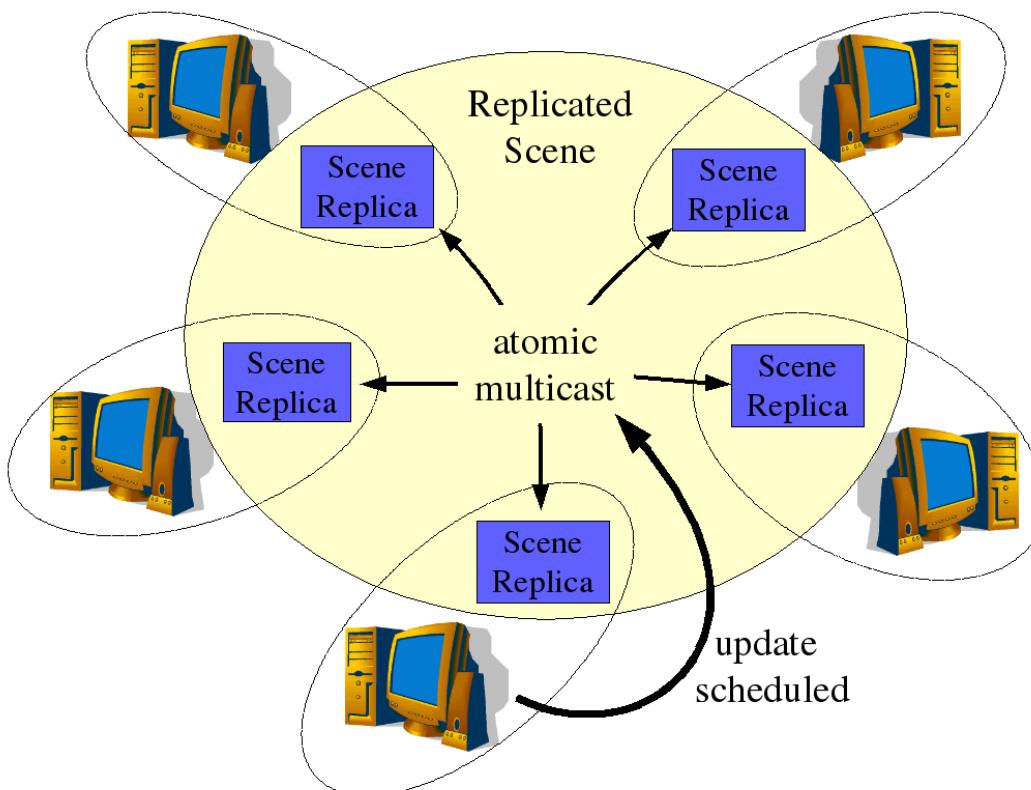
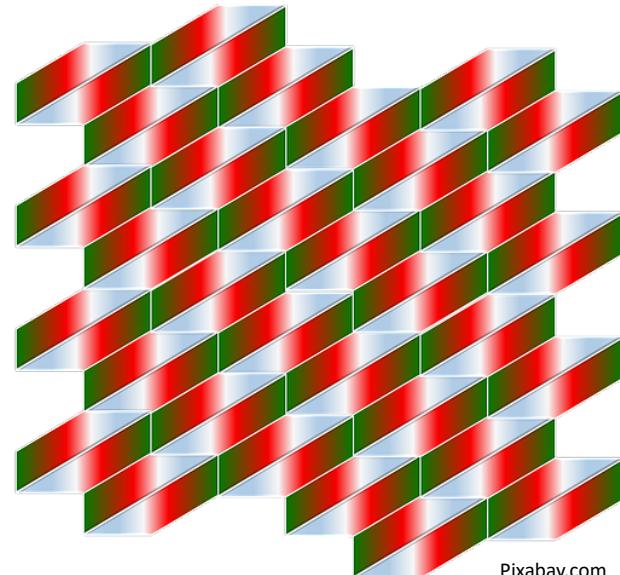


# Replication



# Outline

- Data replication
- Replication patterns
- Chain replication
- Gossiping



Pixabay.com



Pixabay.com

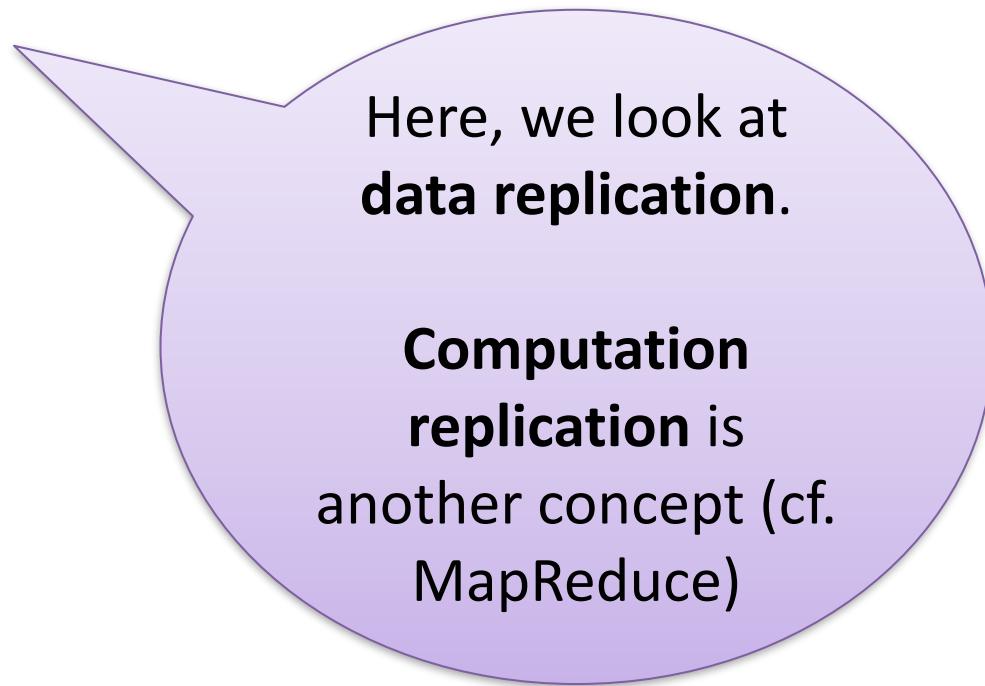
# DATA REPLICATION

# *Why Replicate?*

- Performance
- High availability
  - Fault tolerance

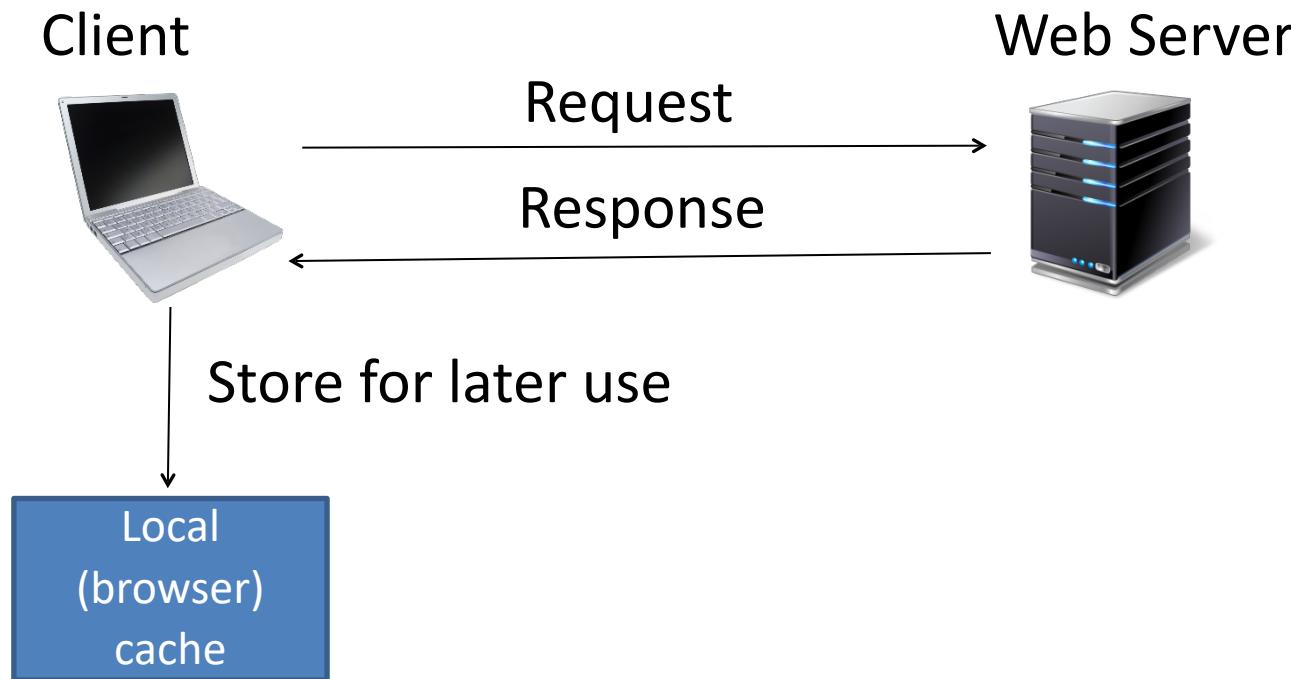
# *Why Replicate?*

- Performance
- High availability
  - Fault tolerance



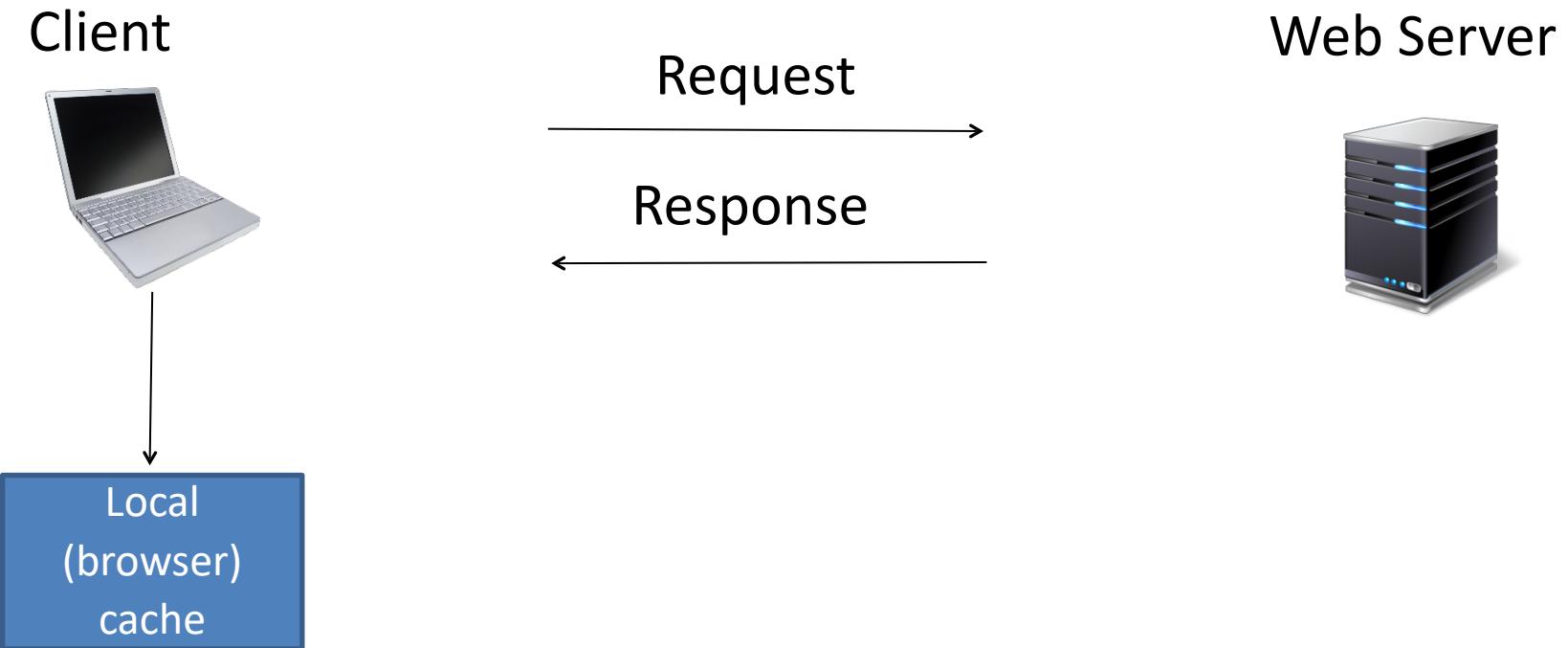
# Performance

## Caching data at browsers and proxy servers



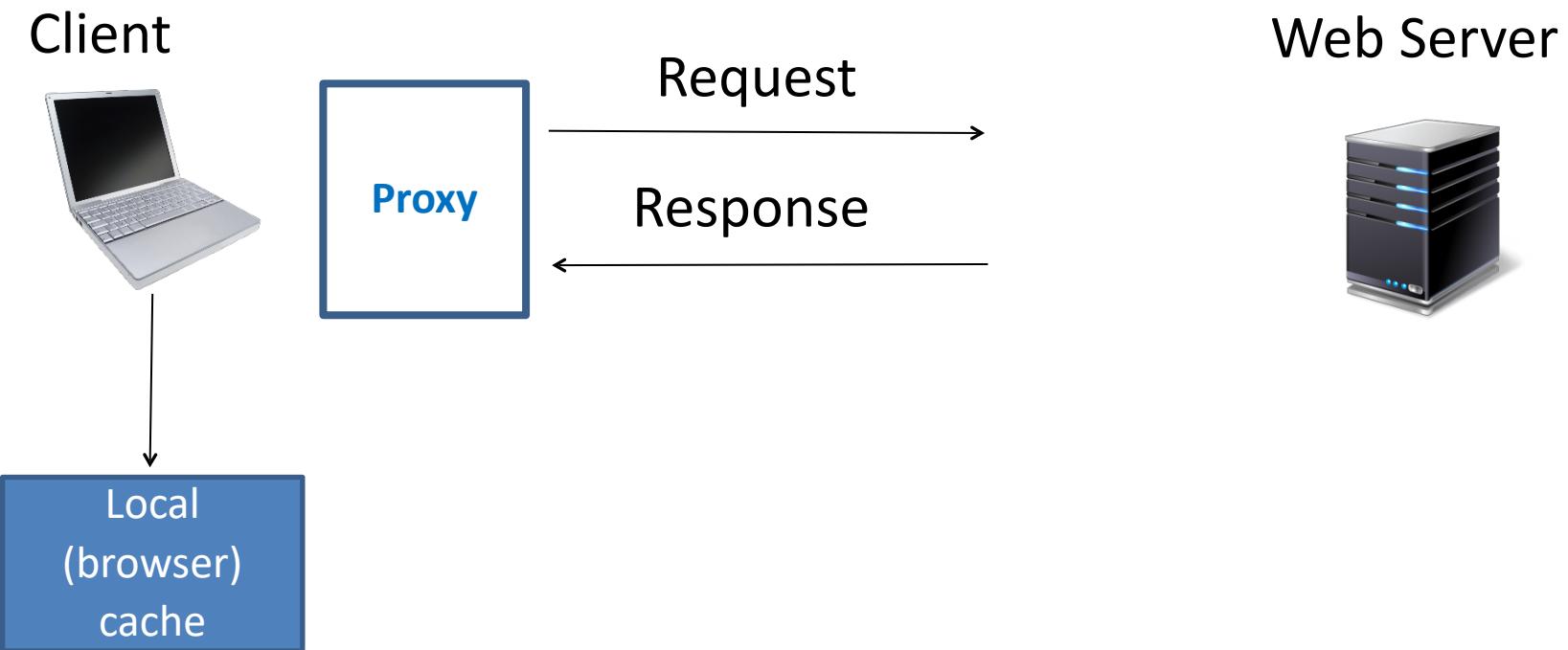
# Performance

## Caching data at browsers and proxy servers



# Performance

## Caching data at browsers and proxy servers



# Caching data at a proxy

Client



Proxy

Local  
(browser)  
cache

Configure Proxy Access to the Internet

No proxy

Auto-detect proxy settings for this network

Use system proxy settings

Manual proxy configuration

HTTP Proxy  Port

Also use this proxy for FTP and HTTPS

HTTPS Proxy  Port

FTP Proxy  Port

SOCKS Host  Port

SOCKS v4  SOCKS v5

Automatic proxy configuration URL  Reload

No proxy for

Example: .mozilla.org, .net.nz, 192.168.1.0/24

Connections to localhost, 127.0.0.1, and ::1 are never proxied.

