# High Throughput Replication with Integrated Membership Management

Pedro Fouto, Nuno Preguiça, João Leitão

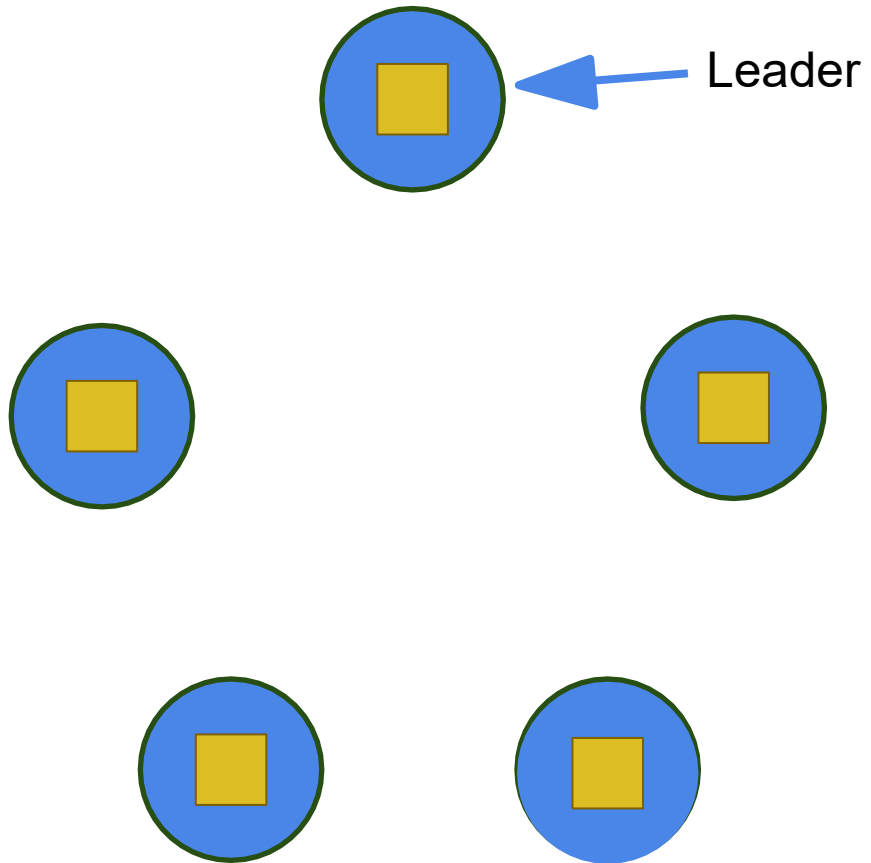Presented By: Aman, Niranjan, Hrishabh

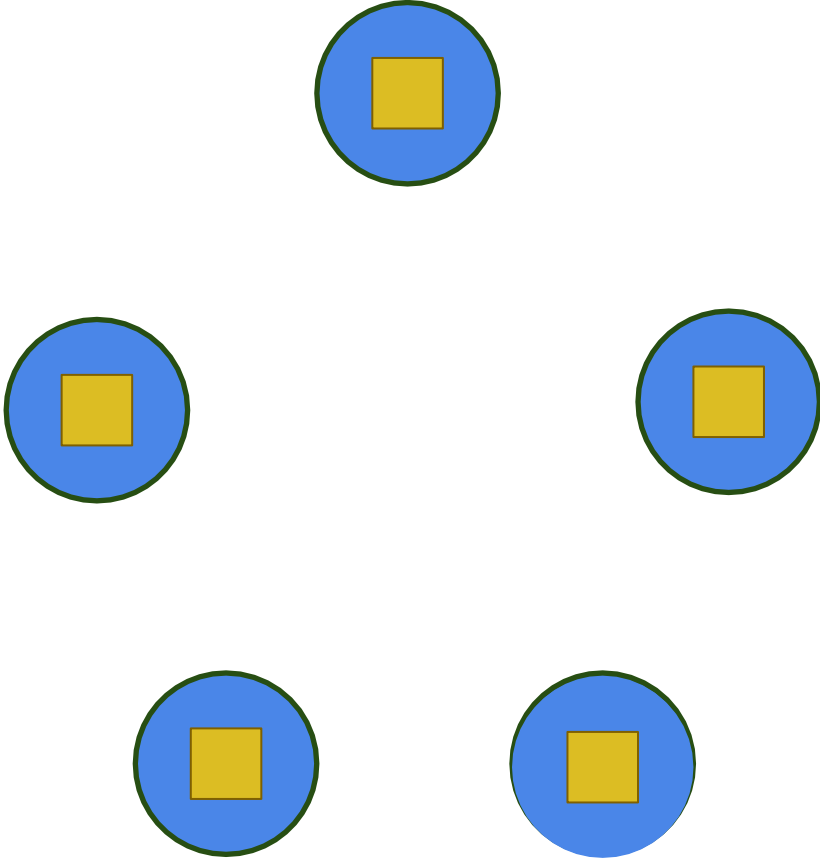17 October, 2025

# Motivation: Consensus and SMR

- Building blocks of numerous practical replication systems

- Their performance is critical

- Many alternatives have been designed

- Two very relevant ones:
    - (Multi-)Paxos
    - Chain Replication

# Motivation: Consensus and SMR

- Building blocks of numerous practical replication systems

- Their performance is critical

- Many alternatives have been designed

- Two very relevant ones:

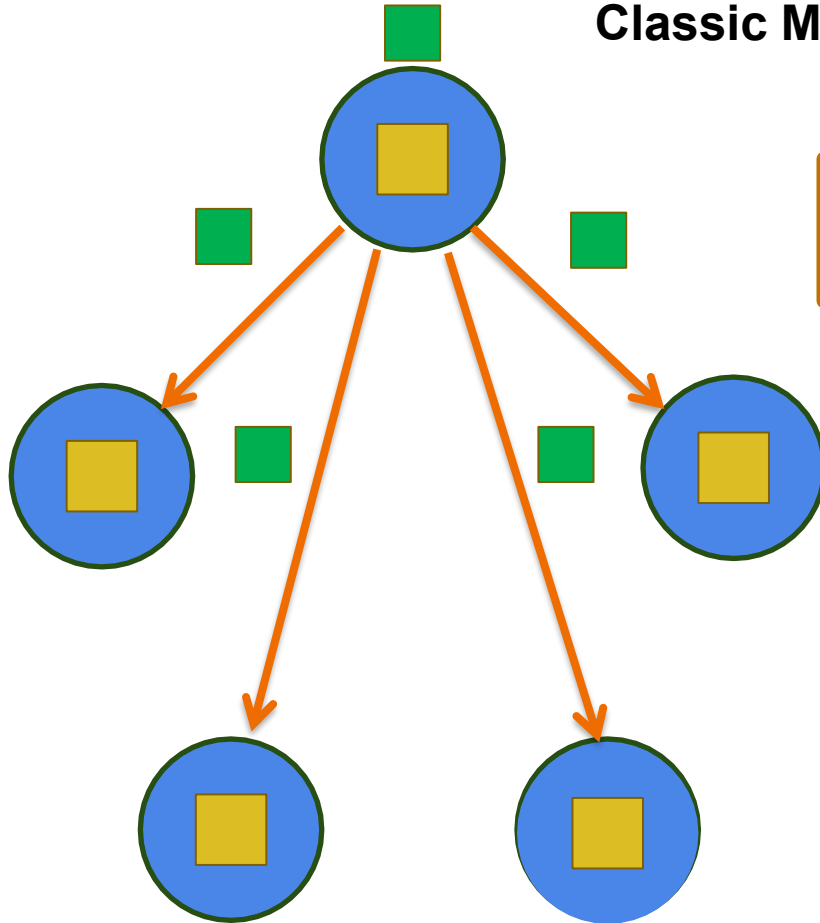  - **(Multi-)Paxos**

  - Chain Replication

# Multi-Paxos



Leader
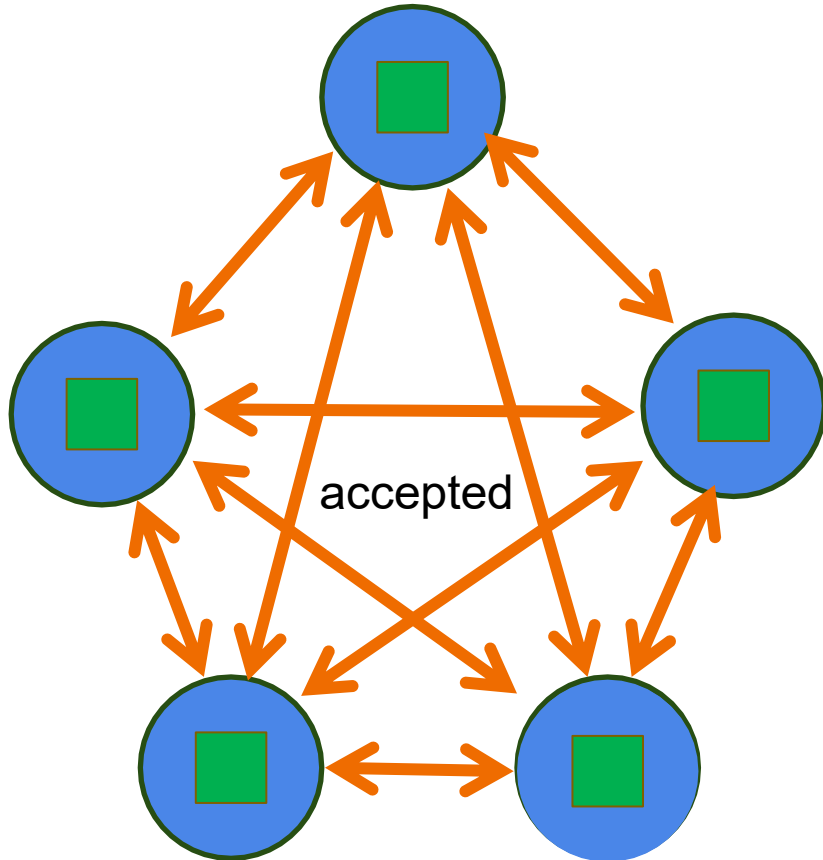
# Multi-Paxos

**Classic Multi-Paxos**

# Multi-Paxos

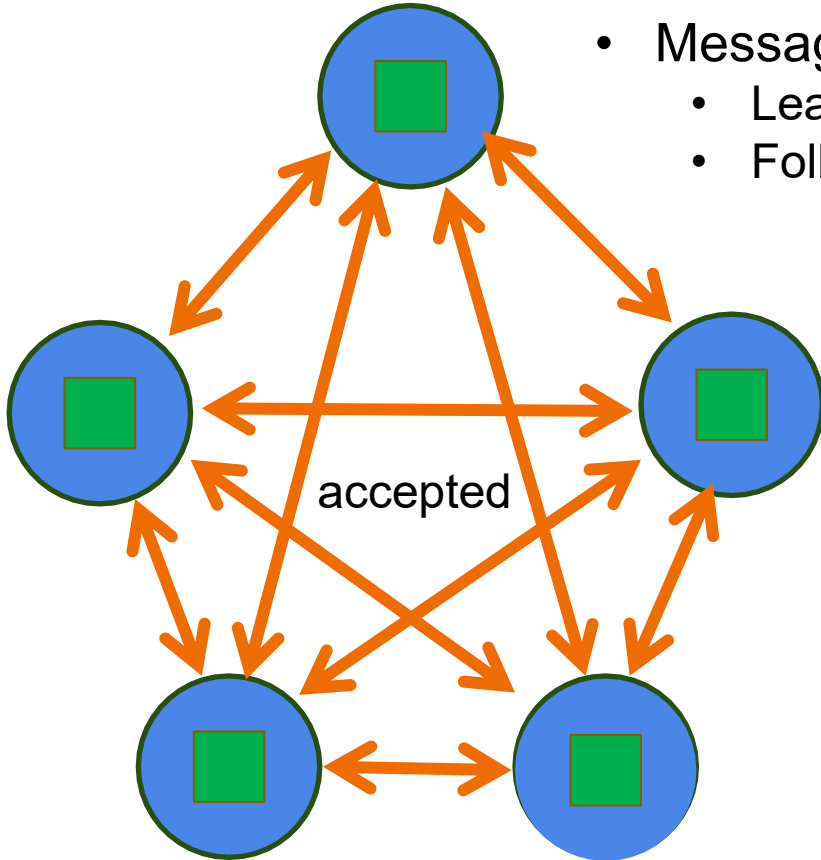**Classic Multi-Paxos**

*Skipping the first phase*

# Multi-Paxos

**Classic Multi-Paxos**



accepted

# Multi-Paxos



**Classic Multi-Paxos**

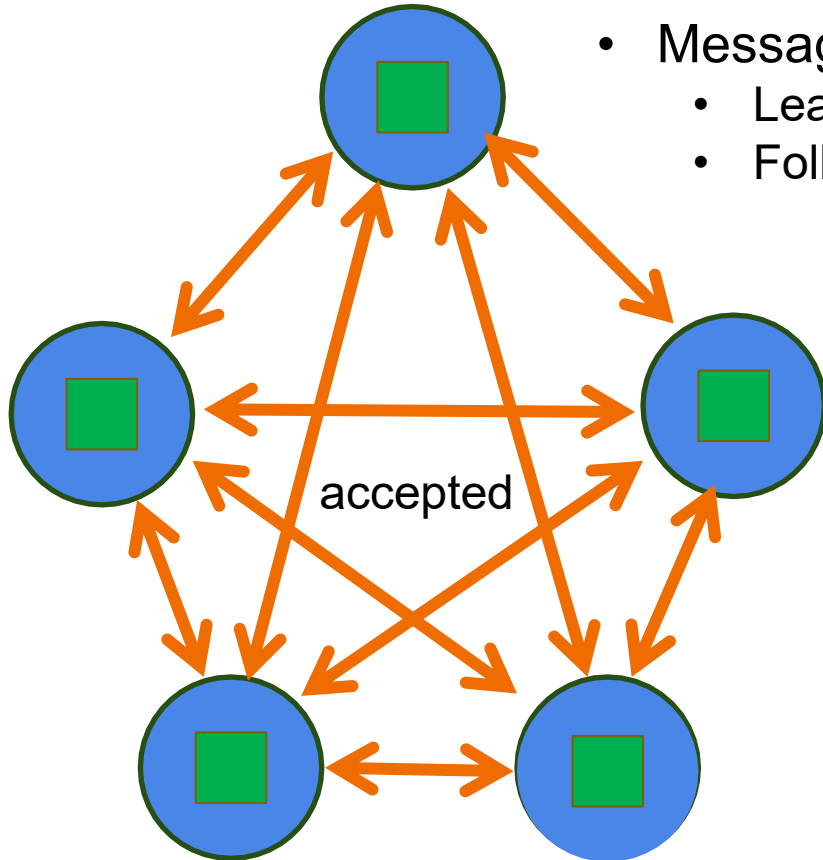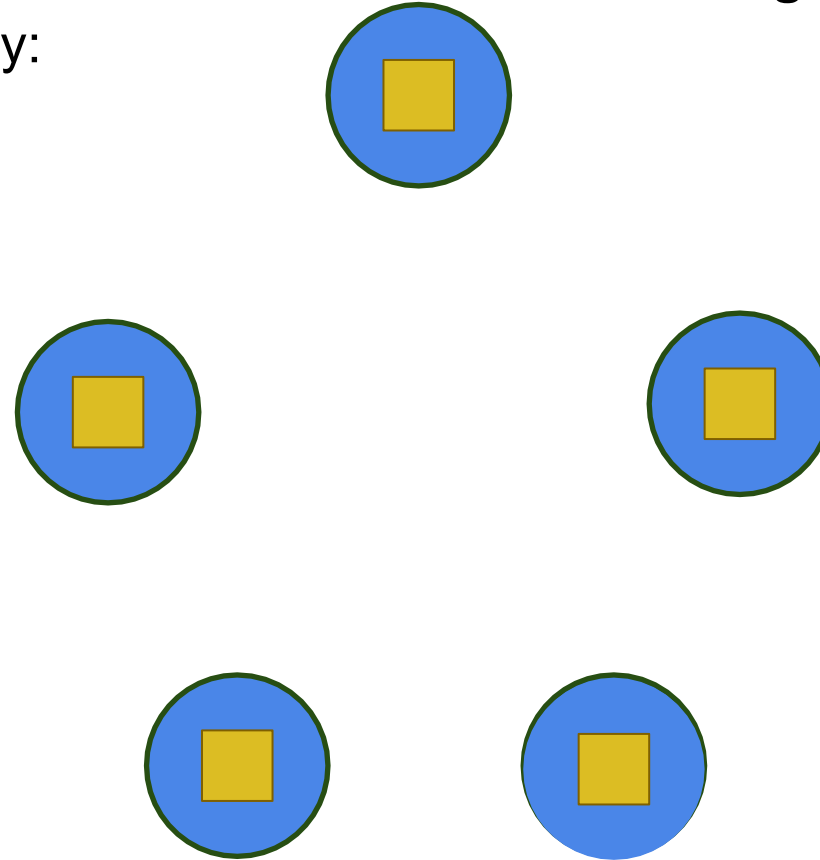- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

accepted

# Multi-Paxos

**Classic Multi-Paxos**

**Distinguished Learner**



- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

accepted

# Multi-Paxos

**Classic Multi-Paxos**

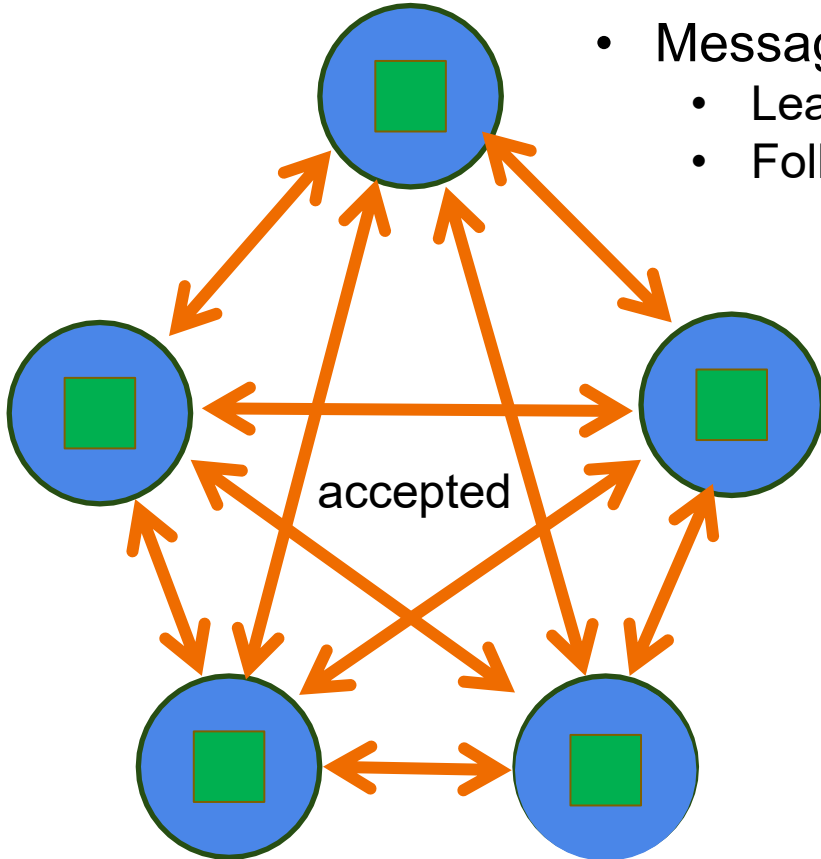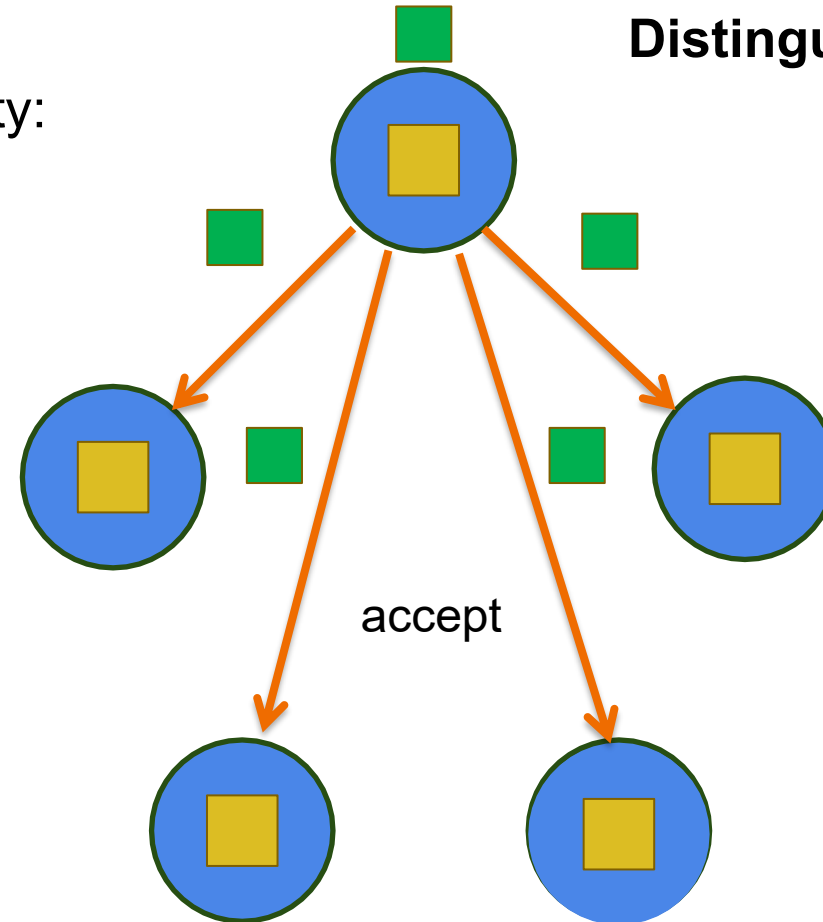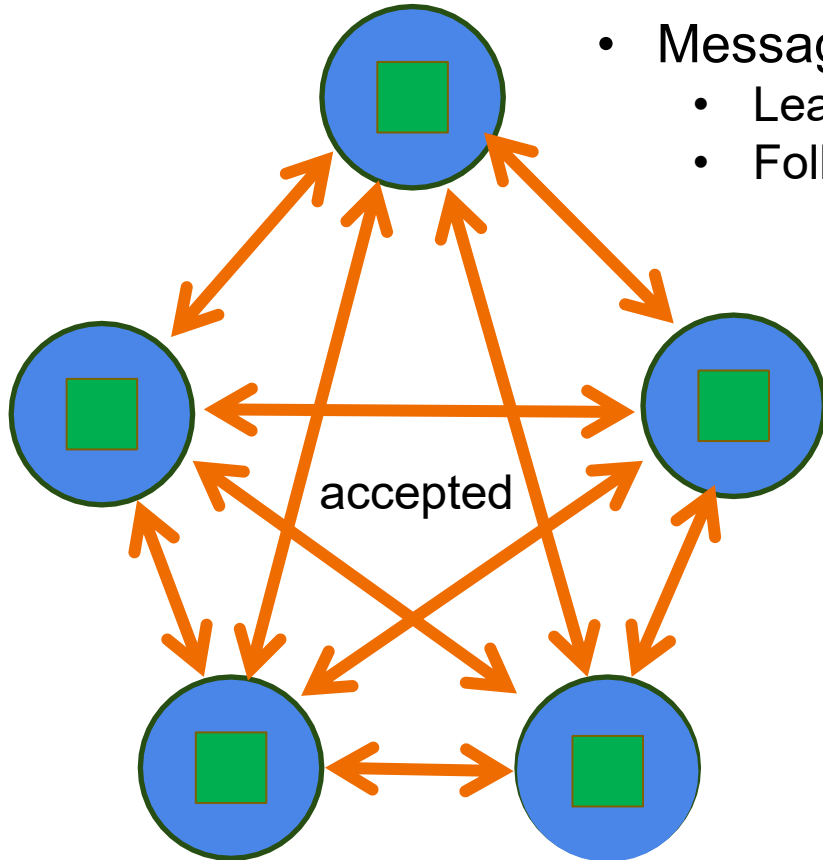- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

