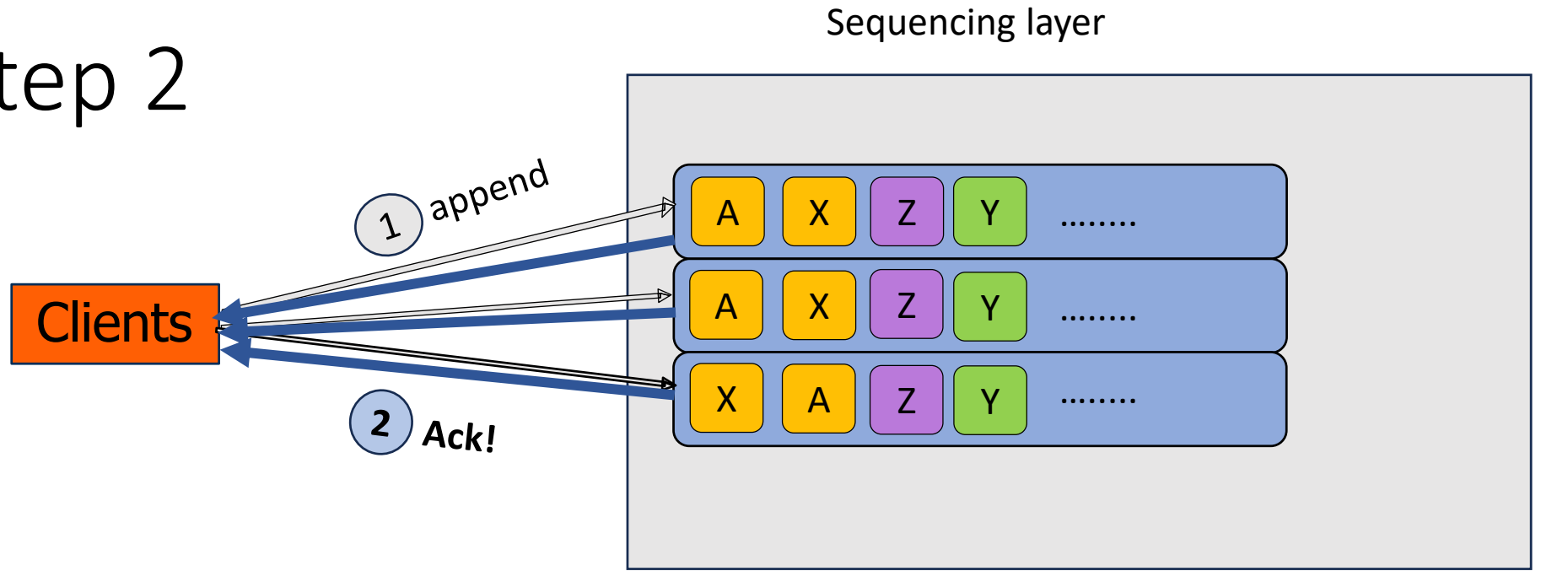
How is Lazy ordering done?

Step by step

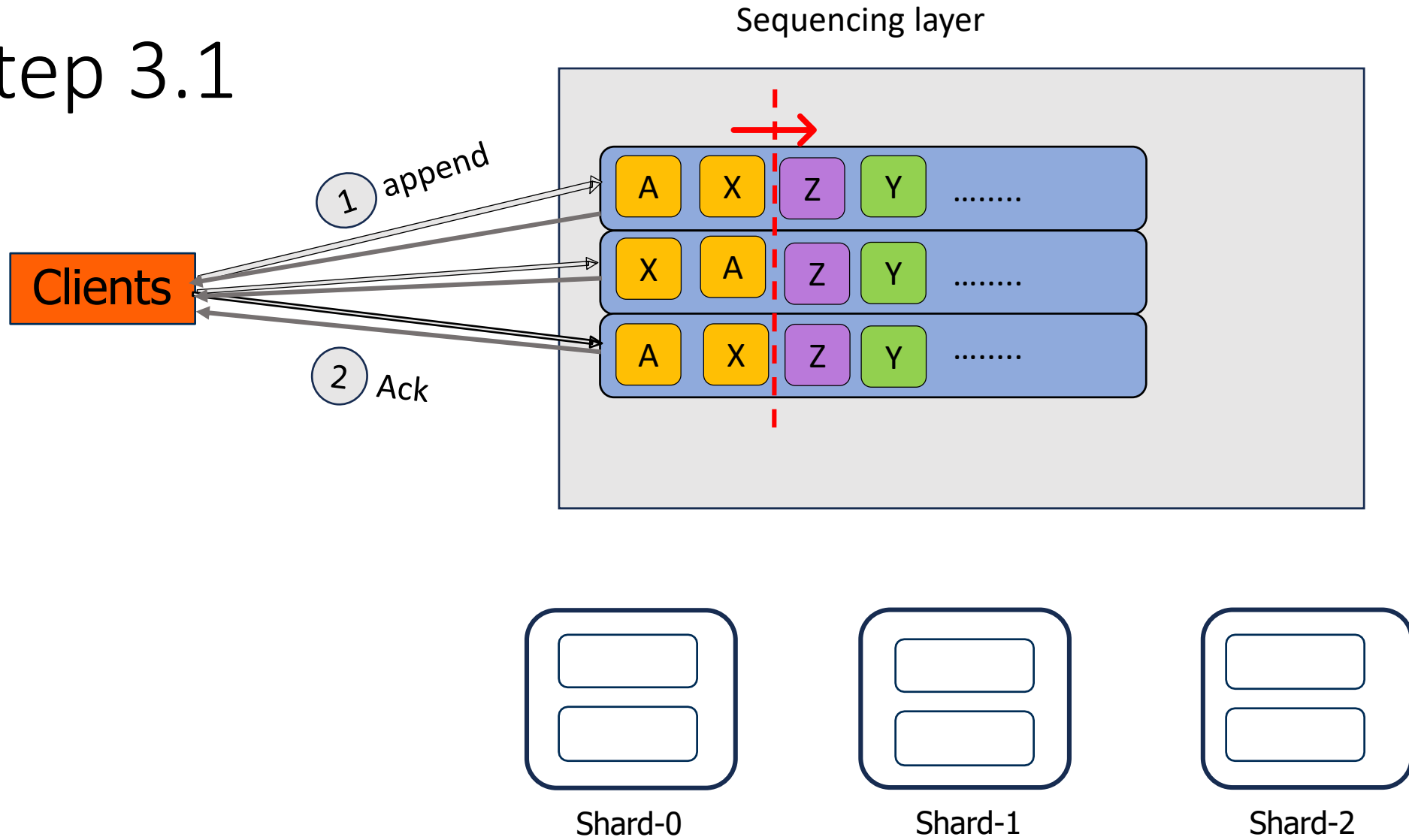
Step 1



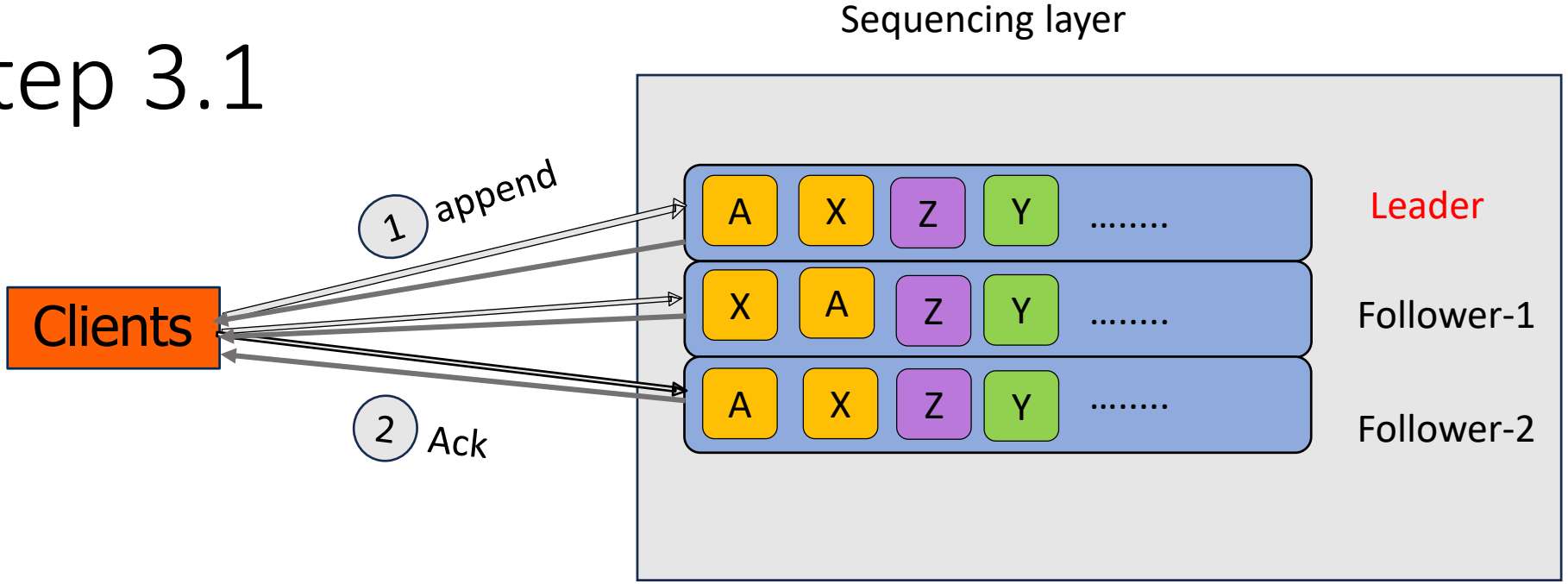
Step 2



Step 3.1



Step 3.1

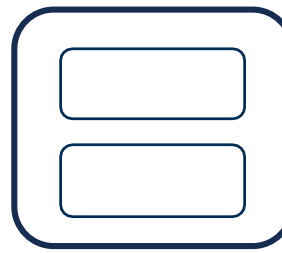


Assuming 3 shards
Mapping: $p \bmod 3$

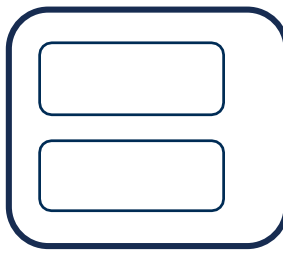
0	A	Shard 0
1	X	Shard 1
2	Z	Shard 2
3	Y	Shard 0