Do not prompt for authentication if password is saved

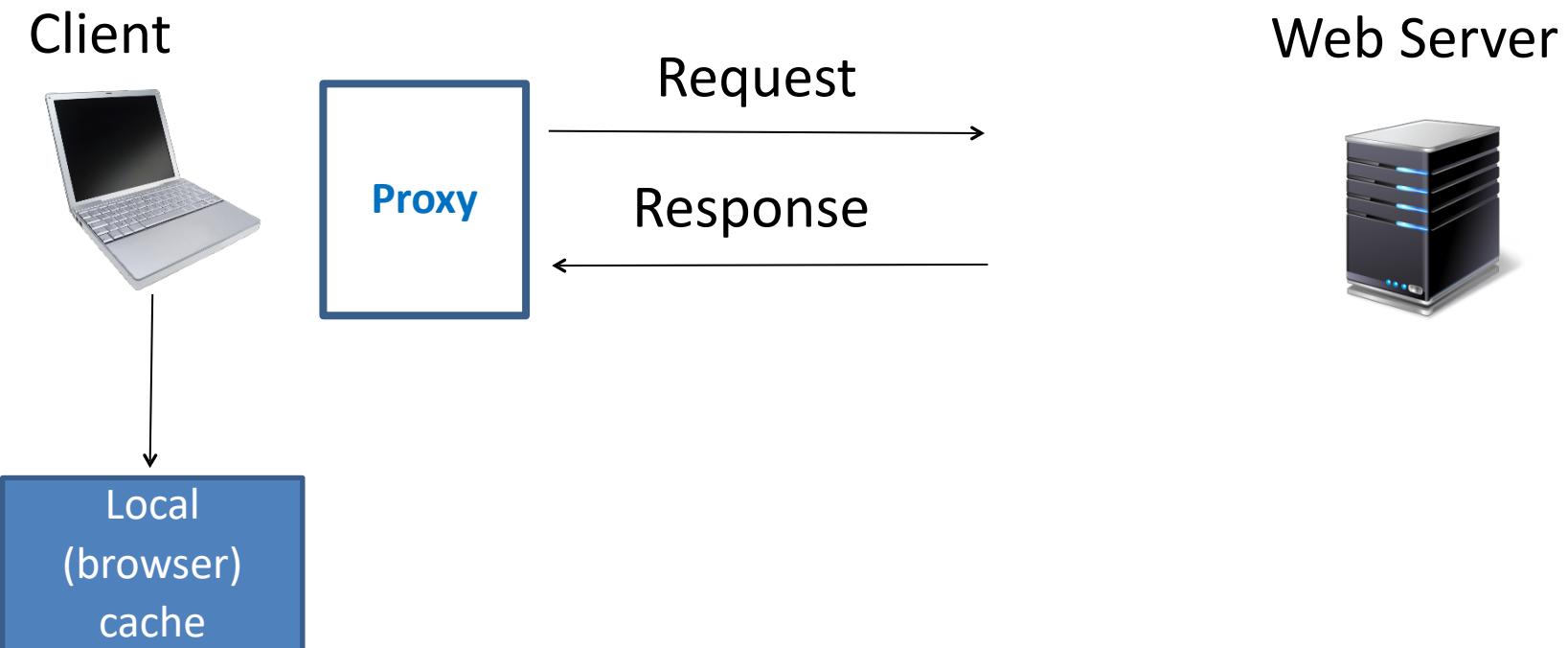
Proxy DNS when using SOCKS v5

Enable DNS over HTTPS

Use Provider

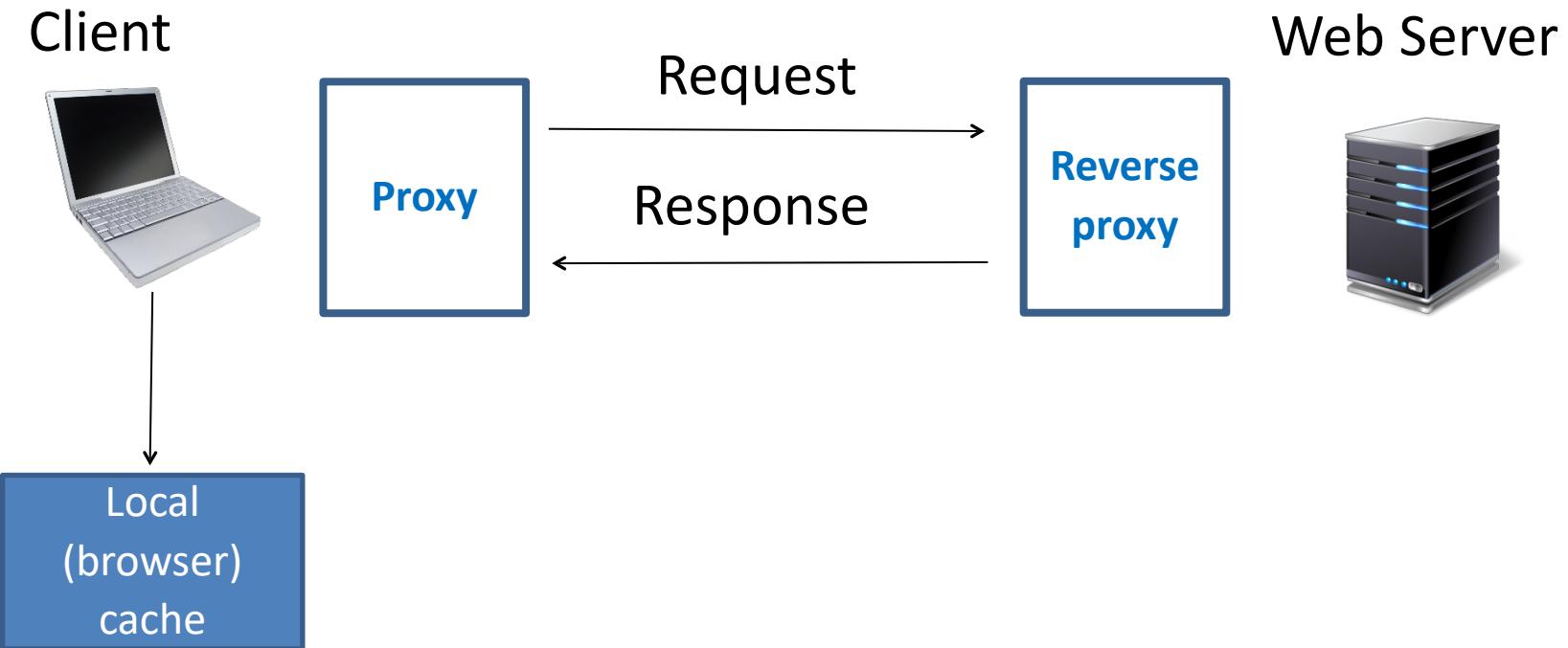
# Performance

## Caching data at browsers and proxy servers



# Performance

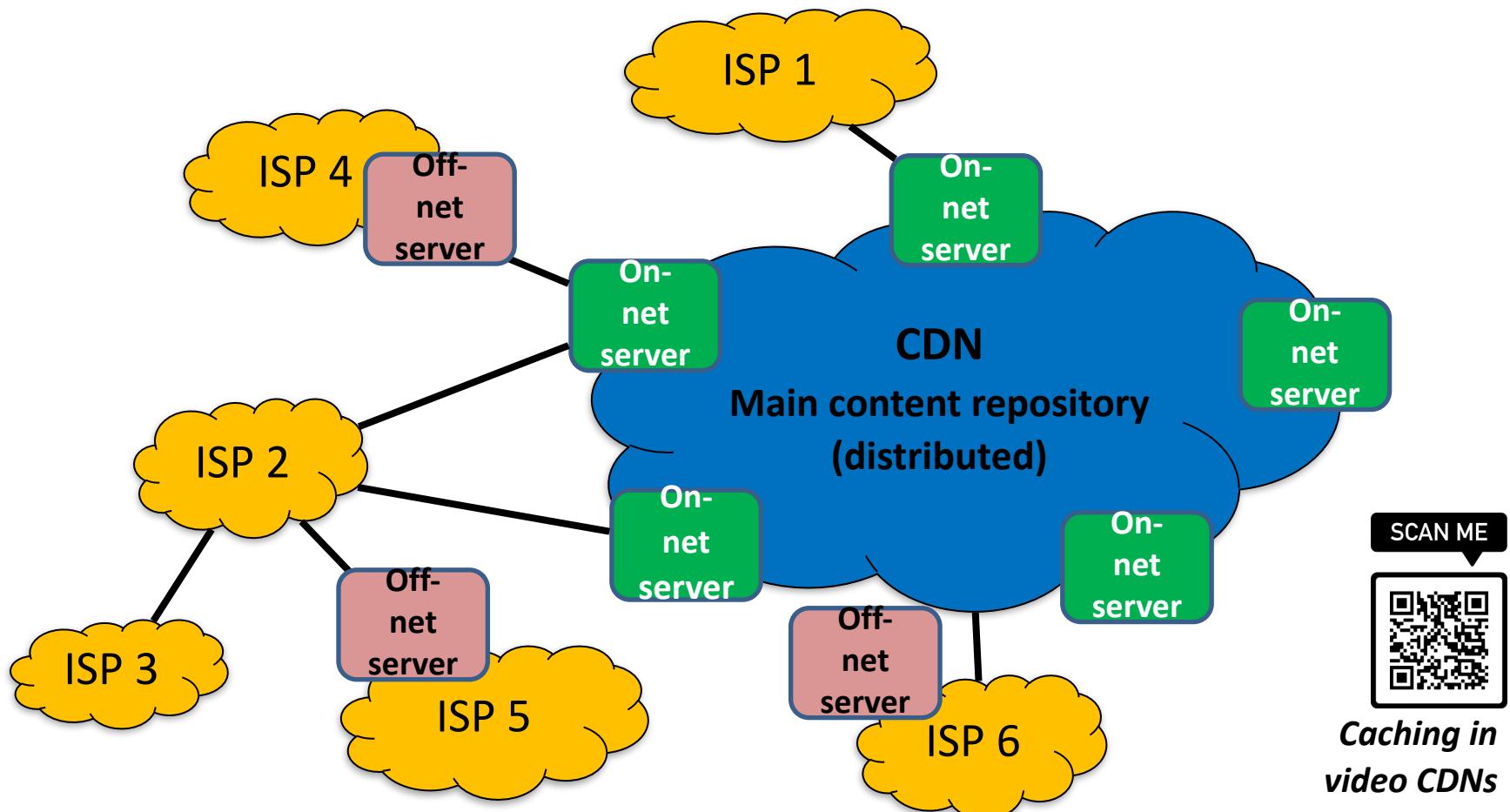
## Caching data at browsers and proxy servers





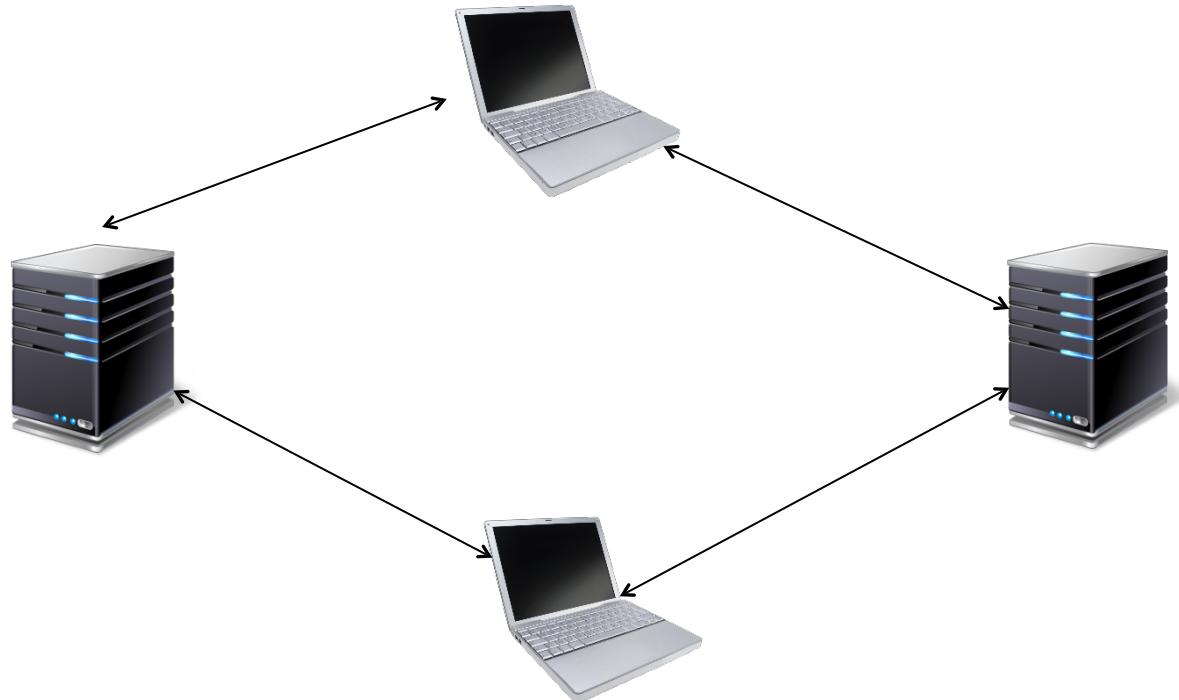
# Typical (video) CDN

(Content Delivery Network)



# High Availability

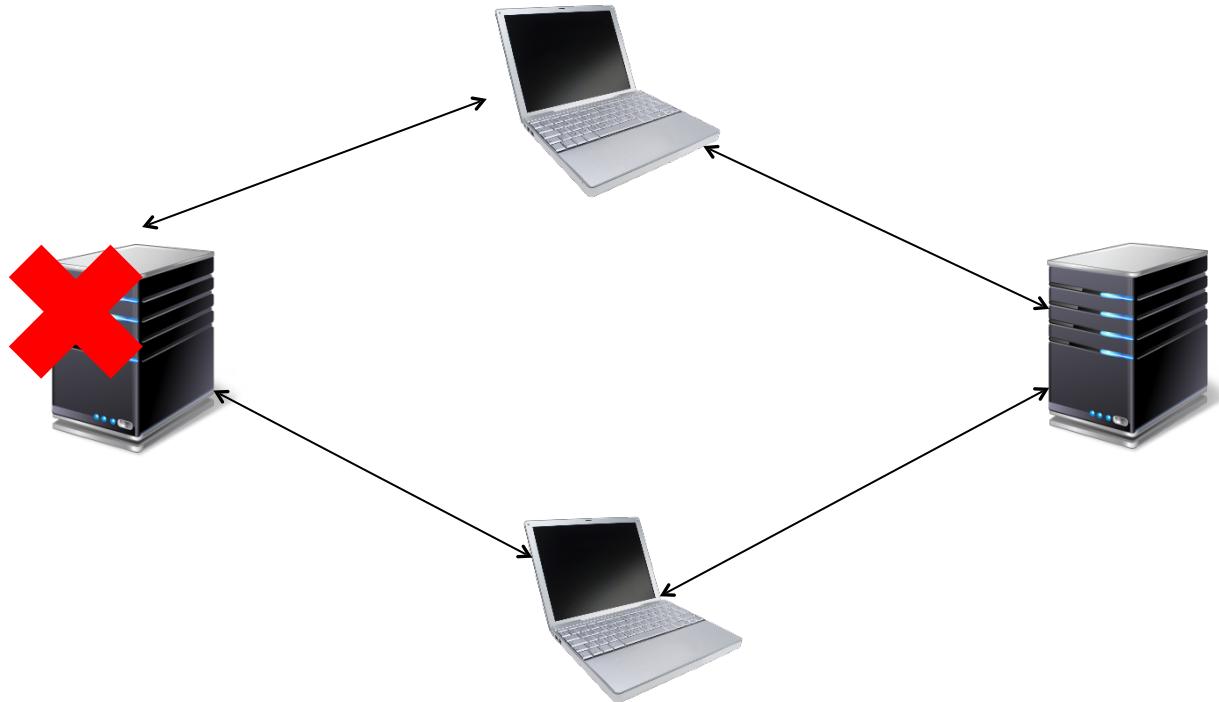
Upon crashes, data offered/retrieved by/from replica



[1] The availability of a service by replicating its nodes grows.