accepted

**Distinguished Learner**

accept

# Multi-Paxos

**Classic Multi-Paxos**

- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

accepted

**Distinguished Learner**

accepted

# Multi-Paxos

**Classic Multi-Paxos**

- Message complexity:
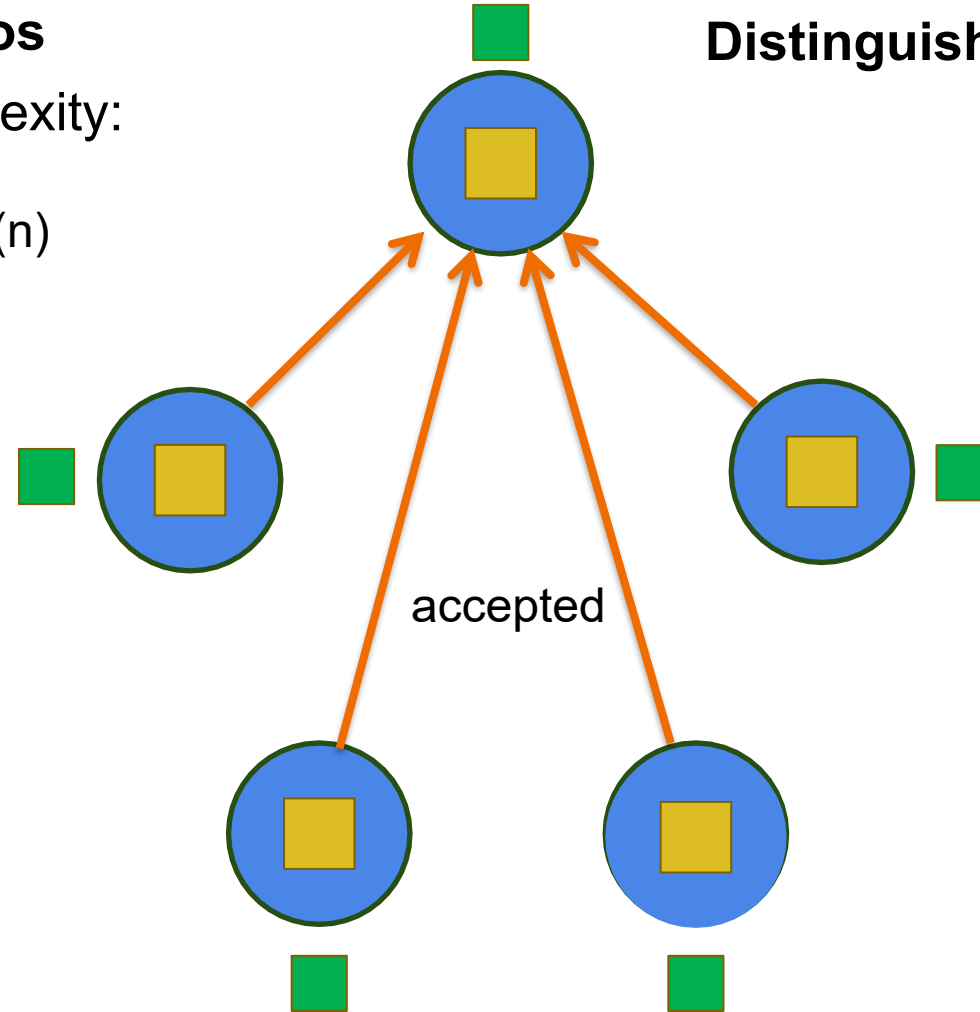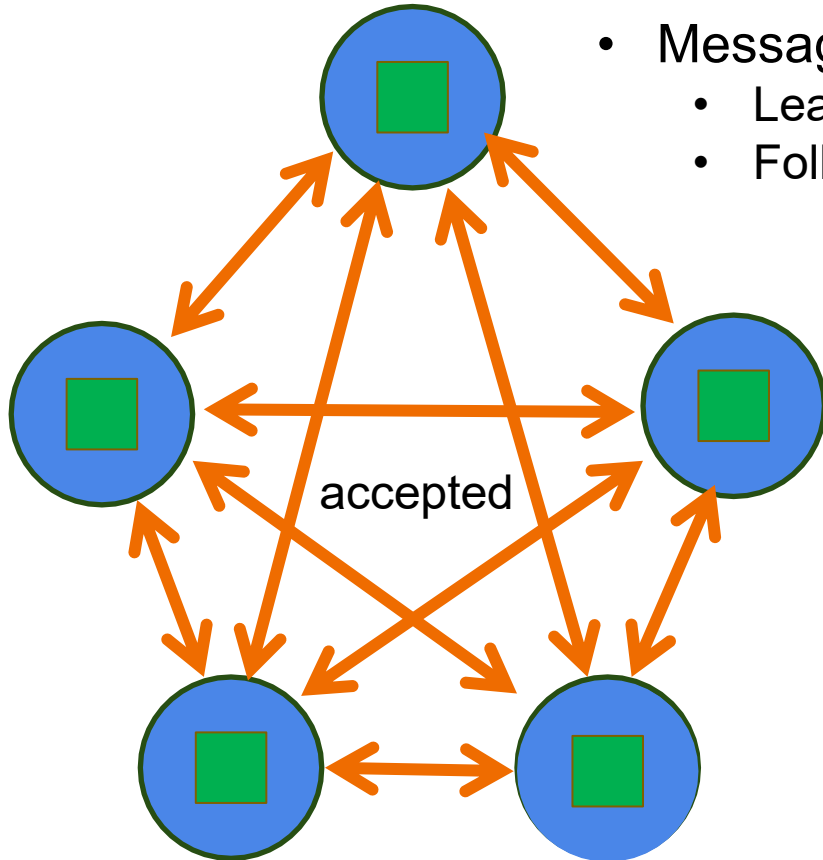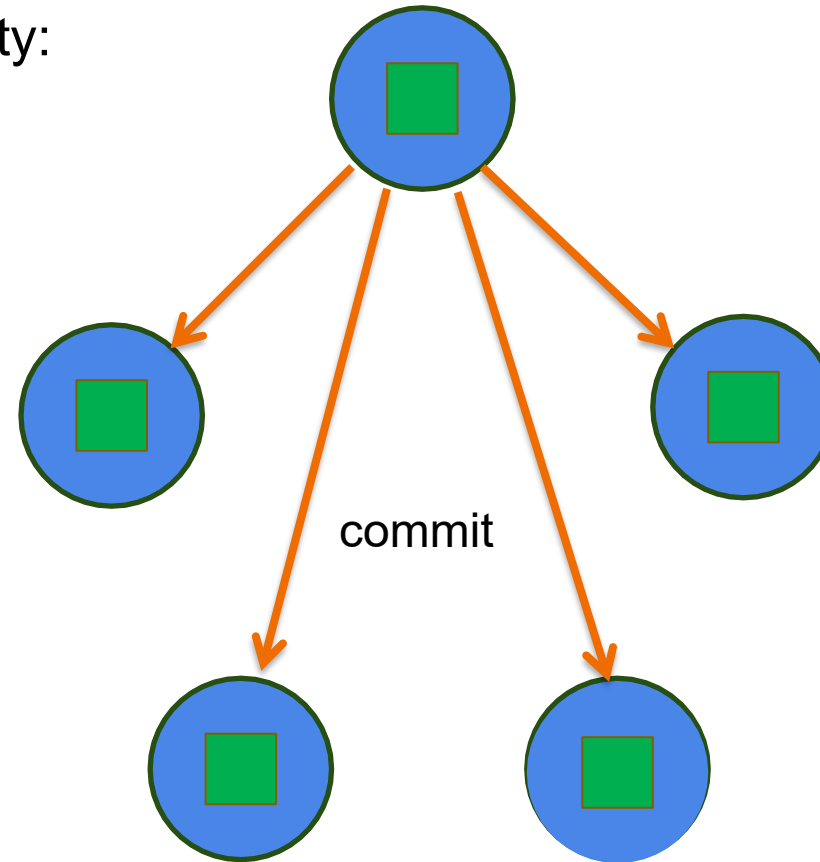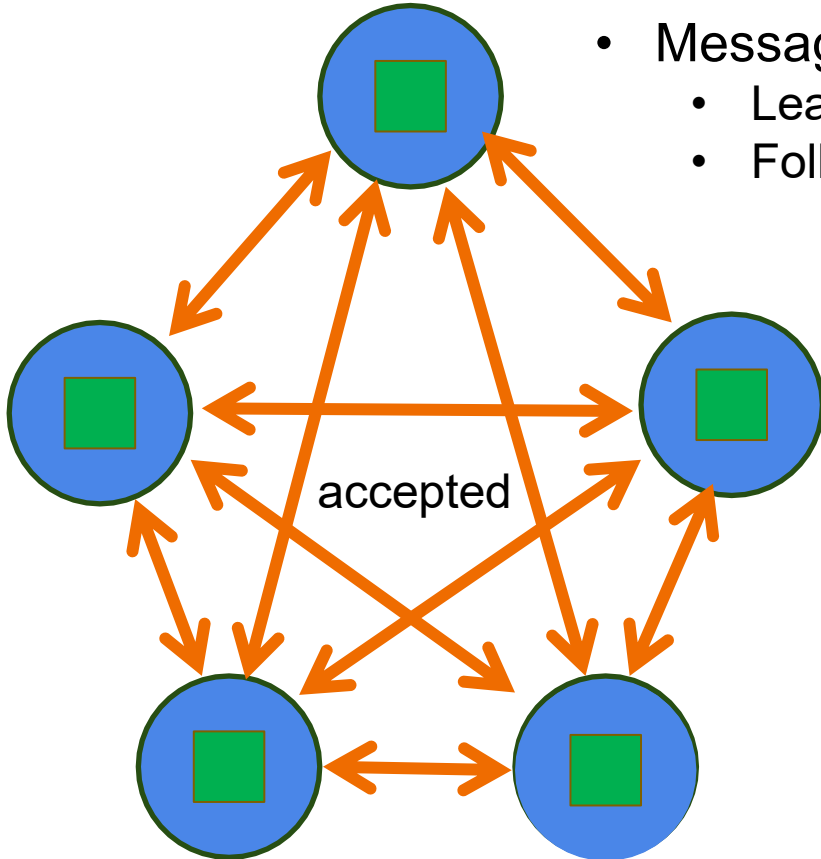  - Leader: O(n)
  - Followers: O(n)

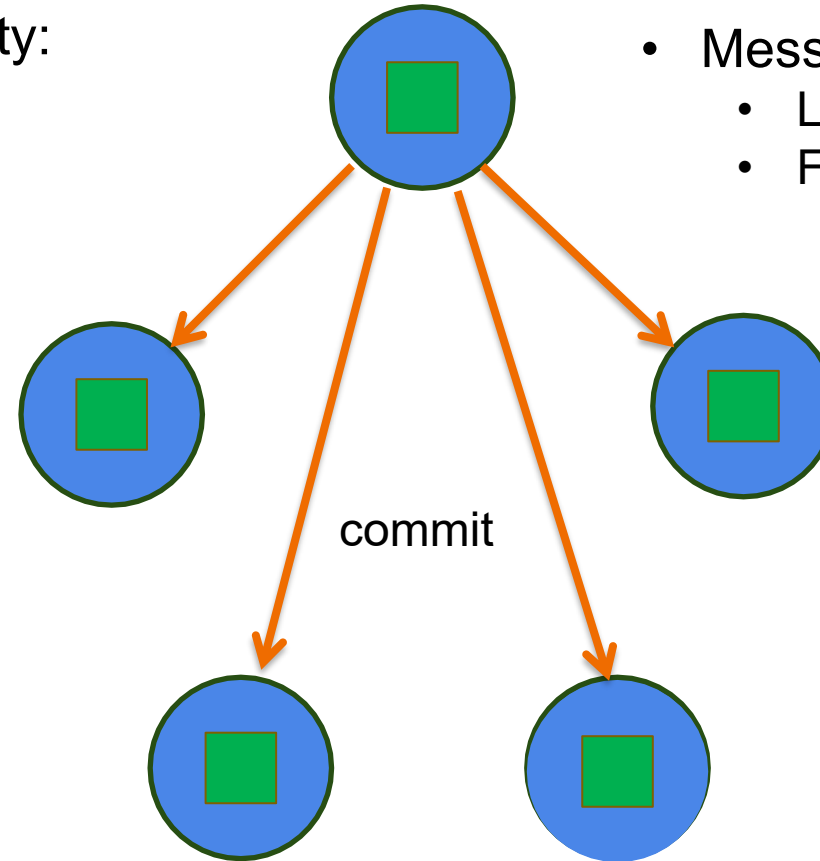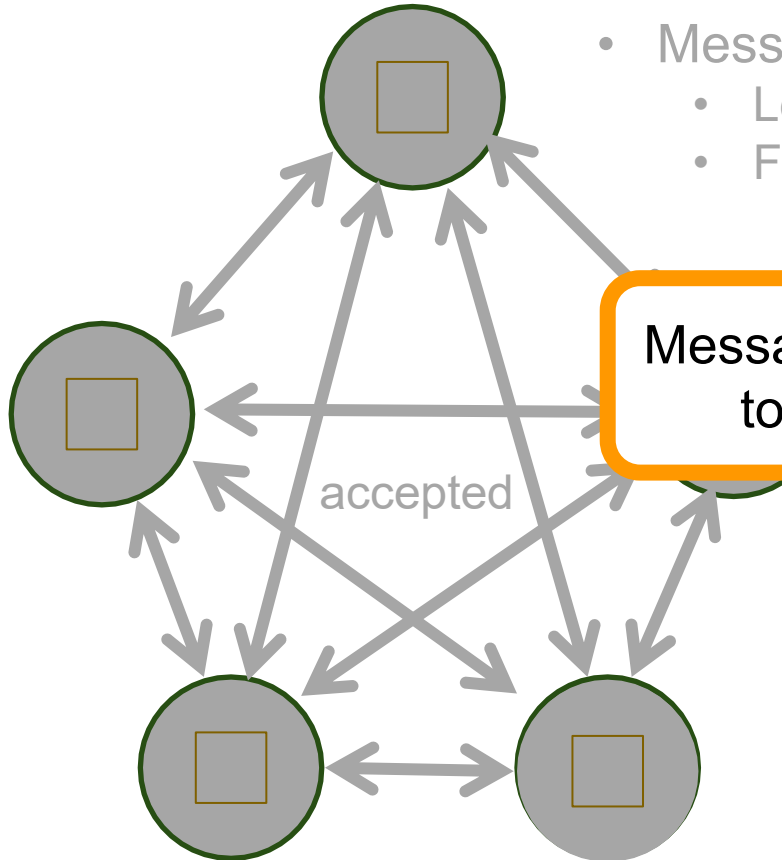**Distinguished Learner**



accepted

commit

# Multi-Paxos

**Classic Multi-Paxos**

- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

**Distinguished Learner**

- Message complexity:
  - Leader: O(n)
  - Followers: O(1)

accepted

commit
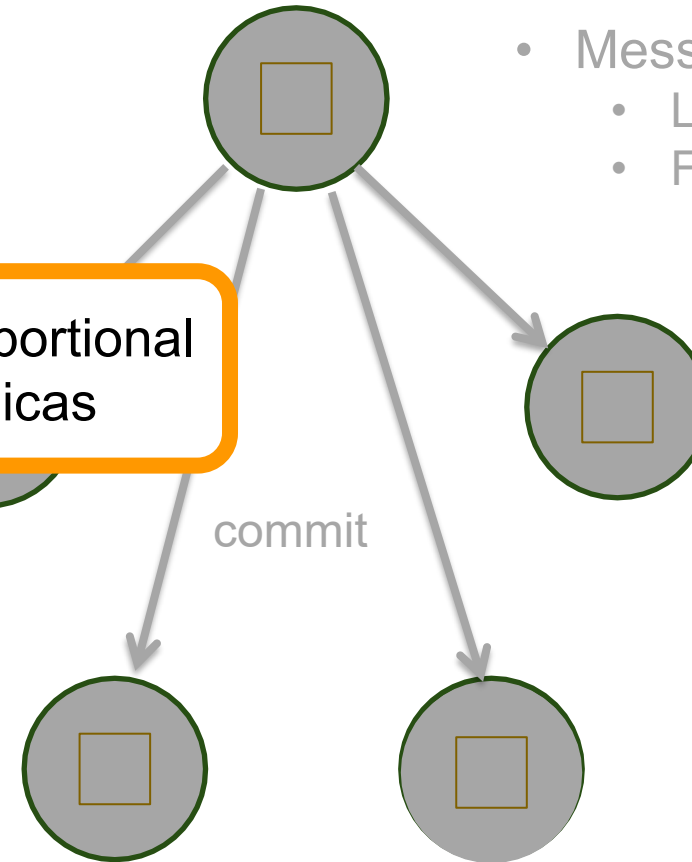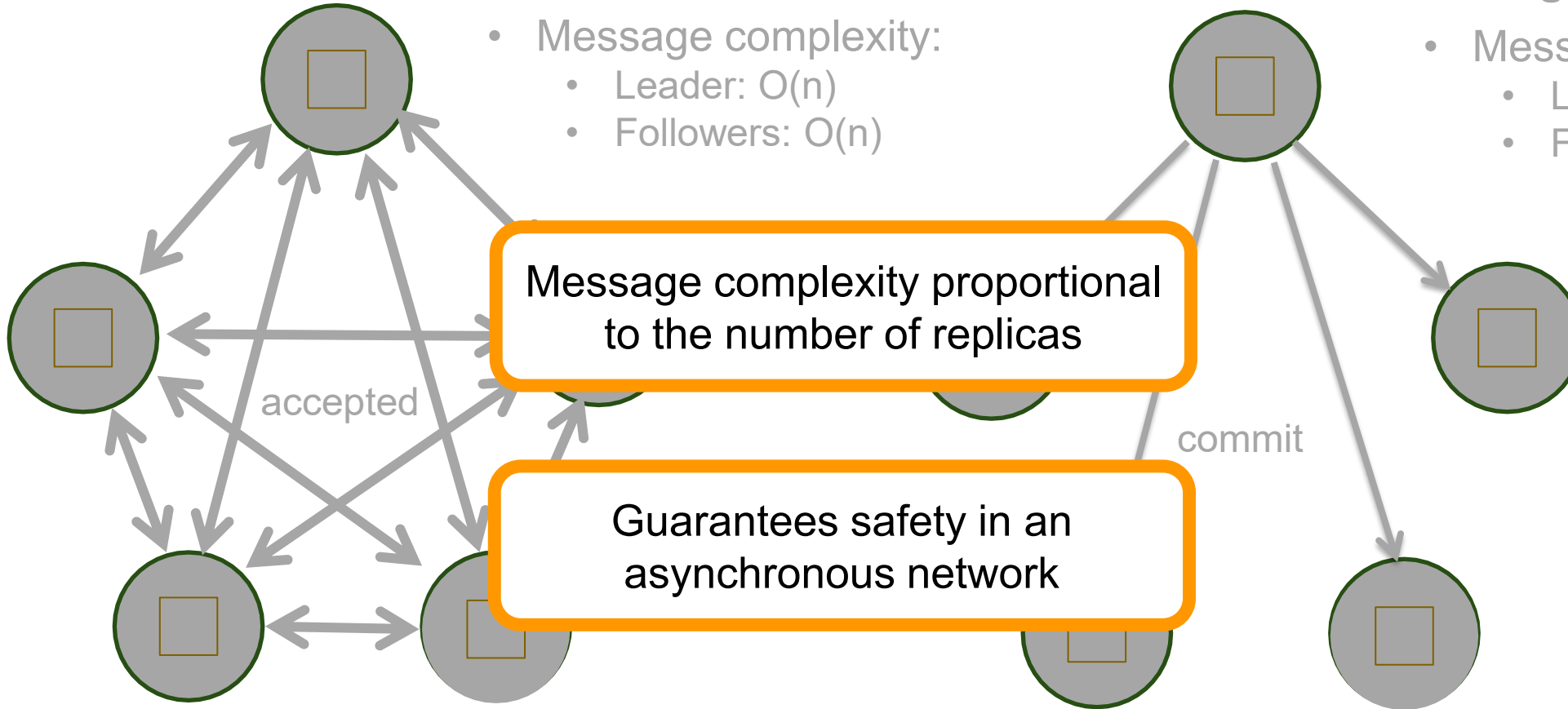
# Multi-Paxos



**Classic Multi-Paxos**

- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

**Distinguished Learner**

- Message complexity:
  - Leader: O(n)
  - Followers: O(1)

accepted

Message complexity proportional
to the number of replicas

commit

15

# Multi-Paxos



**Classic Multi-Paxos**

- Message complexity:
  - Leader: O(n)
  - Followers: O(n)

**Distinguished Learner**

- Message complexity:
  - Leader: O(n)
  - Followers: O(1)

accepted

commit

Message complexity proportional to the number of replicas
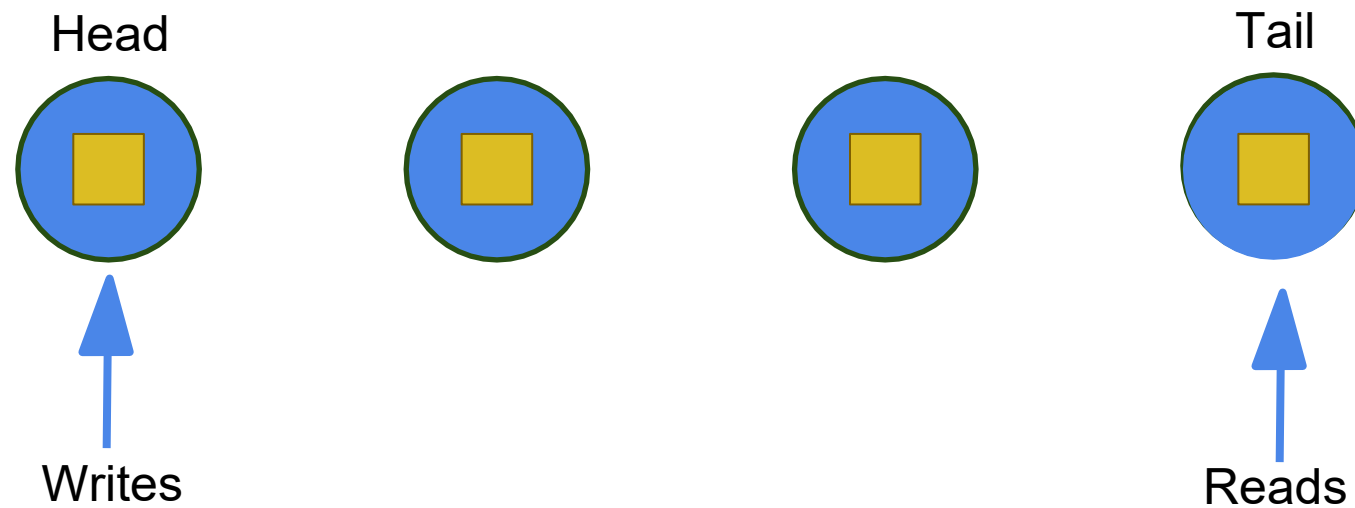
Guarantees safety in an asynchronous network

16

# Motivation: Consensus and SMR

- Building blocks of numerous practical replication systems

- Their performance is critical

- Many alternatives have been designed

- Two very relevant ones:
  - (Multi-)Paxos
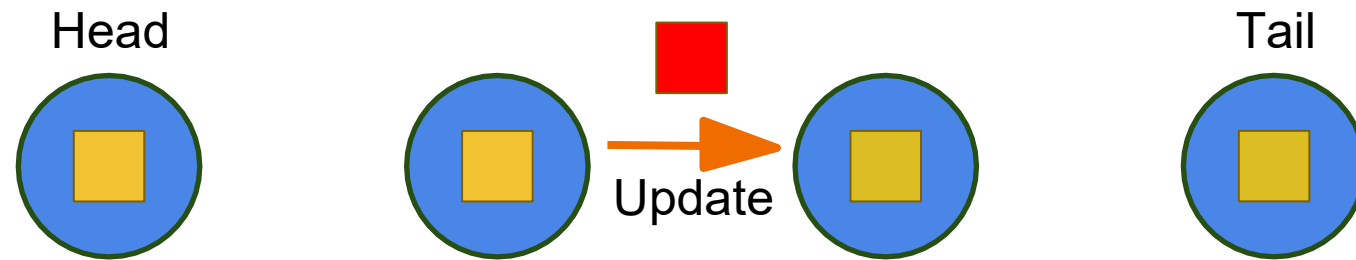  - **Chain Replication**

# Chain Replication (CR)



Head

Tail

Writes

Reads

Any object.
Eg: Key-Value Store: X = 1
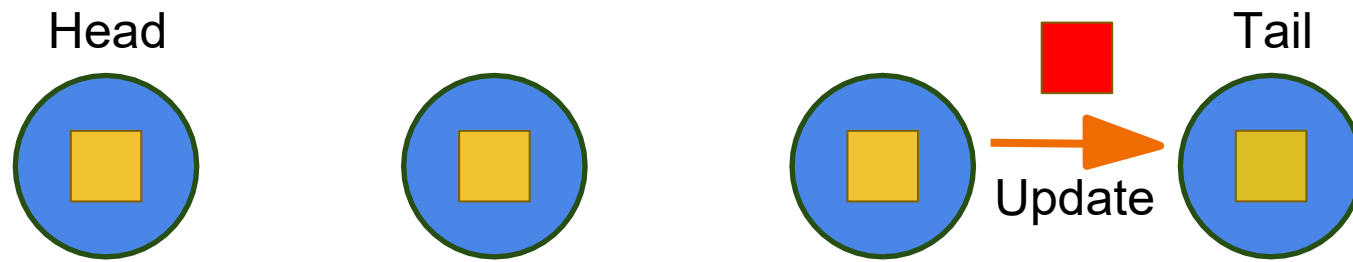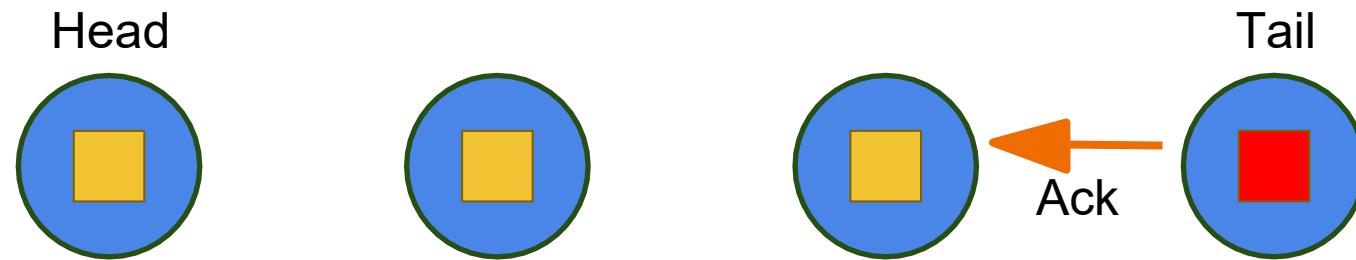
# Chain Replication (CR)