Shard-0

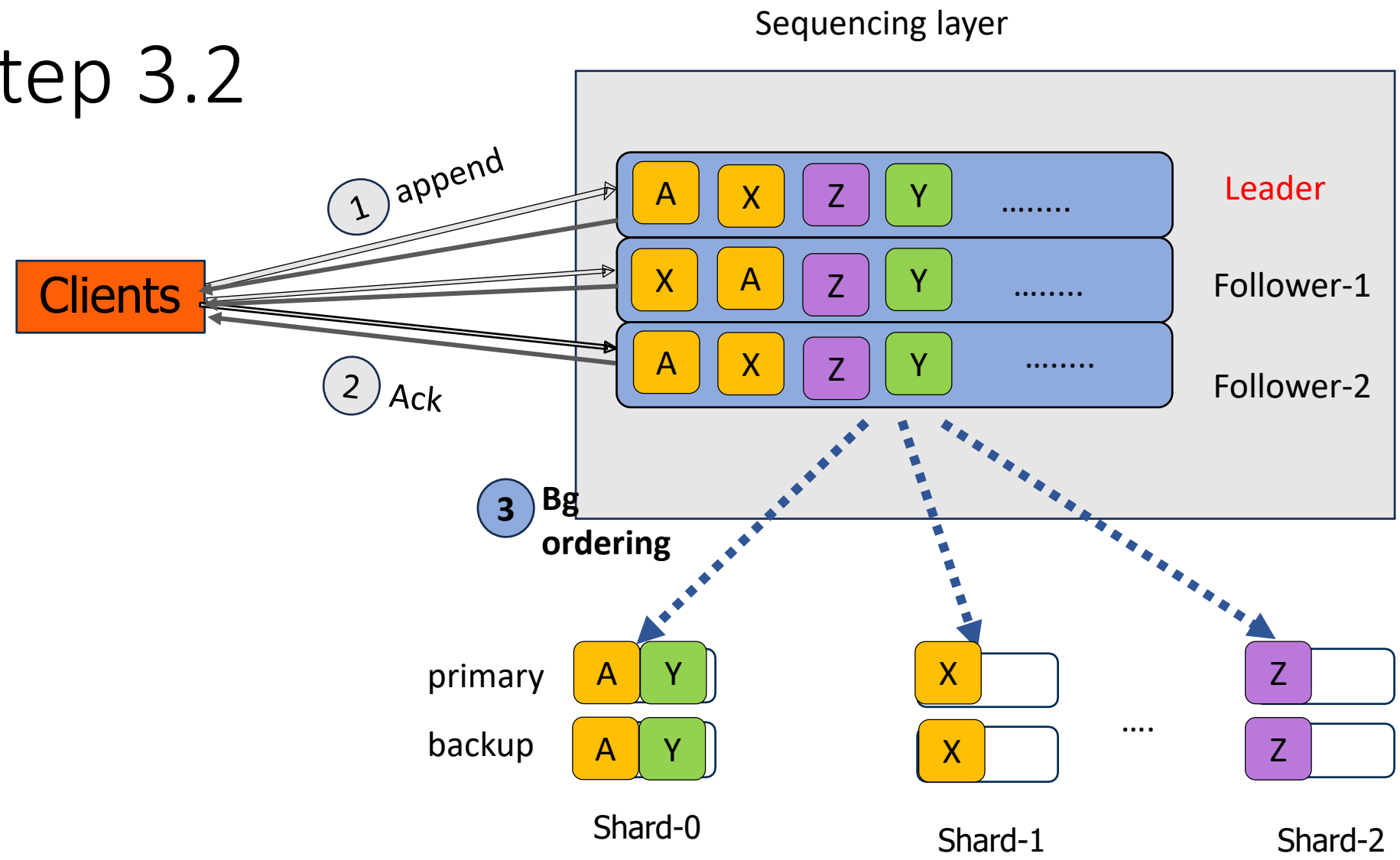


Shard-1

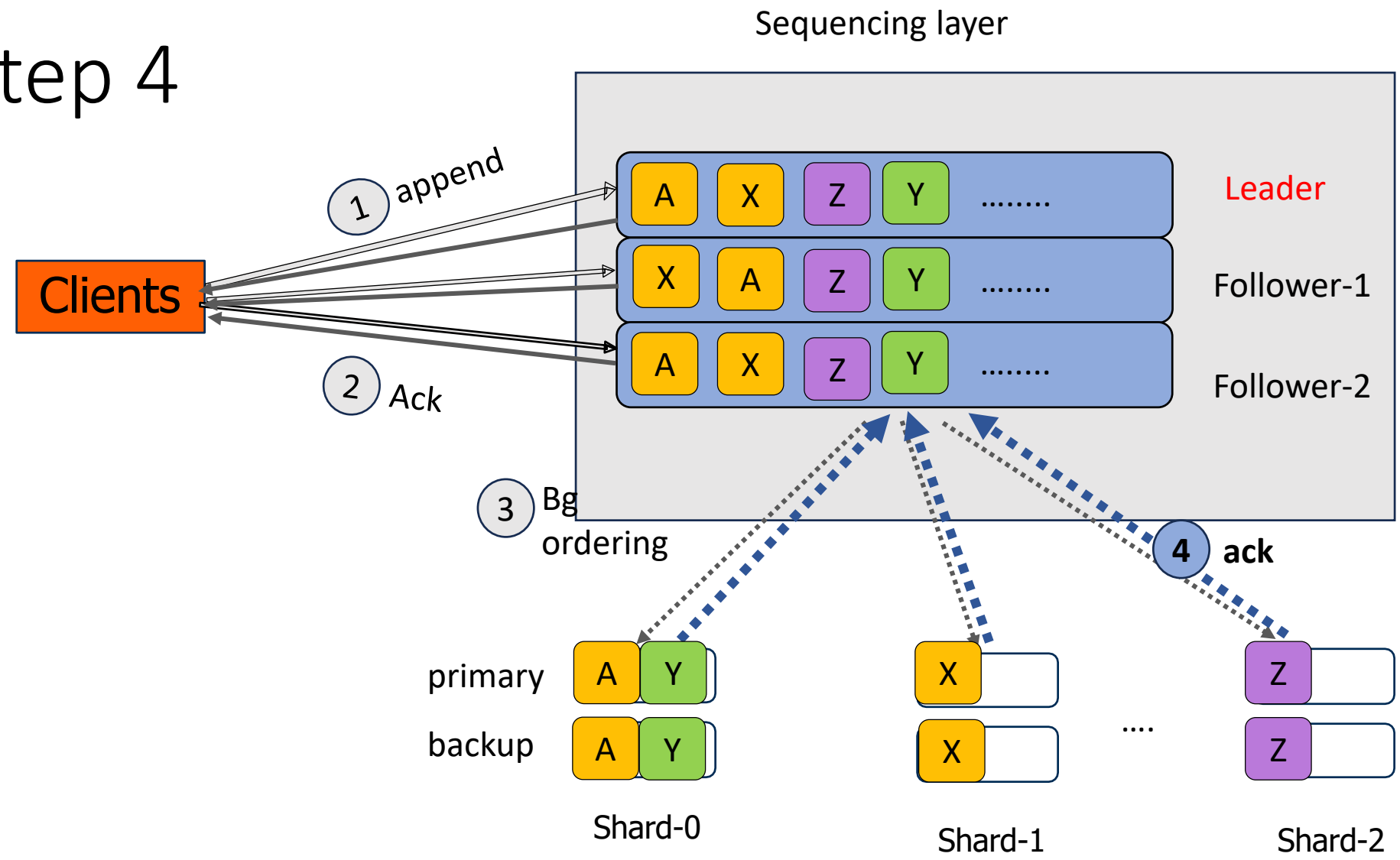


Shard-2

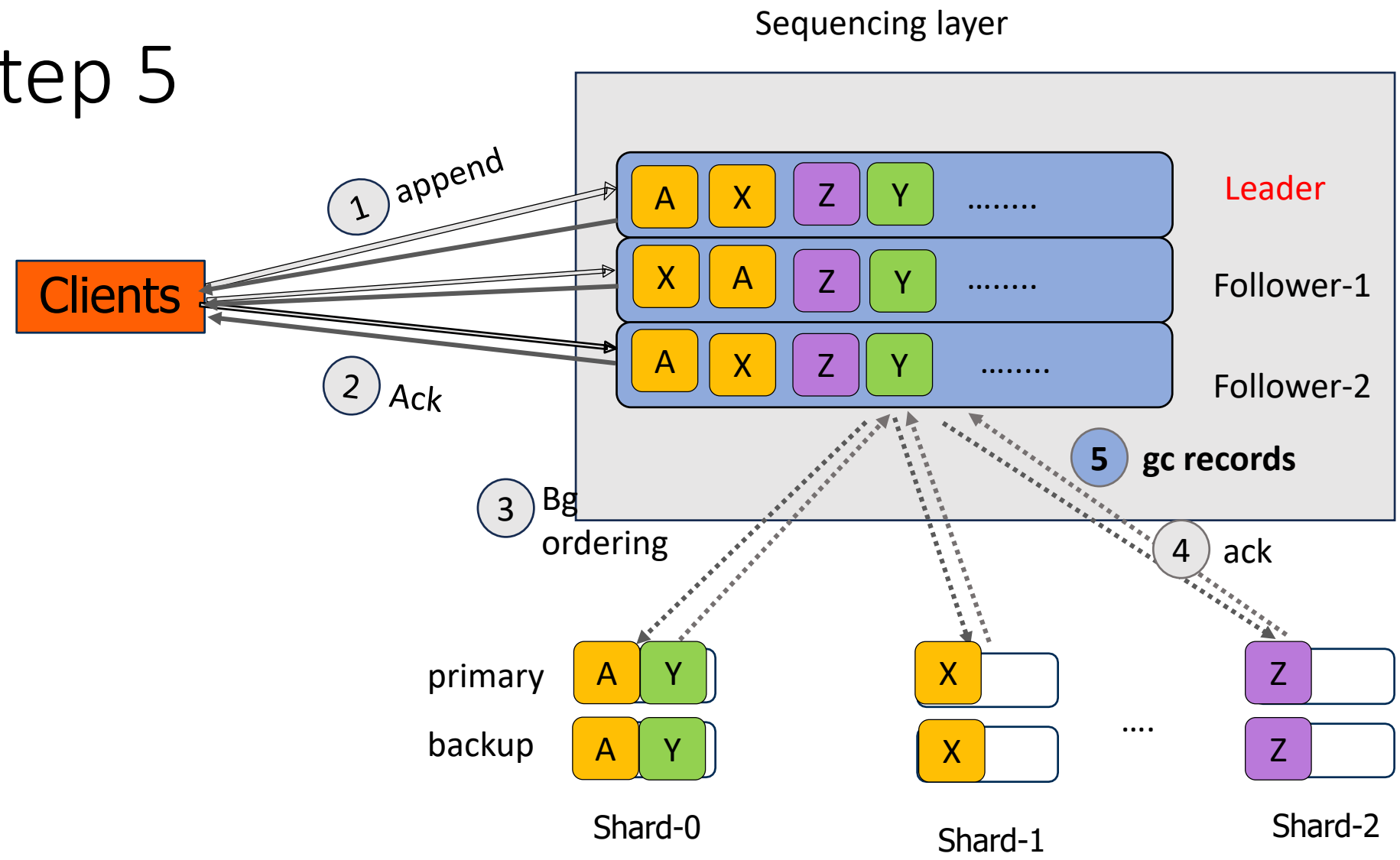
Step 3.2



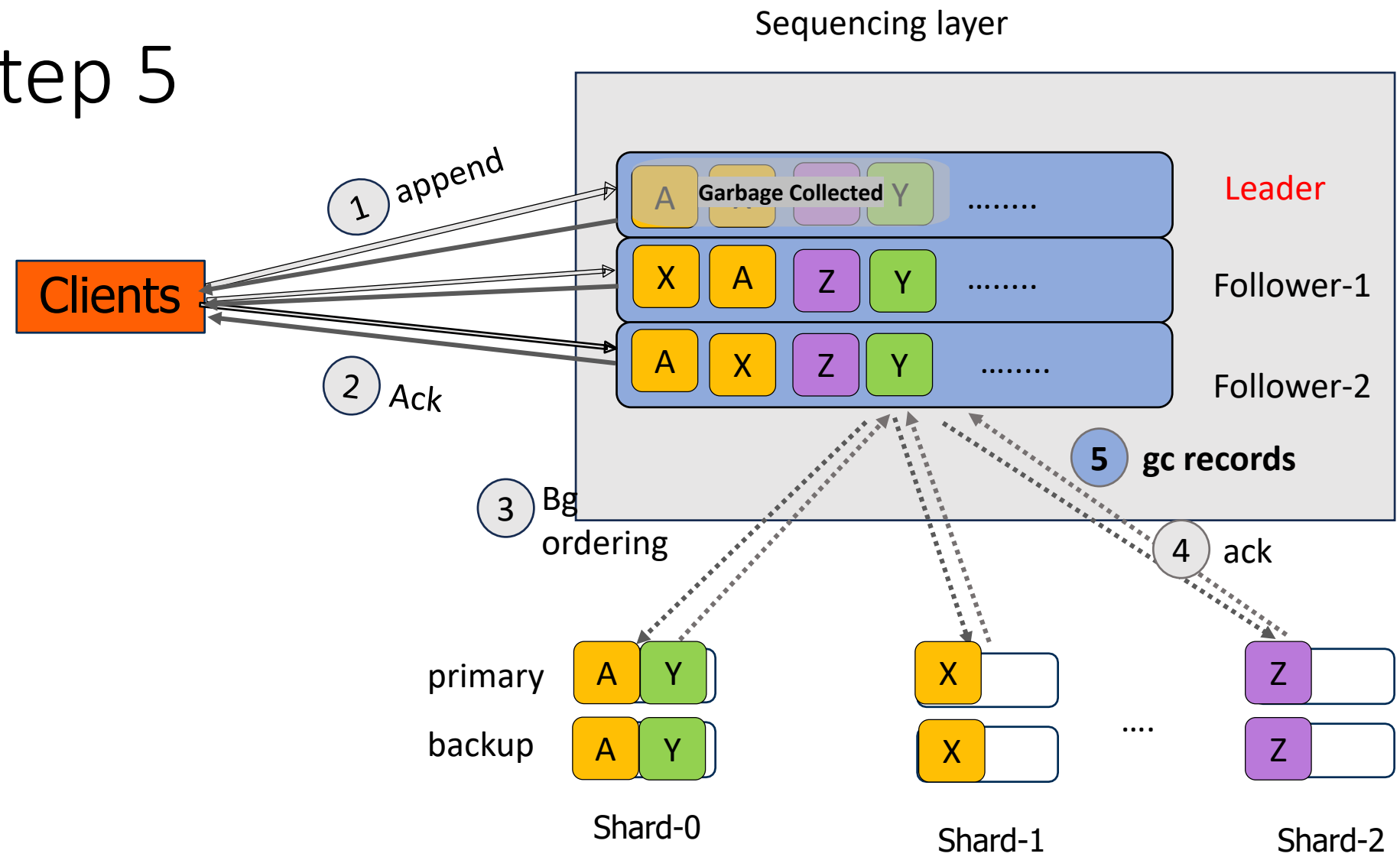
Step 4



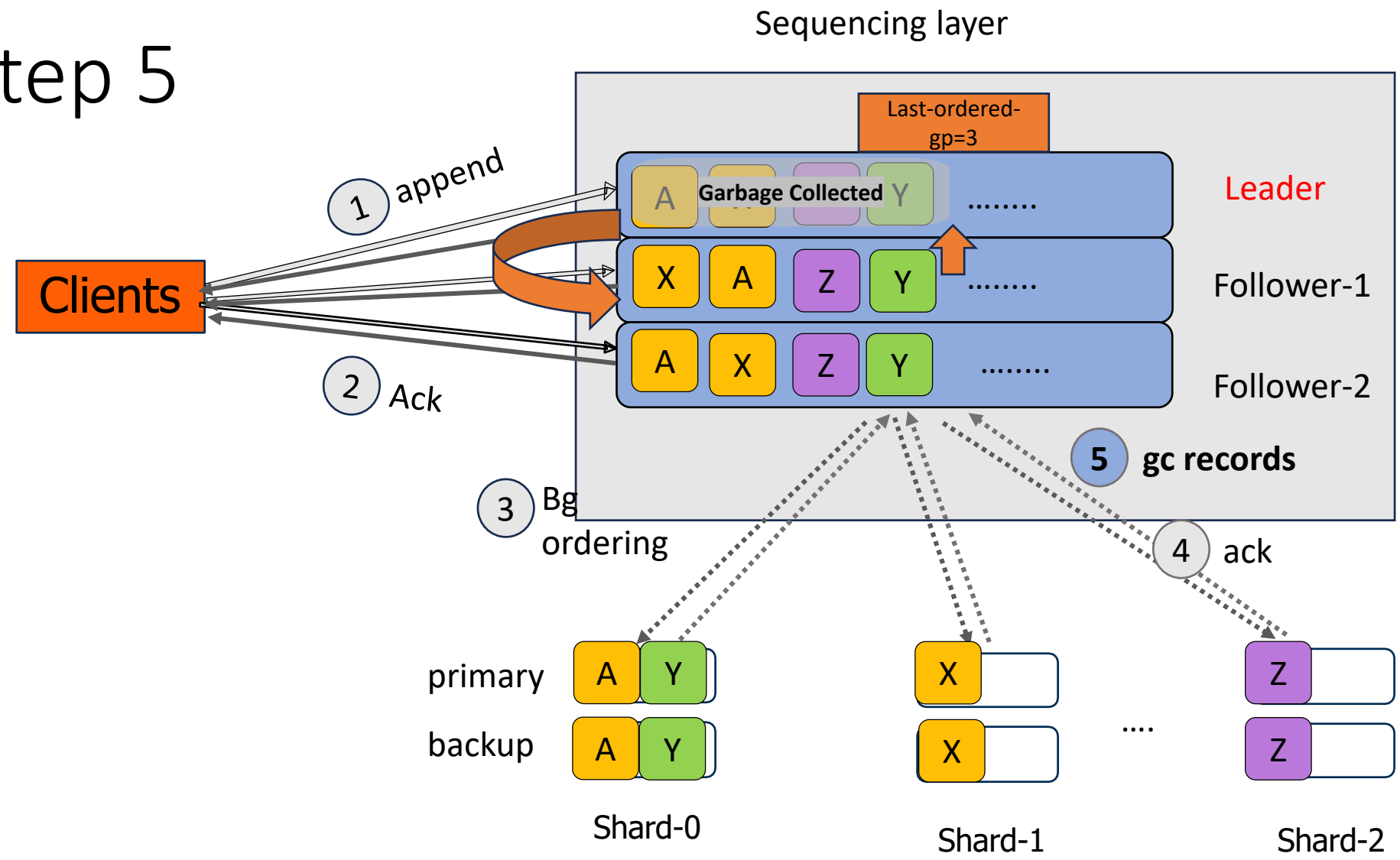
Step 5



Step 5

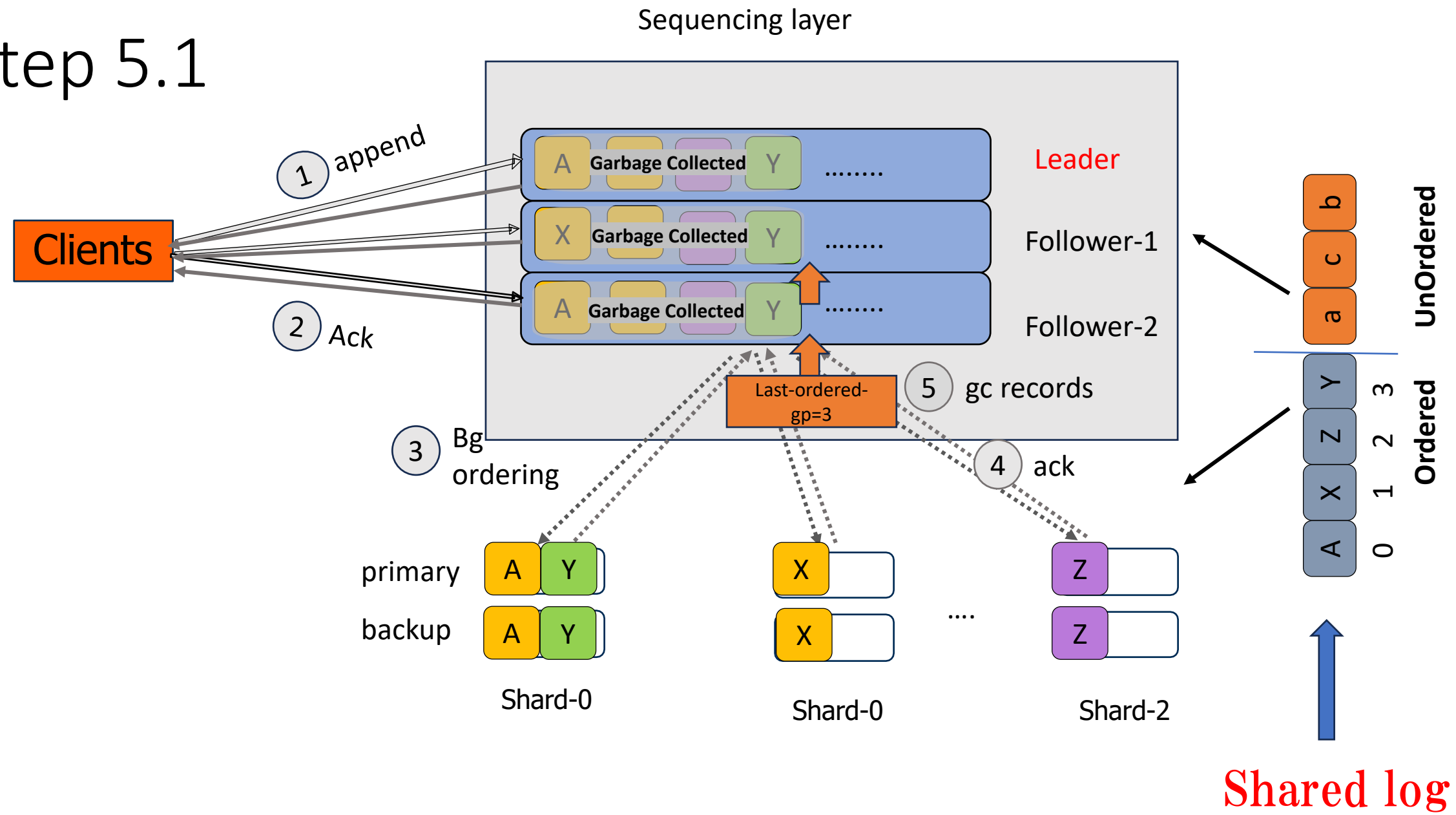


Step 5



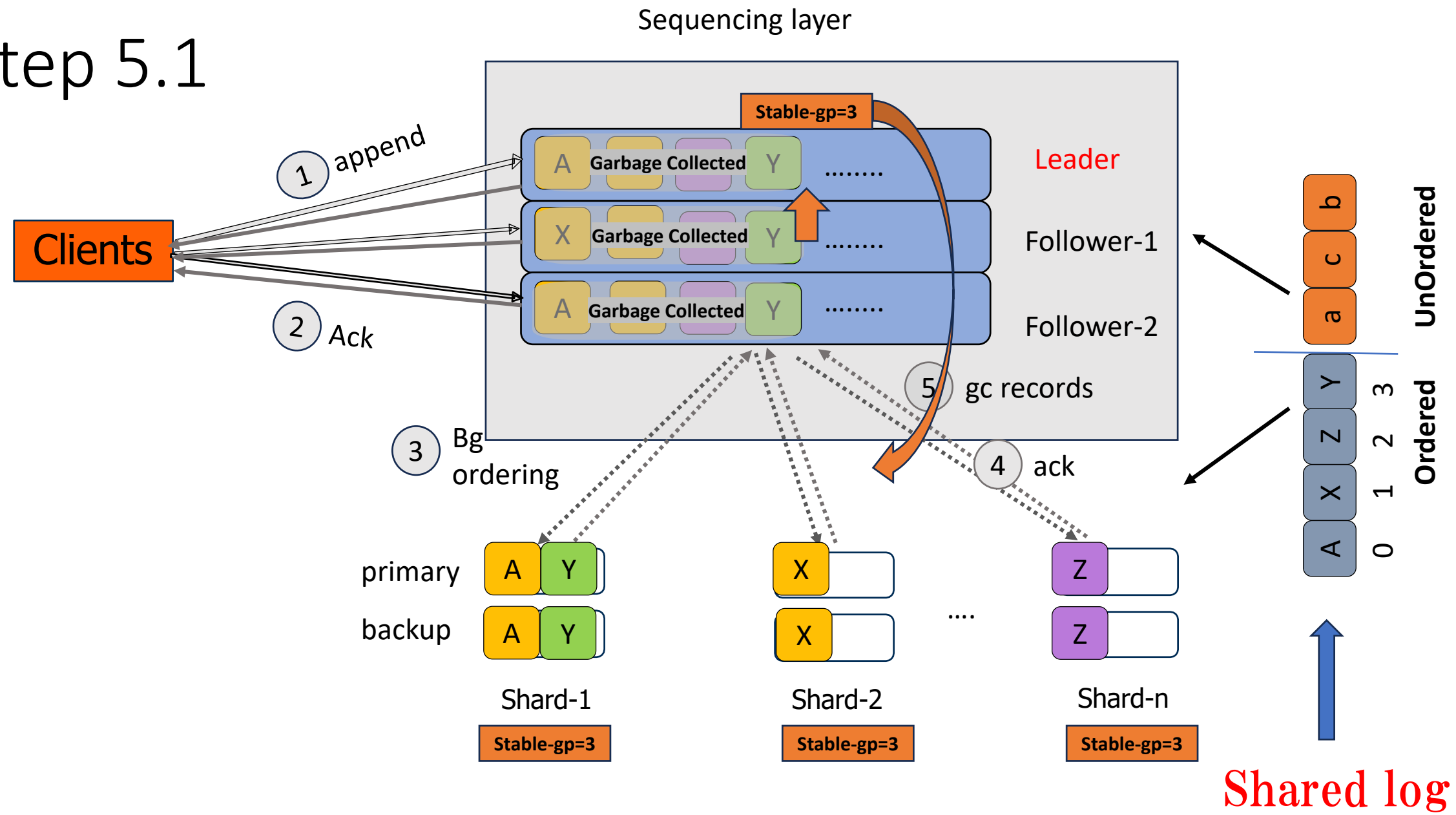
Record

Step 5.1