# High Availability

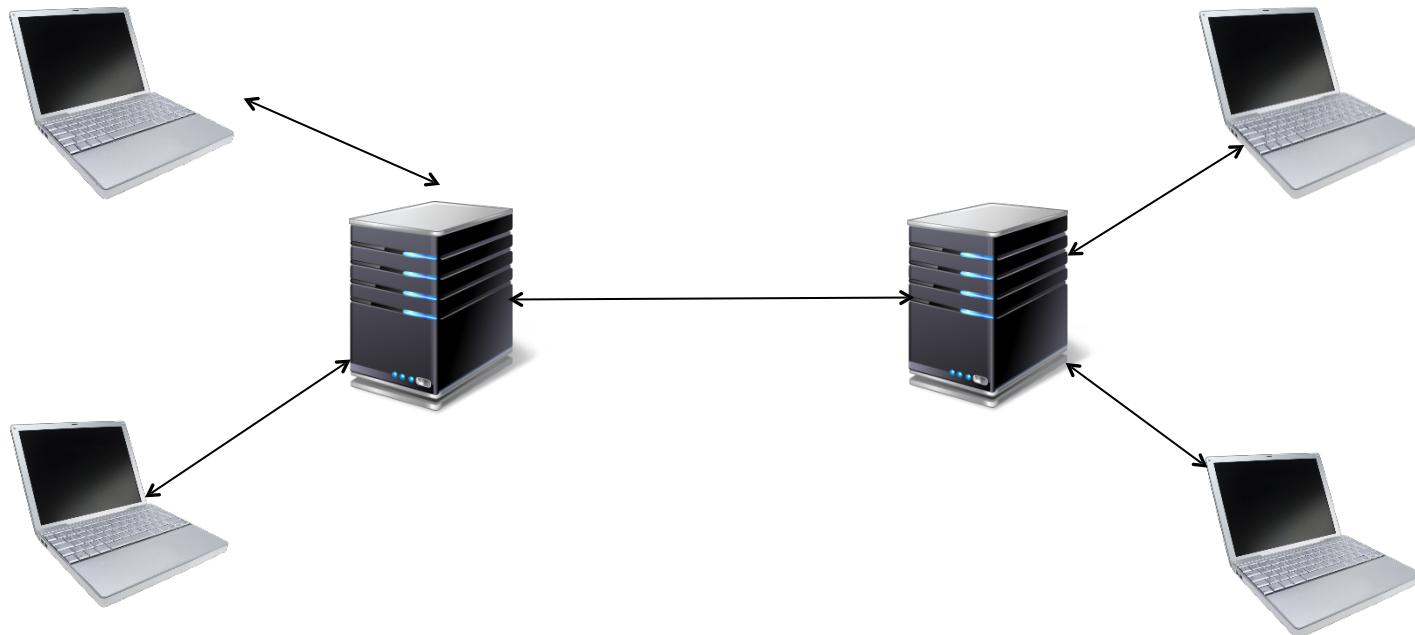
Upon crashes, data offered/retrieved by/from replica



[1] The availability of a service by replicating its nodes grows.

# High Availability: Network Partitioning

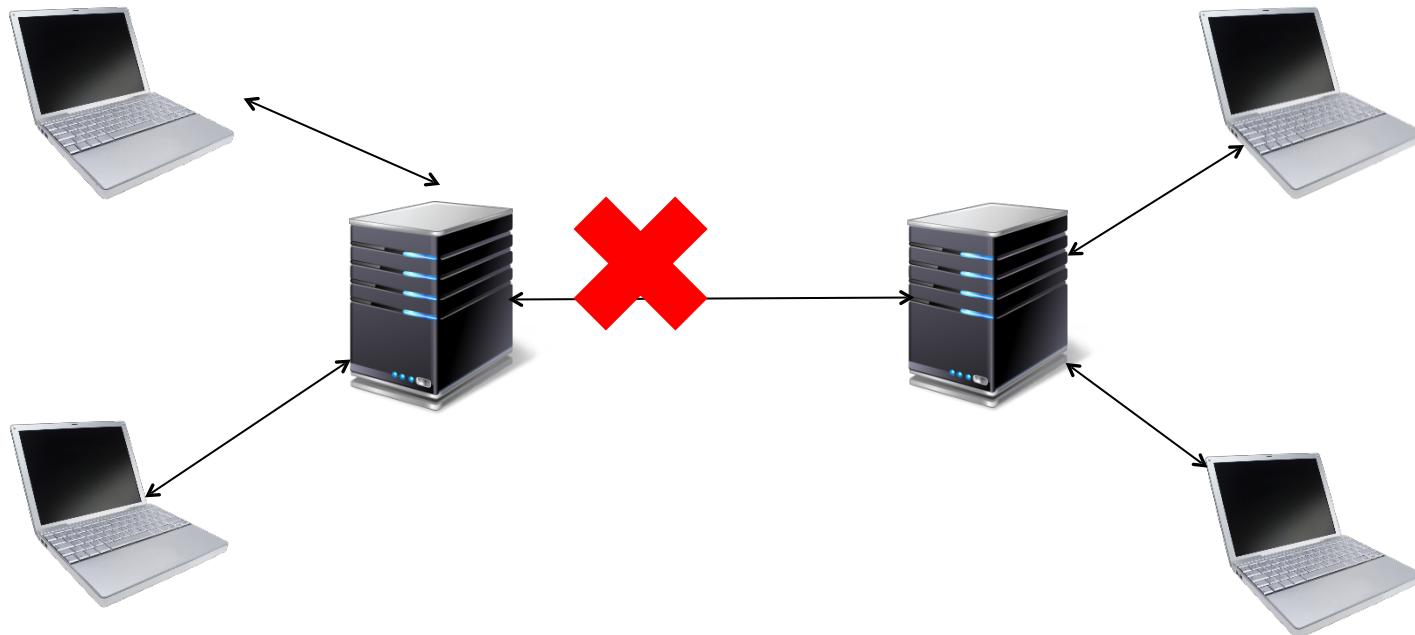
Upon network partitioning, data available to clients from within partition



**Partition tolerance:** A system that continues to operate in face of network partitions

# High Availability: Network Partitioning

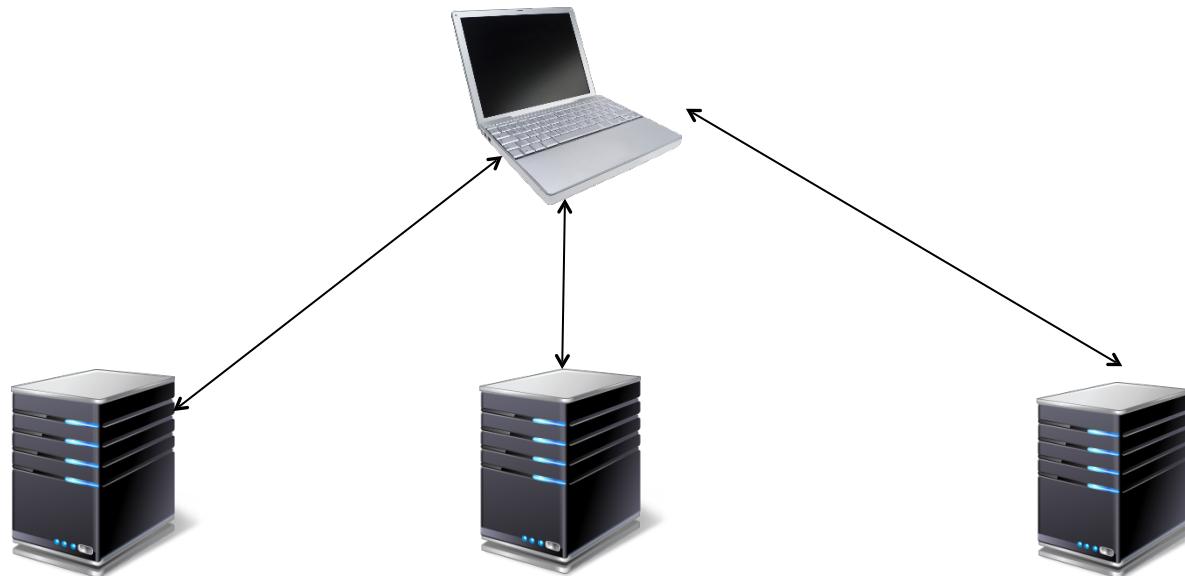
Upon network partitioning, data available to clients from within partition



**Partition tolerance:** A system that continues to operate in face of network partitions

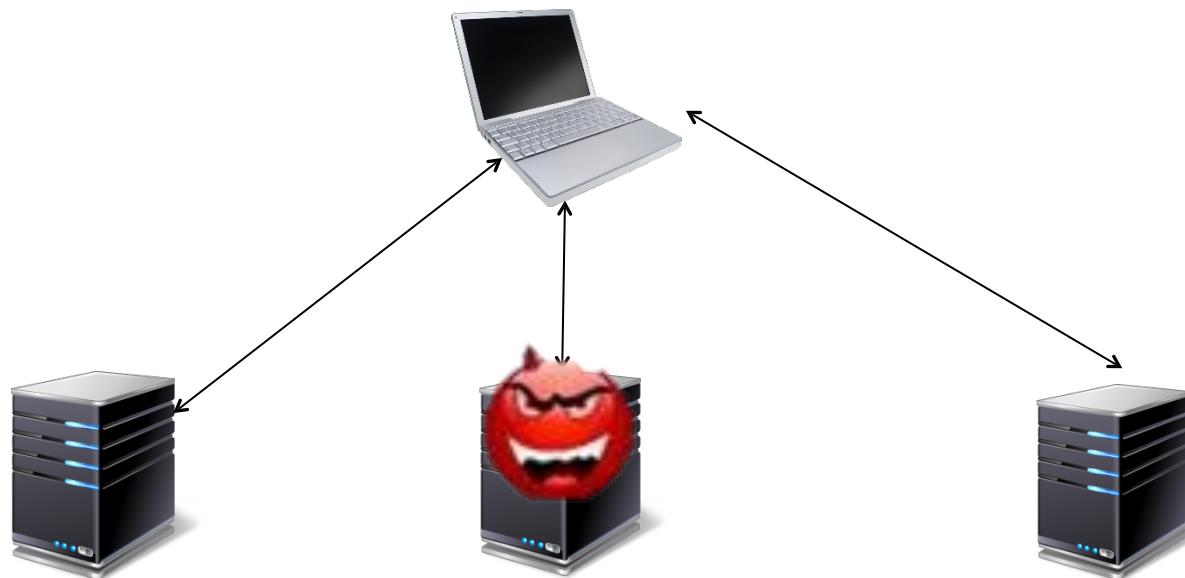
# High Availability: Fault Tolerance

Providing reliable service in face of faulty nodes  
(not just crashes, but also arbitrary failures!)



# High Availability: Fault Tolerance

Providing reliable service in face of faulty nodes  
(not just crashes, but also arbitrary failures!)



# “Cost” of Replication

- Not just cost of **storing** additional copies of data
- Cost to **keep replicas up to date** in face of updates
- *How to deal with **stale** (out-of-date) data at replicas?*

# Self-study Questions

- What is the “cache hierarchy” for web data? Think of web pages with varying complexity of content; understand where each content element is cached.
- Is this cache hierarchy the same for any distributed data processing system?
- Find some popular examples of open-source caching technology and understand their position in the cache hierarchy.
- Think of alternatives for recovering a crashed replica – how does it catch up to the current state?
- Determine conflicts when issuing reads and writes concurrently to replicas.