Head

Tail

Write
Request

# Chain Replication (CR)

# Chain Replication (CR)
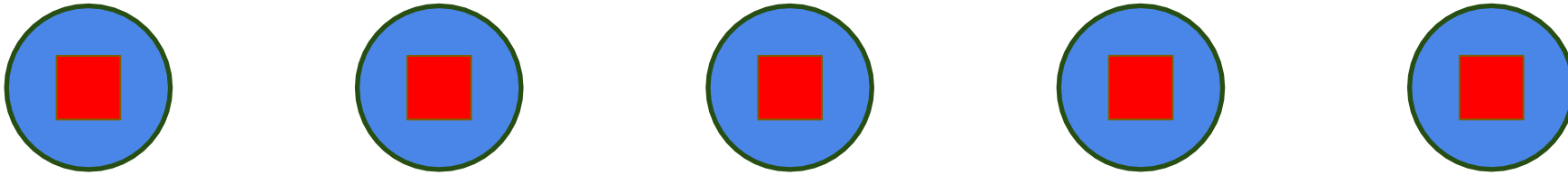


Head

Tail

Update

# Chain Replication (CR)



Head

Tail

Update

# Chain Replication (CR)



Head

Tail

Ack

# Chain Replication (CR)



Head

Tail

Ack

# Chain Replication (CR)



Head

Tail

Ack

# Chain Replication (CR)

Head

Tail

Write Ack
sent to Client

# Chain Replication (CR)

Head

Tail

- Message complexity:
  - All replicas: O(1)

# Chain Replication (CR)

Head

Tail



- Message complexity:
  - All replicas: O(1)

However, it has its limitations…

# Motivation: Consensus and SMR

- Building blocks of numerous practical replication systems

- Their performance is critical

- Many alternatives have been designed

- Two very relevant ones:
  - (Multi-)Paxos
  - Chain Replication

- Two aspects are often overlooked:
  - **Performant linearizable reads**
  - Membership management and reconfiguration

# Motivation: Linearizable Reads

- Some existing solutions assume a synchronous model (e.g. Chain Replication)

- Dealing with asynchrony is complicated. One can:
  - Relax consistency (e.g. ZooKeeper)
  - Add extra (costly) steps to write operations
  - Synchronize with other replicas when reading

# Motivation: Linearizable Reads

# Motivation: Linearizable Reads
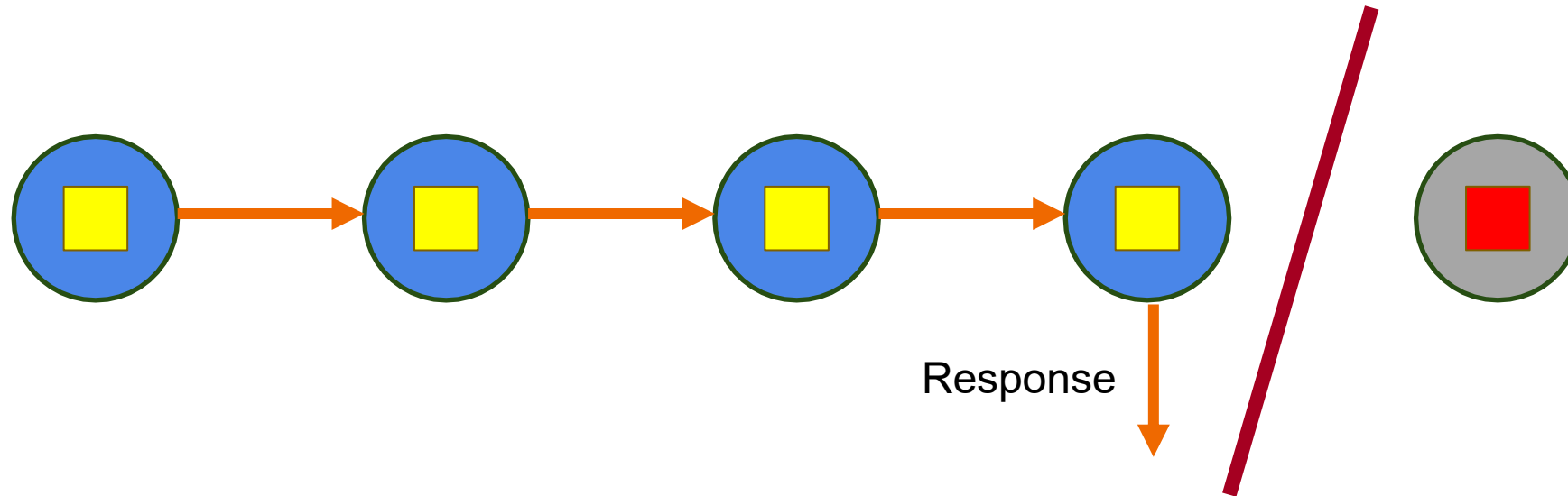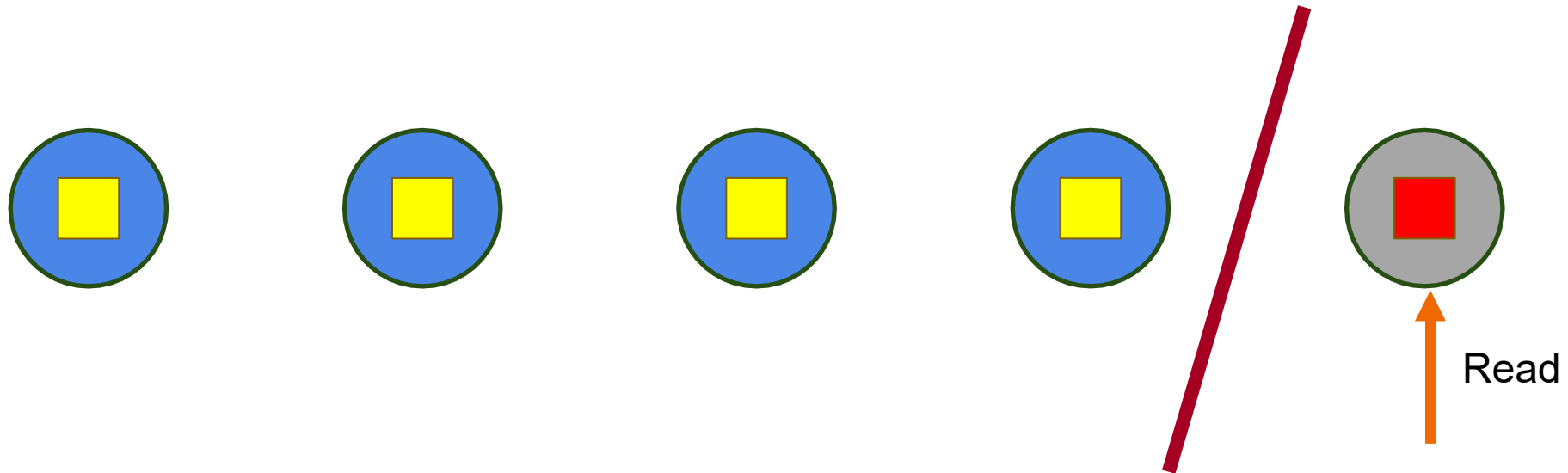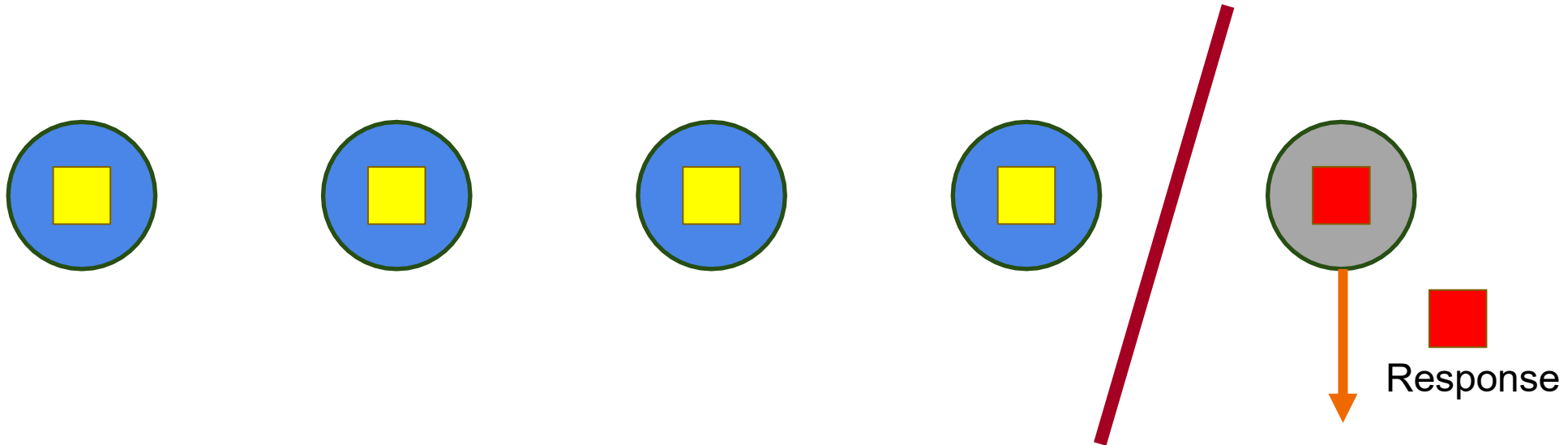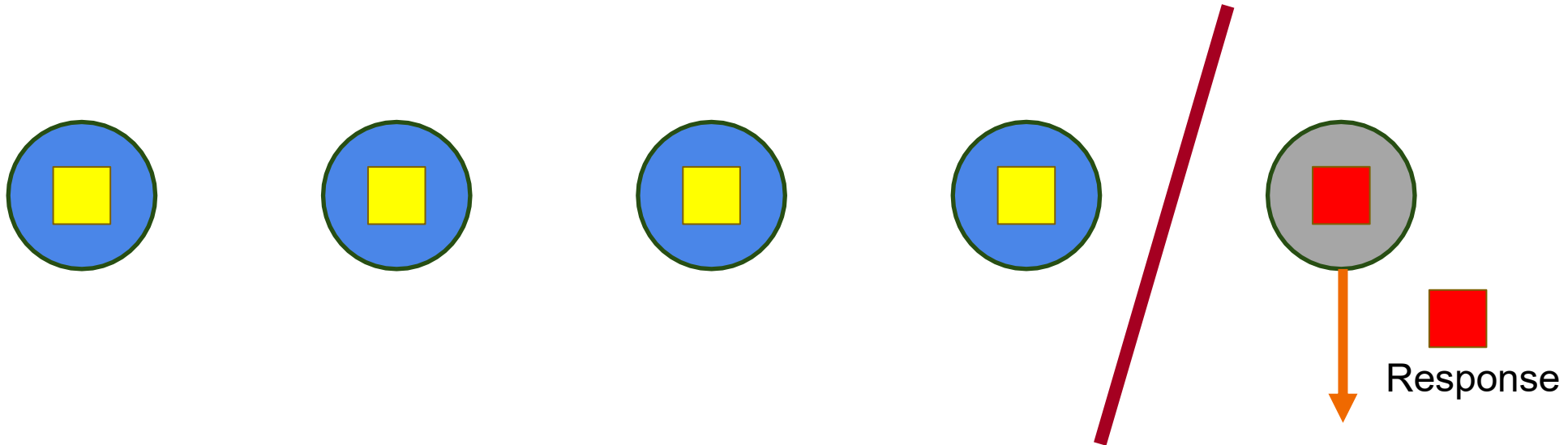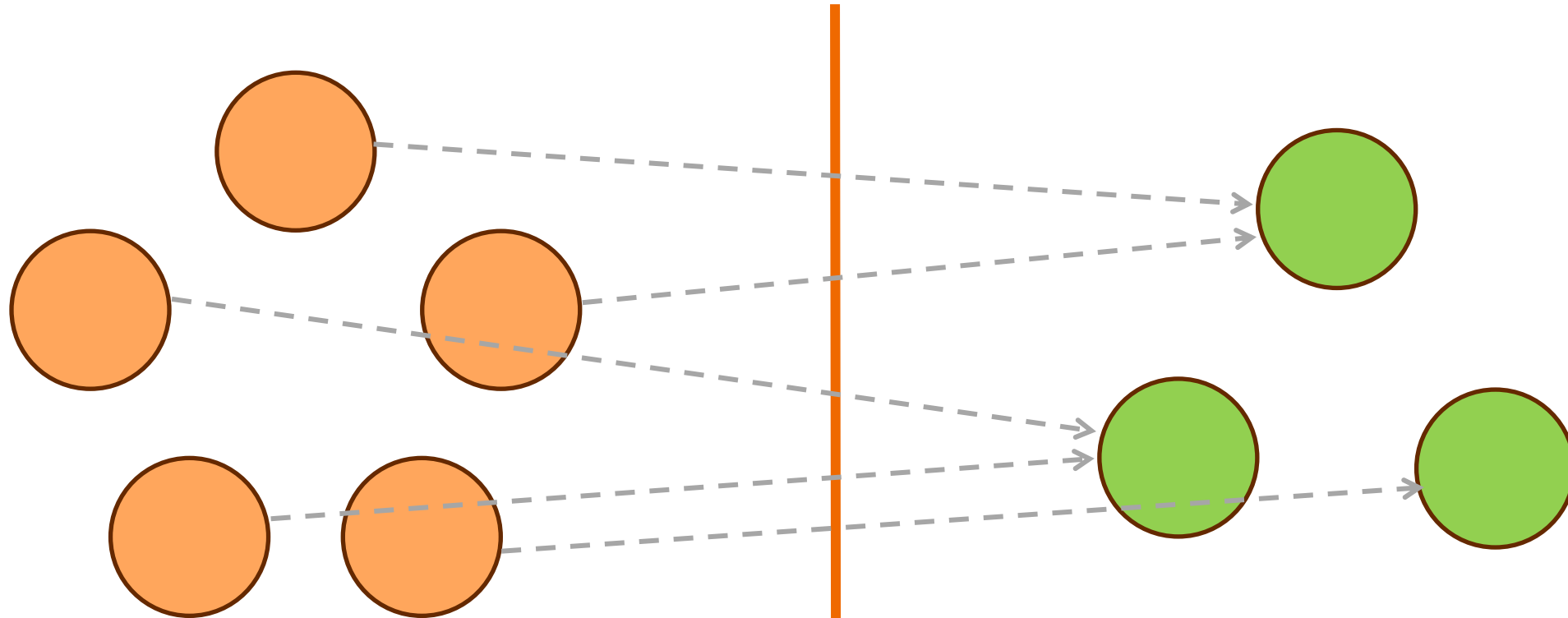
Read

# Motivation: Linearizable Reads



Response

# Motivation: Linearizable Reads



Network Partition

# Motivation: Linearizable Reads



New Tail

# Motivation: Linearizable Reads

# Motivation: Linearizable Reads



Response

# Motivation: Linearizable Reads



Read

# Motivation: Linearizable Reads



Response

# Motivation: Linearizable Reads

Response

**Breaks linearizability in an asynchronous network model**

# Motivation: Linearizable Reads

- Some existing solutions assume a synchronous model (e.g. Chain Replication)

- Dealing with asynchrony is complicated. One can:
  - **Relax consistency** (e.g. ZooKeeper)
  - Add **extra (costly) steps** to write operations
  - **Synchronize** with other replicas when reading

# Motivation: Consensus and SMR

- Building blocks of numerous practical replication systems

- Their performance is critical

- Many alternatives have been designed

- Two very relevant ones:
  - (Multi-)Paxos
  - Chain Replication

- Two aspects are often overlooked:
  - Performant linearizable reads
  - **Membership management and reconfiguration**

# Motivation : Membership Management

- Most consensus solutions overlook membership management

- Often assume that using an external coordinator service (e.g. ZooKeeper) is trivial and the best solution

- This is not the case:

  ▪ Fault-tolerance becomes complex

  ▪ Complex (and redundant) integration with consensus

  ▪ More vulnerable to partial network partitions

# Motivation : Membership Management



- Your consensus solution:
  - Fault-tolerance: 2

- External Coordinator (e.g. ZooKeeper)
  - Fault-tolerance: 1

# Motivation : Membership Management



- Your consensus solution:
  - Fault-tolerance: 2

- External Coordinator (e.g. ZooKeeper)
  - Fault-tolerance: 2

# Motivation : Membership Management

Replicas: 10

Fault-tolerance: 2

- Your consensus solution:
  - Fault-tolerance: 2
- External Coordinator (e.g. ZooKeeper)
  - Fault-tolerance: 2

# Motivation : Membership Management

# Motivation : Membership Management



Consensus round to decide removal

# Motivation : Membership Management



Asynchronous design propagation to consensus replicas

# Motivation : Membership Management



Replicas cannot change view individually!

# Motivation : Membership Management



*Another redundant consensus round is required*

# Motivation : Membership Management



Partition between coordinator and consensus replicas

# Motivation : Membership Management



Correct replicas will be removed

# Motivation : Membership Management



Partition between consensus replicas

# Motivation : Membership Management



Leader

System will halt until partition is healed

55

# Proposal : ChainPaxos

**Novel consensus algorithm:**

- Combining the best properties of Multi-Paxos and Chain Replication

  - Correction in an asynchronous network

  - Constant message complexity

- Going beyond existing solutions:

  - Maximizing throughput of both read and write operations

  - Providing local linearizable reads in any replica

  - Integrated reconfiguration and fault-tolerance

# State of ChainPaxos Nodes

**1**

chain : *array of nodes*
self : *node*
$c_{nextok}$ : *node*
$c_{sleader}$ : *node*
marked : *set of node*

**2**

$np_{leader}$ : *int*
inst : *map int × PaxosInst*

**3**

submitted : *set of requests*
pending : *set of requests*

**4**

$max_{ack}$ : *int*
$max_{acpt}$ : *int*
amLeader : *bool*

# ChainPaxos : Write Path

Leader (regular Multi-Paxos election)

# ChainPaxos : Write Path



: New value to be written

# ChainPaxos : Write Path

Leader

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# ChainPaxos : Write Path

Encapsulate multiple

Multi-Paxos messages

accept
+
accept ok

63

# ChainPaxos : Write Path

accept
+
accept ok ×2

# ChainPaxos : Write Path

Operation is "locked-in"

Reply

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# ChainPaxos : Write Path



Need to garbage collect + execute on the first replicas

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# ChainPaxos : Write Path

Ack is "piggybacked" along

with next write*

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# ChainPaxos : Write Path

# Local Linearizable Reads

**Requirements for Linearizability:**

- The result of a read must contain **all writes that completed** before it started

- The result of a read must contain the result of **all reads that completed** before it started

**Challenge:**

- Read from **any replica**

- **No extra communication** steps

# Local Linearizable Reads

# Local Linearizable Reads



Read
Operation

# Local Linearizable Reads



Read Operation

Response must:
- Contain all completed writes
- Contain all completed reads

79

# Local Linearizable Reads



Read
Operation

Response must:
- Contain all completed writes
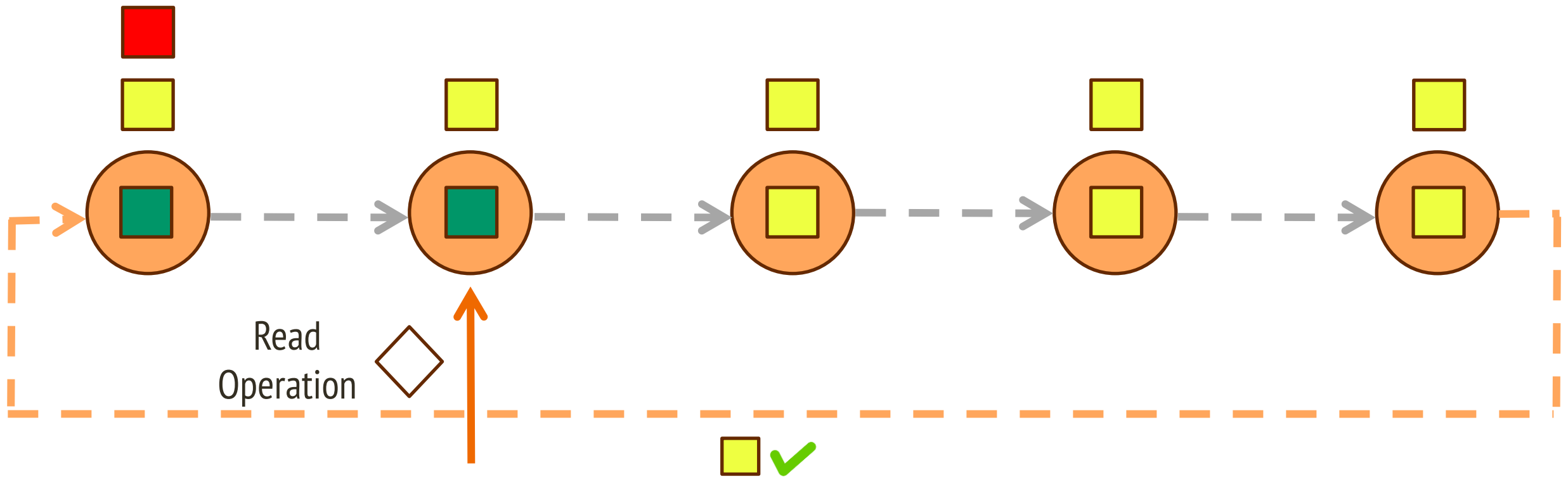- Contain all completed reads

# Local Linearizable Reads

Response must:
- Contain all completed writes (🟩)
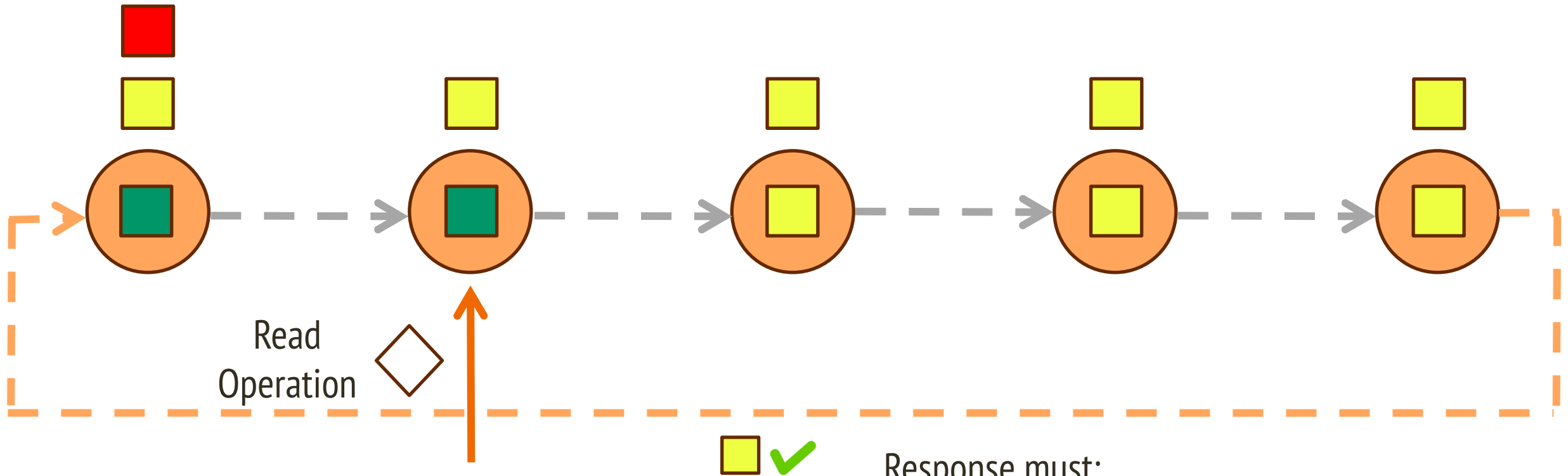- Contain all completed reads