Record

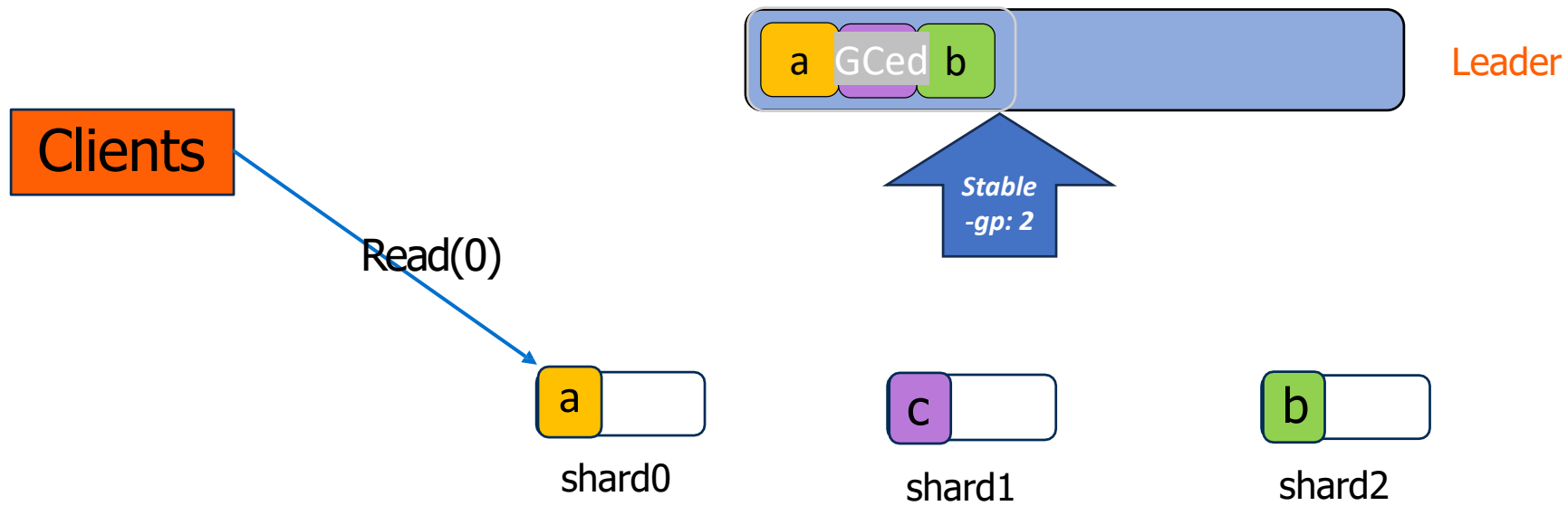
Step 5.1



How does Erwin Read work?

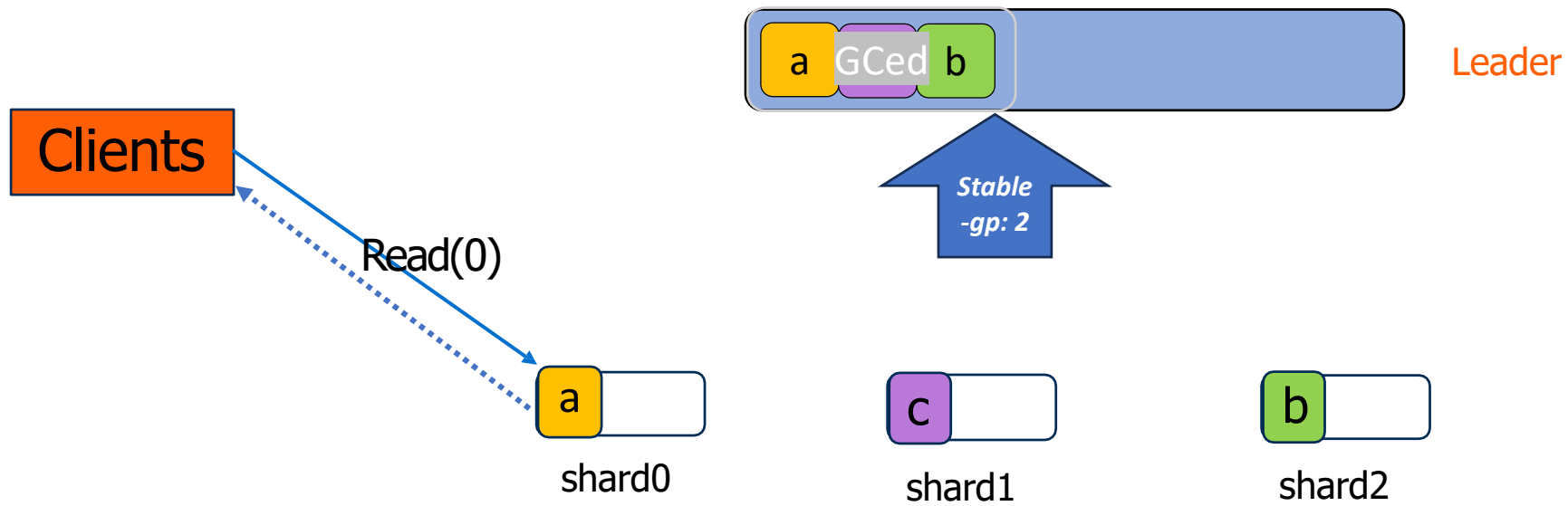
Step by step

Erwin: Read



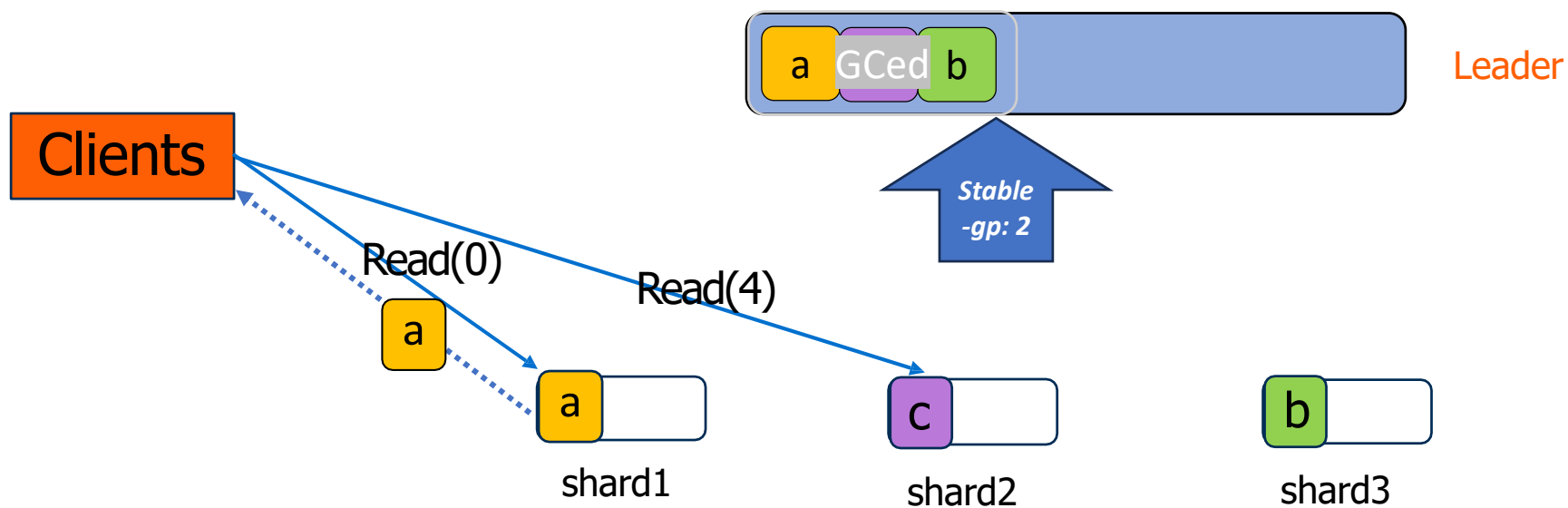
- Reading ordered position (fast read): entry returned directly

Erwin: Read



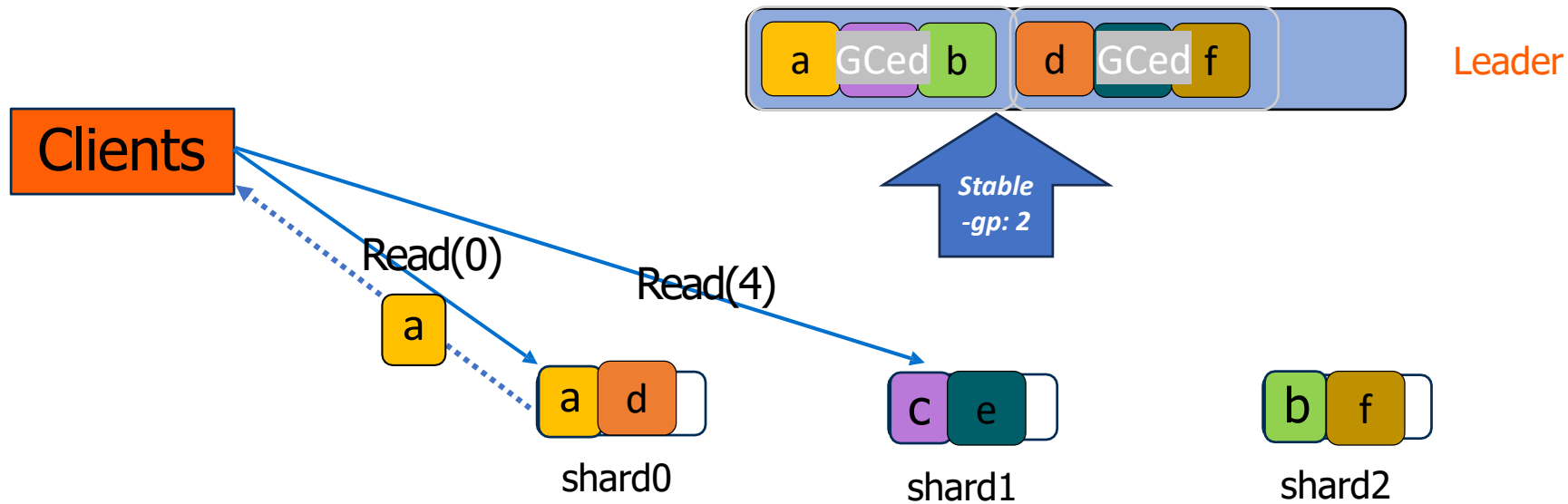
- Reading ordered position (fast read): entry returned directly