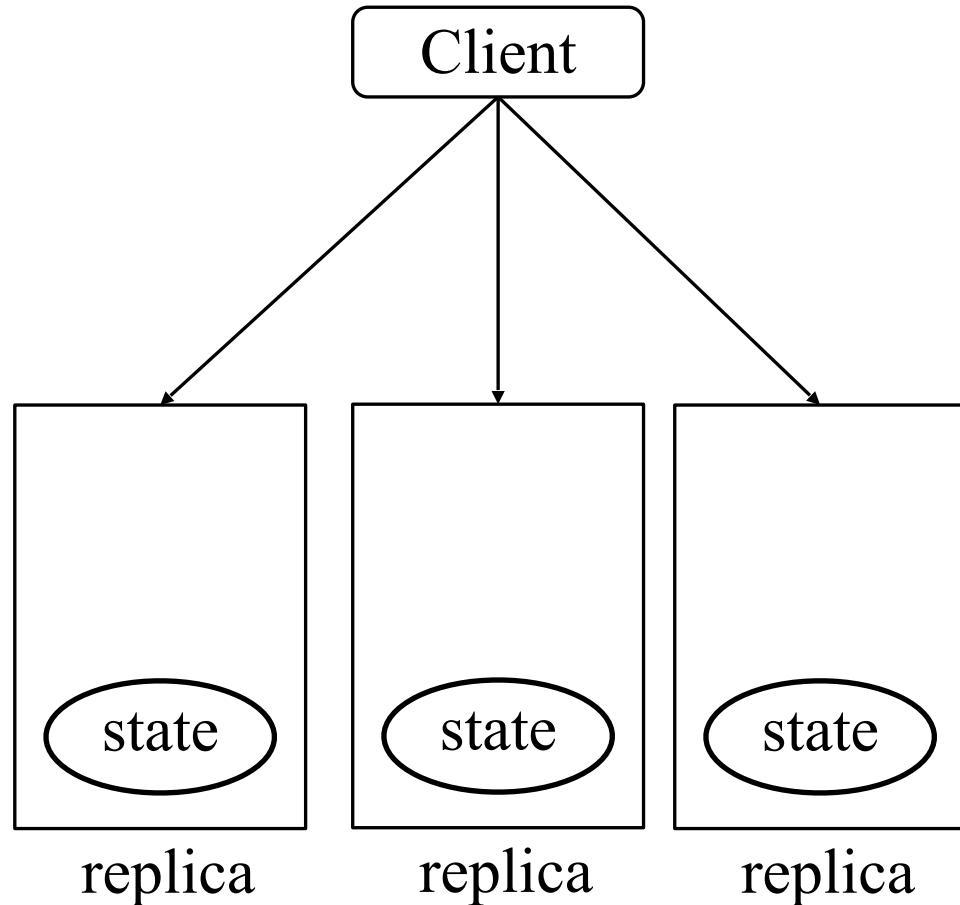


Pixabay.com

# REPLICATION PATTERNS WITH UPDATE PROCESSING

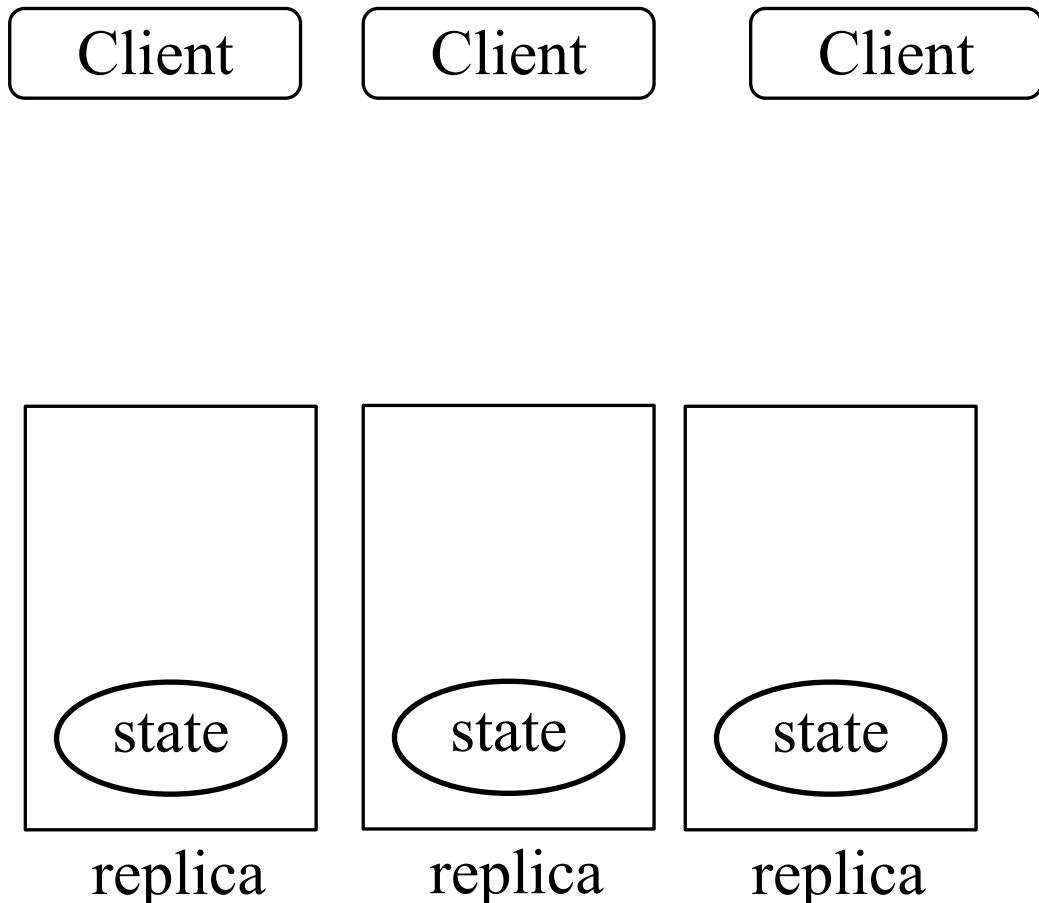
# Replication Pattern

- **Active replication**
- Requests can be reads or writes
- Replicas can crash



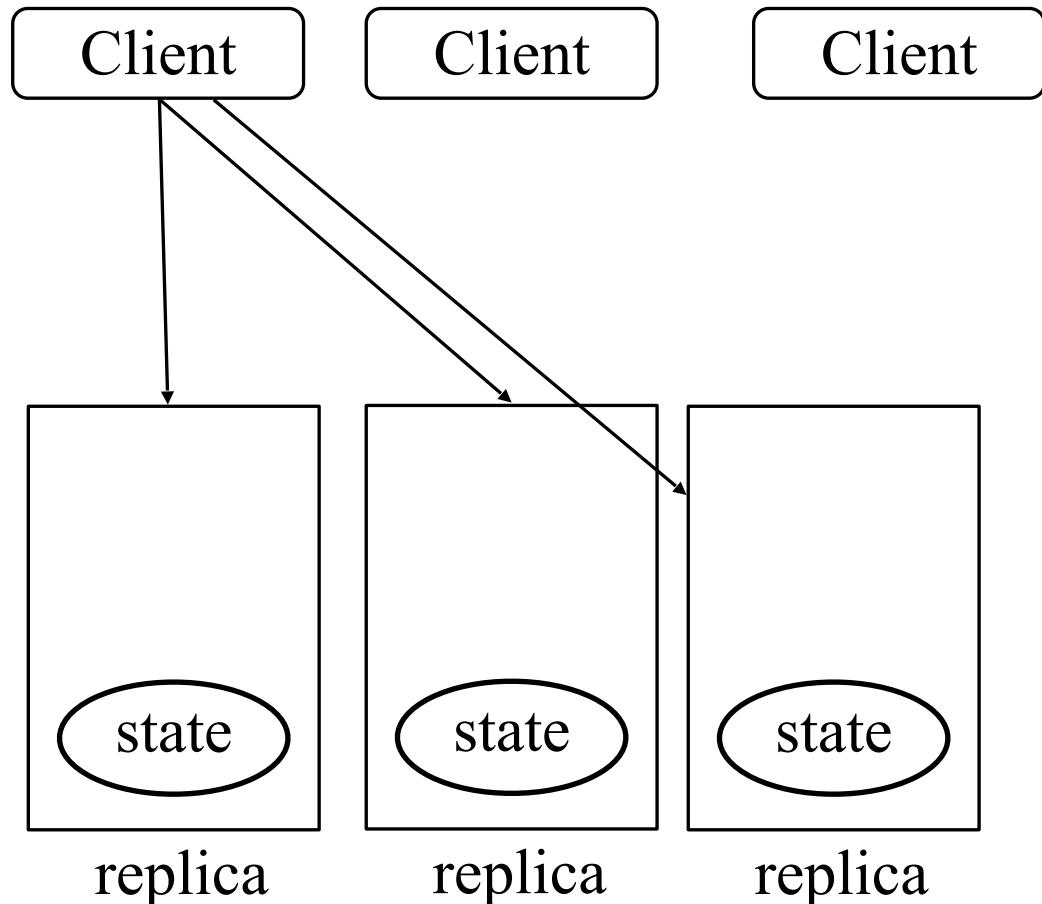
# Replication Pattern

- **Active replication**
- Clients send requests to every replica



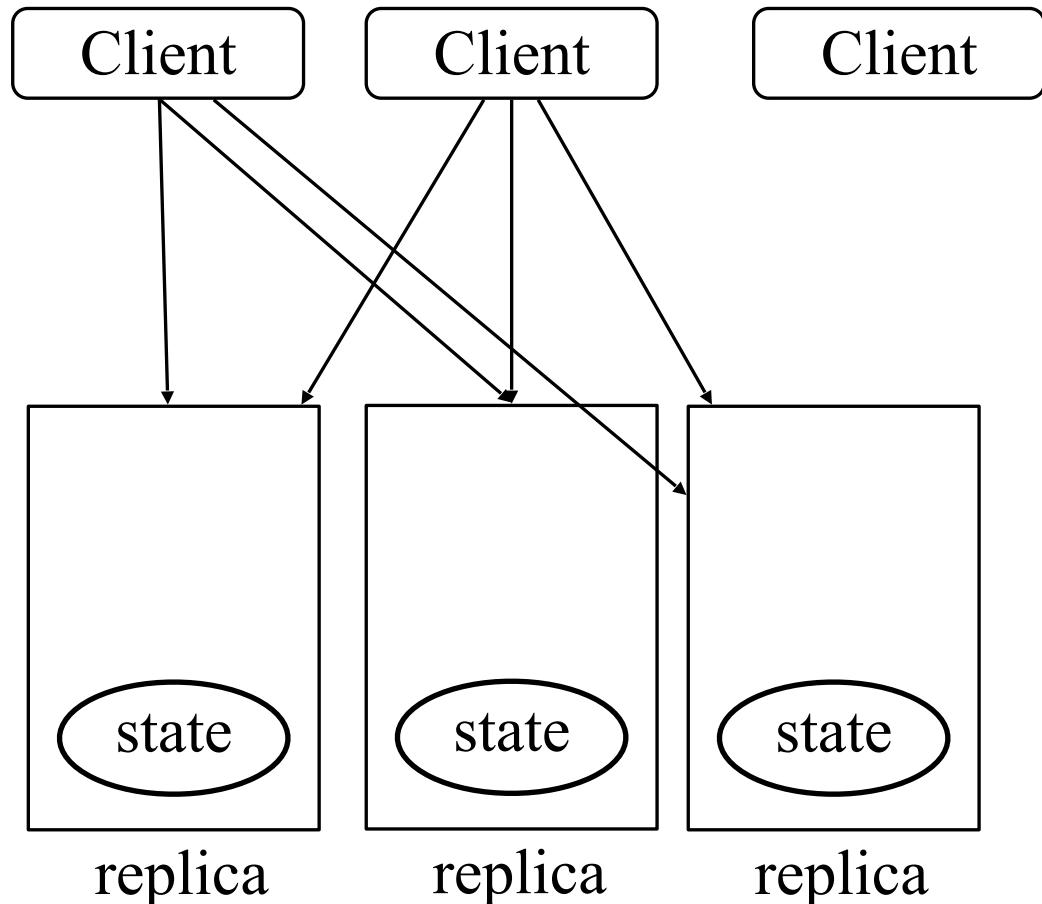
# Replication Pattern

- **Active replication**
- Clients send requests to every replica

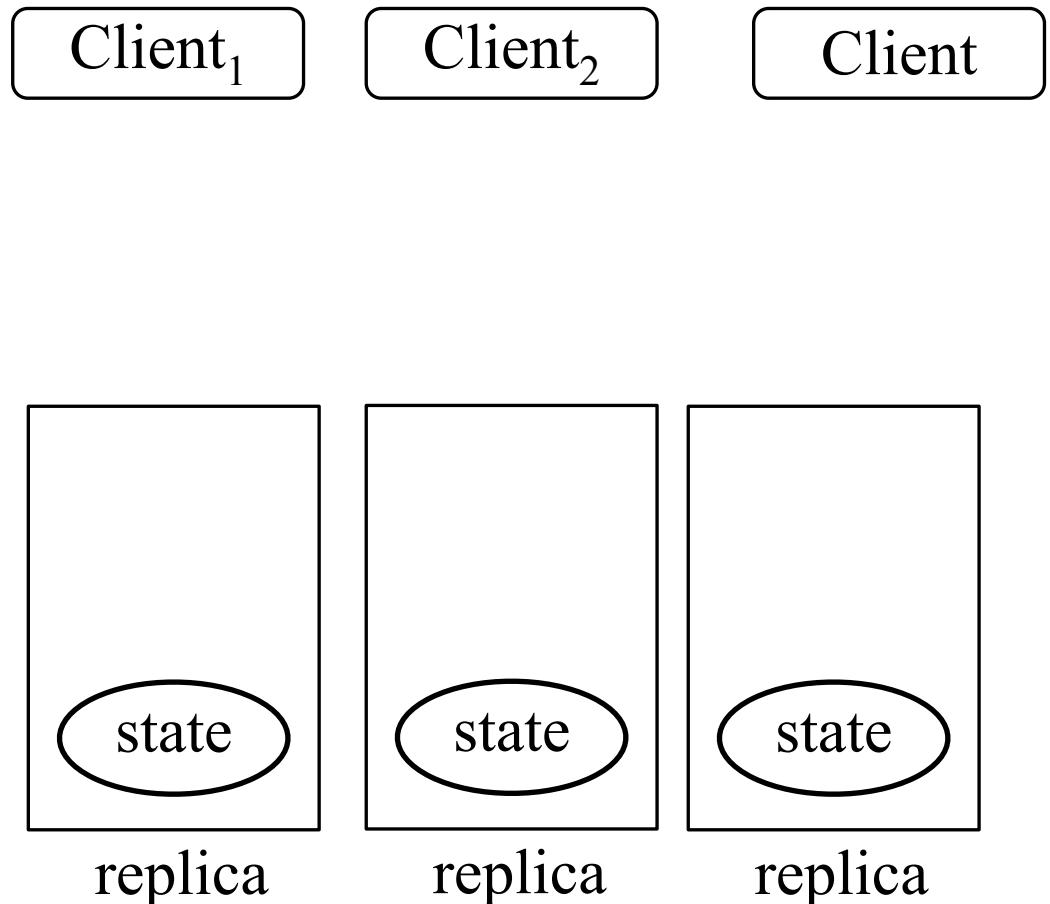


# Replication Pattern

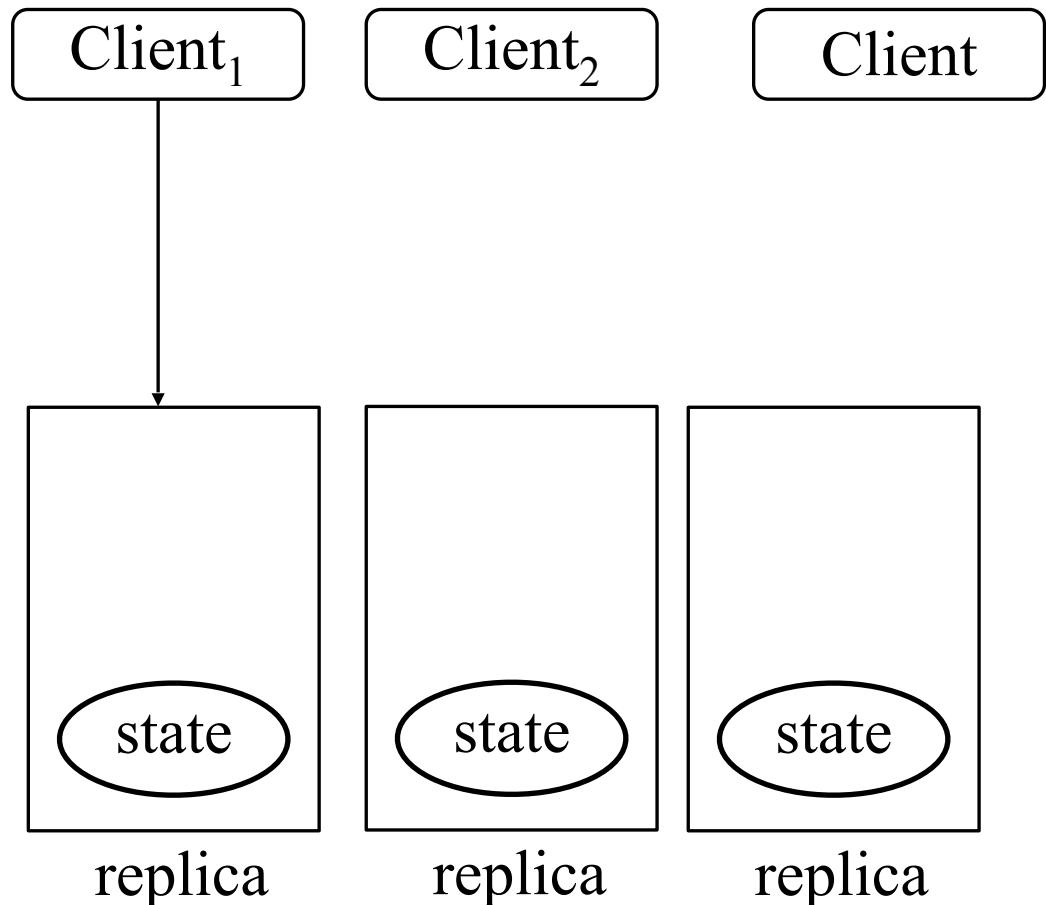
- **Active replication**
- Clients send requests to every replica