# Local Linearizable Reads



Response must:
- Contain all completed writes (🟩)
- Contain all completed reads

# Local Linearizable Reads

Attach read to next
unseen operation (green)



: Waiting Read

# Local Linearizable Reads

# Local Linearizable Reads

Read is now attached
to green



Wait until green ack goes

around chain

# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads
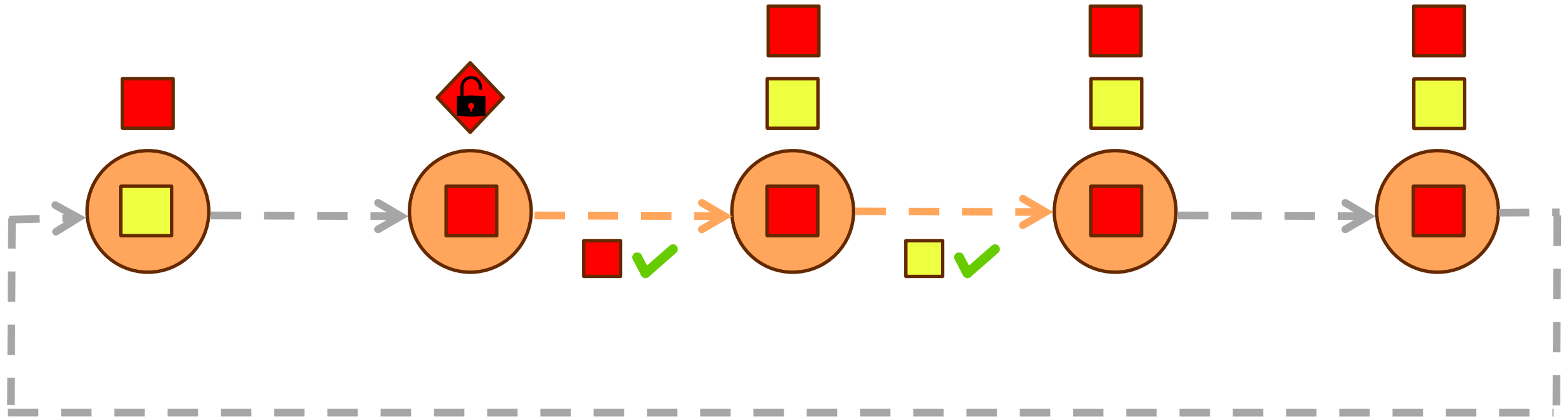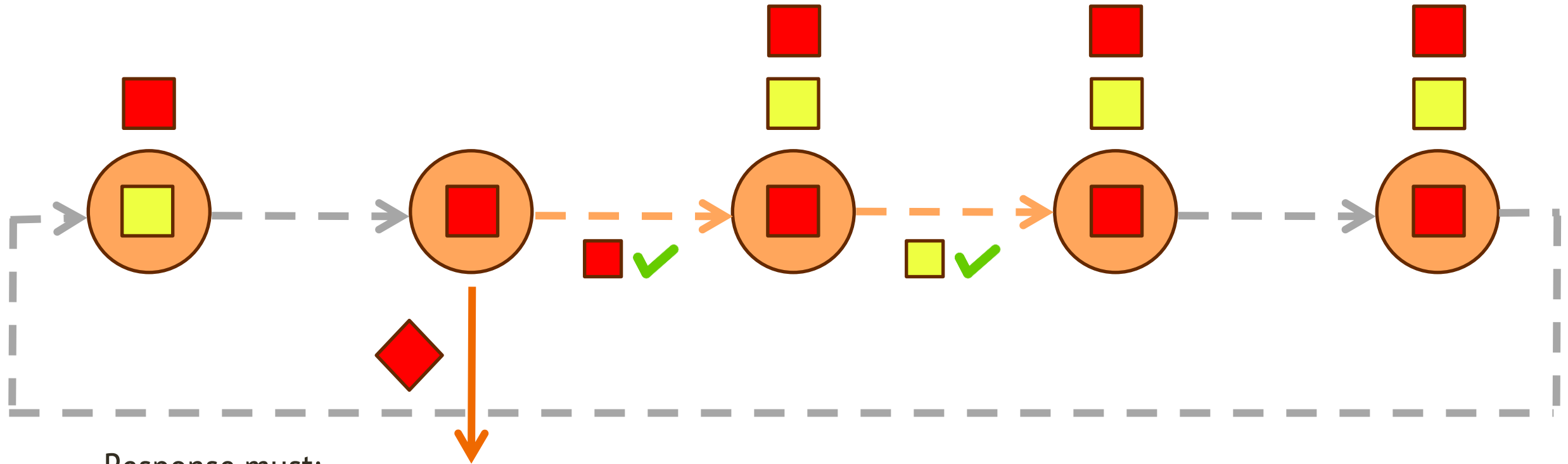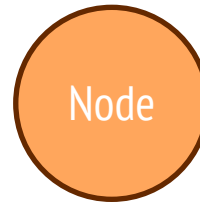
# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads

Acknowledge of green
will unlock the read

# Local Linearizable Reads



Response must:
- Contain all completed writes (■)
- Contain all completed reads

Respond with current state

# Local Linearizable Reads

Stronger than required ( ▢ )

Best we can do without extra communication

Respond with current state

Response must:
- Contain all completed writes ( ▢ )
- Contain all completed reads

93

# Local Linearizable Reads



: New value to be written

# Local Linearizable Reads

# Local Linearizable Reads



Read Operation ◇

✓

# Local Linearizable Reads



Read
Operation ◇

Response must:
- Contain all completed writes
- Contain all completed reads ( □ )

# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads

Red operation will "push" every

pending operation and ack

# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads

# Local Linearizable Reads



Response must:
- Contain all completed writes
- Contain all completed reads (☐)

105

# Local Linearizable Reads



Response must:

- Contain all completed writes
- Contain all completed reads ( ⬜ )

# Local Linearizable Reads

Response must:
- Contain all completed writes
- Contain all completed reads ( □ )

# Local Linearizable Reads: Summary

- Local read operations do not break linearizability

  o Ensures that all previously **completed reads and writes** are visible


- **No additional communication** steps are required

  o More **conservative** than required, but **unavoidable without coordination**


- **Periodic No-Op** messages ensure that reads are replied to during low write workloads

# In Depth: The Log Representation

Log contains the ordered list of all operations seen

Node

| 1 | 2 | 3 | |

# In Depth: The Log Representation

Log contains the ordered list of all operations seen

Node

Each Node/Replica has its own view of the log

| 1 | 2 | 3 | |

# In Depth: The Log Representation

Log contains the ordered

list of all operations seen

Node

Each Node/Replica has its

own view of the log

1 2 3

Each Instance writes to a

new log index

# In Depth: The Log Representation

# In Depth: The Log Representation

maxInst keeps track of

replica's max seen instance

Node

| 1 | 2 | 3 | |

maxInst

maxAck

# In Depth: The Log Representation

maxInst keeps track of
replica's max seen instance

Node

| 1 | 2 | 3 | |

maxInst

maxAck

maxInst is essentially the
length of the replica's log

# In Depth: The Log Representation

maxInst keeps track of replica's max seen instance

maxInst is essentially the length of the replica's log

Node

| 1 | 2 | 3 | |

maxInst

maxAck

maxAck keeps track of replica's max seen acknowledgment

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: The Log Representation

W1

1

maxInst

maxAck

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: The Log Representation

# In Depth: Writes



W1

maxInst maxAck
maxInst maxAck
maxInst maxAck
maxInst maxAck
maxInst maxAck

125

# In Depth: Writes

# In Depth: Writes



127

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes

# In Depth: Writes
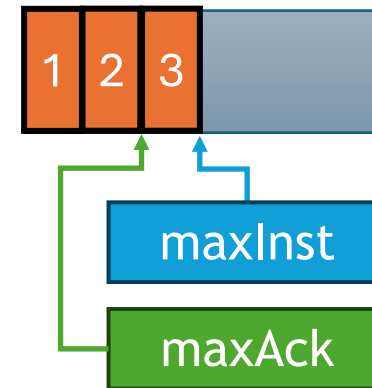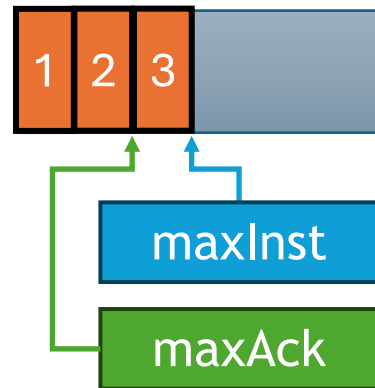
# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

138

# In Depth: Reads

# In Depth: Reads

Next *"unseen instance"* - W2

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

# In Depth: Reads

*"**ack**" of next "**unseen instance**" -* W2 ✓

W3

W3    W3    W3    W3    🔒

W2    W2    W3    W3    W3

W3 ✓

| 1 | 2 | 3 |
|---|---|---|

maxInst

maxAck

148

# In Depth: Reads

# In Depth: Reads

Reads wait for

*"**ack**" of next "**unseen instance**"*

# In Depth: Reads

Why acks of next unseen instance?

# In Depth: Reads

Why acks of next unseen instance?

Why not immediately reply?

# In Depth: Reads

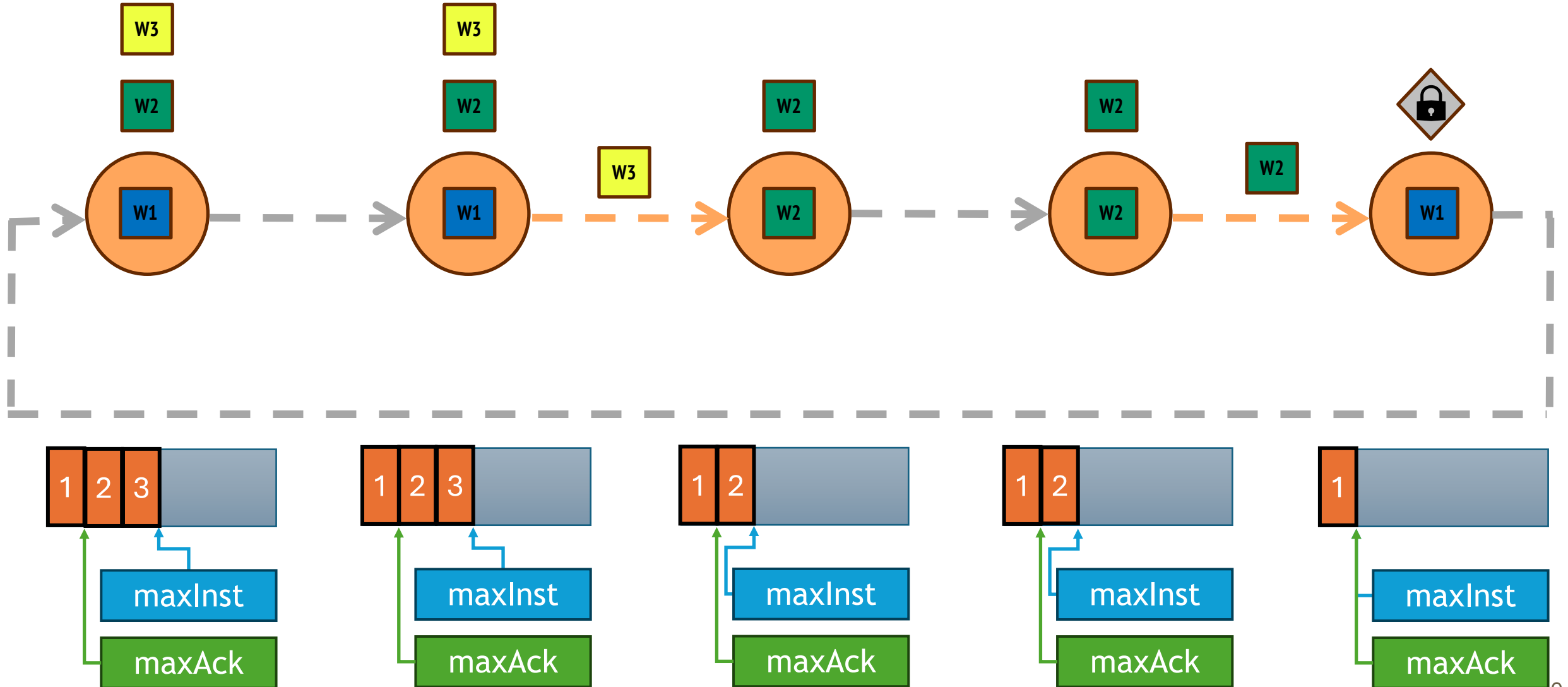Why acks of next unseen instance?

Why not immediately reply?

Replicas cannot be sure whether other replicas have replied to new writes

# In Depth: Why not immediately reply?

# In Depth: Why not immediately reply?

# In Depth: Why not immediately reply?



Not Linearizable ✗

# In Depth: Reads

Why acks of next unseen instance?

Why not immediately reply?

# In Depth: Reads

Why acks of next unseen instance?

~~Why not immediately reply?~~

Why not next unseen instance?

Again, replicas cannot be sure whether other replicas have replied to new writes

# In Depth: Why not next unseen instance?
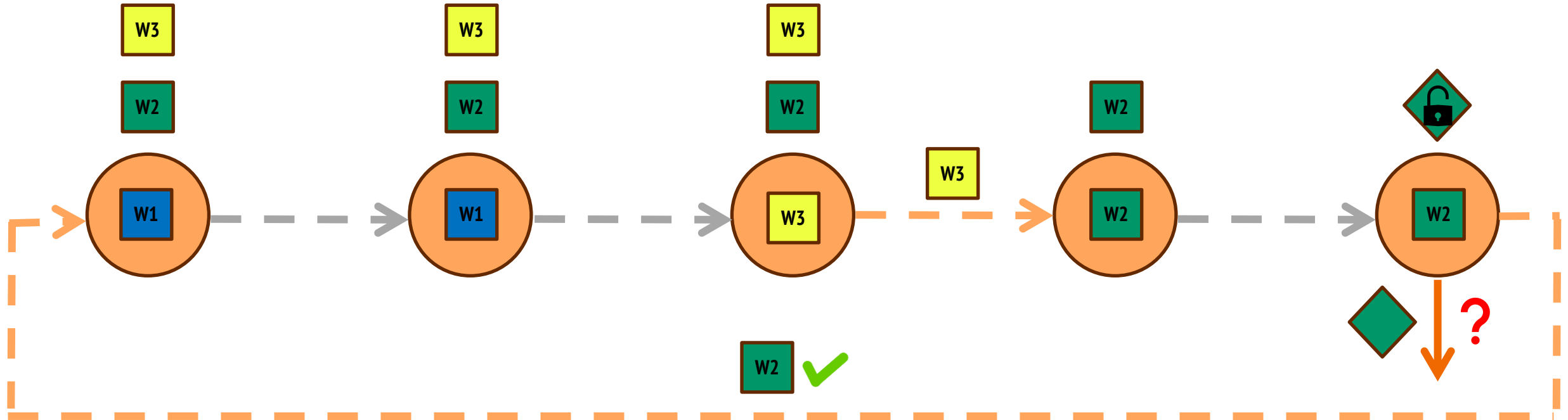
# In Depth: Why not next unseen instance?
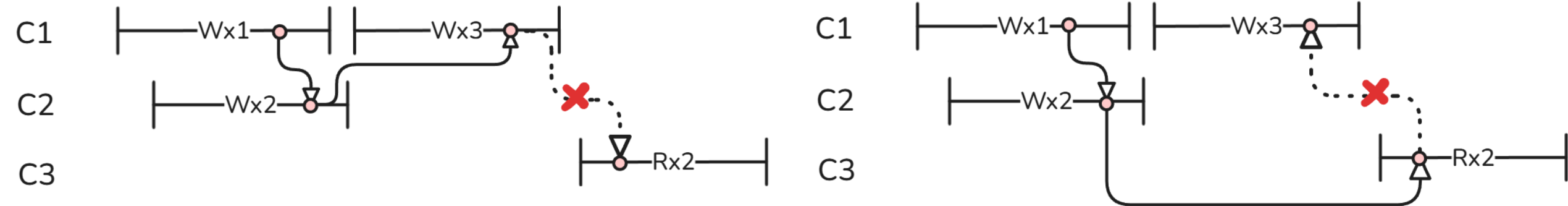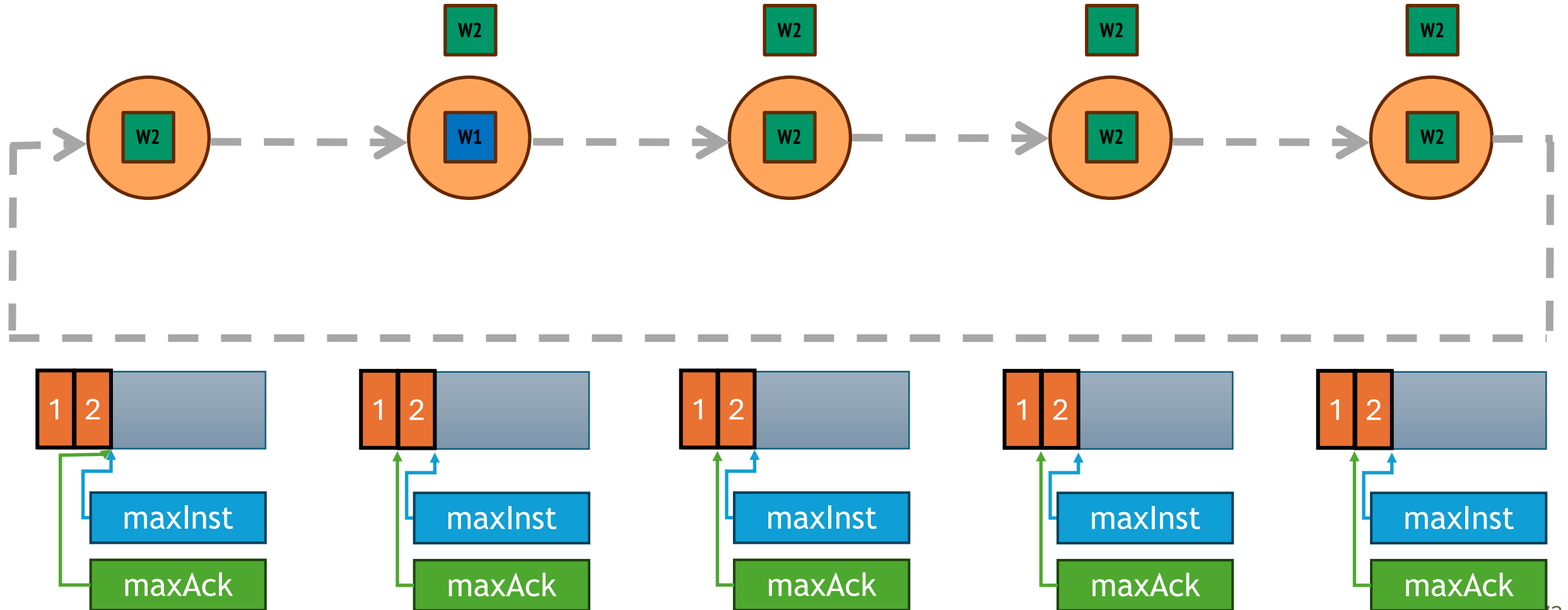
# In Depth: Why not next unseen instance?

# In Depth: Why not next unseen instance?

Consider a slightly different

situation of the system ...

# In Depth: Why not next unseen instance?

# In Depth: Why not next unseen instance?



Differences

# In Depth: Why not next unseen instance?

# In Depth: Why not next unseen instance?

# In Depth: Why not next unseen instance?

# In Depth: Why not next unseen instance?



Not Linearizable ✘

# In Depth: Reads

Why acks of next unseen instance?

~~Why not immediately reply?~~

~~Why not next unseen instance?~~

# In Depth: Reads

Why acks of next unseen instance?

~~Why not immediately reply?~~

~~Why not next unseen instance?~~

Why not ack of max seen instance?