Erwin: Read



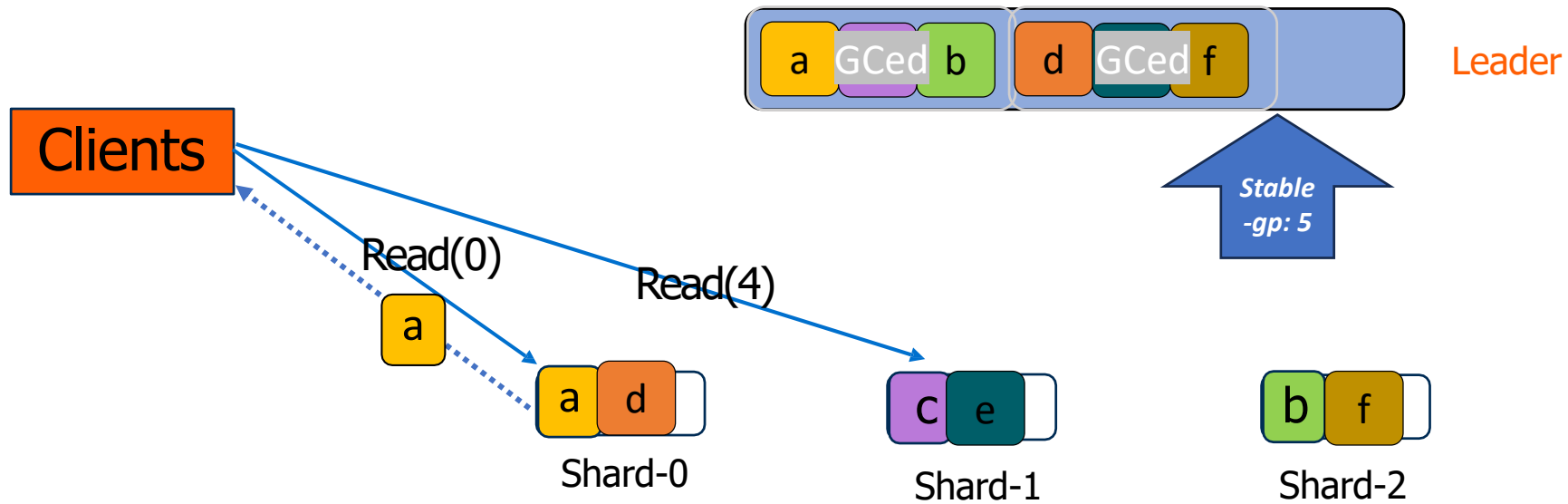
- Reading ordered position (fast read): entry returned directly
- Reading unordered position (slow read): must wait until ***stable-gp*** is advanced to the read position

Erwin: Read



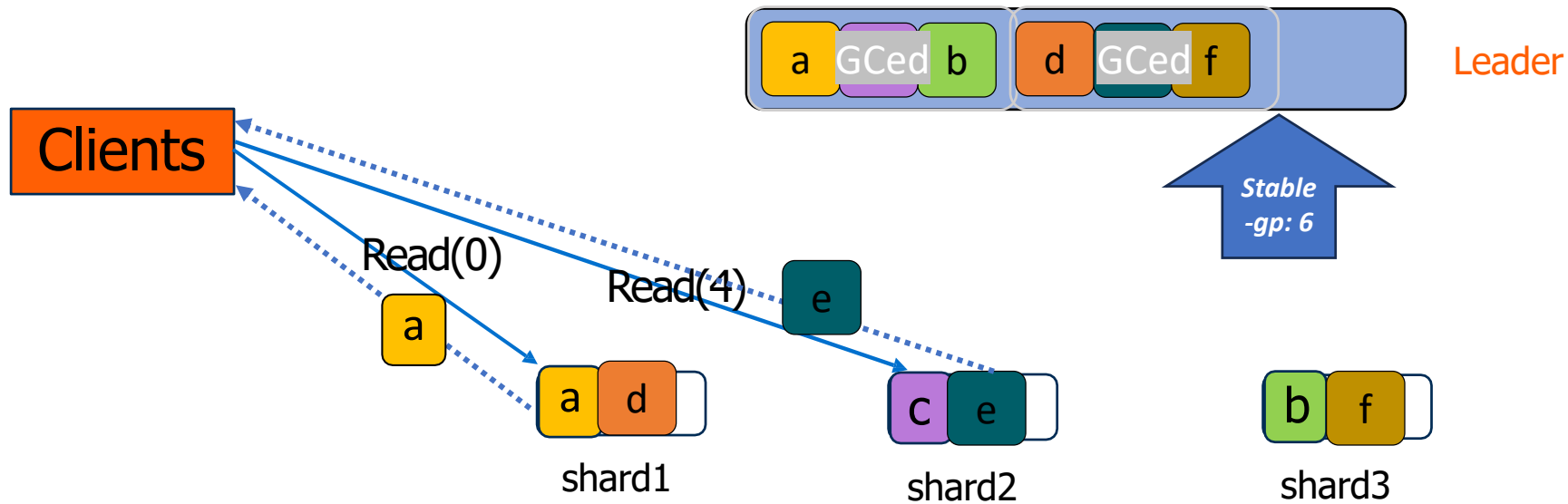
- Reading ordered position (fast read): entry returned directly
- Reading unordered position (slow read): must wait until **stable-gp** is advanced to the read position

Erwin: Read



- Reading ordered position (fast read): entry returned directly
- Reading unordered position (slow read): must wait until ***stable-gp*** is advanced to the read position

Erwin: Read



- Reading ordered position (fast read): entry returned directly
- Reading unordered position (slow read): must wait until ***stable-gp*** is advanced to the read position

Problem with Erwin Blackbox

Large Record --> Sequencing layer will be saturated

Solution

Split Record --> into Data and
Metadata (Record_Id, Shard_Id)

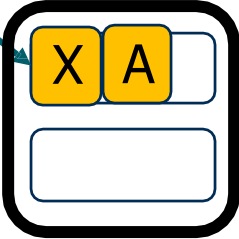


Record

Clients

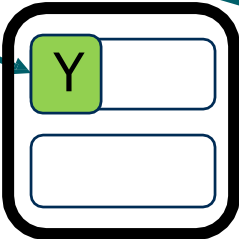
X

A



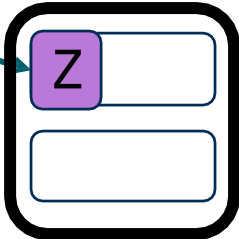
shard1

Y



shard2

Z



shard3

Primary

Backup



Record

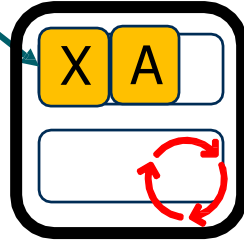
Clients

Problem:

Require coordination within
a shard

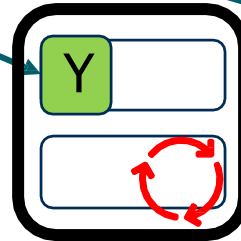
X

A



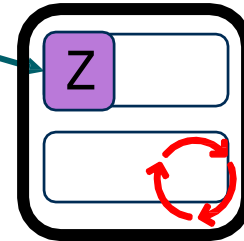
shard1

Y



shard2

Z



shard3

Primary

Backup



Record

Clients

Problem:

Require coordination within a shard

Solution:

Send record to all replicas in the shard

X

X

A

A

X

A

A

X

shard1

Y

Y

Y

Y

shard2

Z

Z

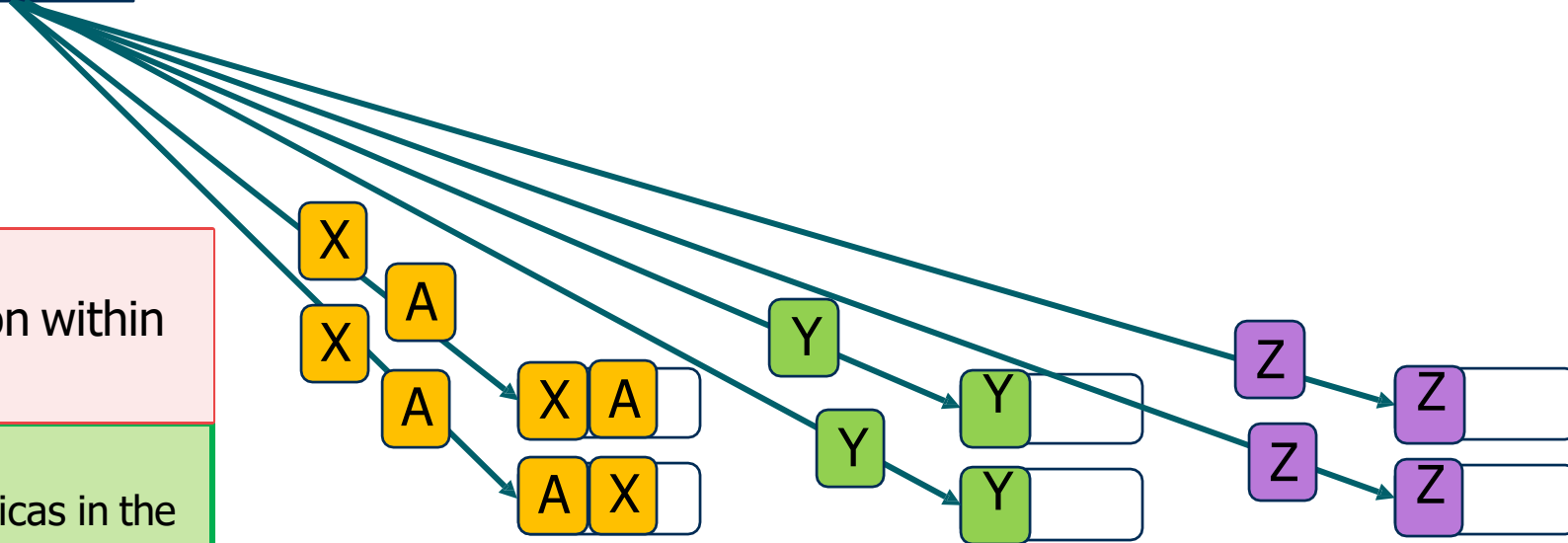
Z

Z

shard3

Primary

Backup



Motivation

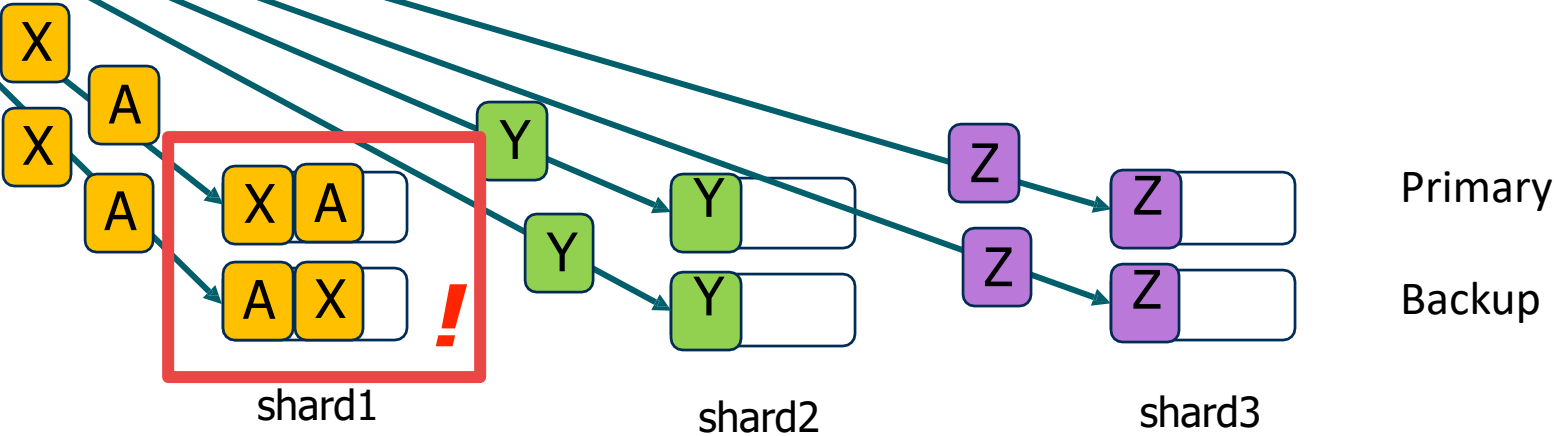
 Record

Problem:
No order across and within shards

Clients

Problem:
Require coordination within a shard

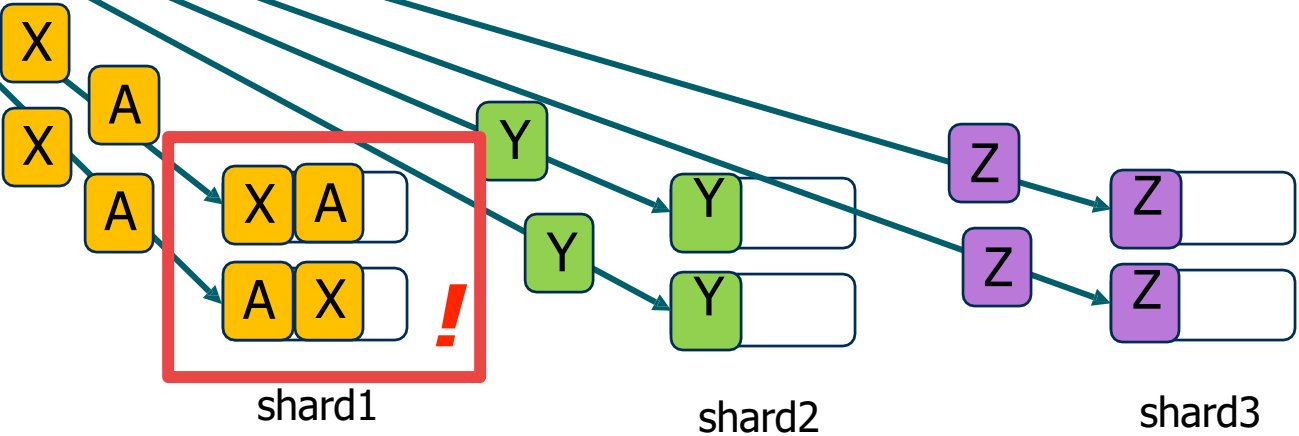
Solution:
Send record to all replicas in the shard