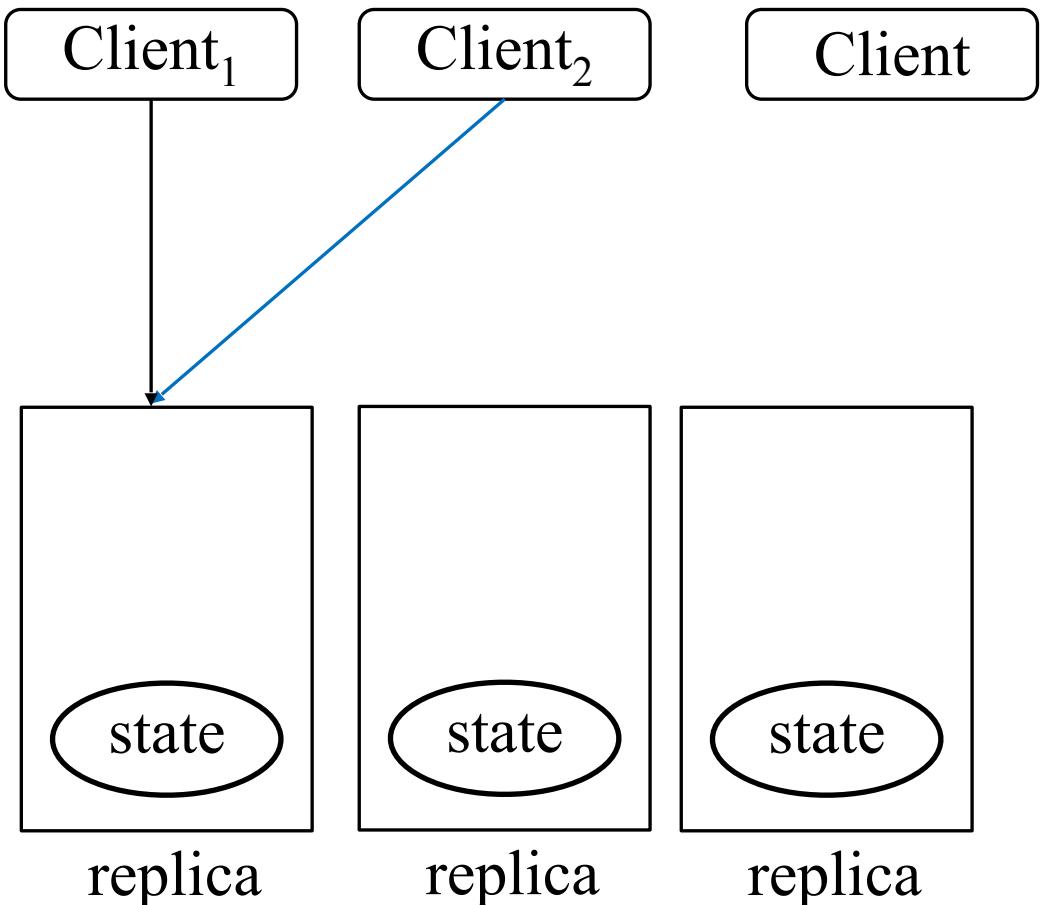
# Replication Pattern



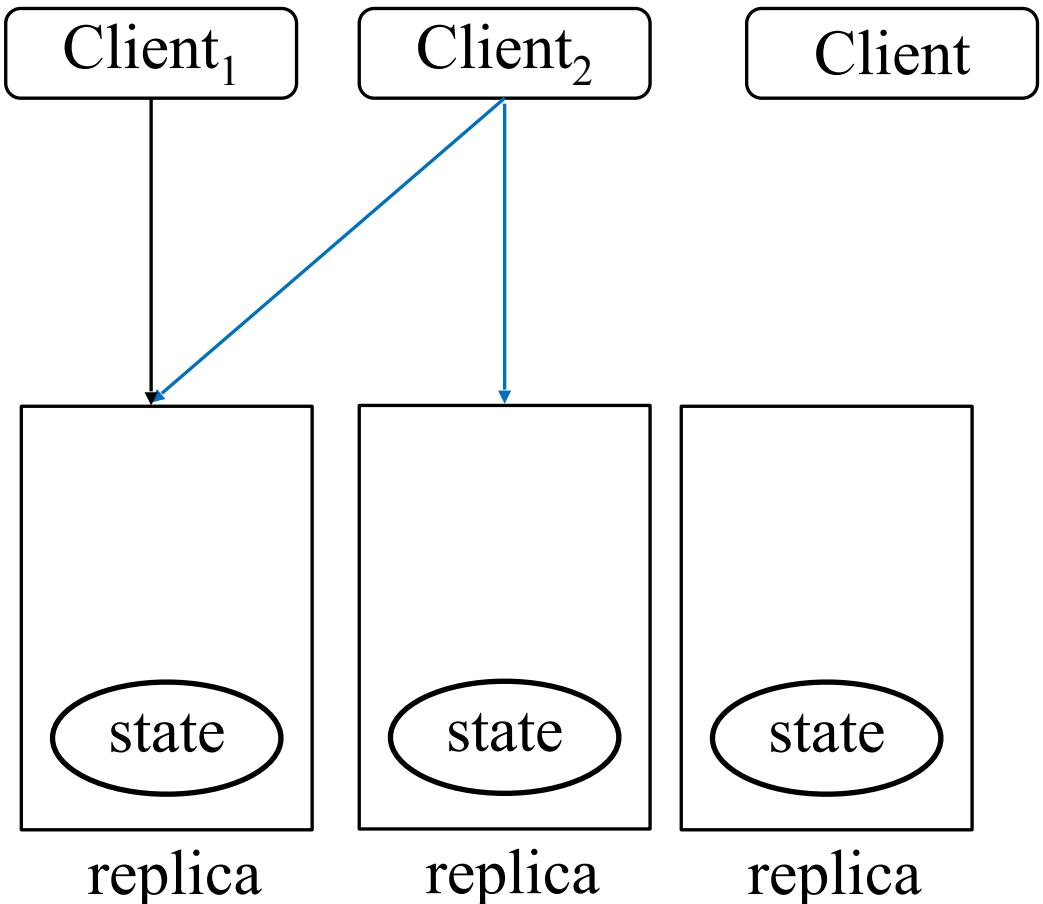
# Replication Pattern



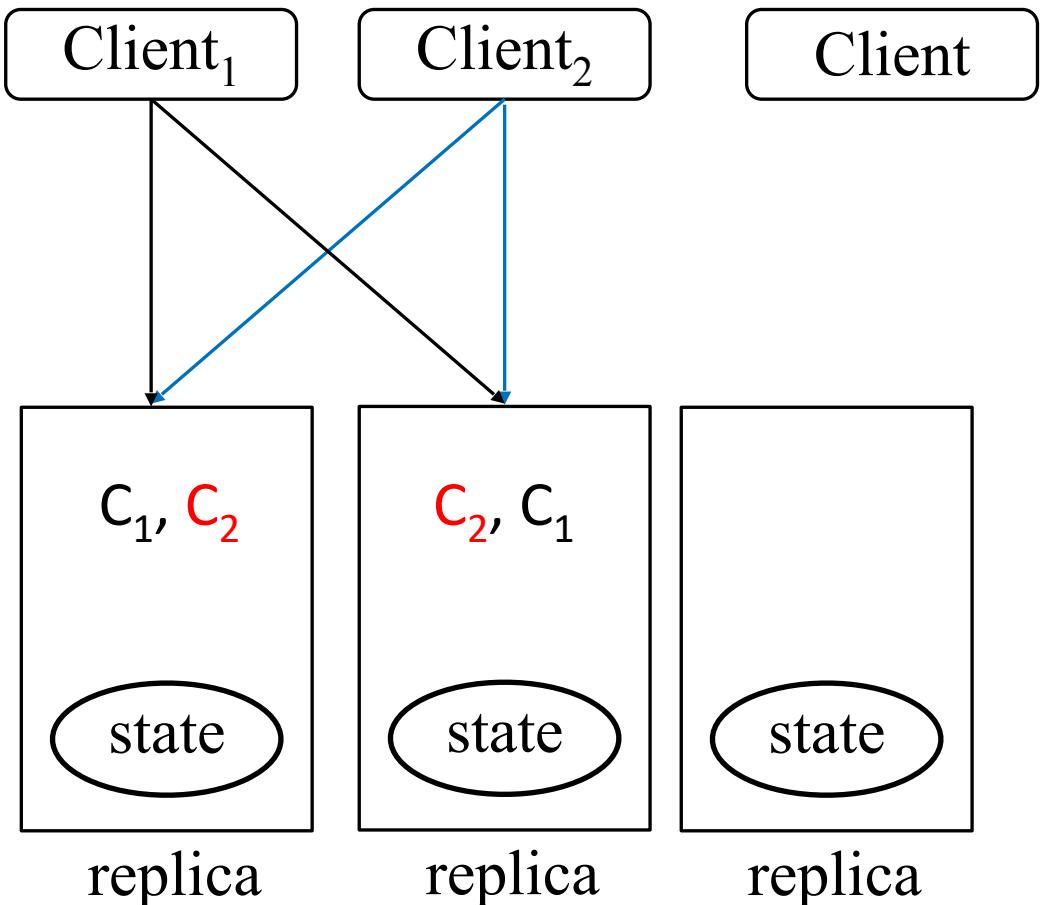
# Replication Pattern



# Replication Pattern

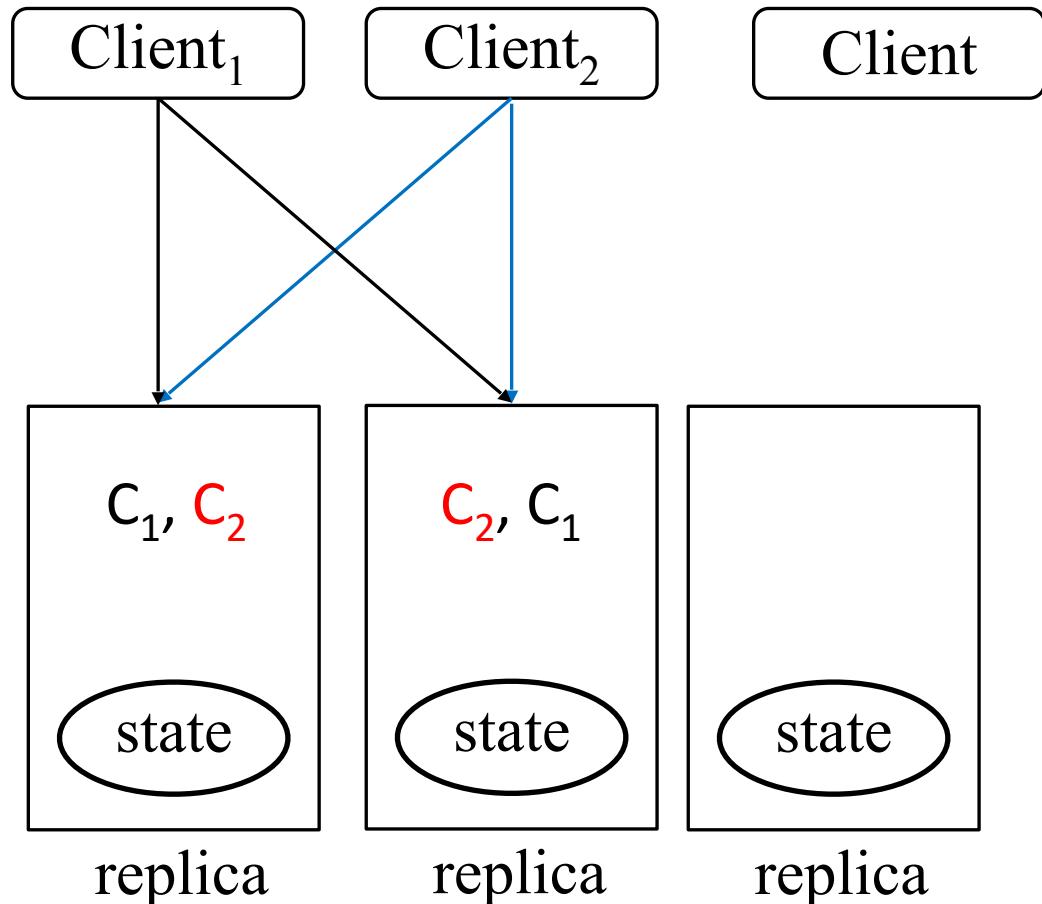


# Replication Pattern



# Replication Pattern

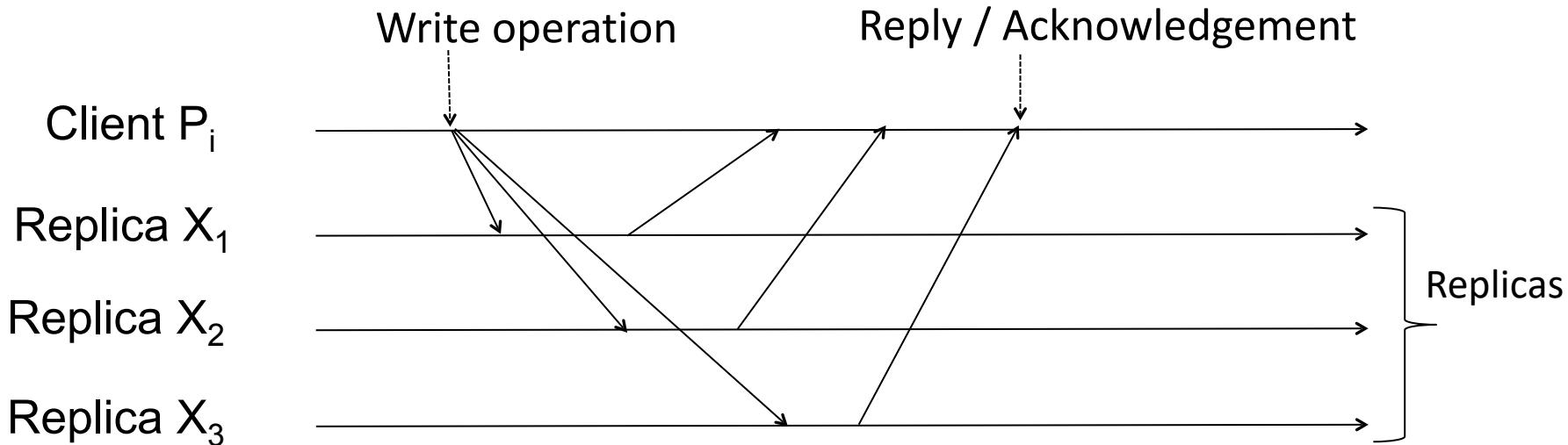
- **Active replication**
- Replicas may diverge



# Active Replication

- Requires **total order broadcast** to guarantee each replica receives
  - all requests (from all clients)
  - in the same order
- Like client-server, “natural” to think about
- Fast to get a response, first result received, unless Byzantine failure assumptions (majority result)

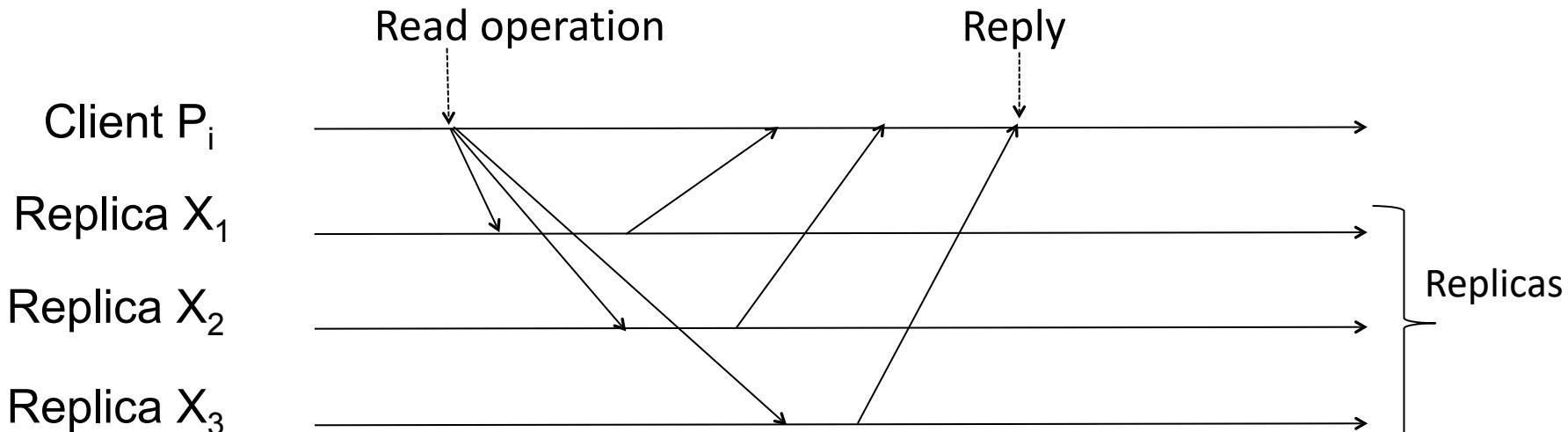
# Active Replication



Configuration service:

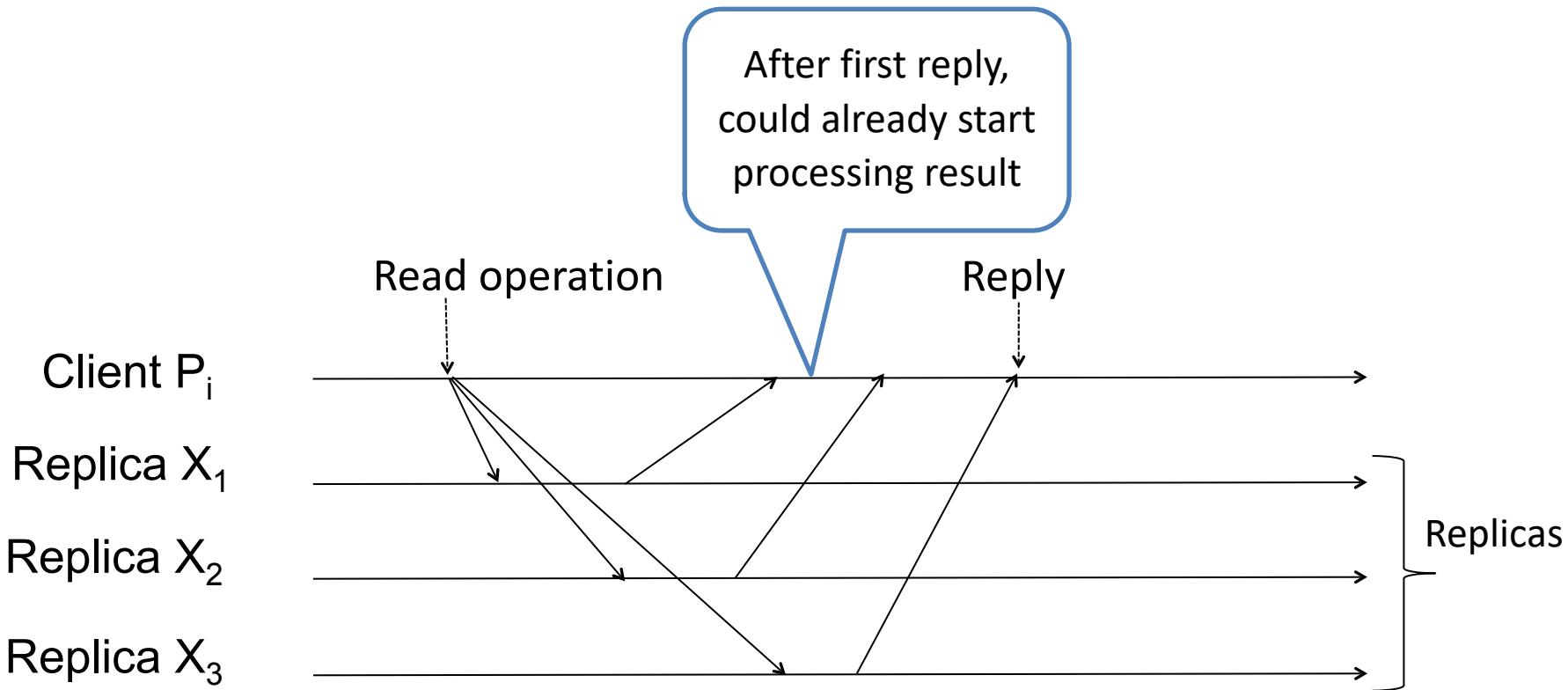
- Failure detection
- Configuration management