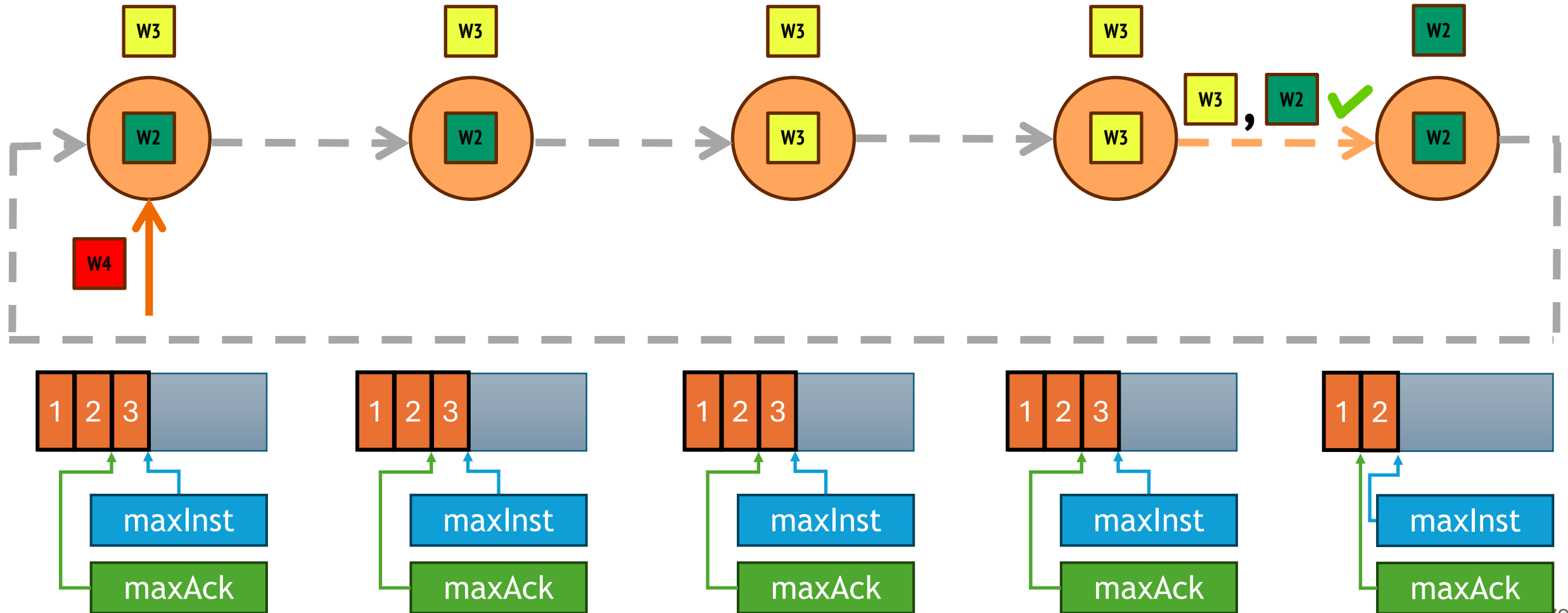Slightly more complicated to show, but it still leads to non-linearizable history
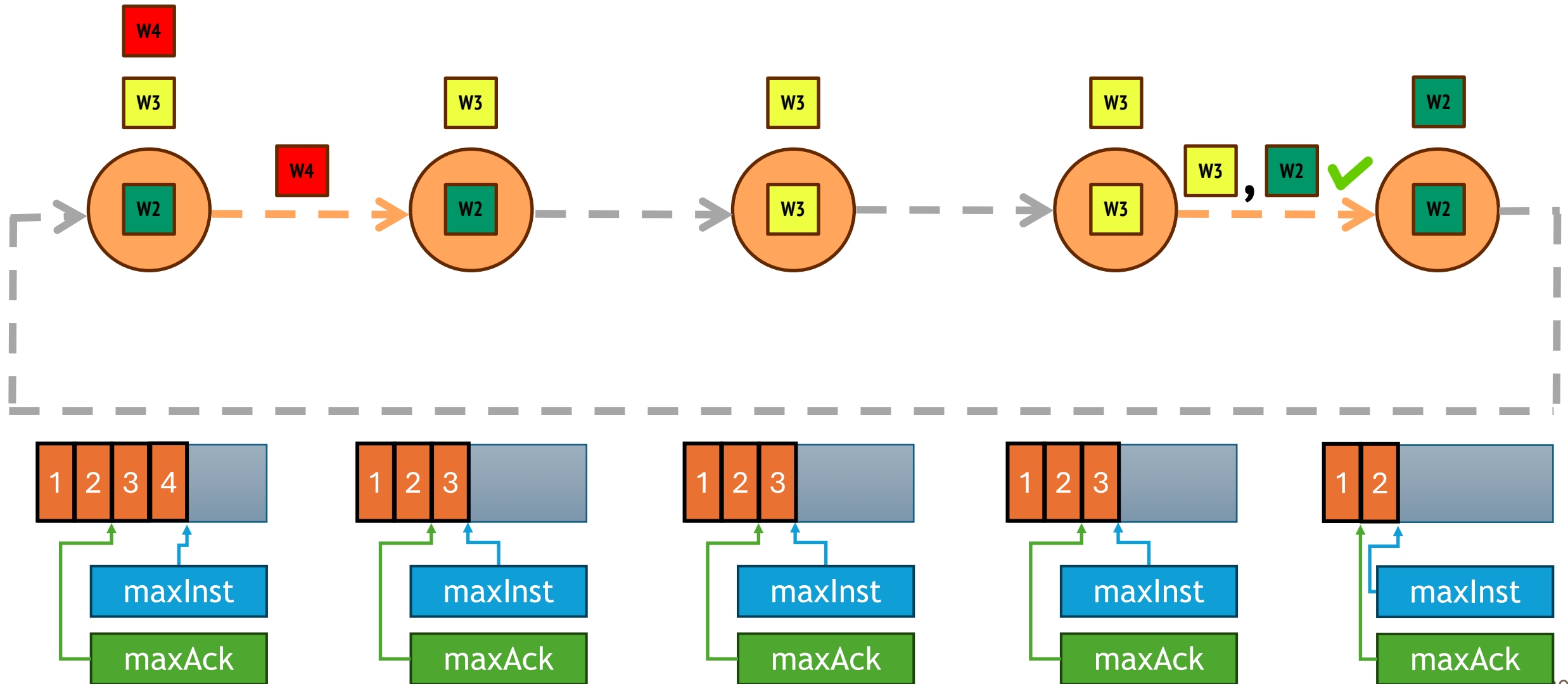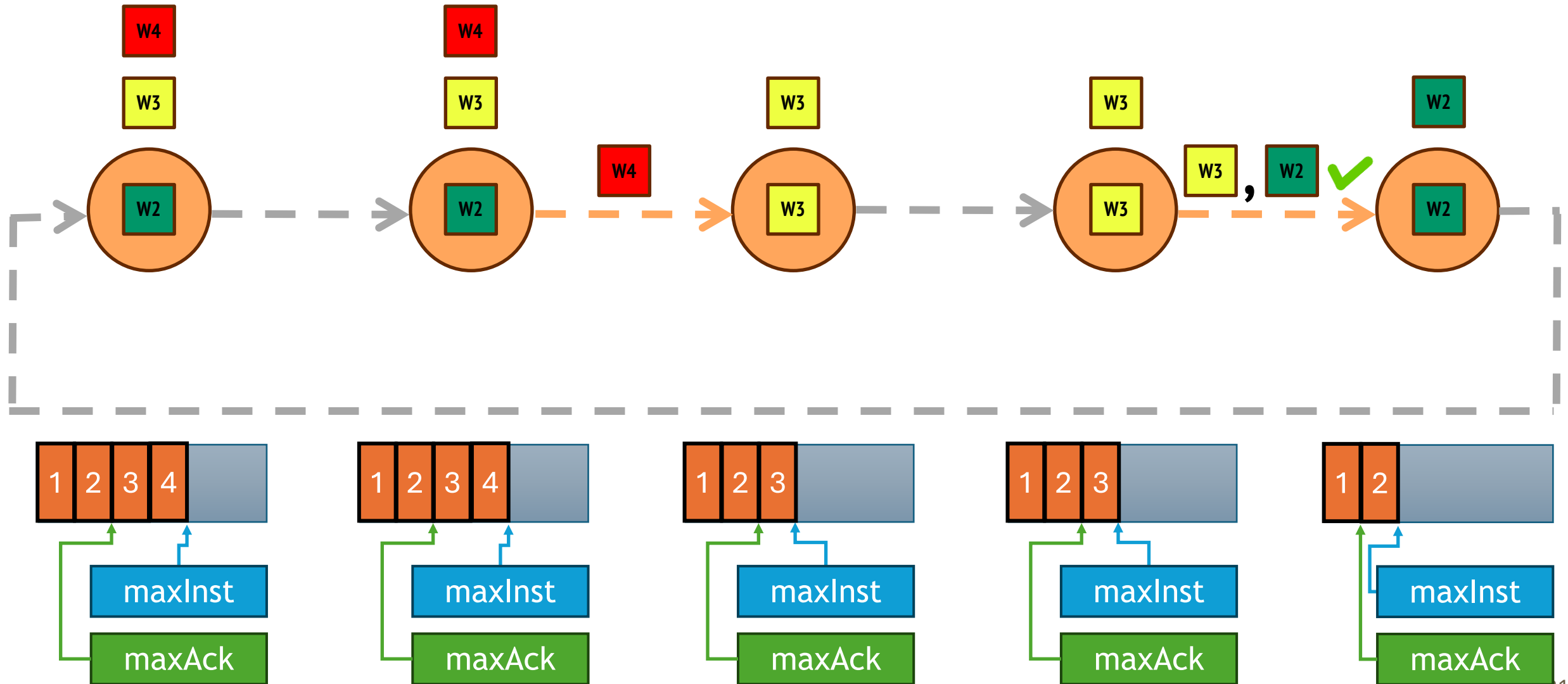
# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?
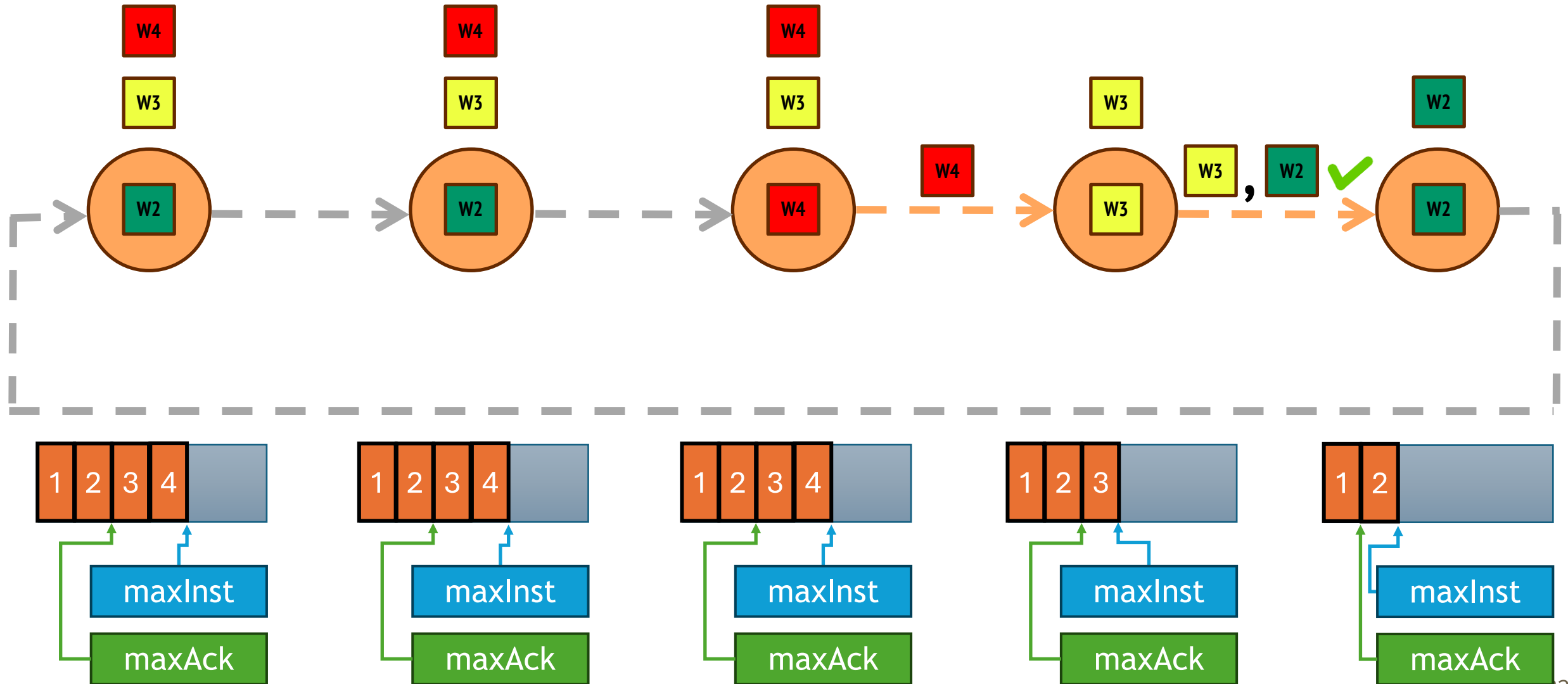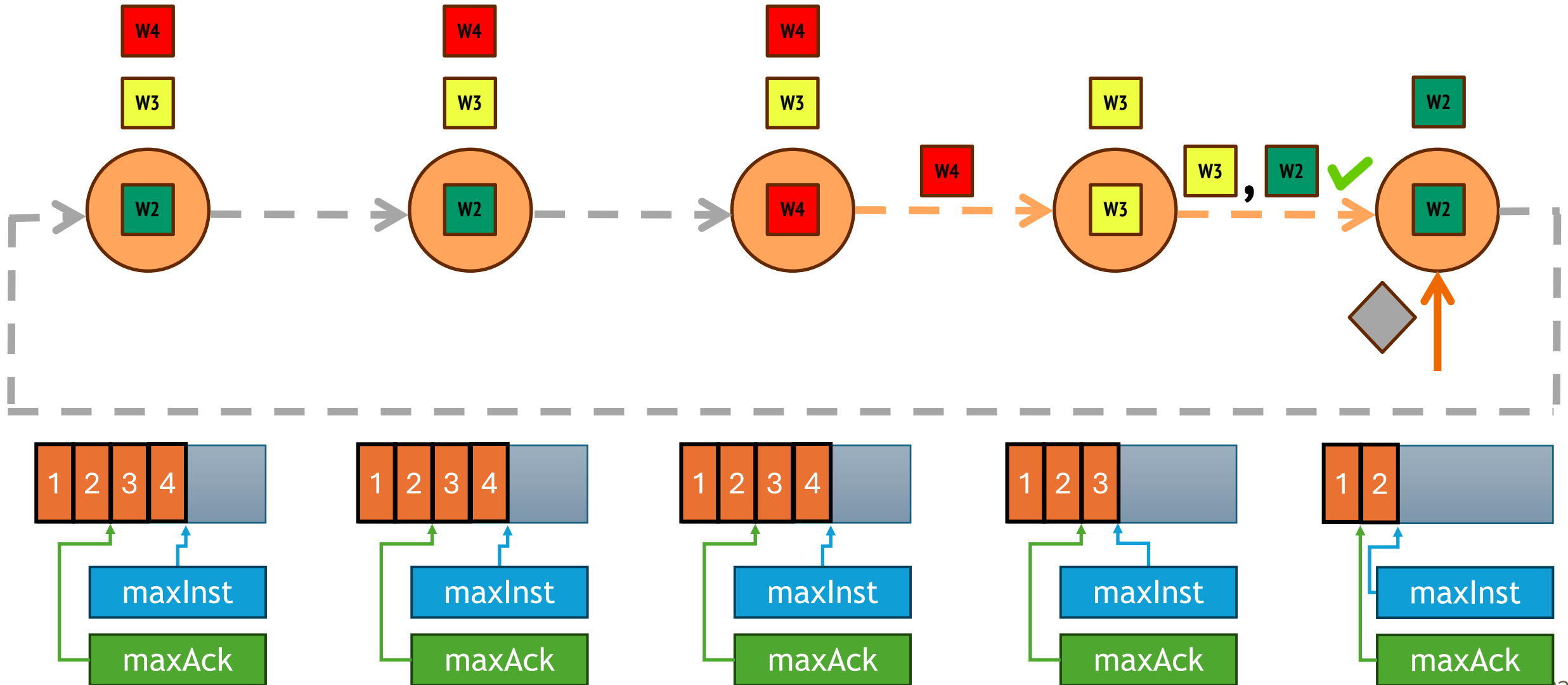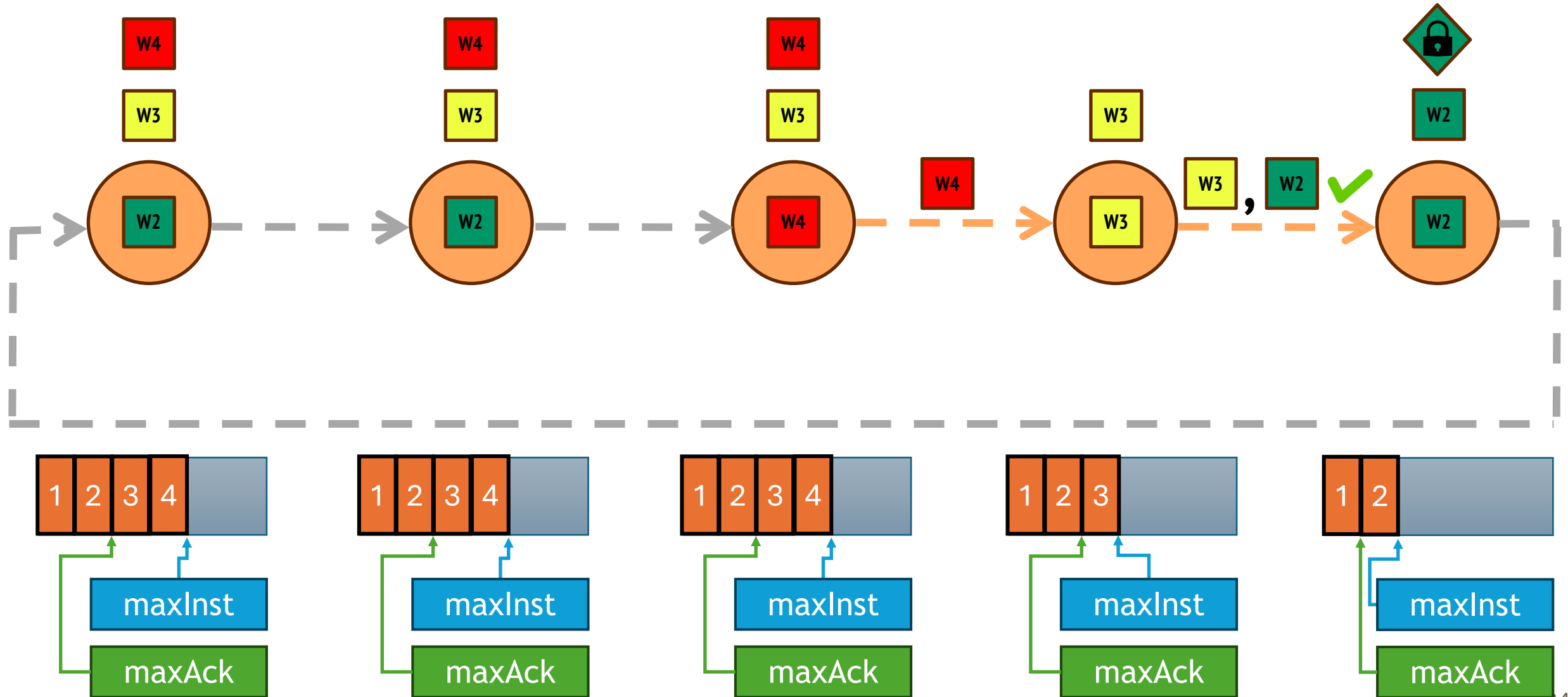
# In Depth: Why not ack of max seen instance?

172

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?
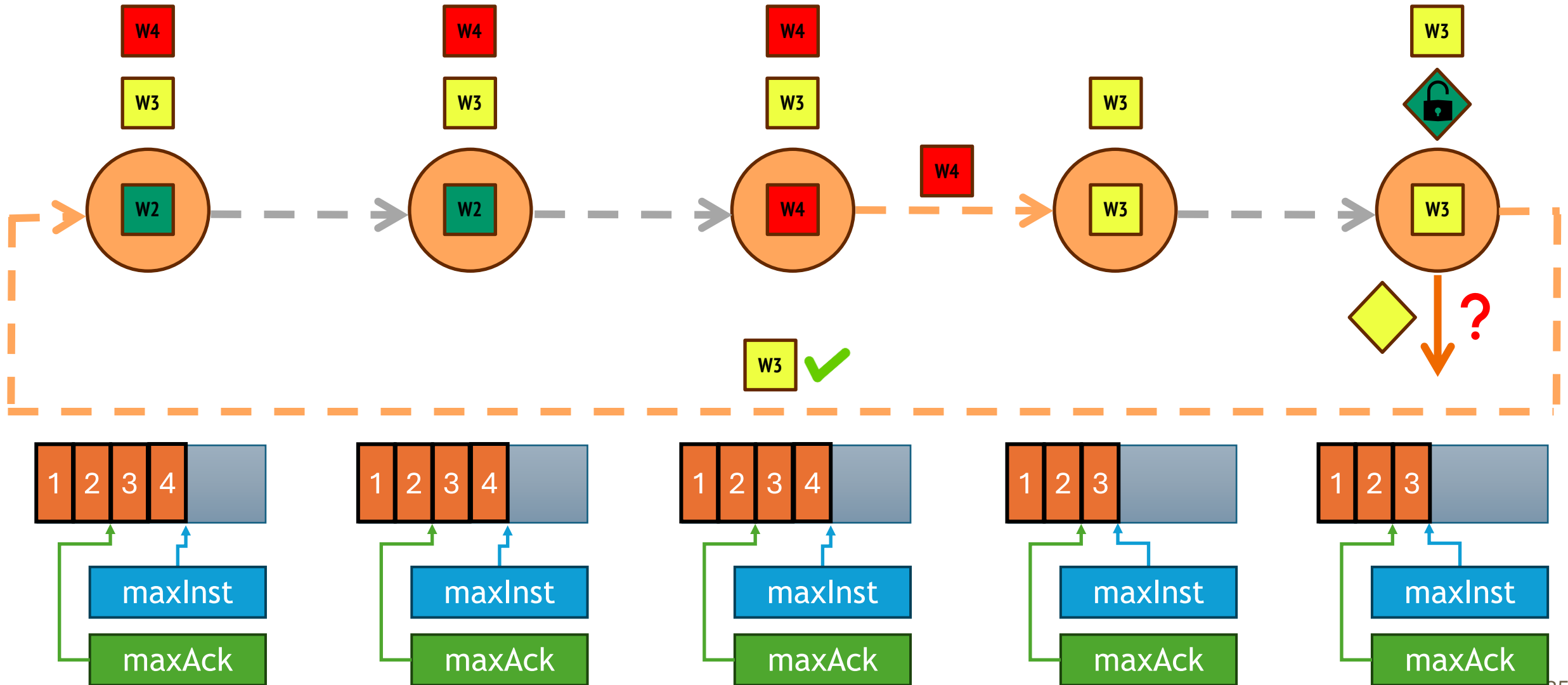
# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?
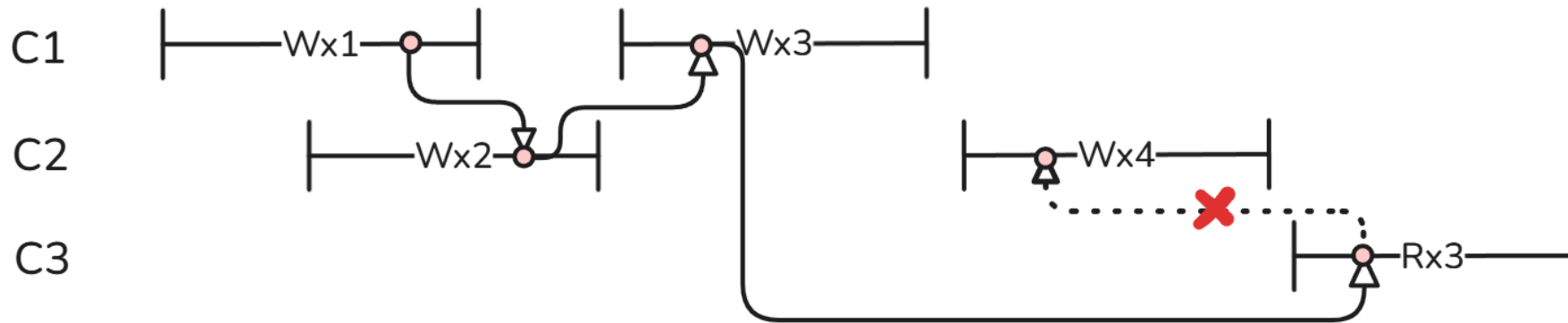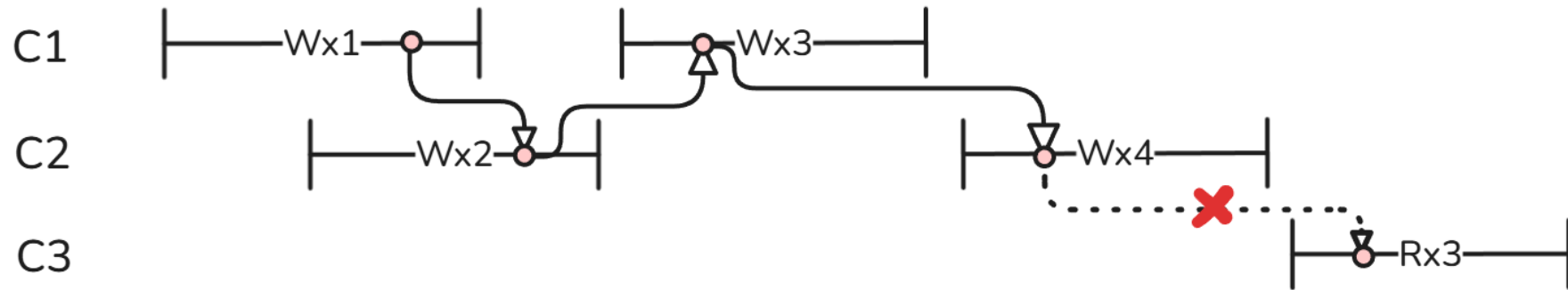
# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?

# In Depth: Why not ack of max seen instance?
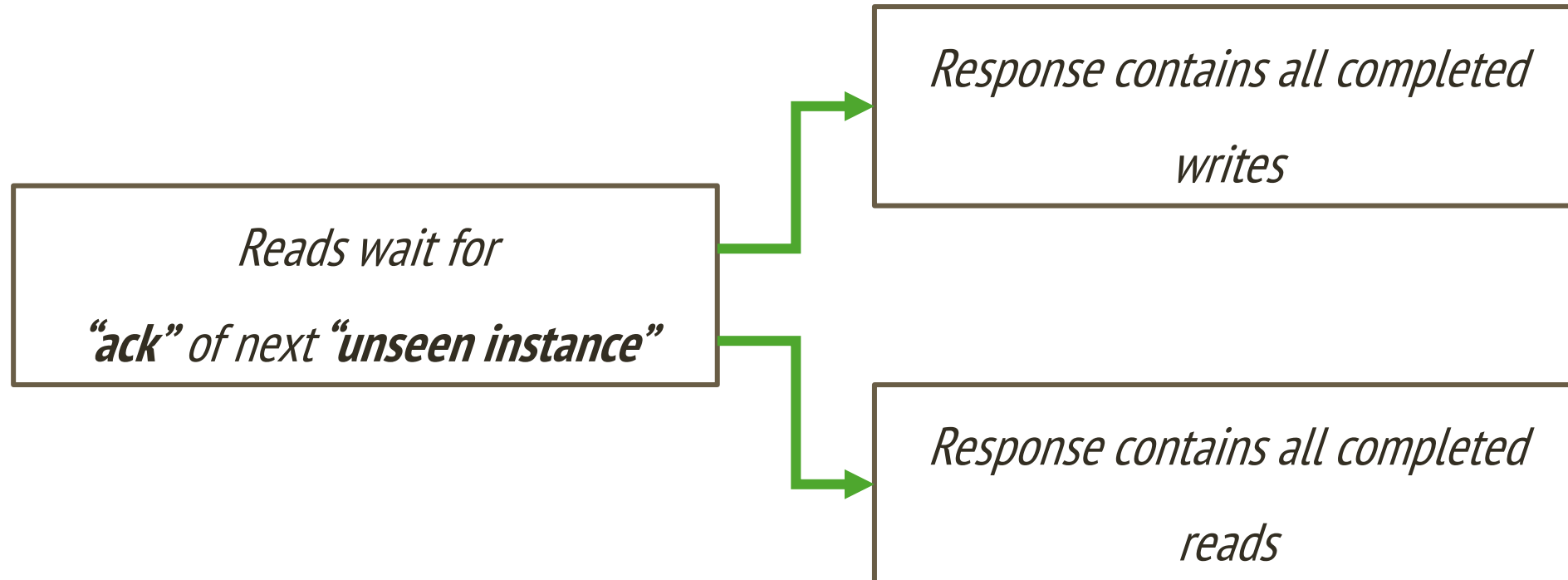


Not Linearizable ✗

# In Depth: Reads

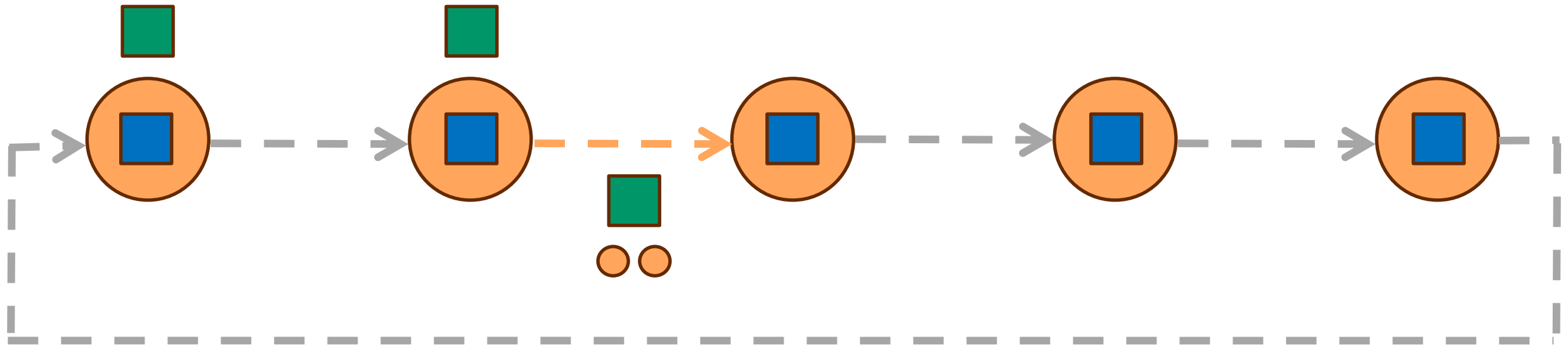Why acks of next unseen instance?