Motivation

 Record

Clients



Problem:
Require coordination within a shard

Solution:
Send record to all replicas in the shard

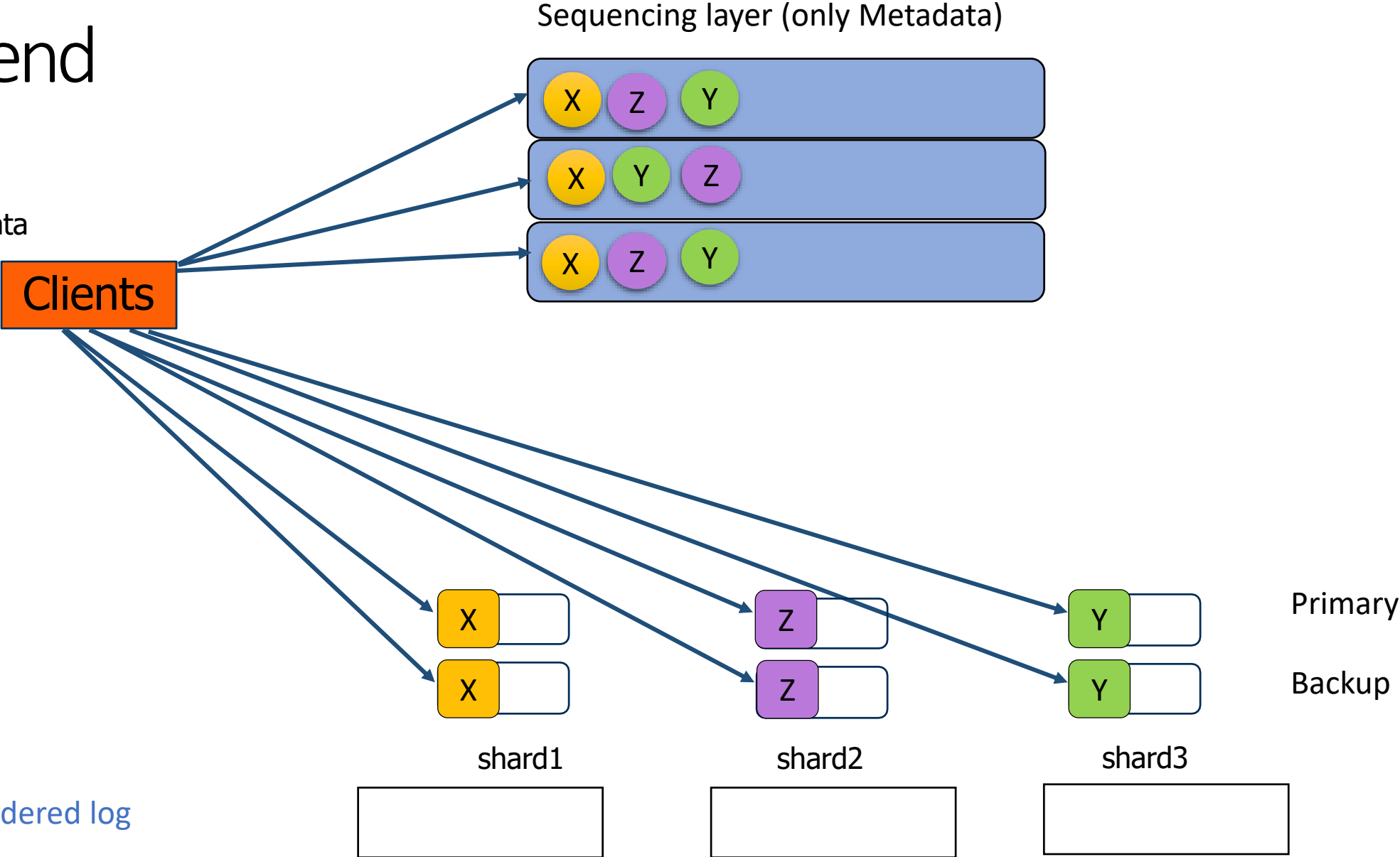
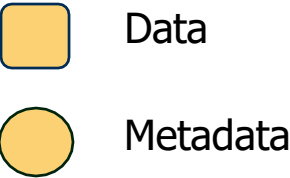
Problem:
No order across and within shards

Solution:
Split the Record into Data and Metadata and then send the metadata to sequencing layer

Erwin-st(scalable throughput)

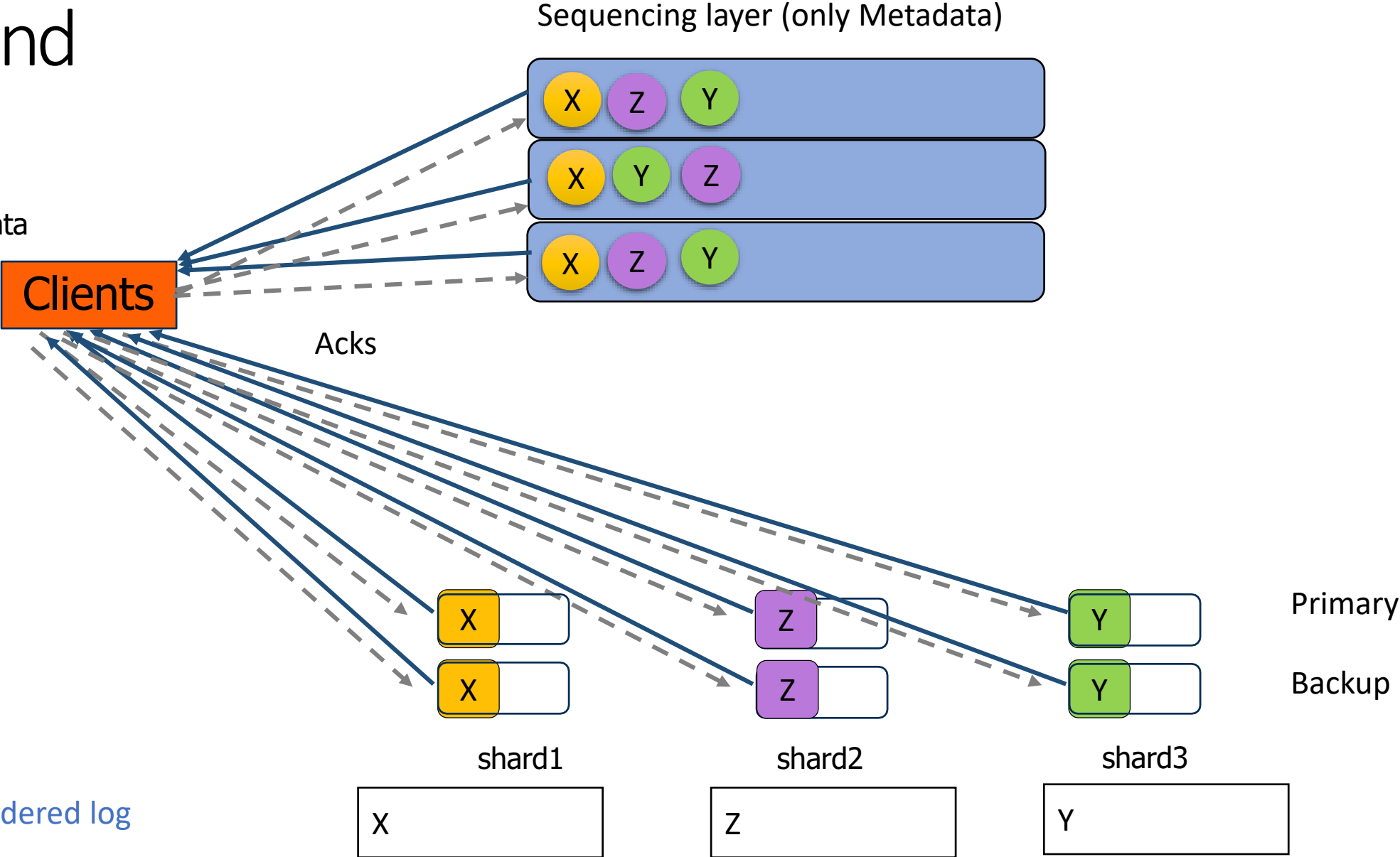
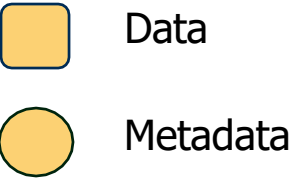
Erwin's Goal: 1-RTT

Append



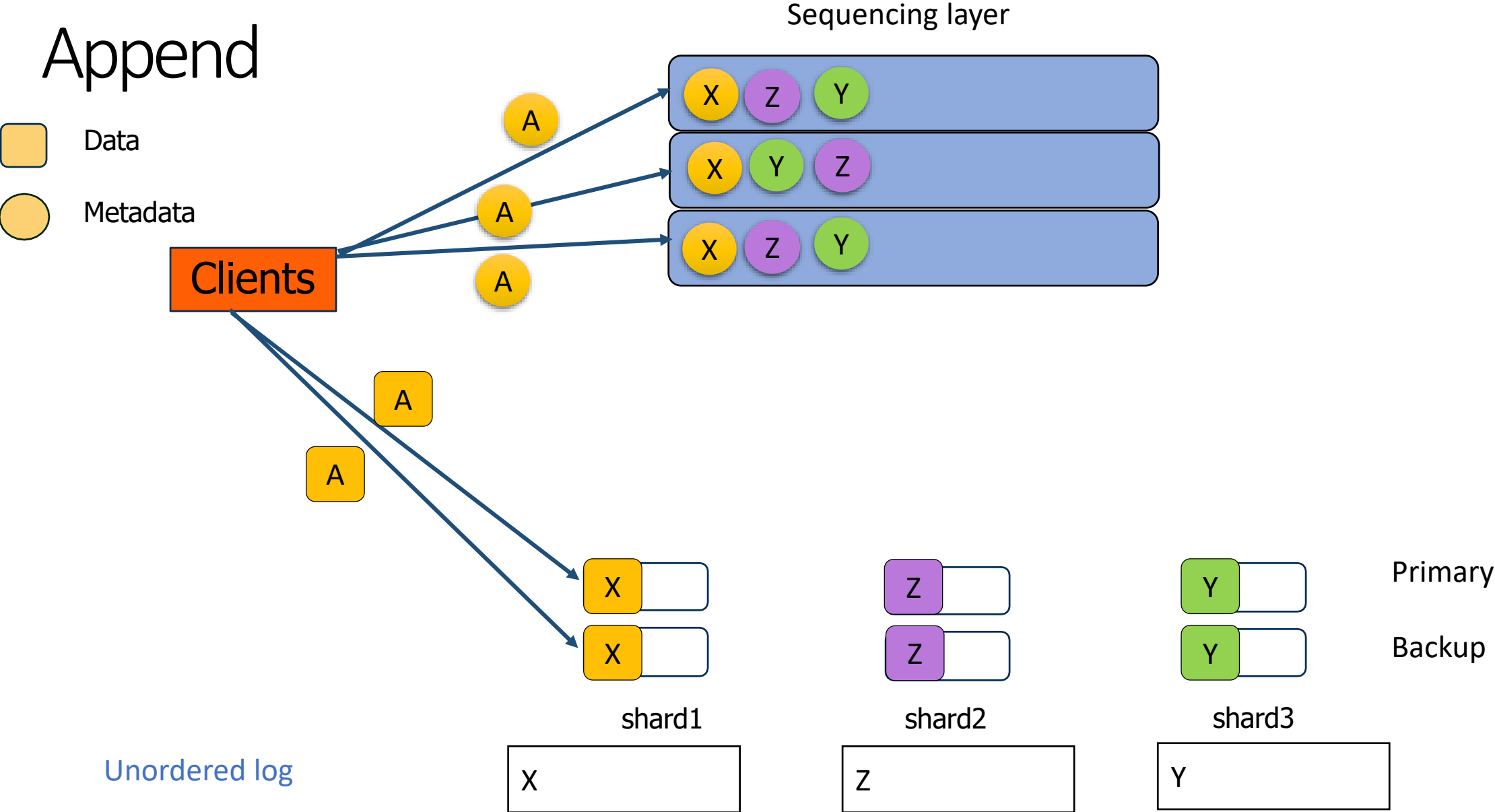
Erwin's Goal: 1-RTT

Append



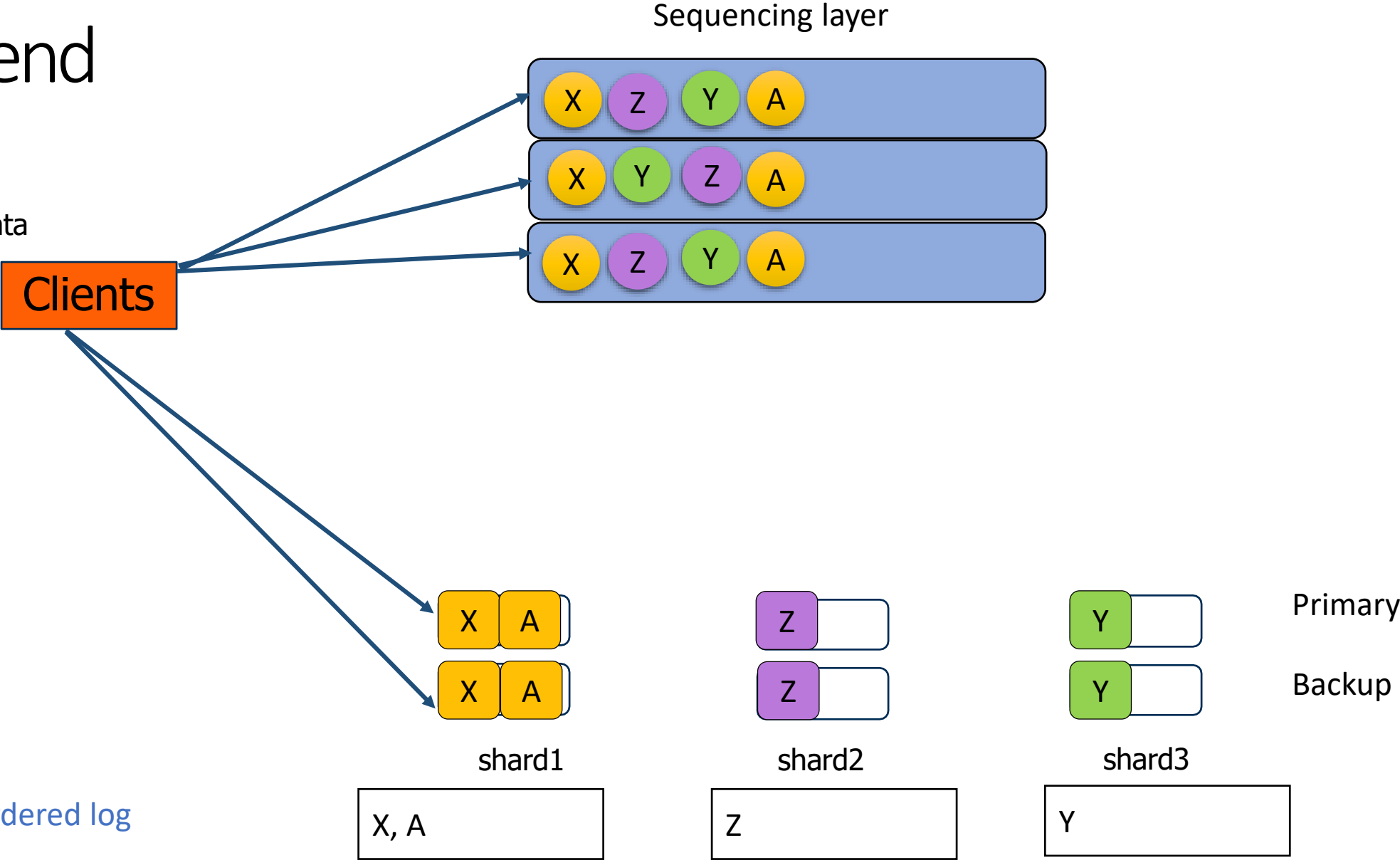
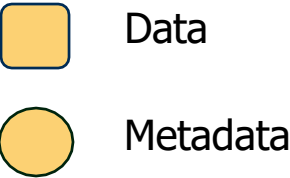
Erwin's Goal: 1-RTT