# Active Replication



Alternative: Read from a quorum

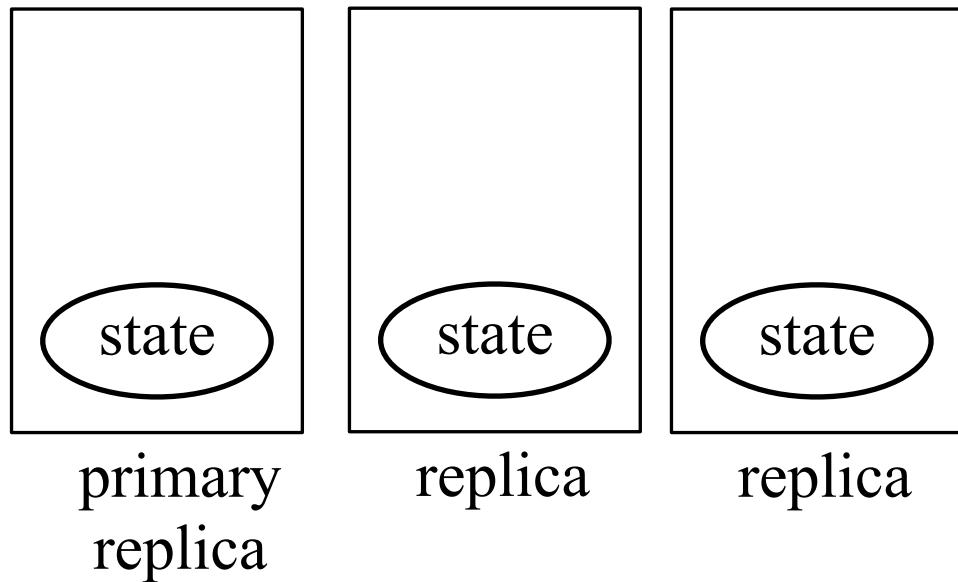
# Active Replication



Alternative: Read from a quorum

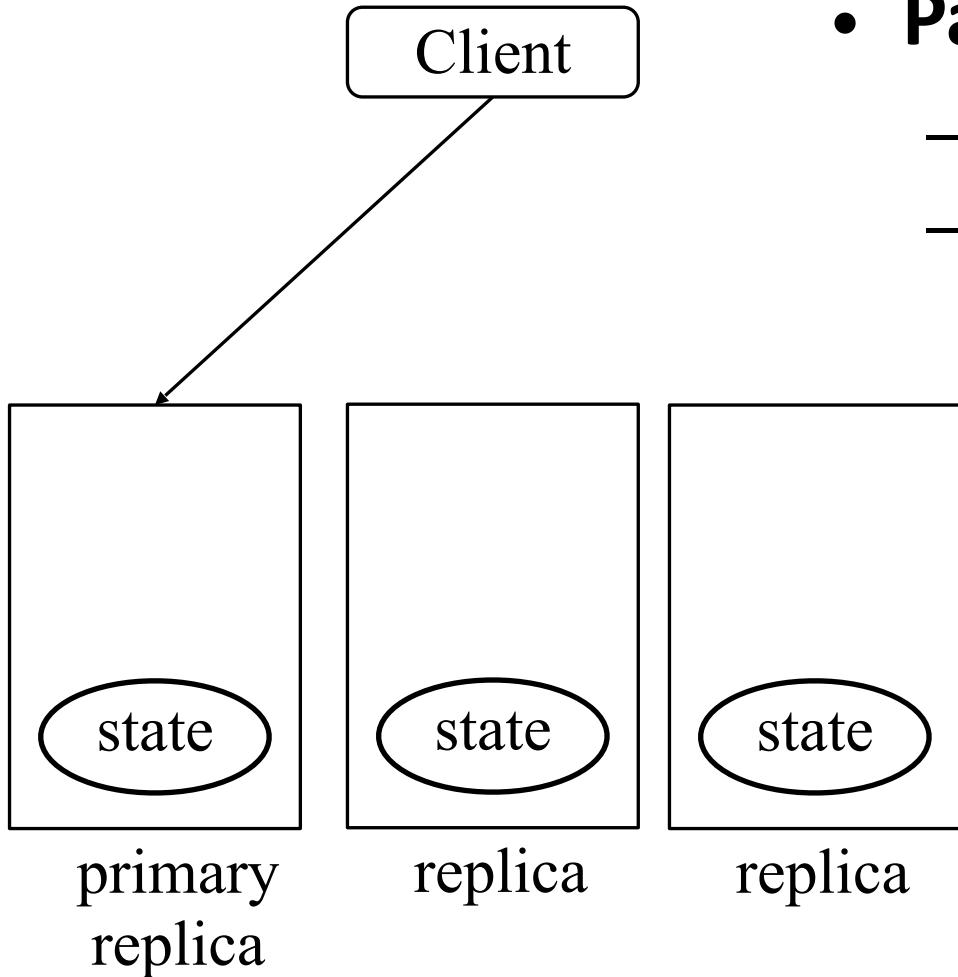
# Replication Pattern

Client



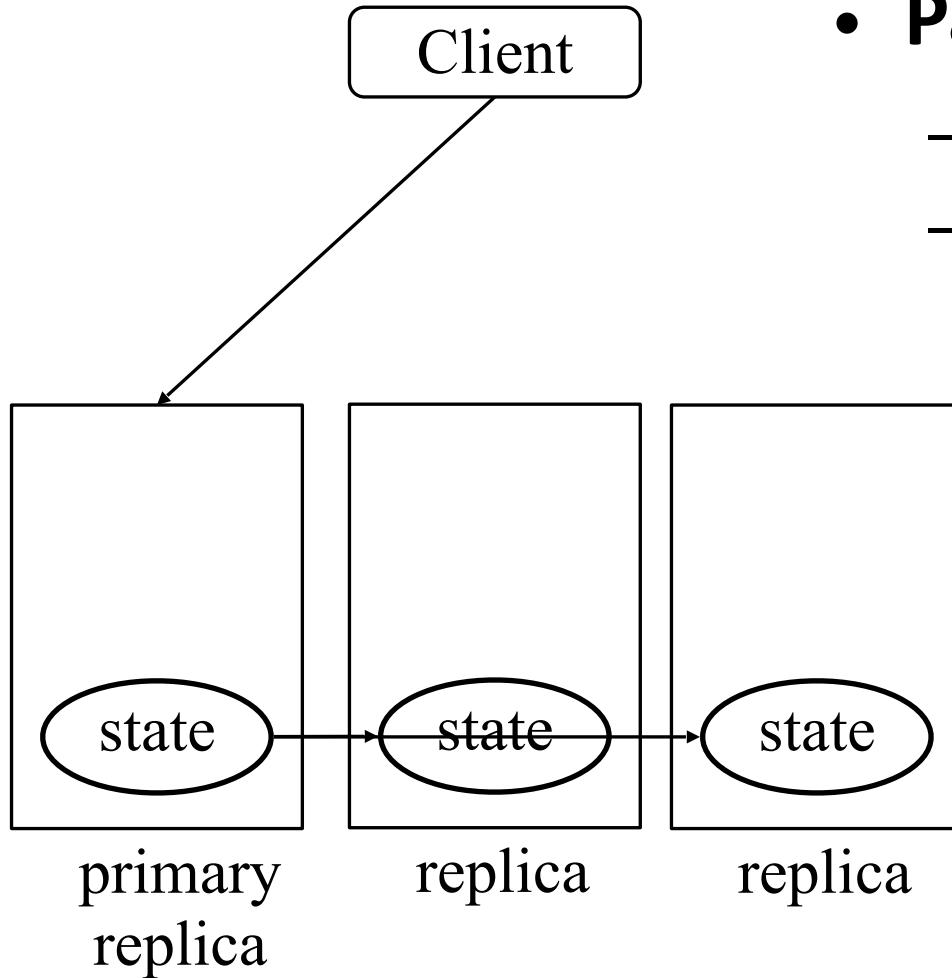
- **Passive replication**
  - Primary-backup
  - Multi-primary

# Replication Pattern



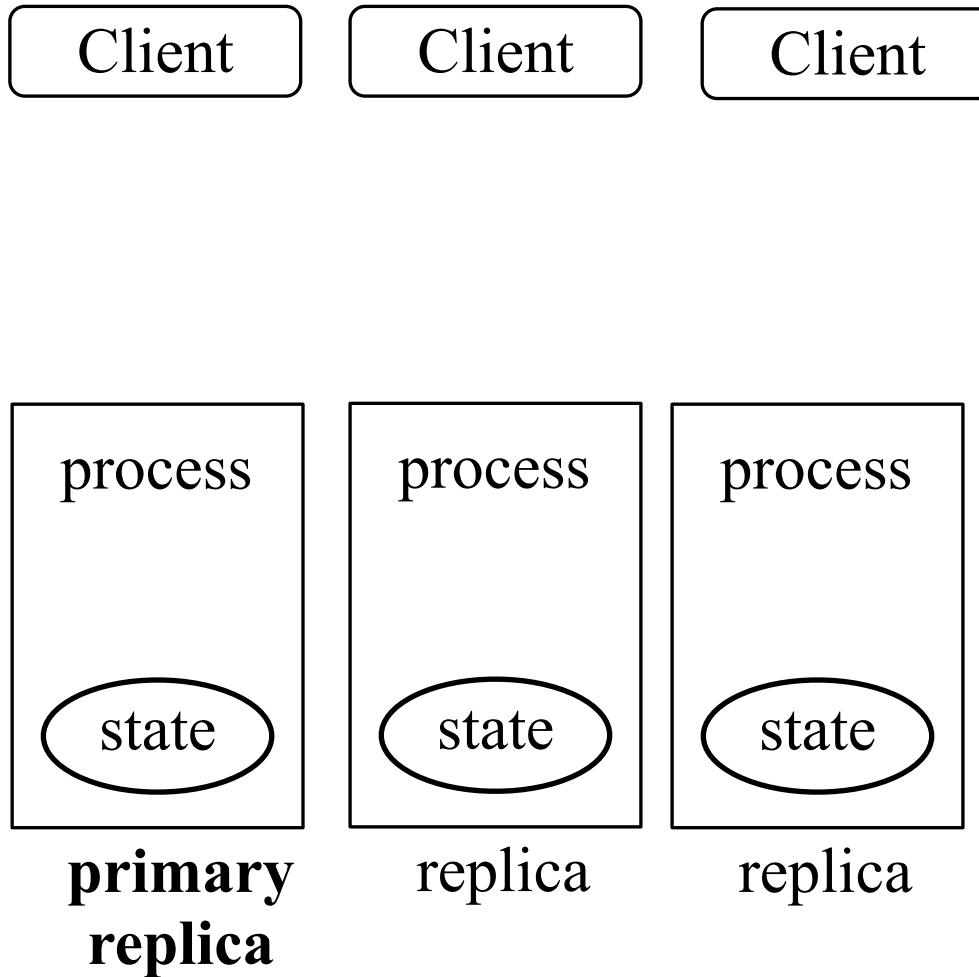
- **Passive replication**
  - Primary-backup
  - Multi-primary