~~Why not immediately reply?~~

~~Why not next unseen instance?~~
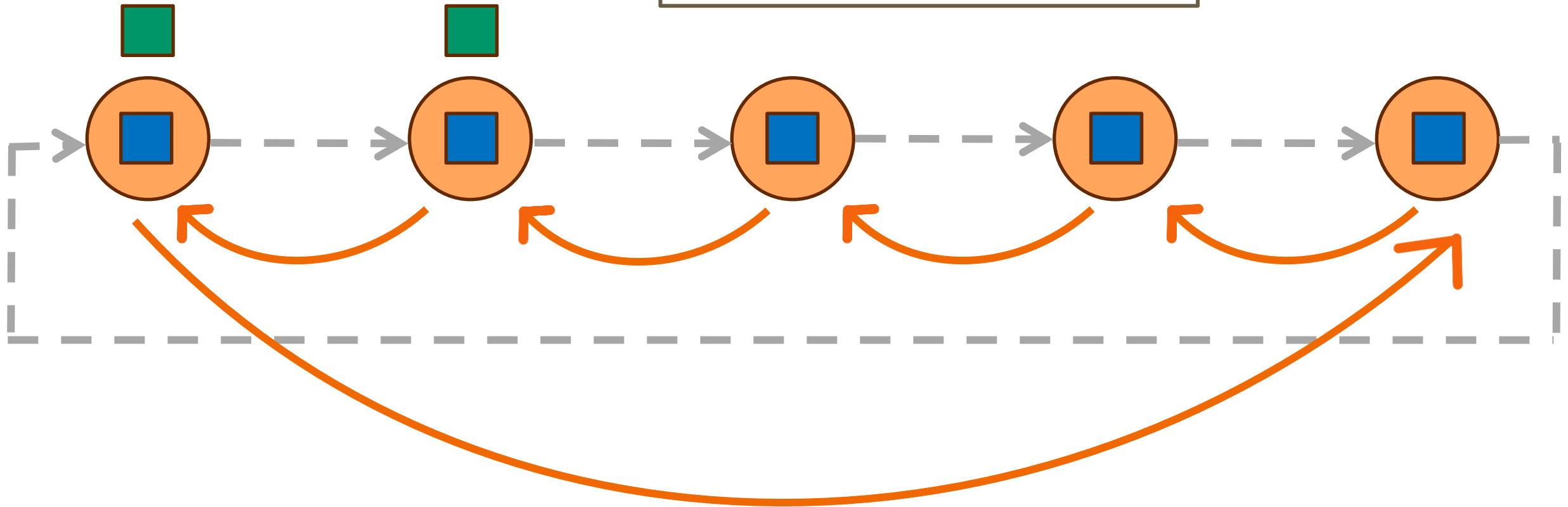
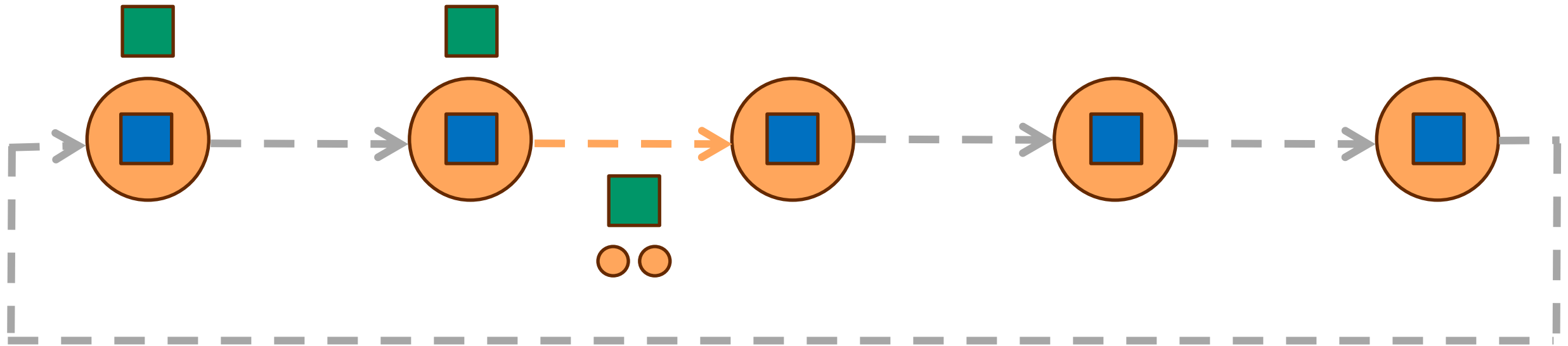~~Why not ack of max seen instance?~~

# In Depth: Reads



Reads wait for

**"ack"** of next **"unseen instance"**

→ Response contains all completed writes

→ Response contains all completed reads

# Membership Management: Removal

# Membership Management: Removal

Replicas send "keep-alive" messages

to previous replica

# Membership Management: Removal

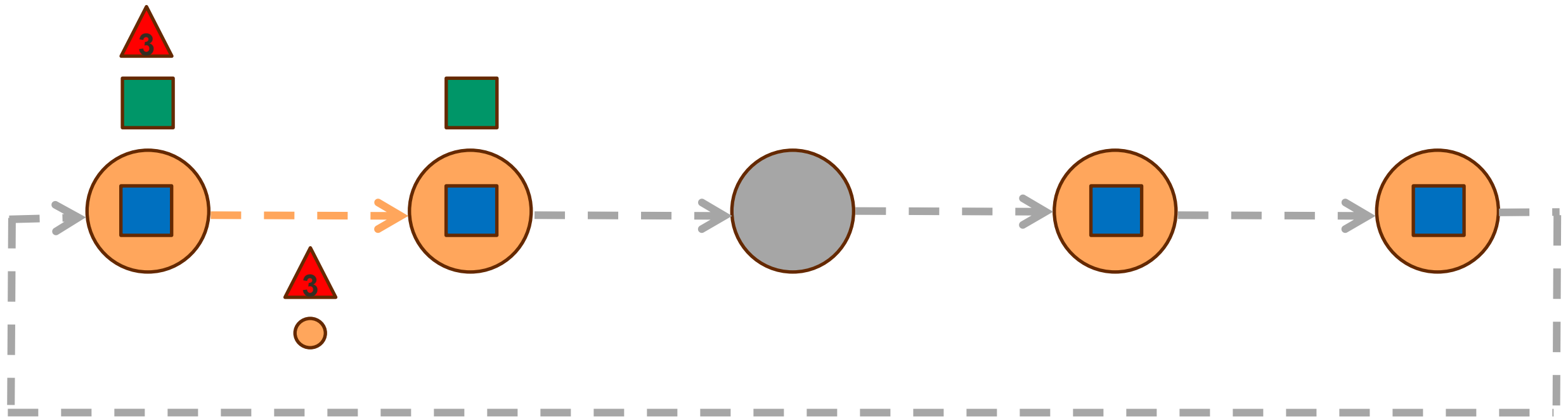# Membership Management: Removal

# Membership Management: Removal
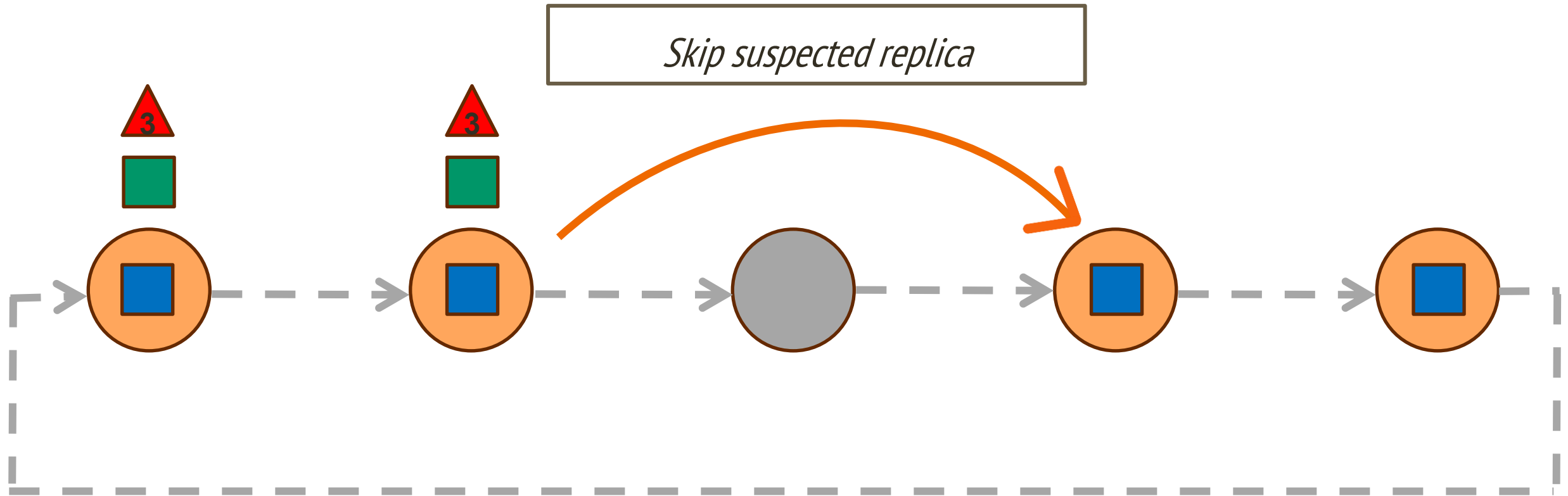
Remove requests are handled like regular write operations
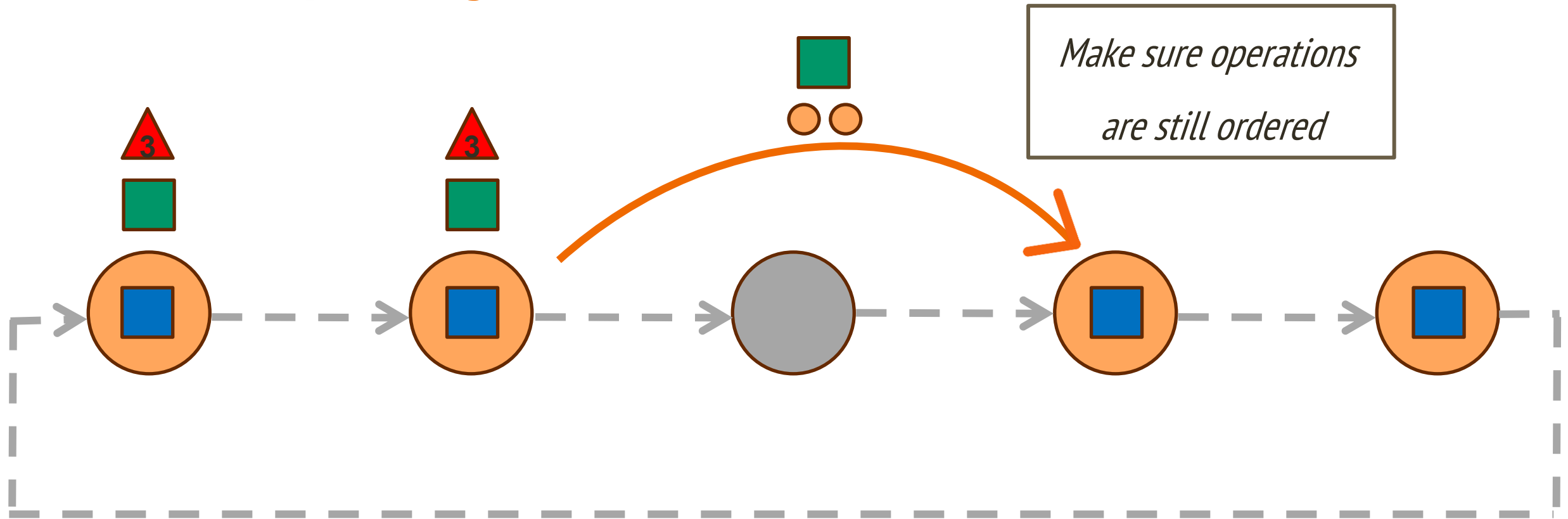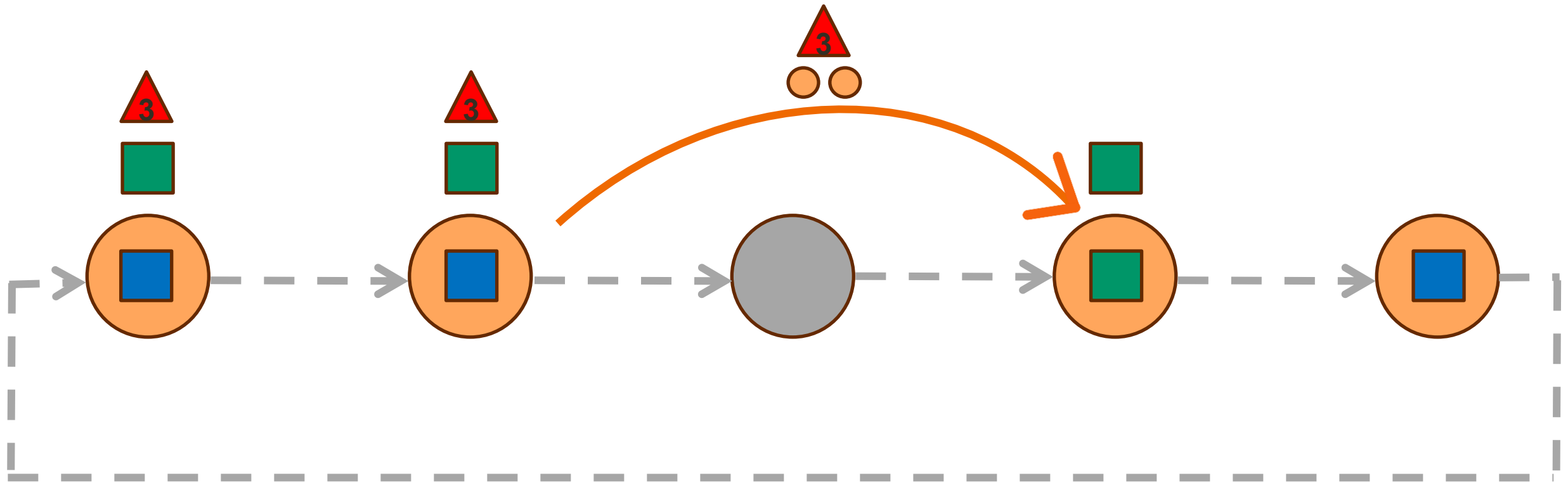


3

⚠ i  : Remove i'th replica request

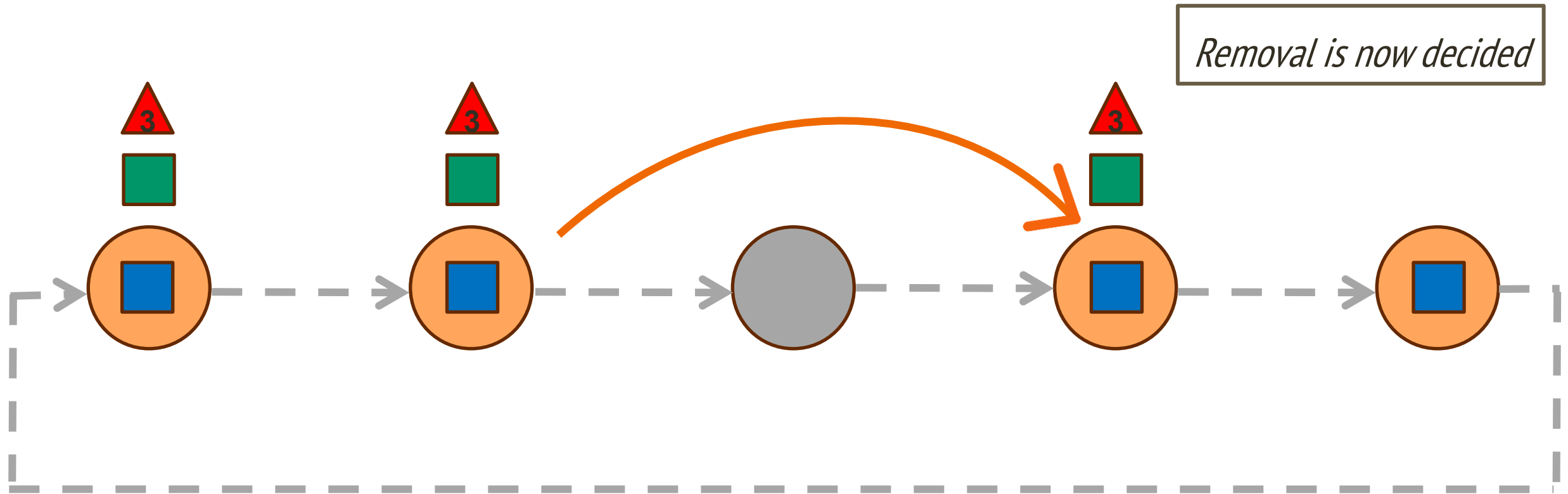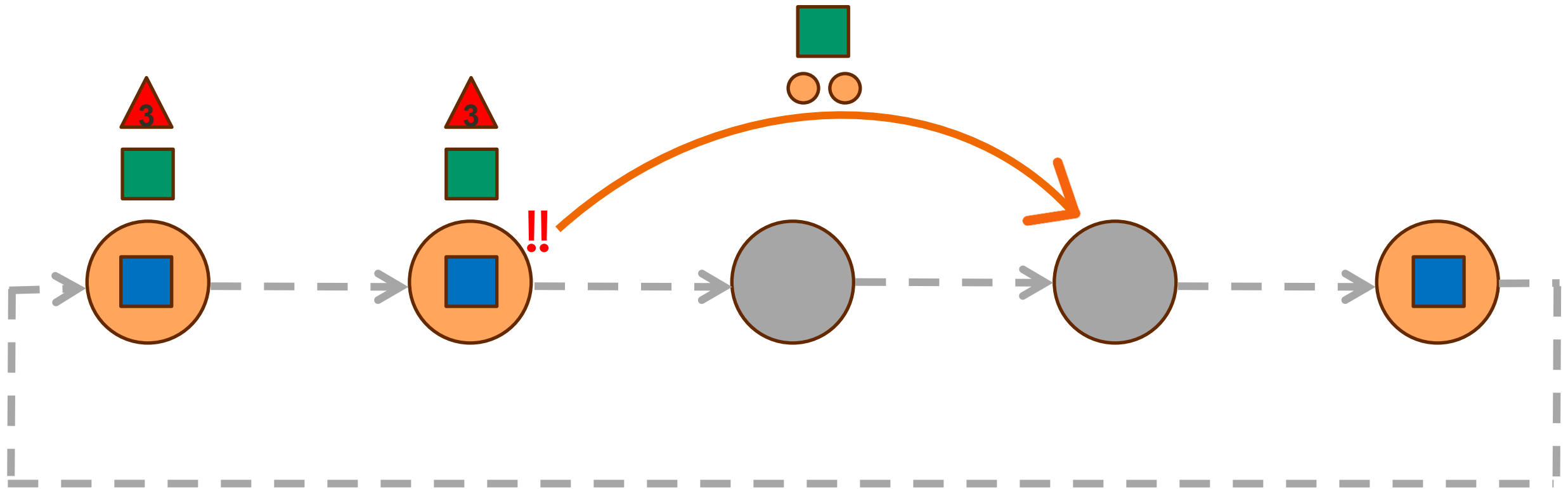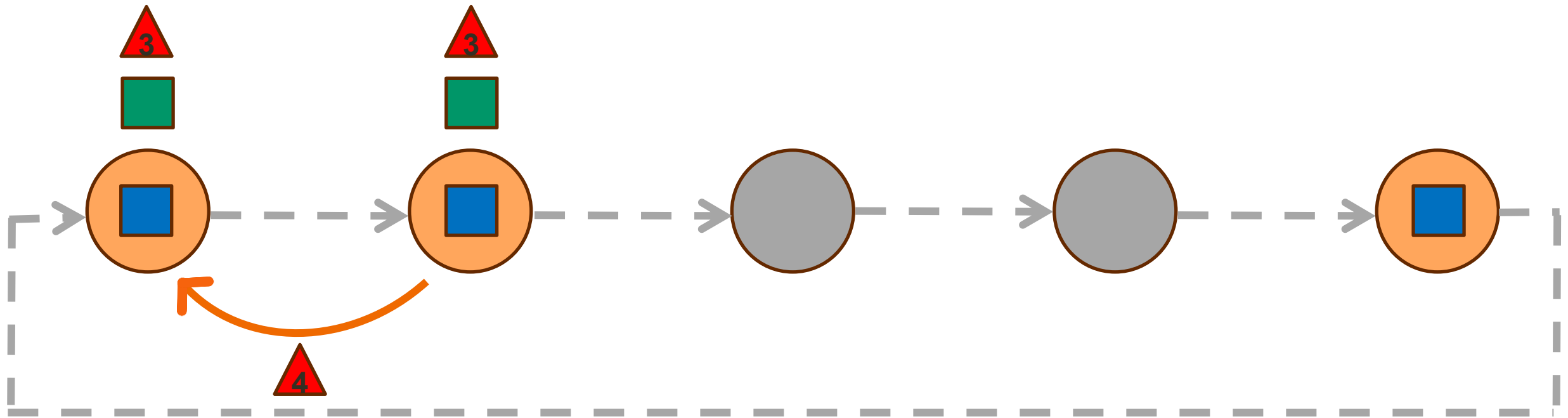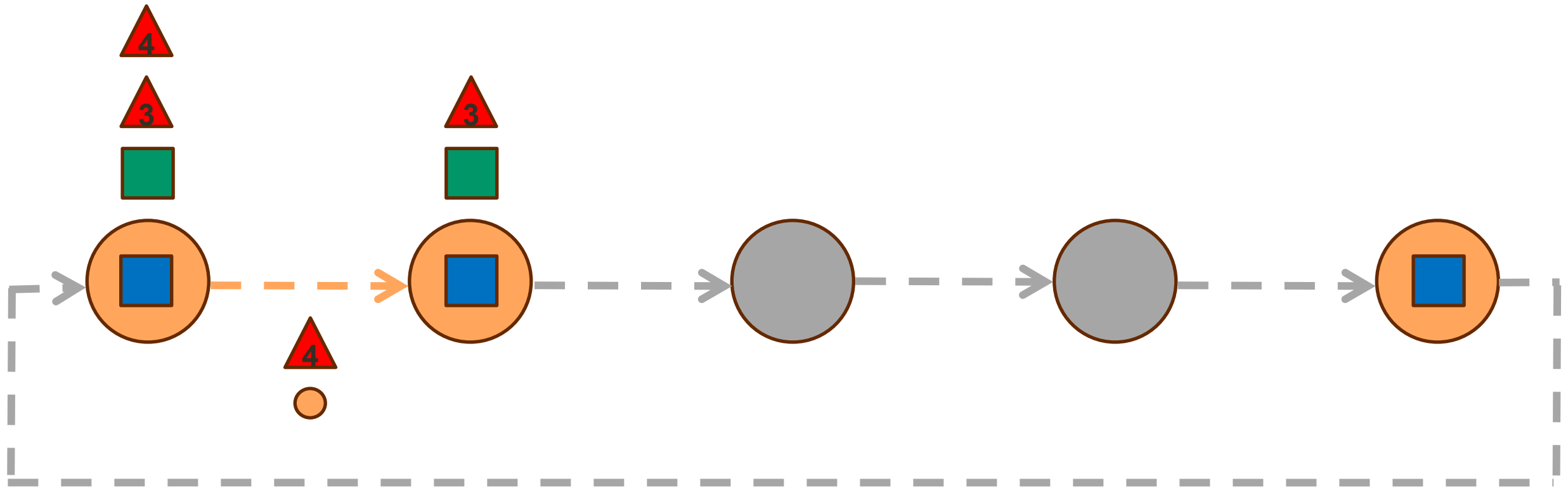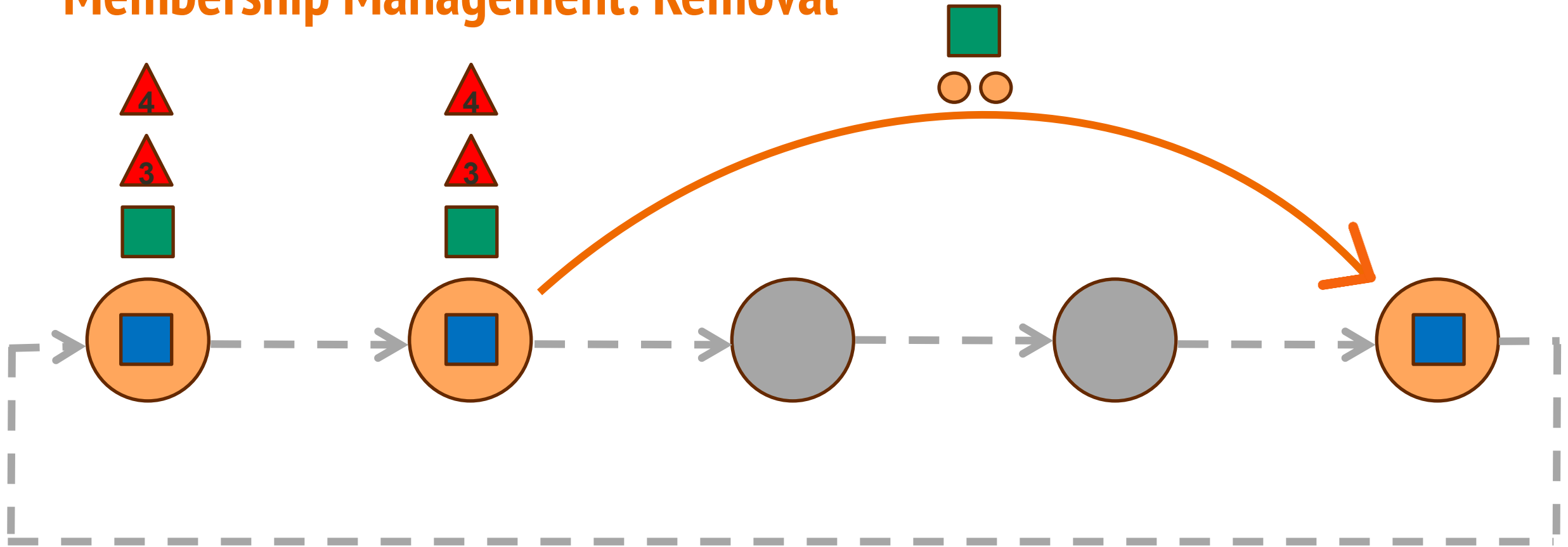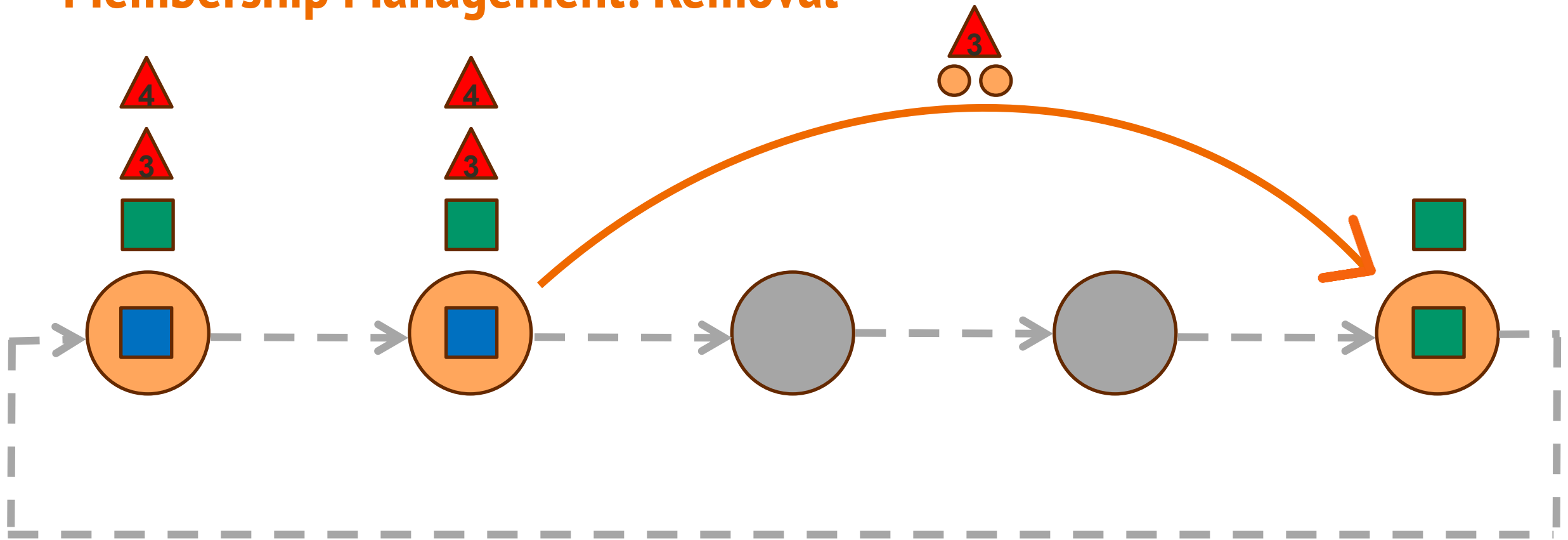# Membership Management: Removal