Append



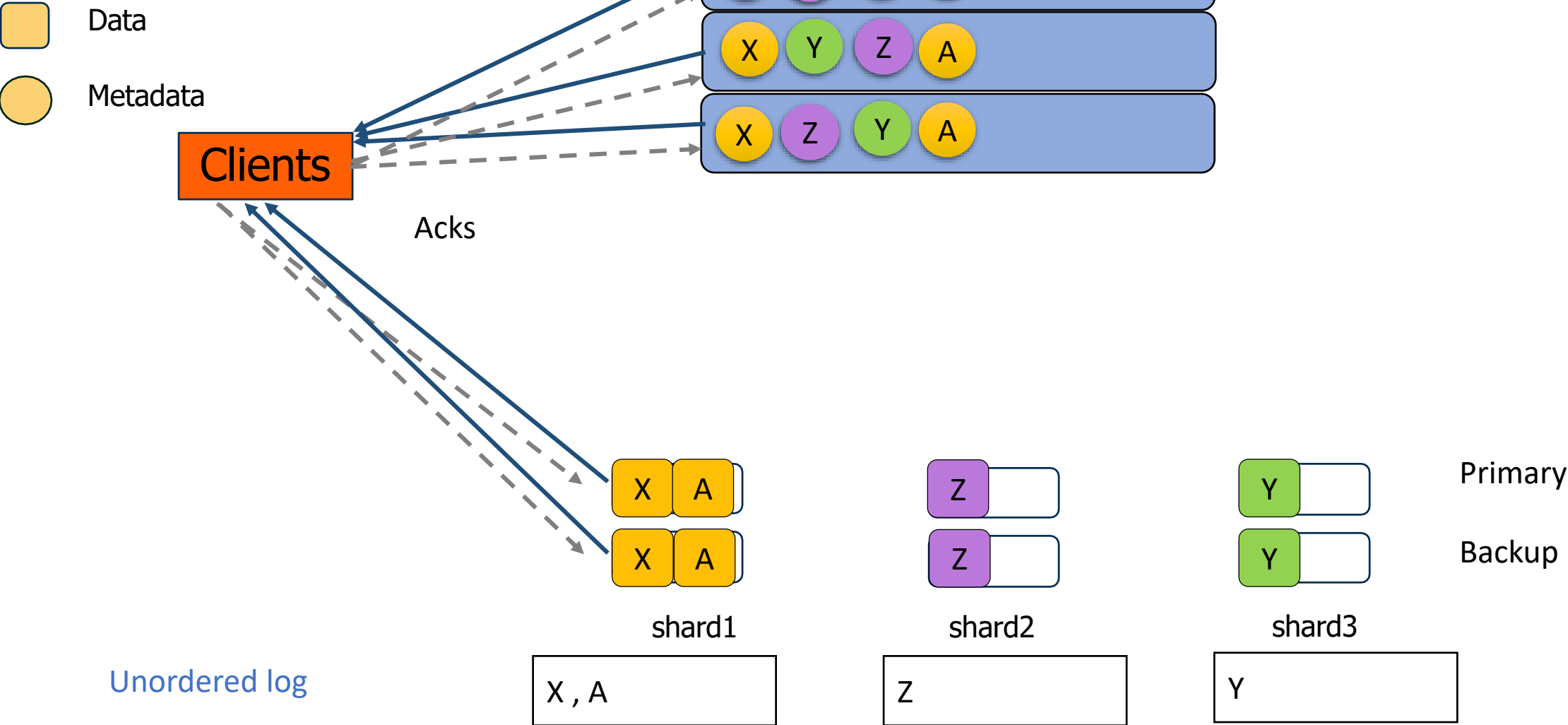
Erwin's Goal: 1-RTT

Append





Erwin's Goal: 1-RTT

Append



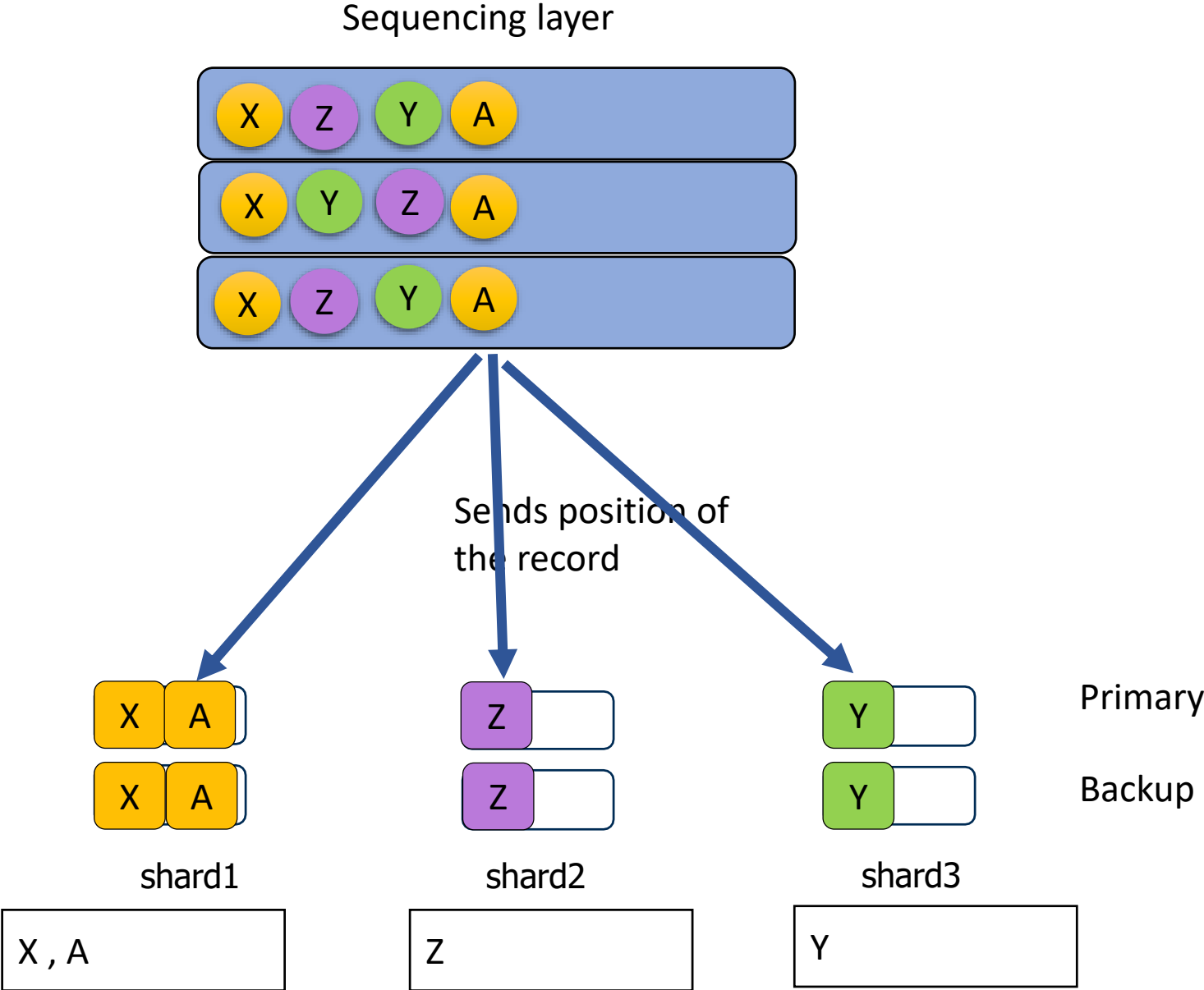
Erwin's Goal: 1-RTT

Append

-  Data
-  Metadata



Clients

Unordered log

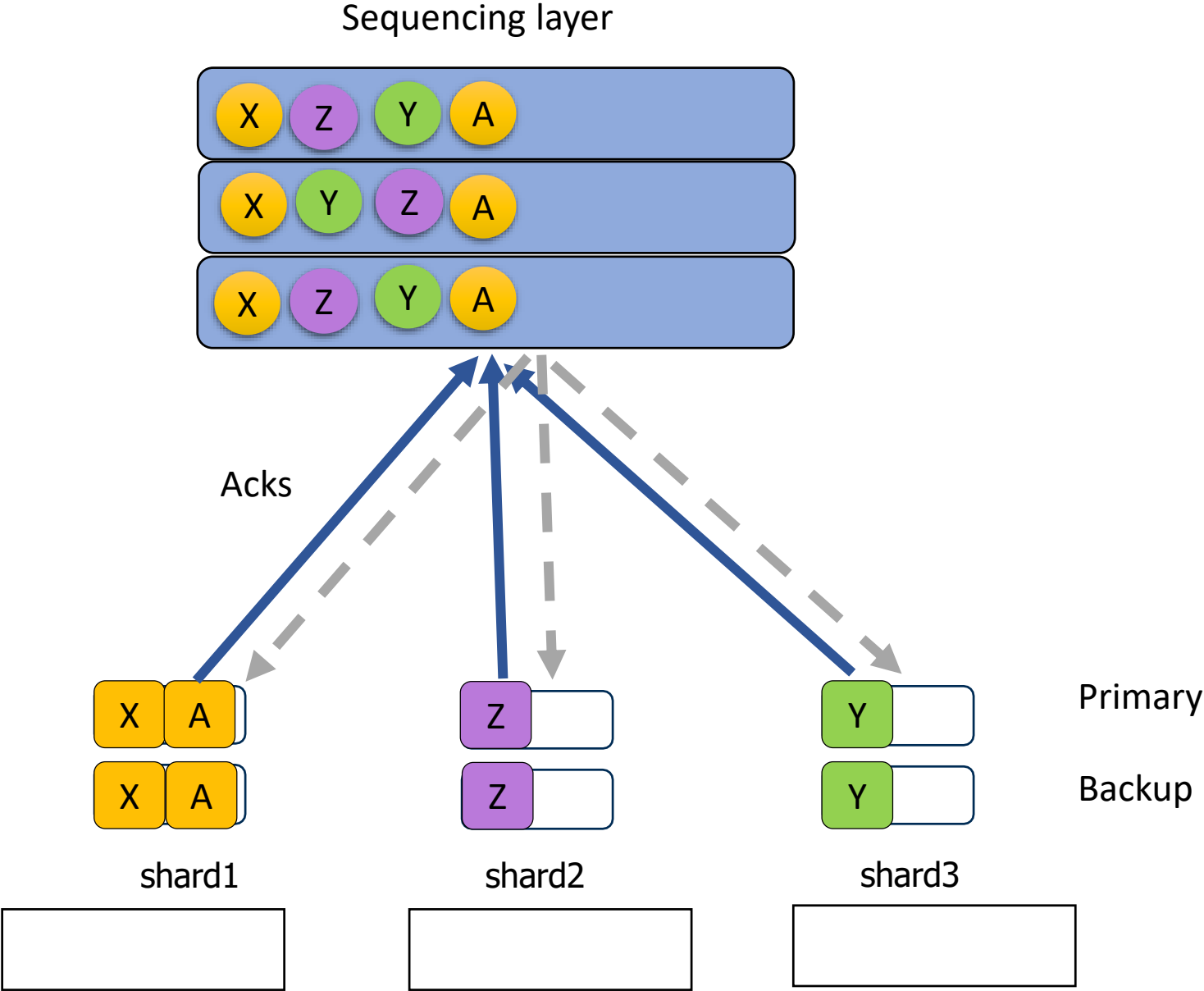


Erwin's Goal: 1-RTT

Append

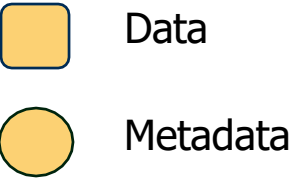
-  Data
-  Metadata

Clients

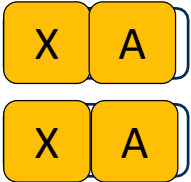
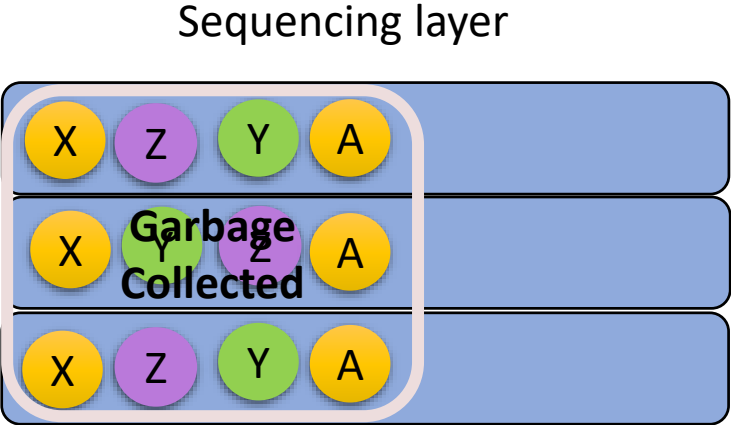


Erwin's Goal: 1-RTT

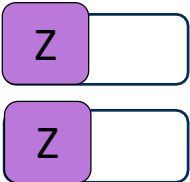
Append



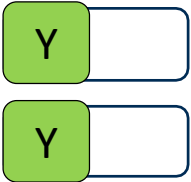
Clients



shard1



shard2



shard3

Primary

Backup

Unordered log



Failure Handling

Failure Handling

- Shard failure is handled within the shard- Paxos/Raft
- Network blips – clients can retry their append requests
- Sequencing layer replica failure- handled using **view and reconfiguration**

Step 1: Failure detection

- Using a control plane which has a zookeeper and a controller
- Every replica in the sequencing layer maintains heartbeat with the zookeeper and if it breaks then it informs the controller about the failed replica

Step 2: Start a new view

- Controller now seals the old **view**
- Any record sent with old view won't be appended
- Any of the other replica is assigned as recovery replica
- Flush records to shards starting from last-ordered-gp+1

Failure Handling Erwin St

- Sequencing replica failure
- Shard Failure

Same as Erwin Blackbox

Failure Handling

- Sequencing replica failure
- Shard Failure

Same as Erwin Blackbox

Client Failures?

Failure Handling

- Sequencing replica failure
- Shard Failure

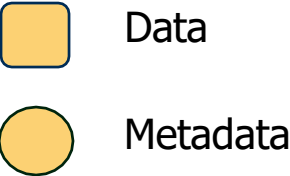
Same as Erwin Blackbox

Client
Failures?