# Replication Pattern

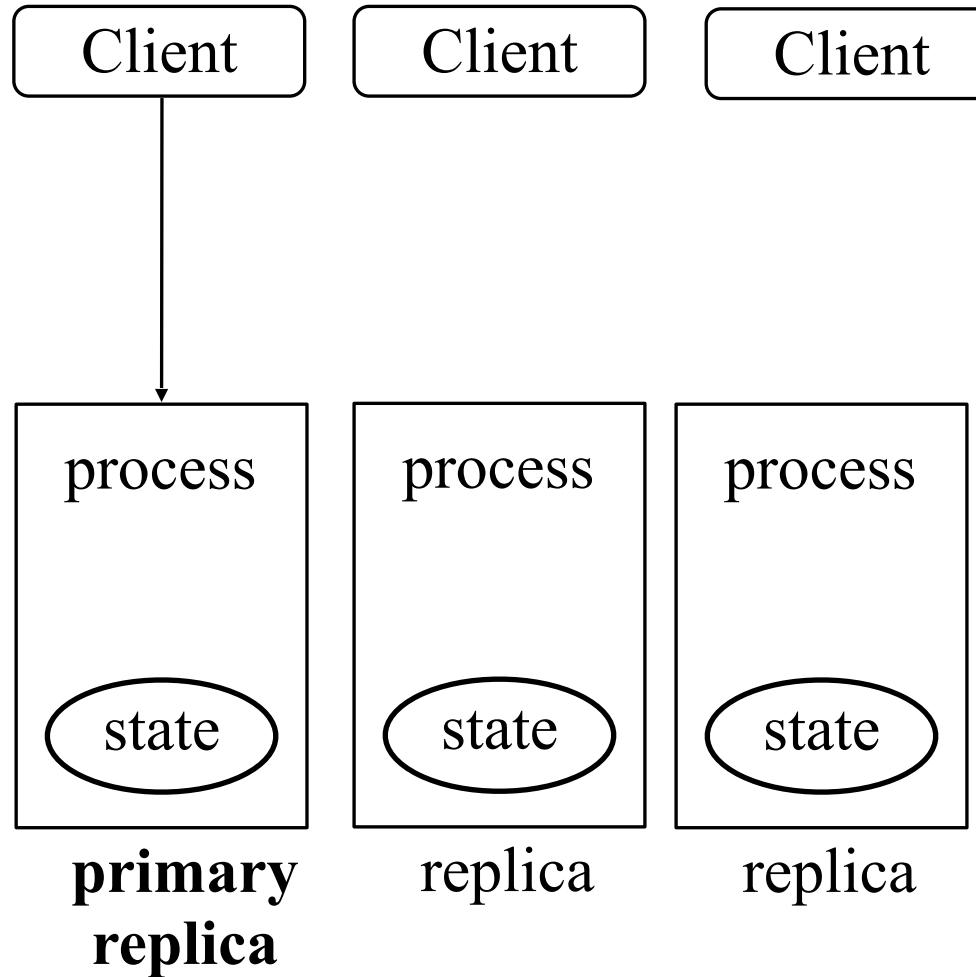


- **Passive replication**
  - Primary-backup
  - Multi-primary

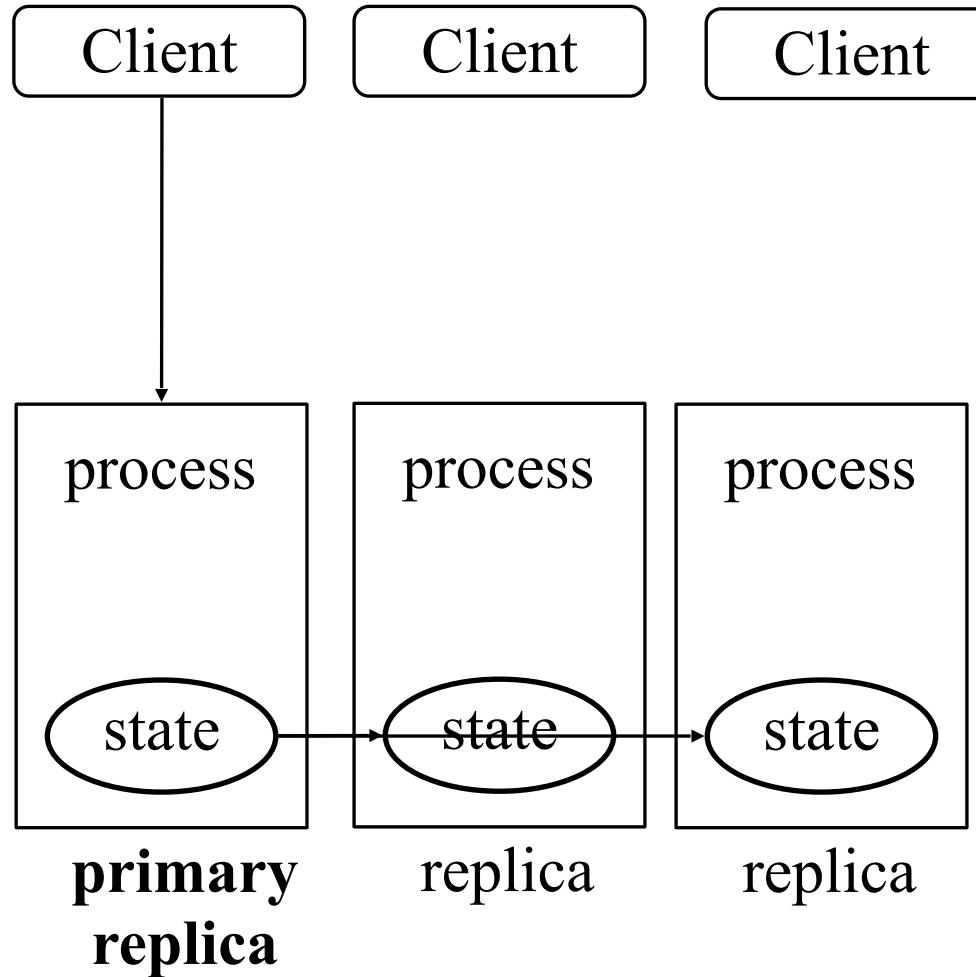
# Replication Pattern



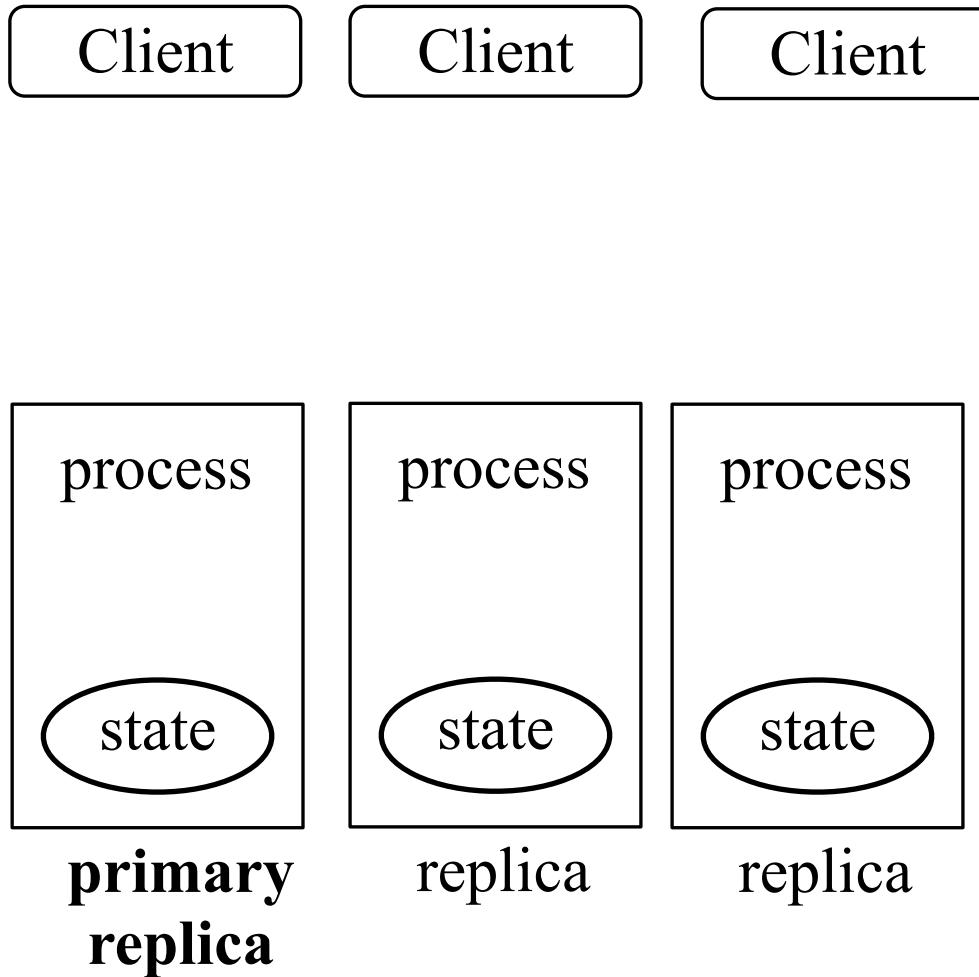
# Replication Pattern



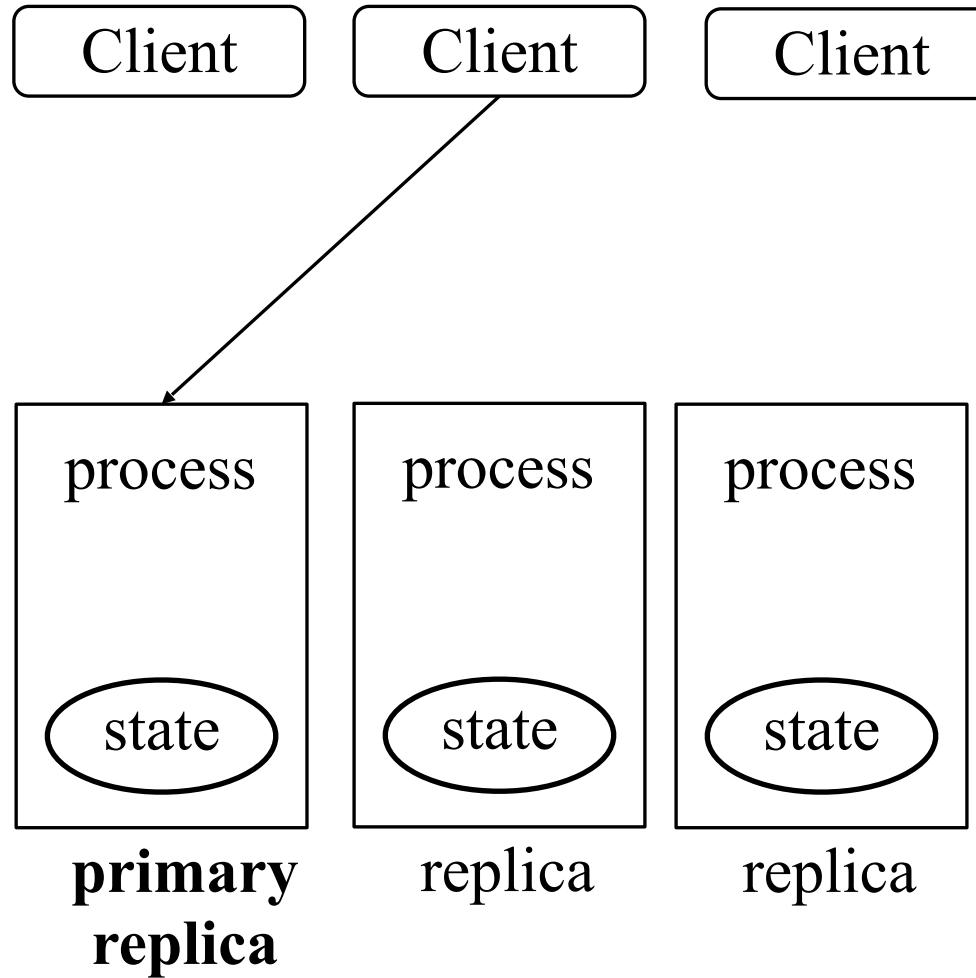
# Replication Pattern



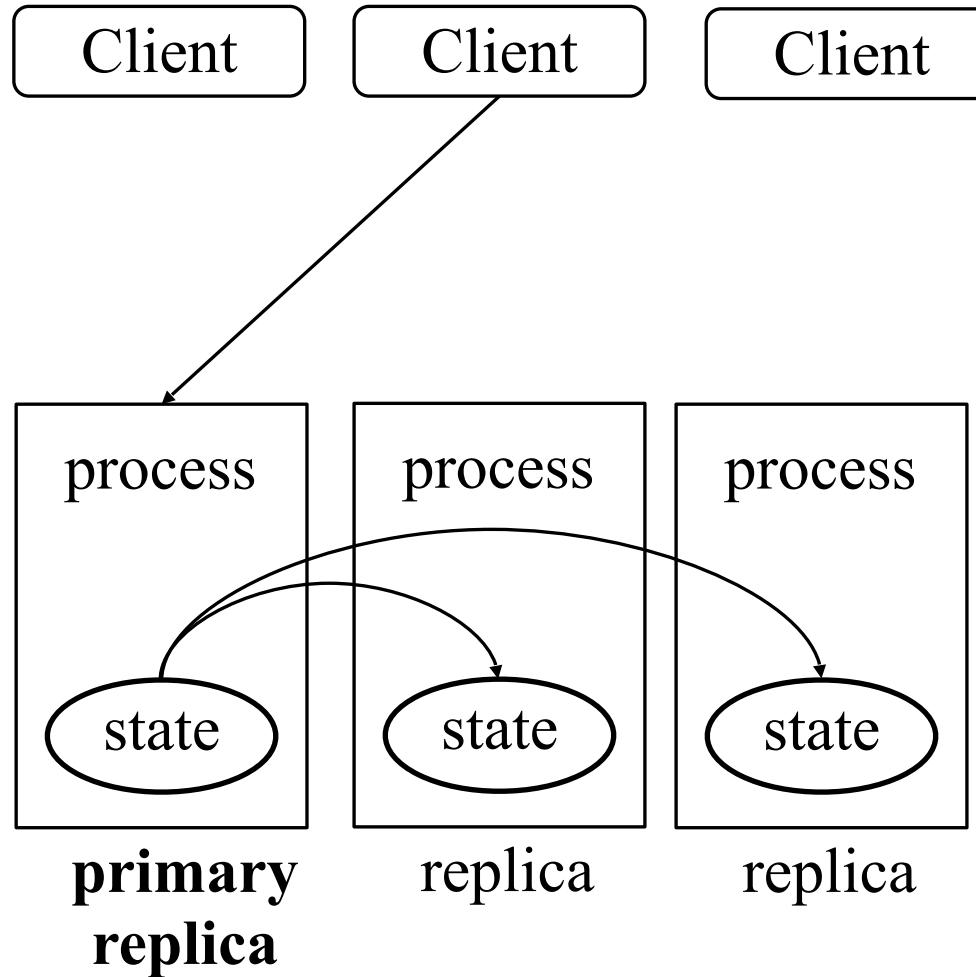
# Replication Pattern



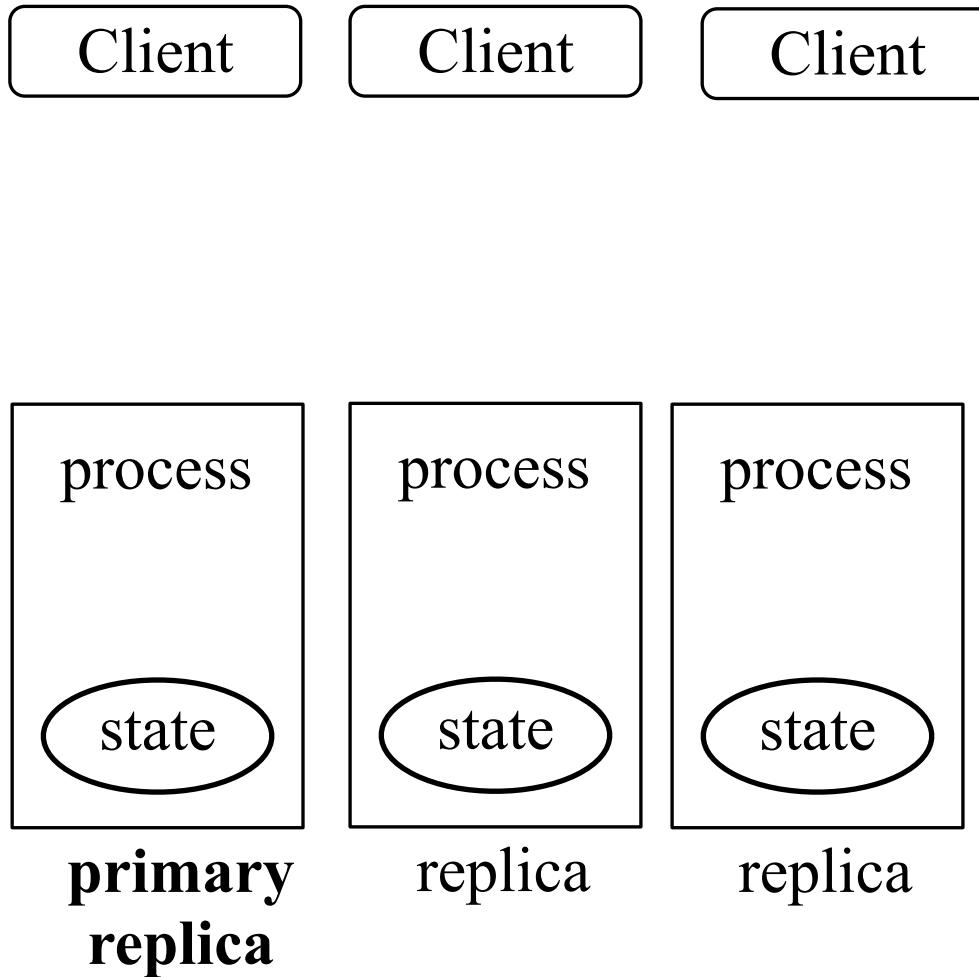
# Replication Pattern



# Replication Pattern



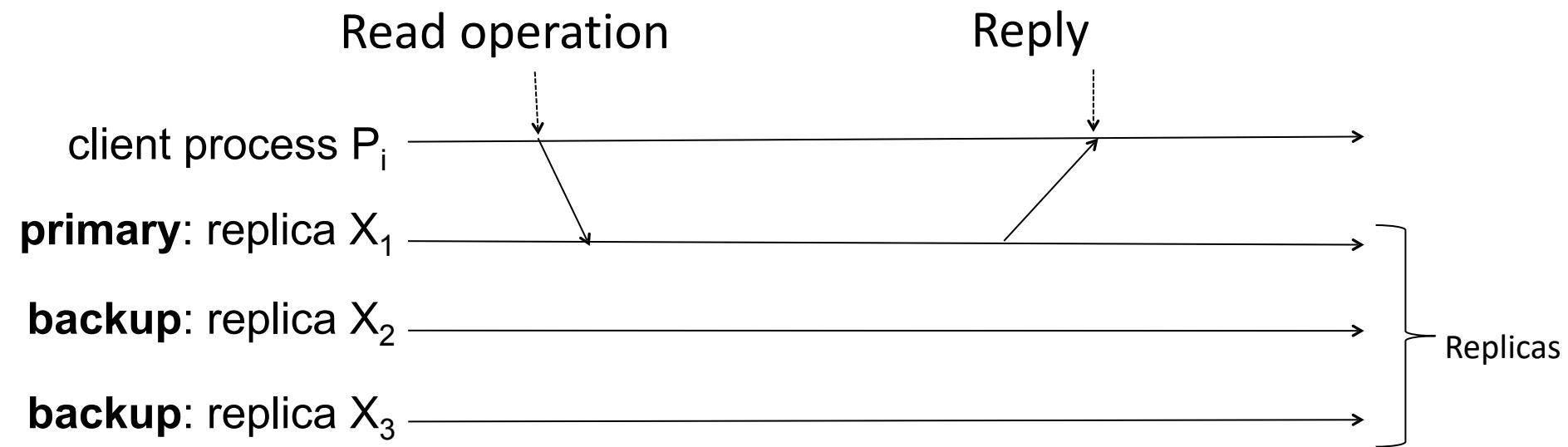
# Replication Pattern



# Passive Replication: Primary-Backup

- A replica is chosen to be the **primary** (leader election)
- **Primary**
  - Receives invocations from clients
  - Executes requests and sends back replies
  - Replicates the state to other replicas
- **Backup**
  - Interacts with primary only
  - Used to replace primary when it crashes (leader election)
- Called **eager replication** if replication is performed **within request boundary** (e.g., before the reply is sent)