# Membership Management: Removal



Skip suspected replica

# Membership Management: Removal



Make sure operations are still ordered

# Membership Management: Removal

# Membership Management: Removal



Removal is now decided

# Membership Management: Removal

# Membership Management: Removal
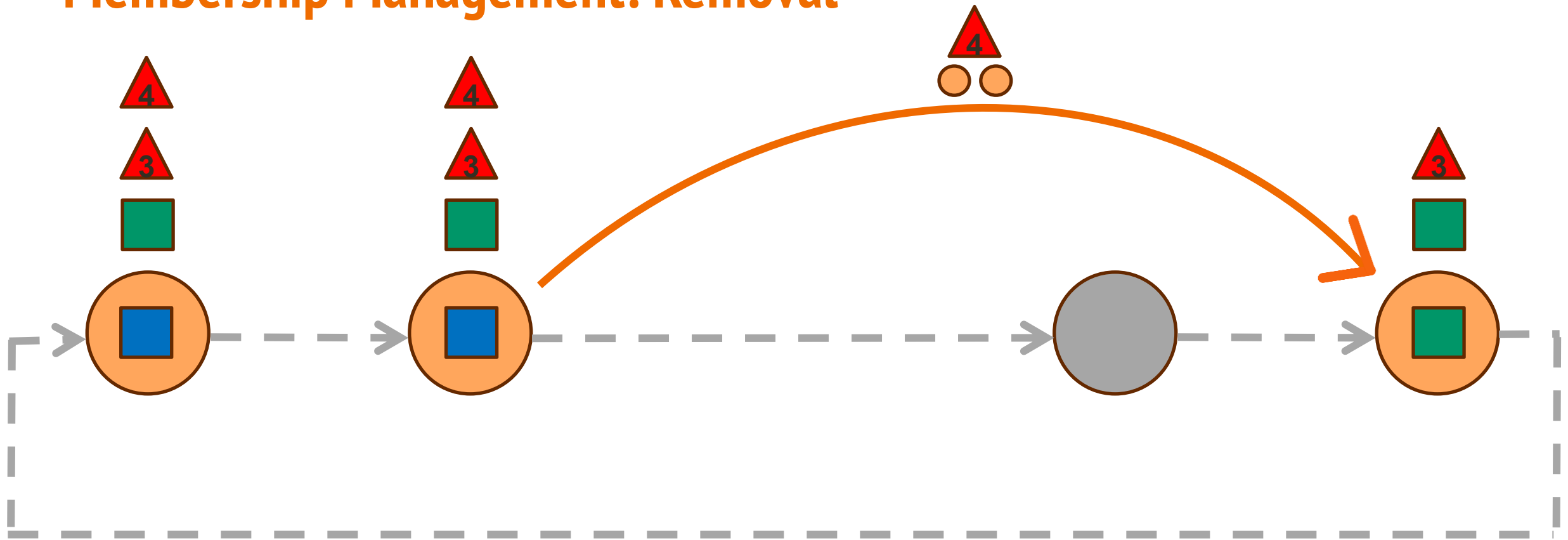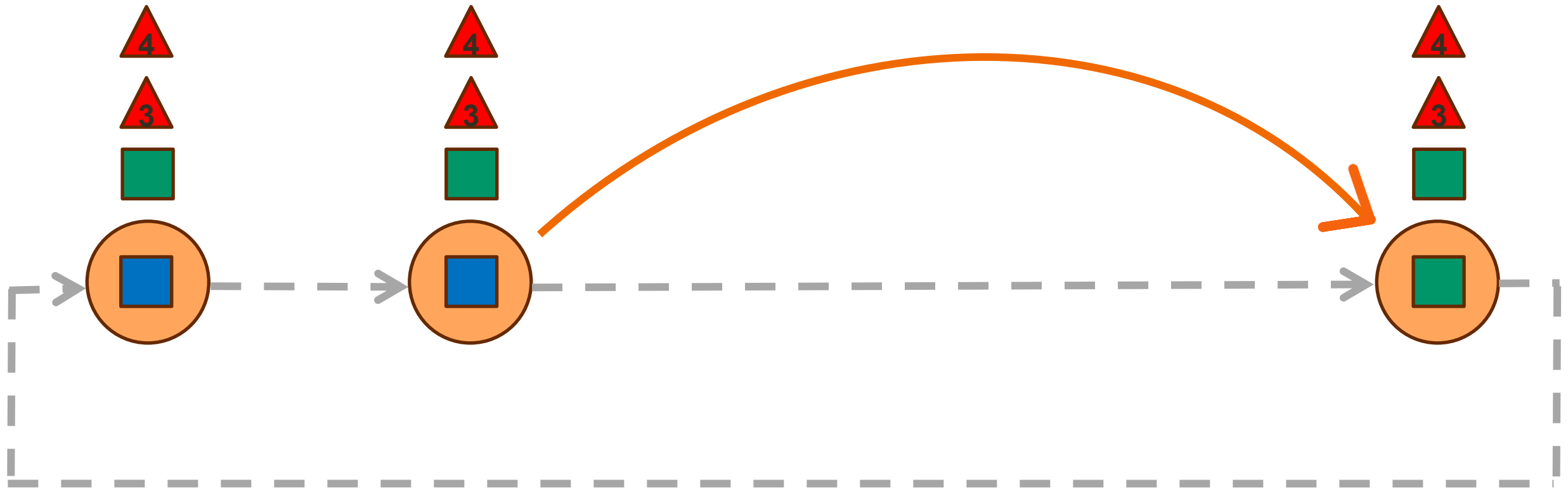
# Membership Management: Removal

# Membership Management: Removal

# Membership Management: Removal
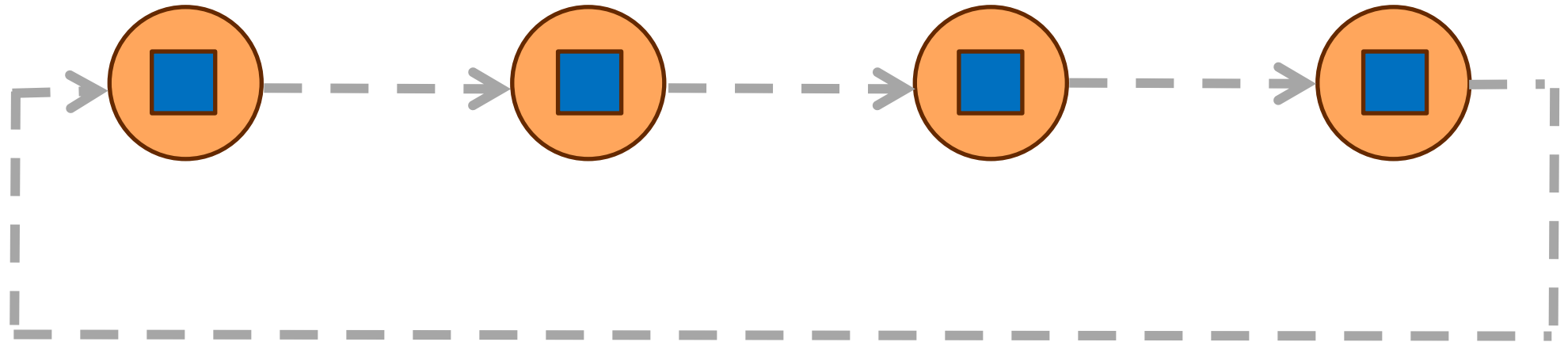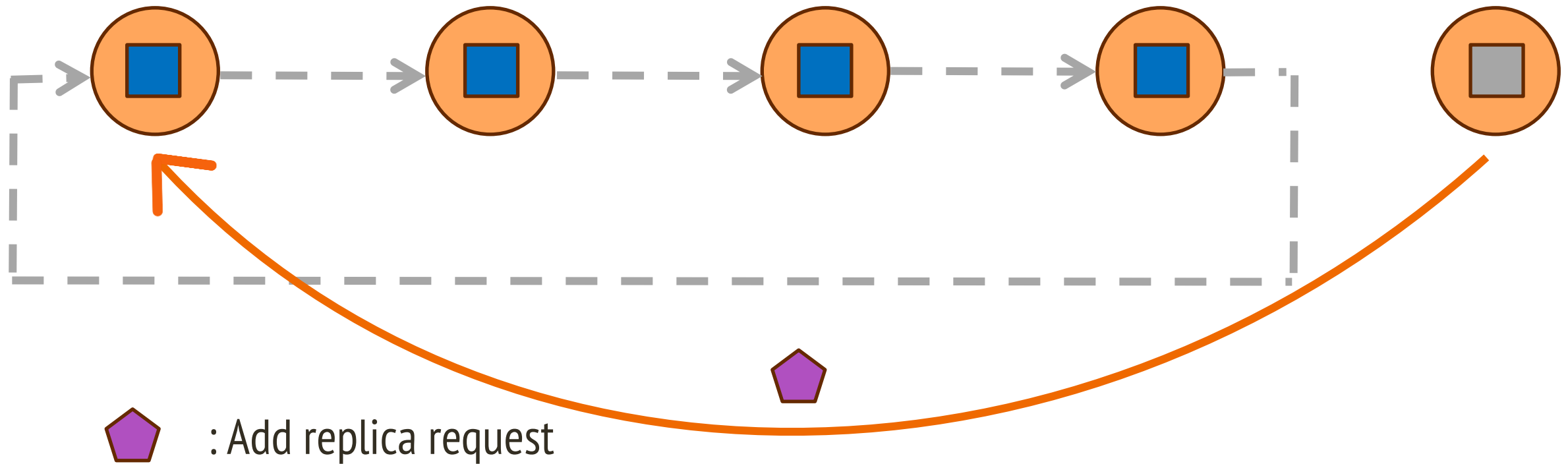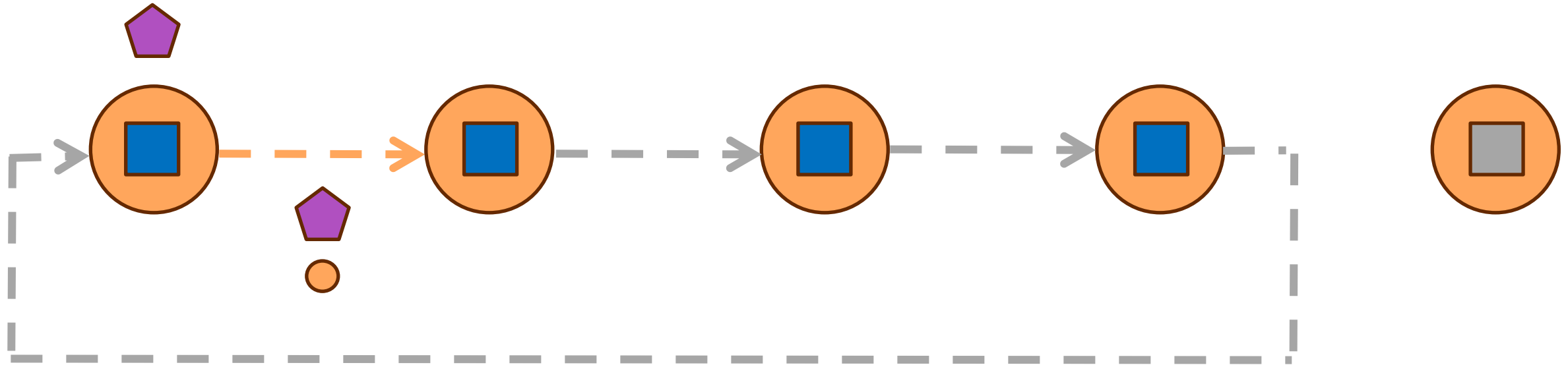
# Membership Management: Removal

# Membership Management: Removal

# Membership Management: Removal

# Membership Management: Addition

# Membership Management: Addition



: Add replica request

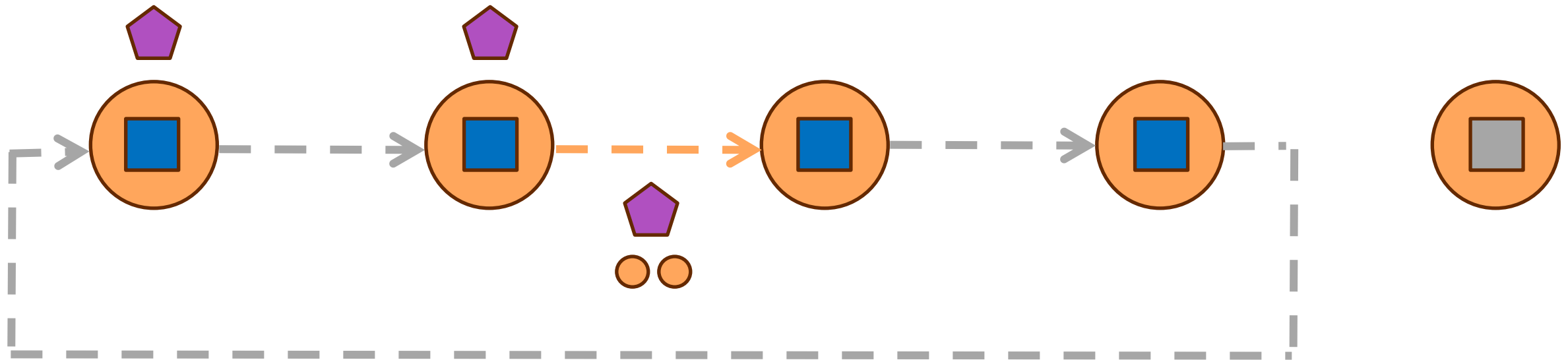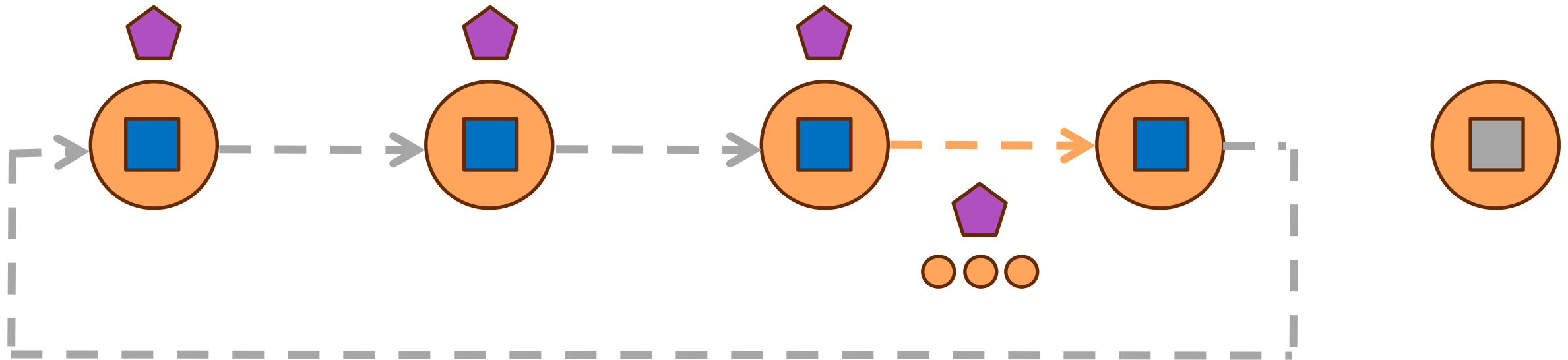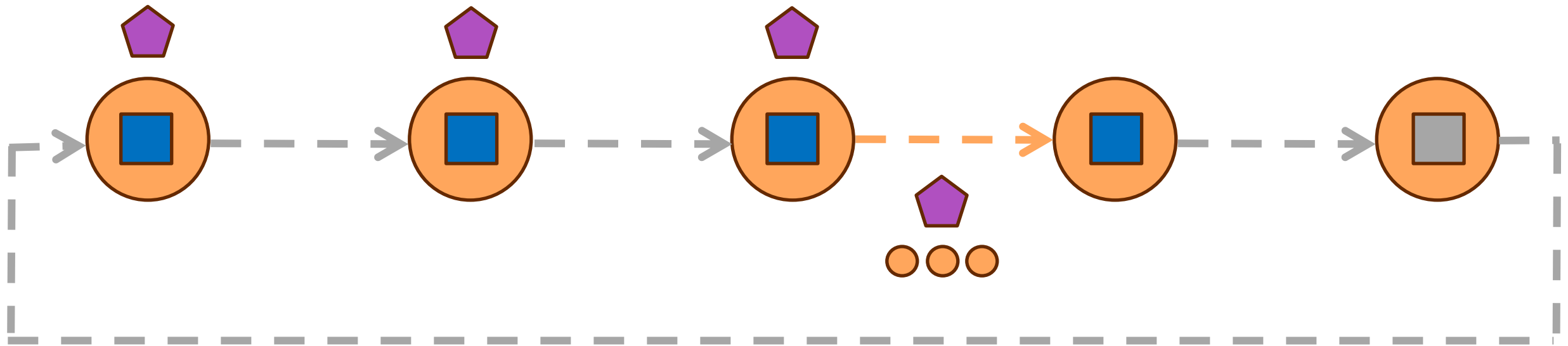# Membership Management: Addition

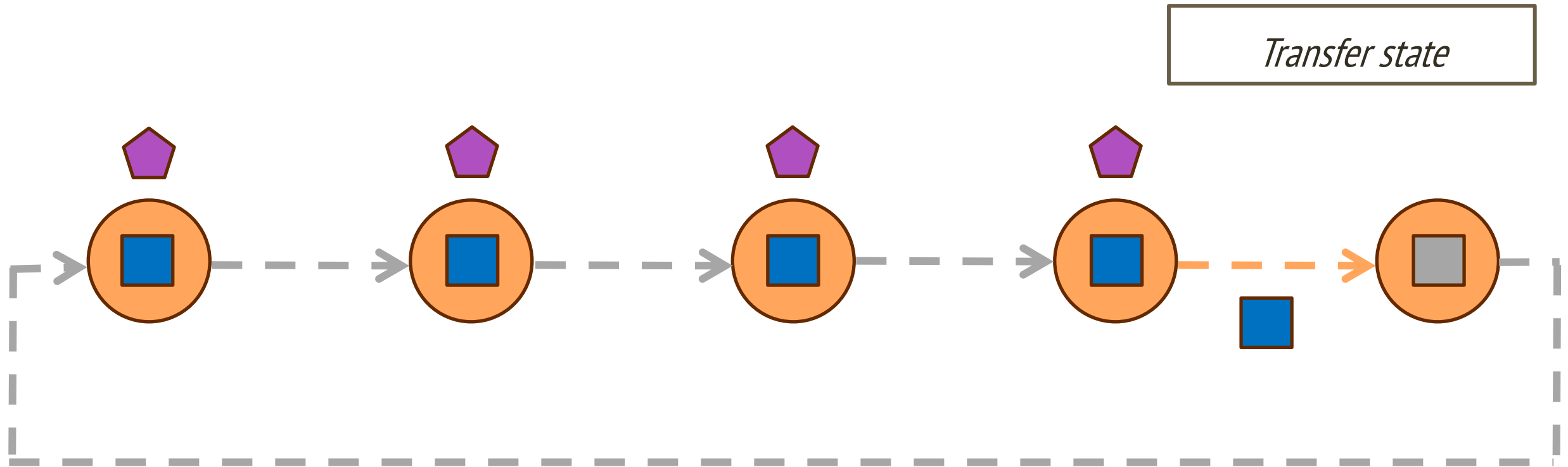# Membership Management: Addition
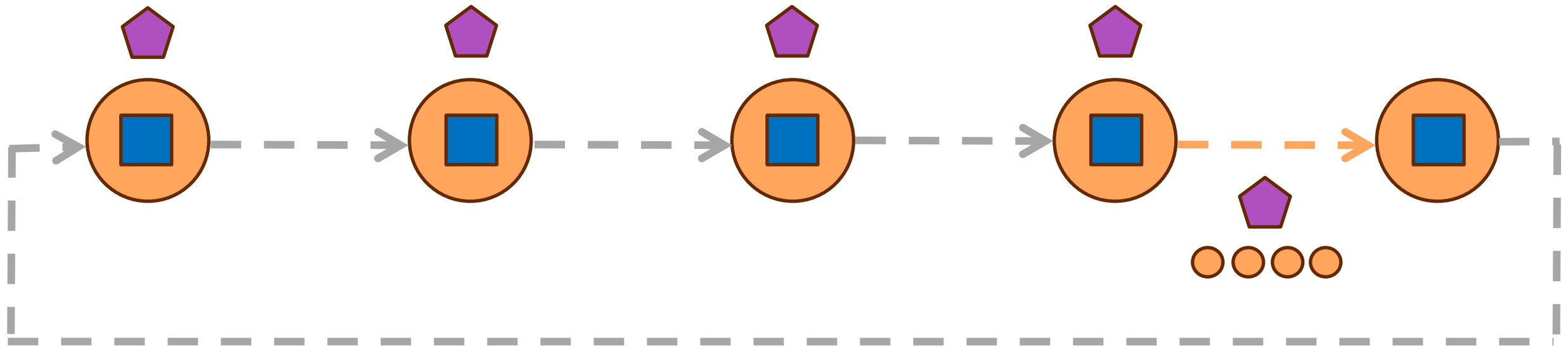
# Membership Management: Addition