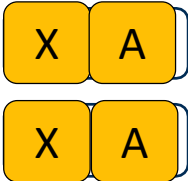
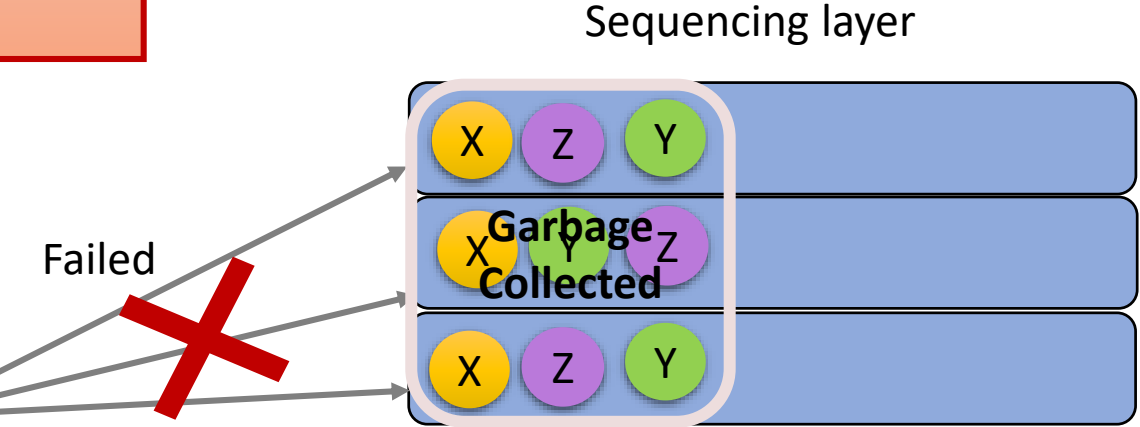


Two Problems

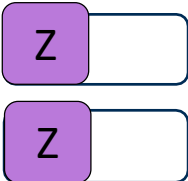
Problem-1



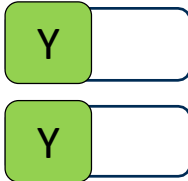
Clients



shard1



shard2

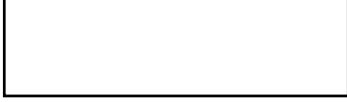


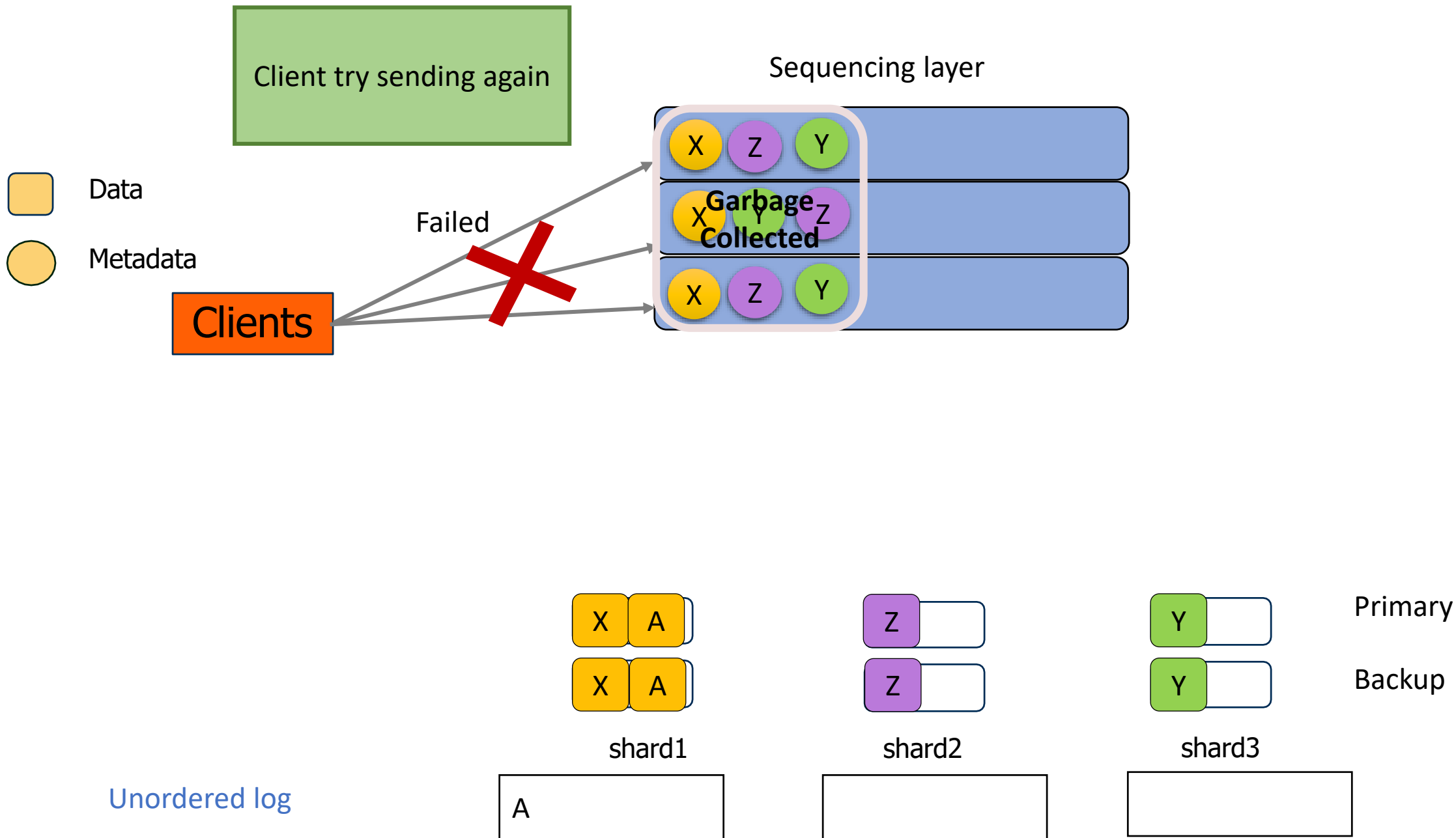
shard3

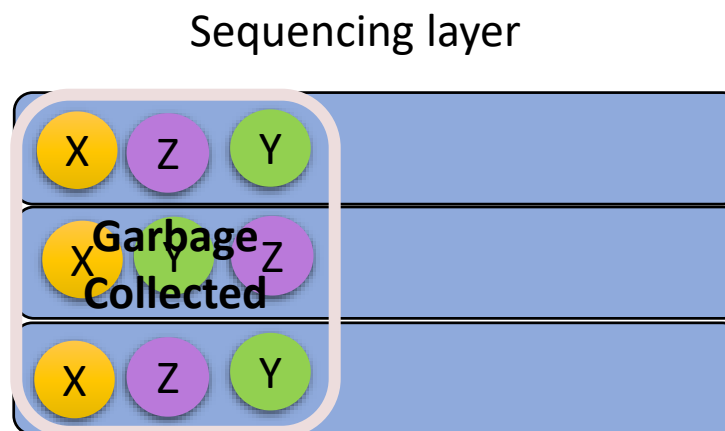
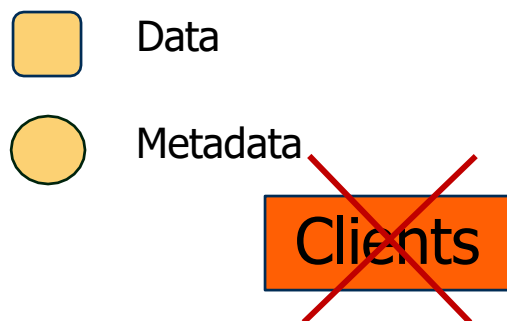
Primary

Backup

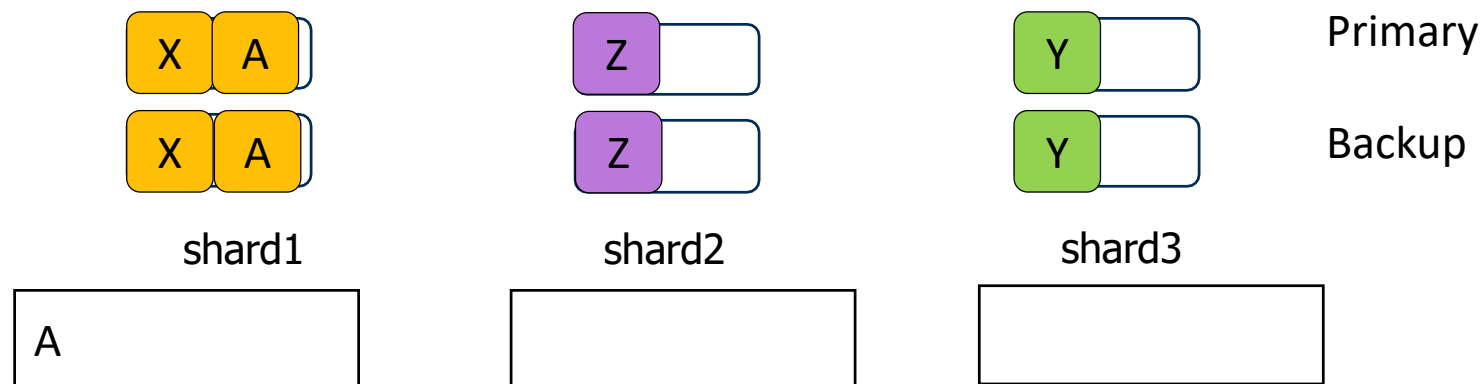
Unordered log

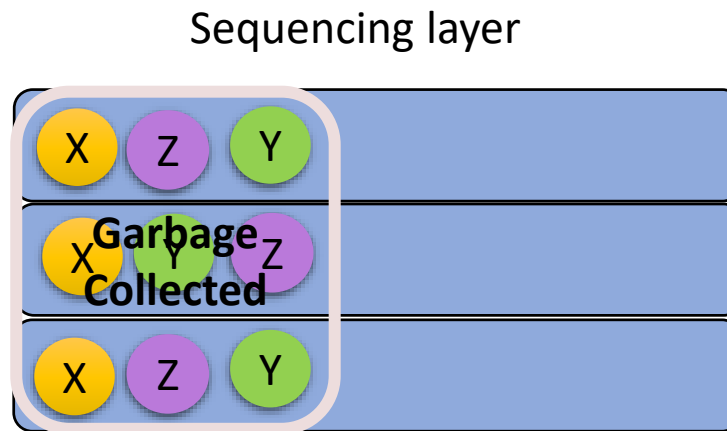
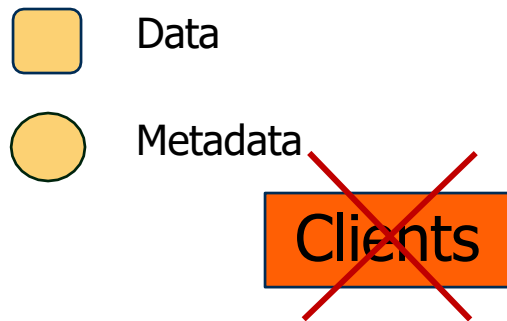




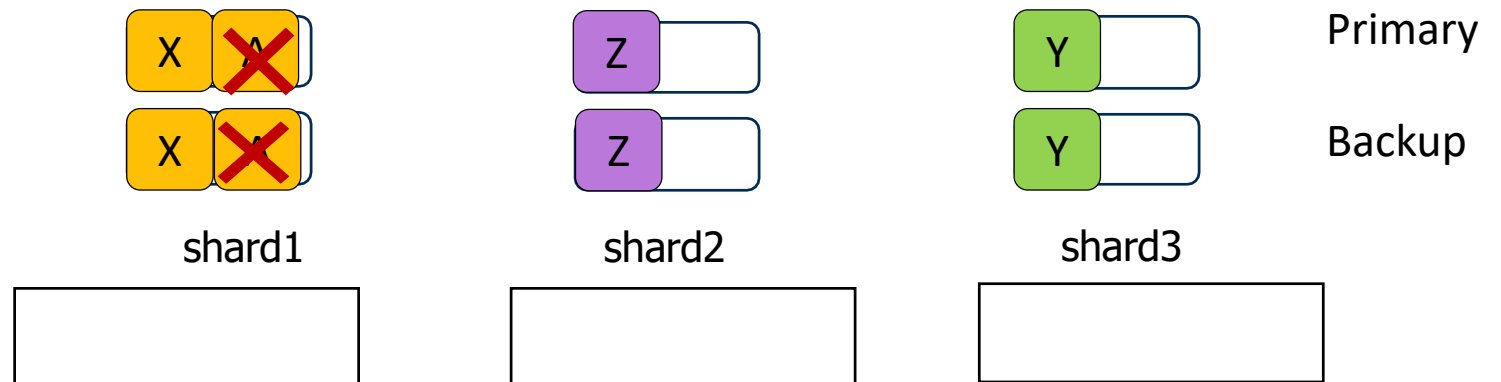


Data is Orphaned

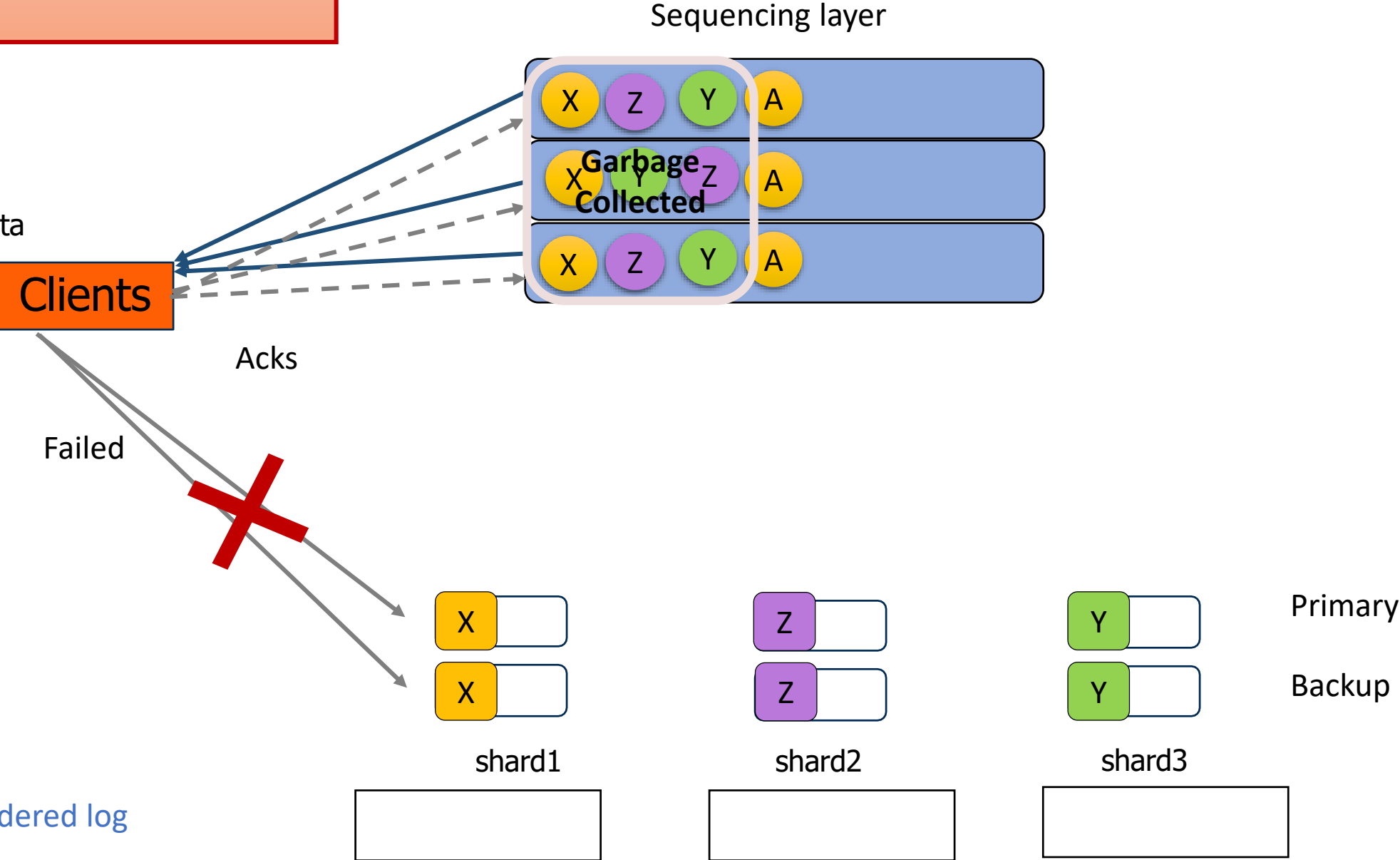
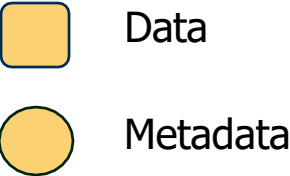






Data is Orphaned



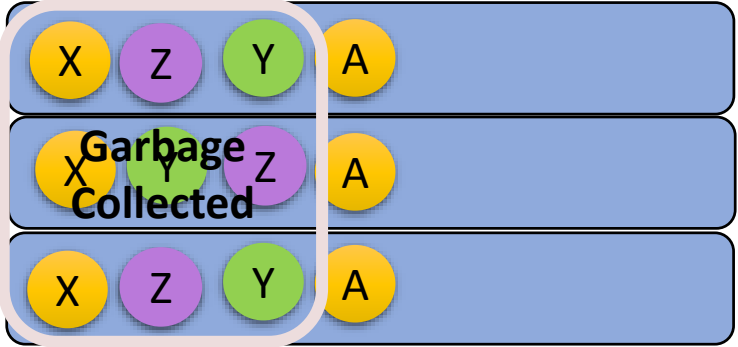
Problem-2



Problem-2

-  Data
-  Metadata

Sequencing layer



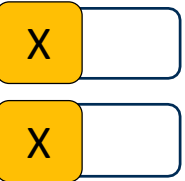
Clients

Failed

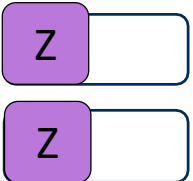


Client try sending again

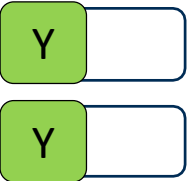
Unordered log



shard1



shard2

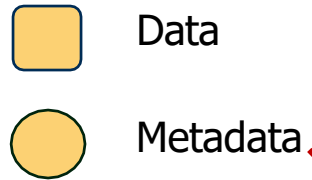


shard3

Primary

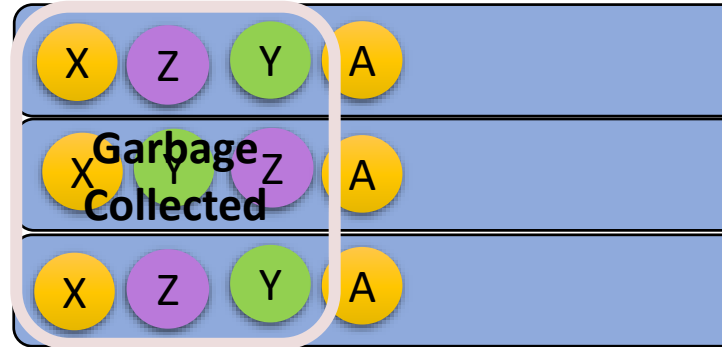
Backup

Problem-2



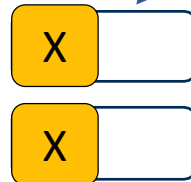
~~Clients~~

Sequencing layer

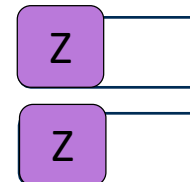


Sends position of the record

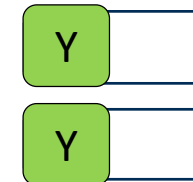
No Data with record_id



shard1



shard2



shard3



Primary

Backup

Unordered log

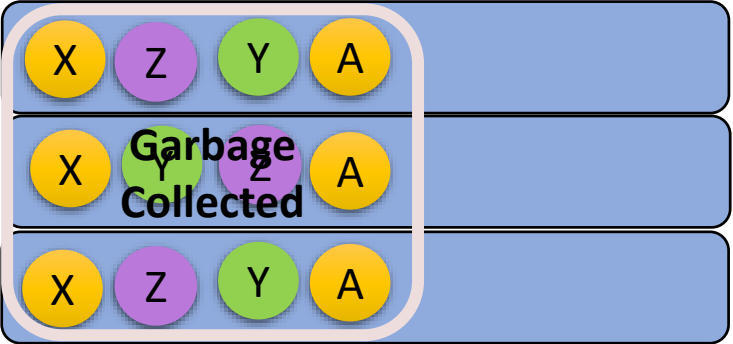


Problem-2

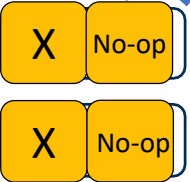
-  Data
-  Metadata

~~Clients~~

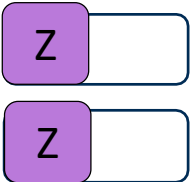
Sequencing layer



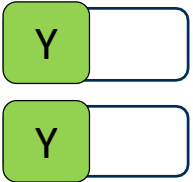
Ack



shard1



shard2





shard3

Primary

Backup

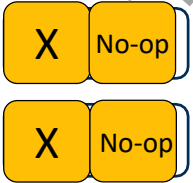
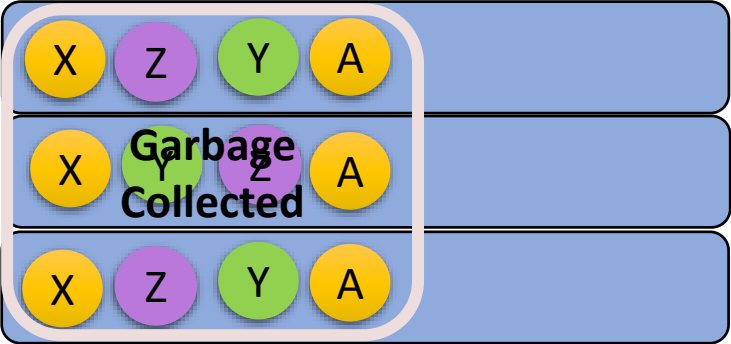
Unordered log

Problem-2

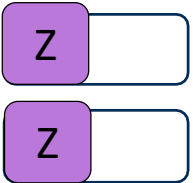
-  Data
-  Metadata

~~Clients~~

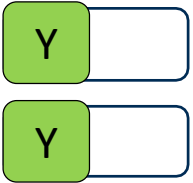
Sequencing layer



shard1



shard2

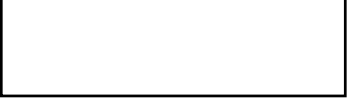


shard3

Primary

Backup

Unordered log



Performance Evaluation

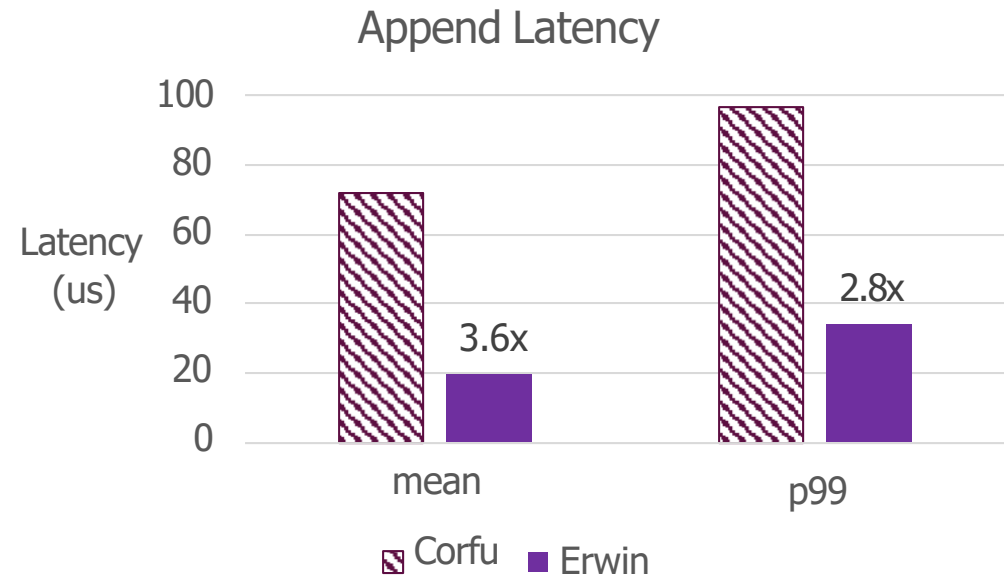
- What's the latency benefit of lazy ordering?
- How do reads perform in LazyLog?
- Do end applications benefit?

What's the Latency Benefit of Lazy Ordering?

Workload: 4KB record append-only
3 replicas per shard with 5 shards

Erwin reduces append latency

- Avg: By 3.6x compared to Corfu
- P99: By 2.8x compared to Corfu



How Do Reads Perform in LazyLog?

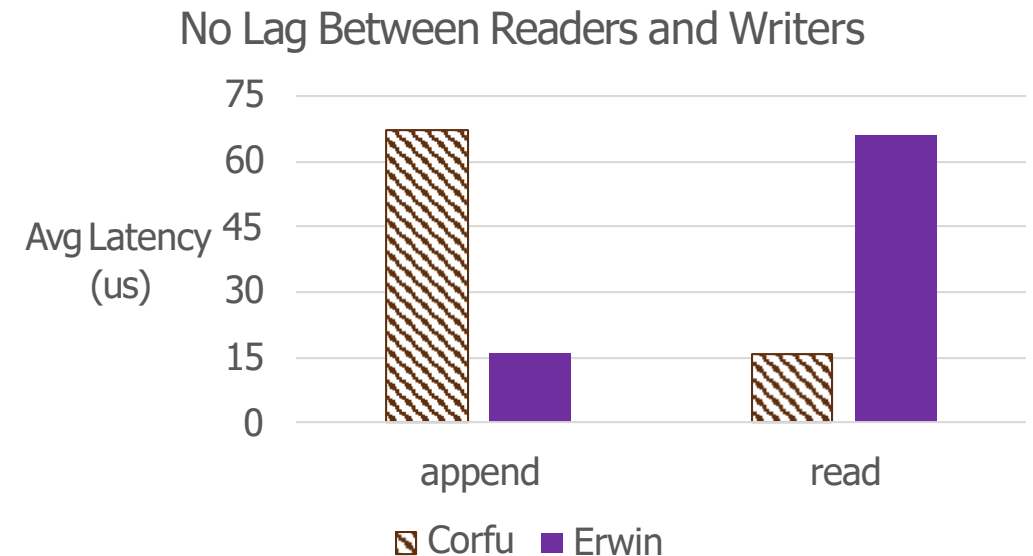
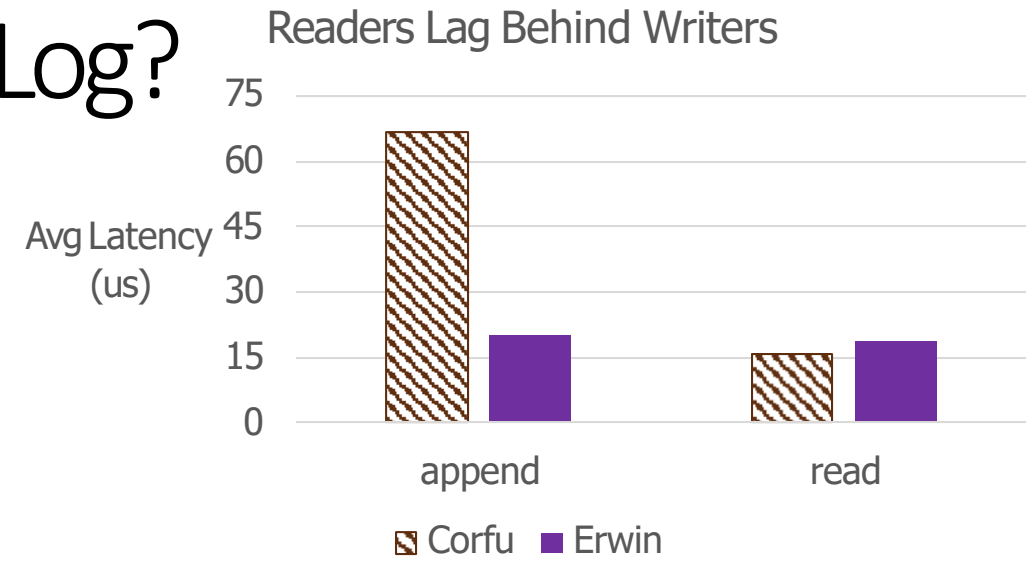
4KB record read after append

For many applications in which reads lag behind writes:

- Erwin achieves low append latency and read latency

In the worst case when there is *no* lag:

- Erwin shifts ordering cost from append to read
- Append +read latency remains the same



Do End Apps Benefit from LazyLog?

Built 3 Apps: *KV Store*, *Audit Log*, and *Journal for stream processing system*

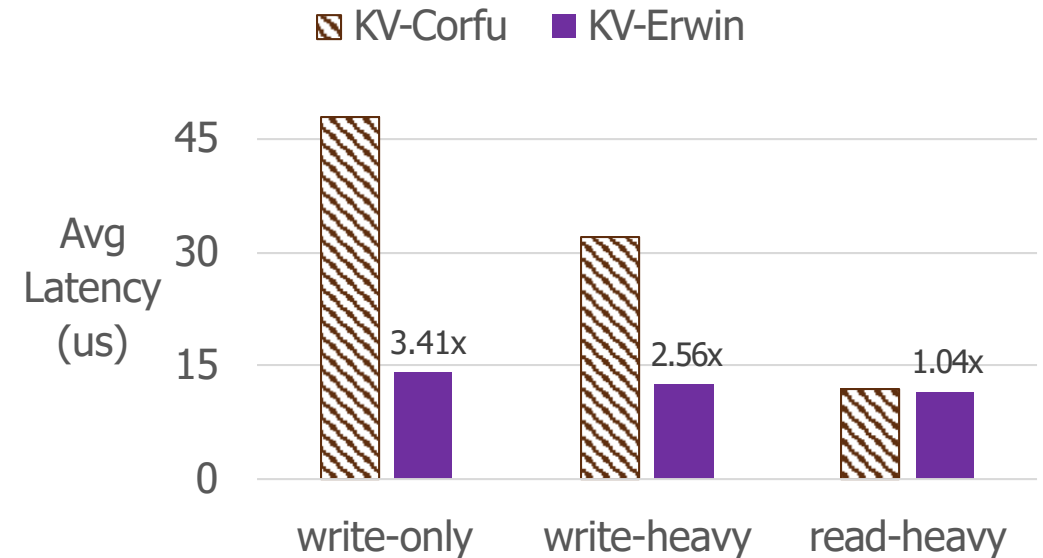
KV Store (decoupled WR-er and RD-er):

Append to log on PUTs

Reader reads log, constructs state, serves GETs

Erwin benefits applications by reducing ingestion latency

- Benefit is more pronounced when shared-log interaction takes significant partition of app request execution



Summary

- Eager-ordering shared logs incur high latencies, impacts app performance
- Eager ordering is not needed for many applications and readers are time-decoupled from writers
- LazyLog – a new shared-log abstraction that defers ordering
- Low ingestion latency with little overhead upon reads
- LazyLog systems deliver benefits for applications