Addition is now decided
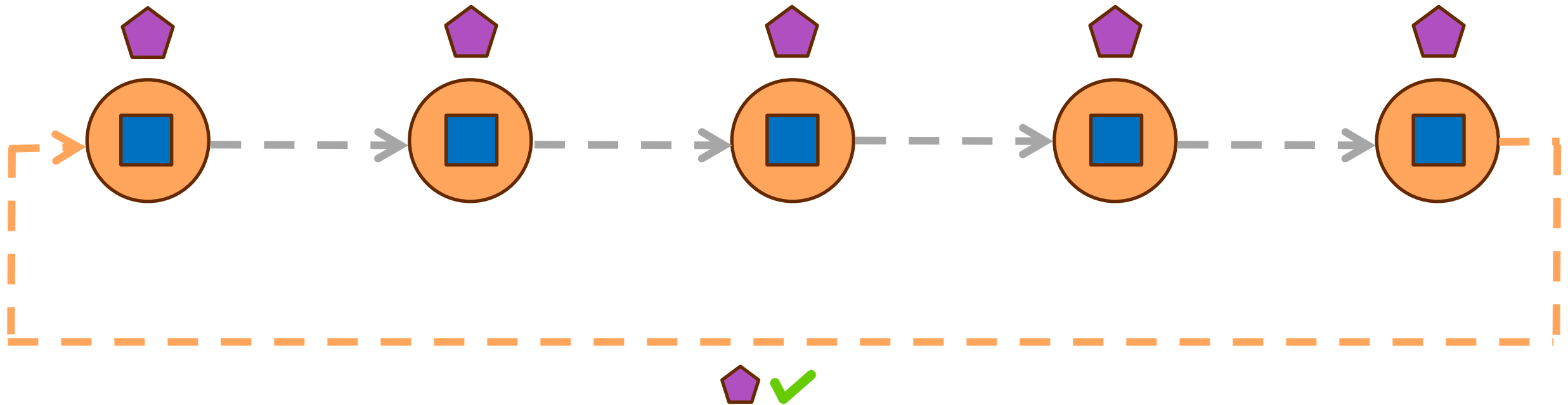


211

# Membership Management: Addition
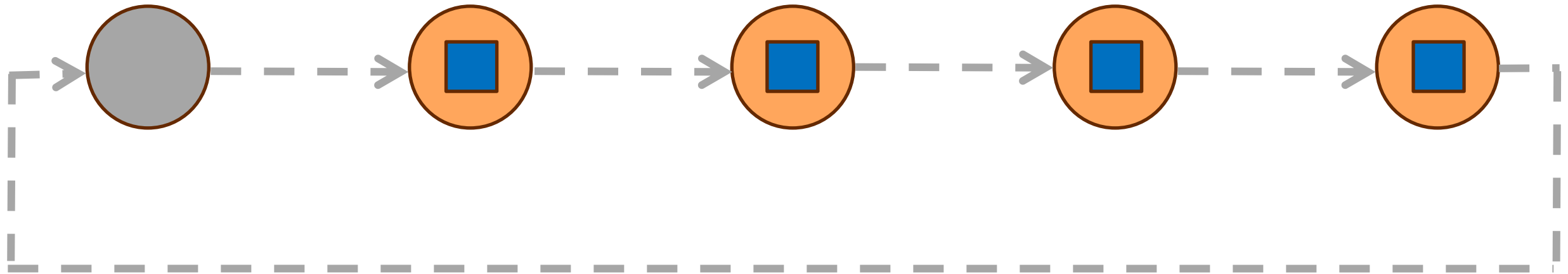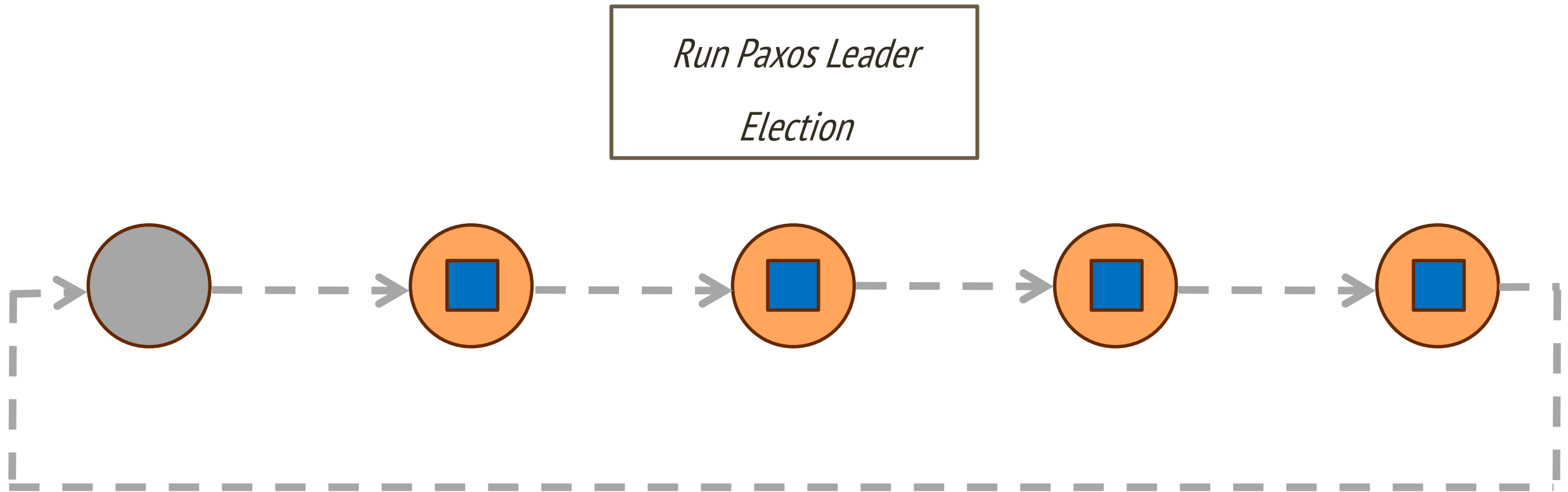
# Membership Management: Addition



Transfer state

213

# Membership Management: Addition

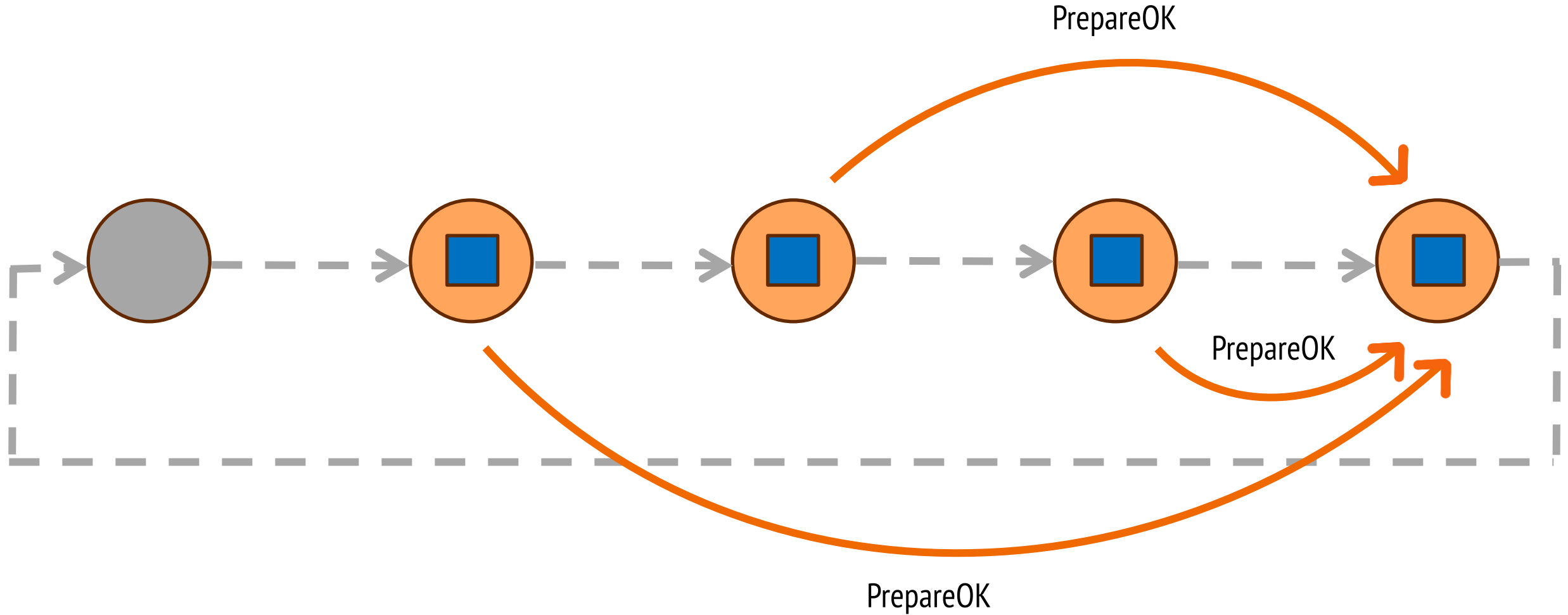# Membership Management: Addition
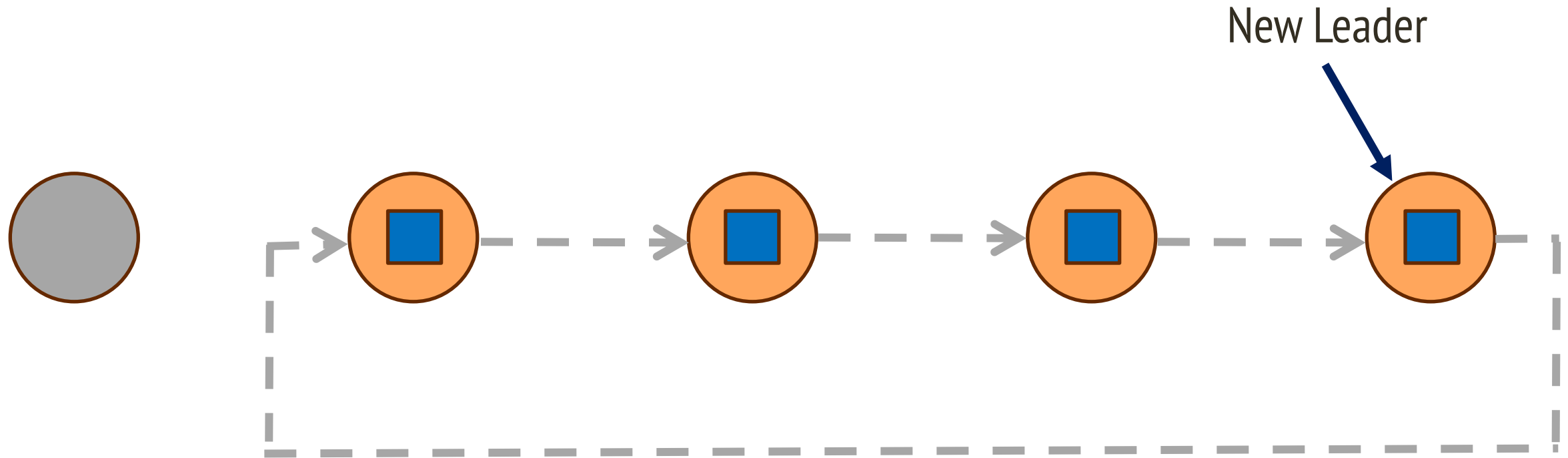
# Leader Failure

# Leader Failure

Run Paxos Leader

Election

# Leader Failure



PrepareMsg

PrepareMsg

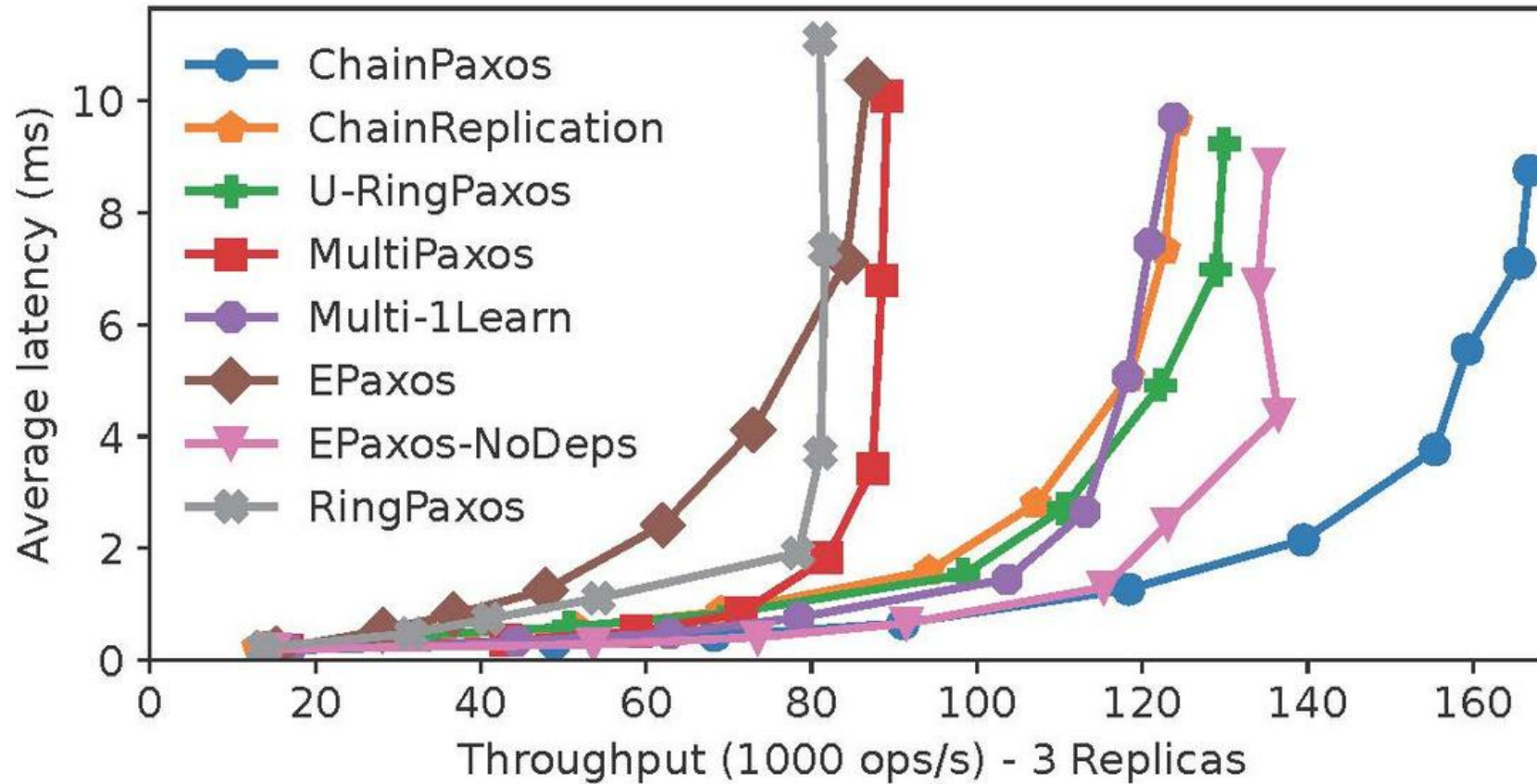PrepareMsg

# Leader Failure
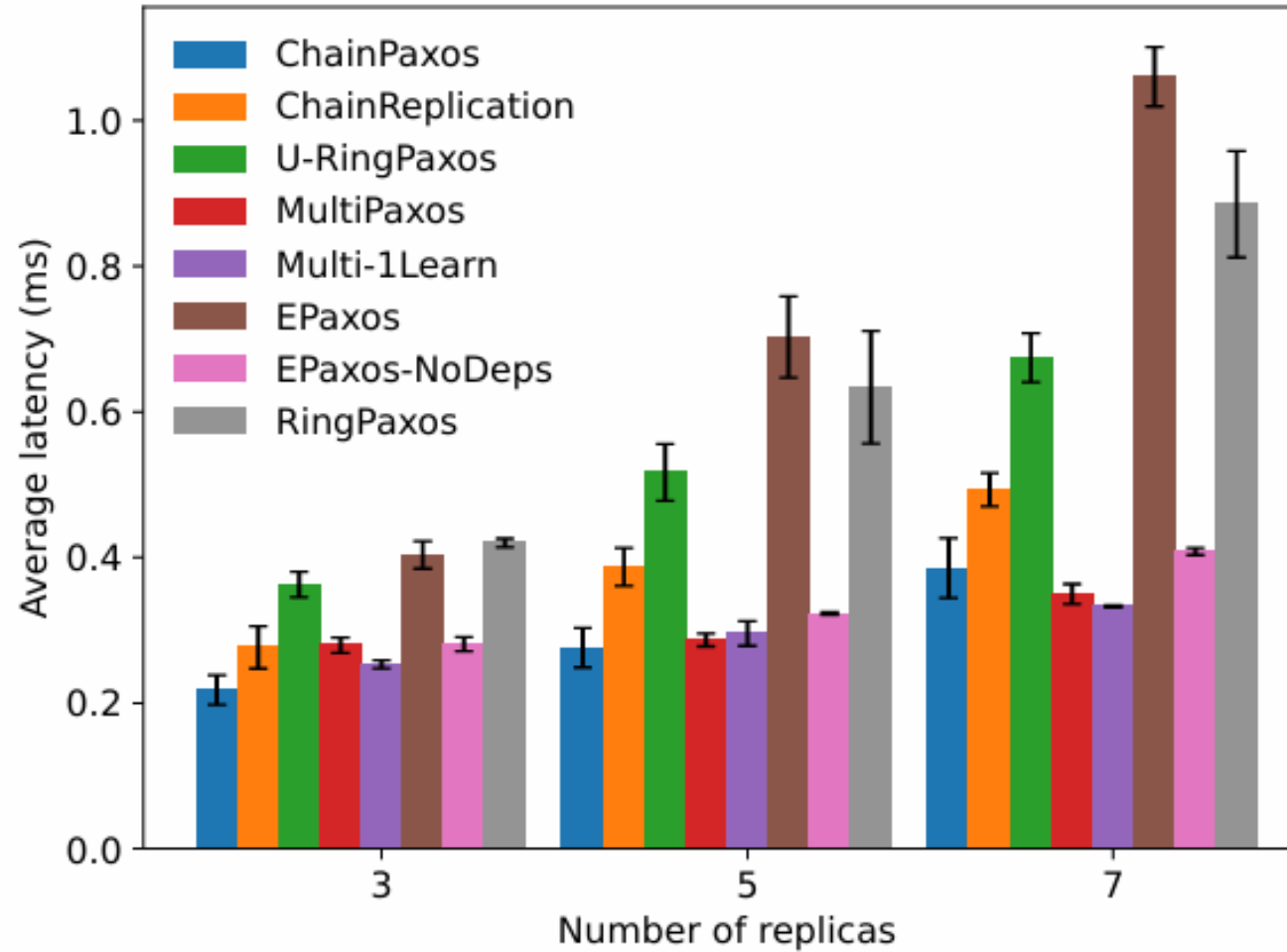


PrepareOK

PrepareOK
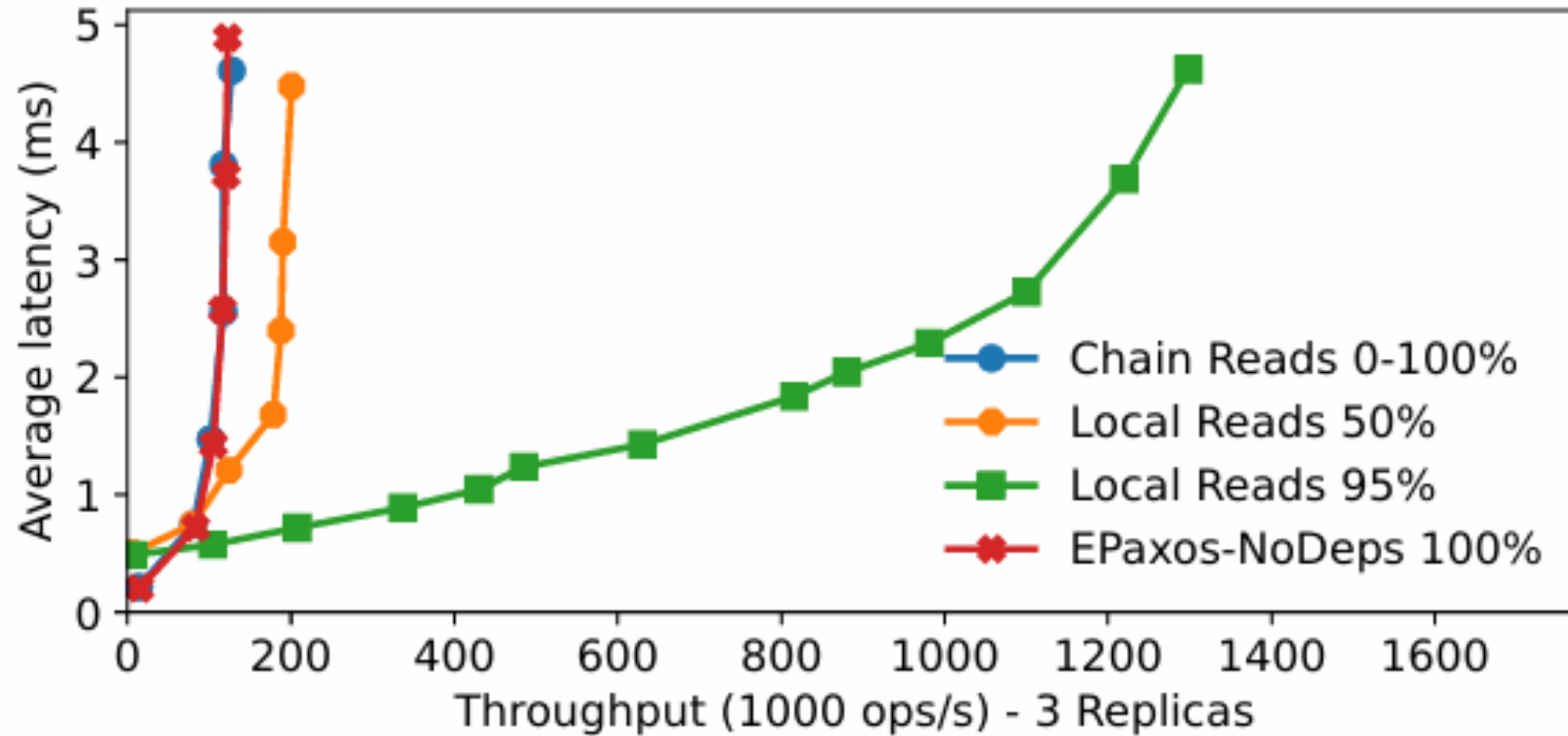
PrepareOK

# Leader Failure

New Leader

# Experiment Results: CPU Bottleneck

# Experiment Results: Latency Test

# Experiment Results: Local Read Test

# References

- Fouto, P., Preguiça, N., & Leitão, J. (2022). High Throughput Replication with Integrated Membership Management. USENIX Annual Technical Conference (USENIX ATC 22), 575–592.

# Thank you