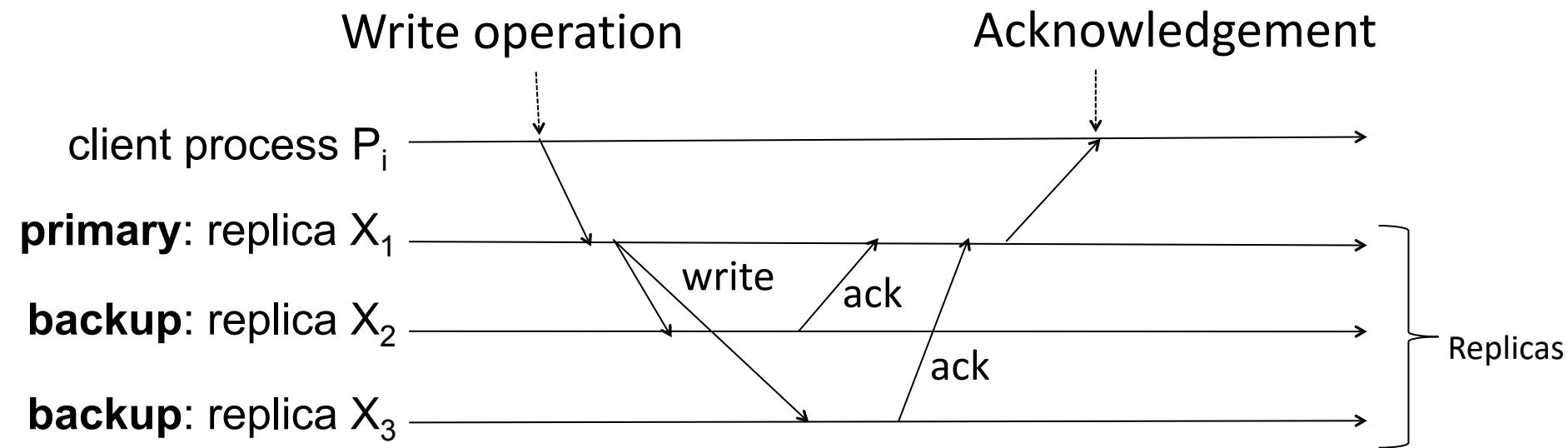
# Primary-Backup Scenario



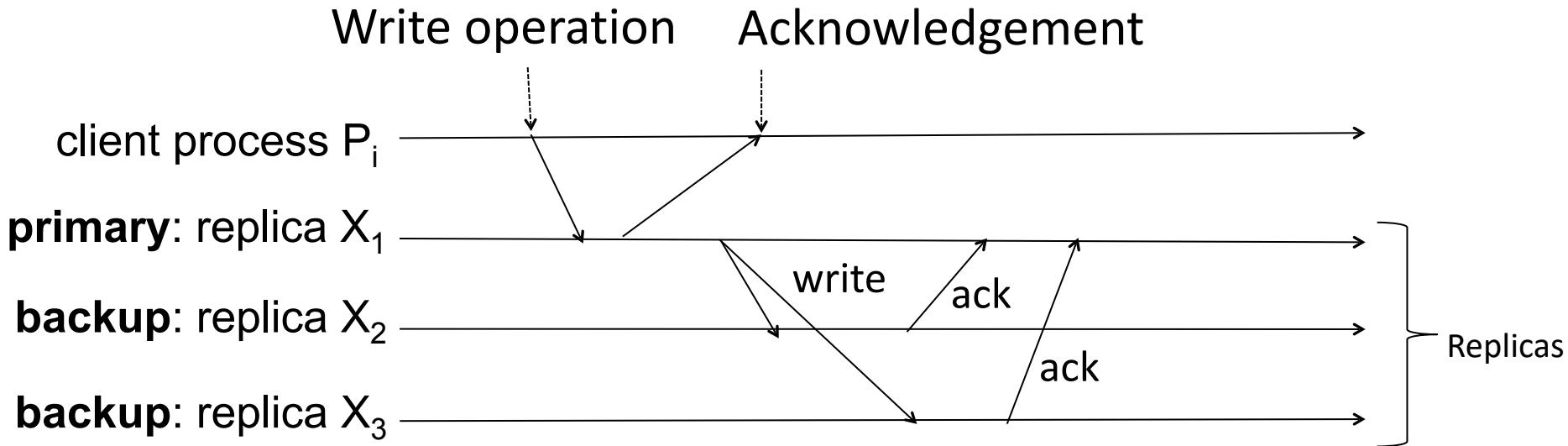
Configuration service:

- Failure detection
- Configuration management

# Primary-Backup Scenario

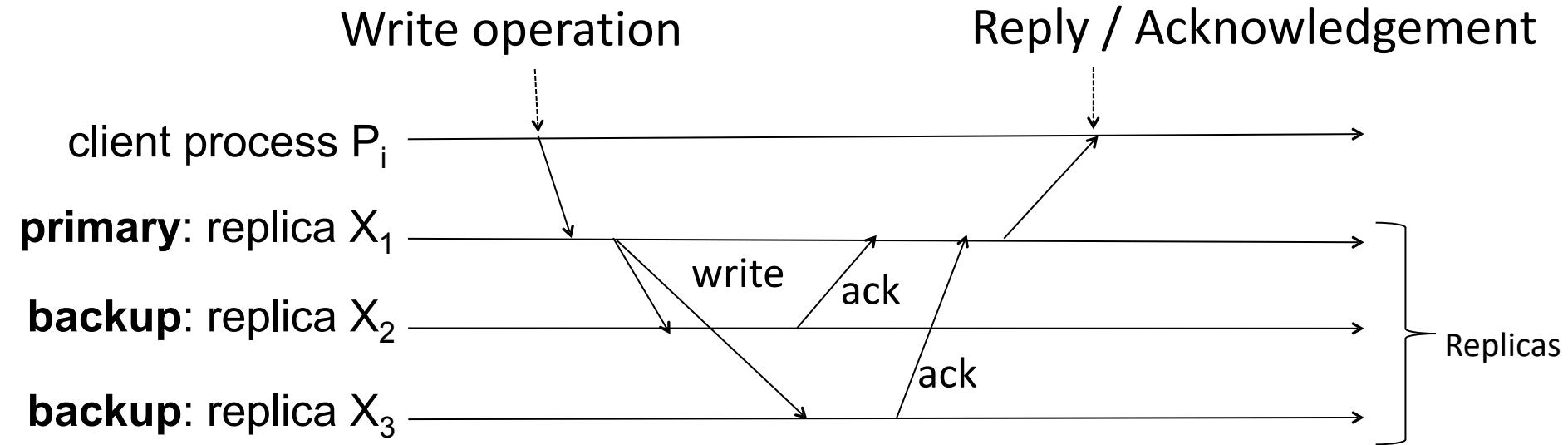


# Primary-Backup Scenario



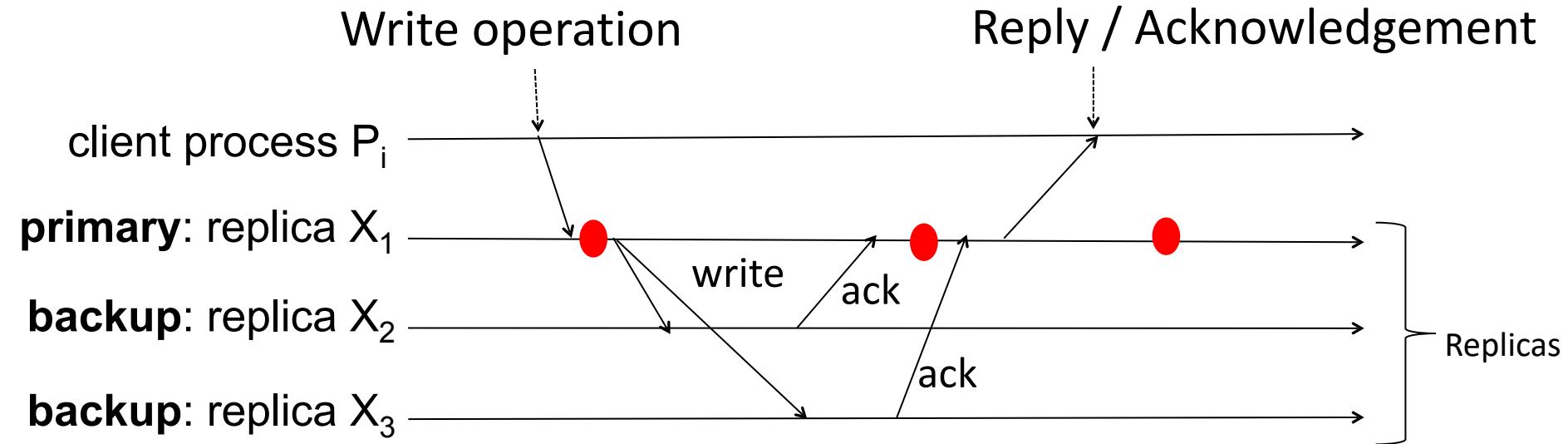
- Writes are propagated asynchronously after primary acknowledges update to client
- Replicas may diverge
- Requires additional mechanism to deal with primary crashing

# Primary-Backup: Presence of Failures



- In all cases, a new **primary is elected** from among the backups

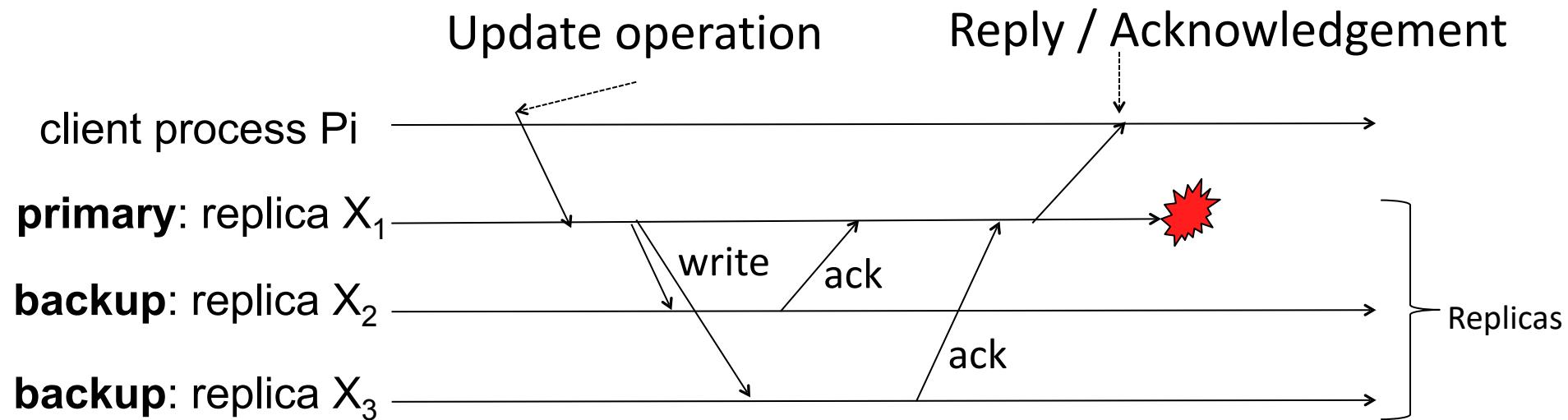
# Primary-Backup: Presence of Failures



- In all cases, a new **primary is elected** from among the backups

# Scenario 1

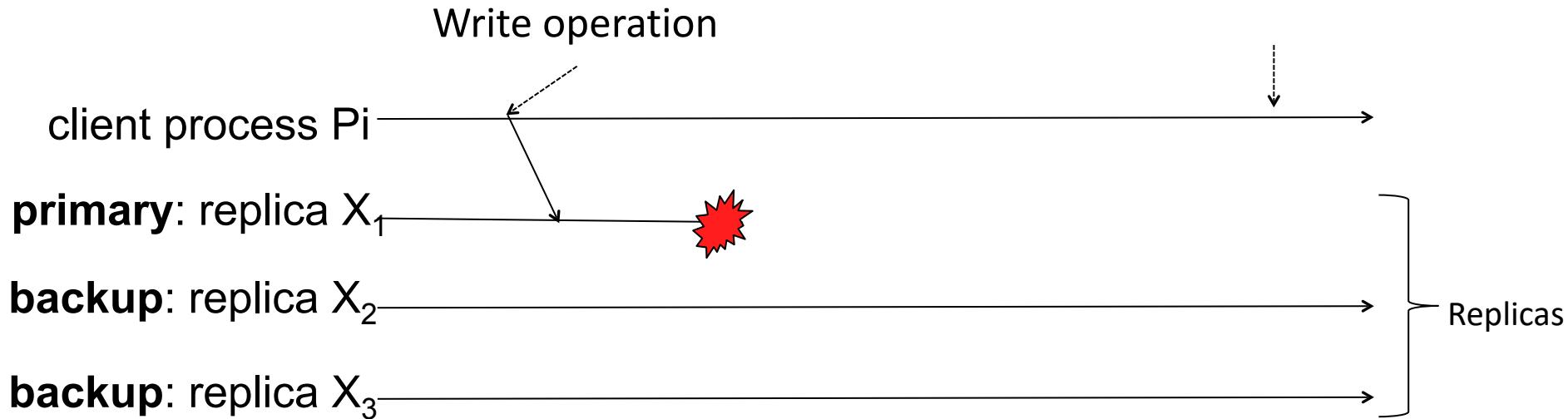
## Primary fails after client receives reply



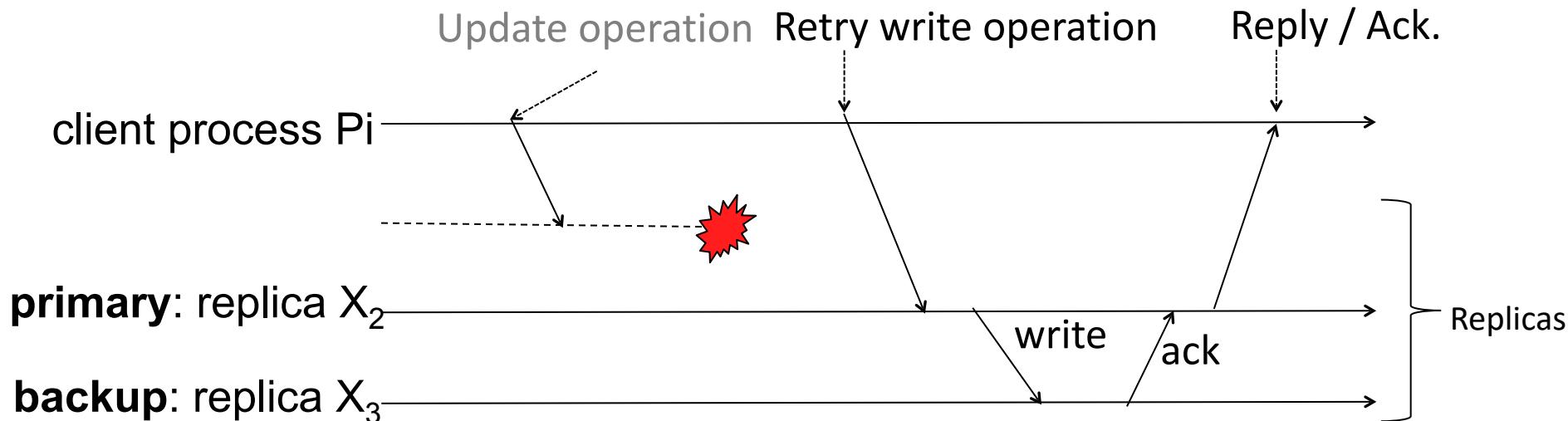
- New primary is elected

# Scenario 2

## Primary fails before propagating updates



# Scenario 2



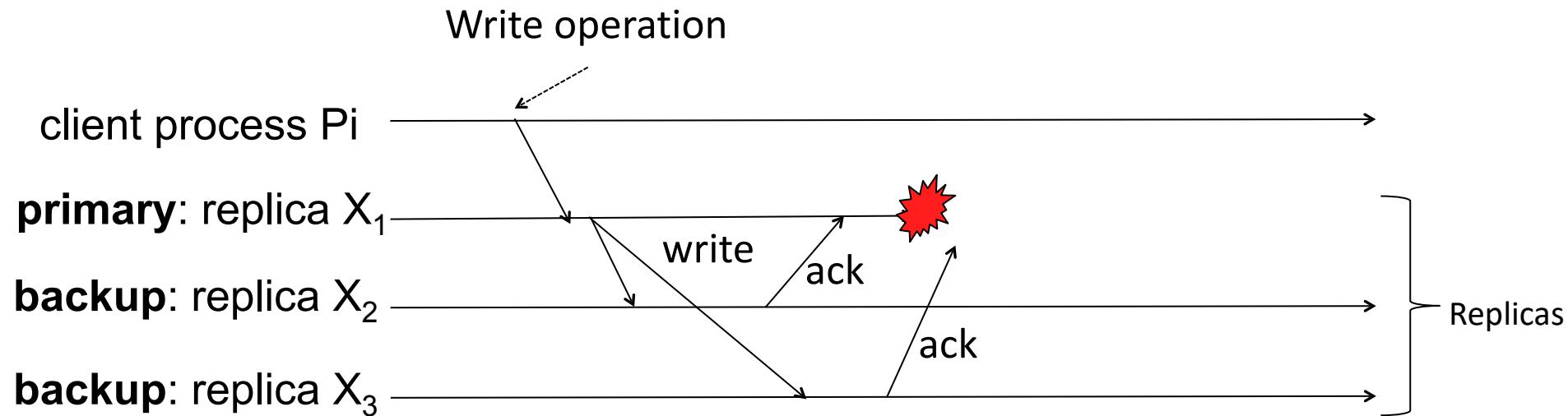
## Configuration service:

- Failure detection
- Leader election
- Configuration management

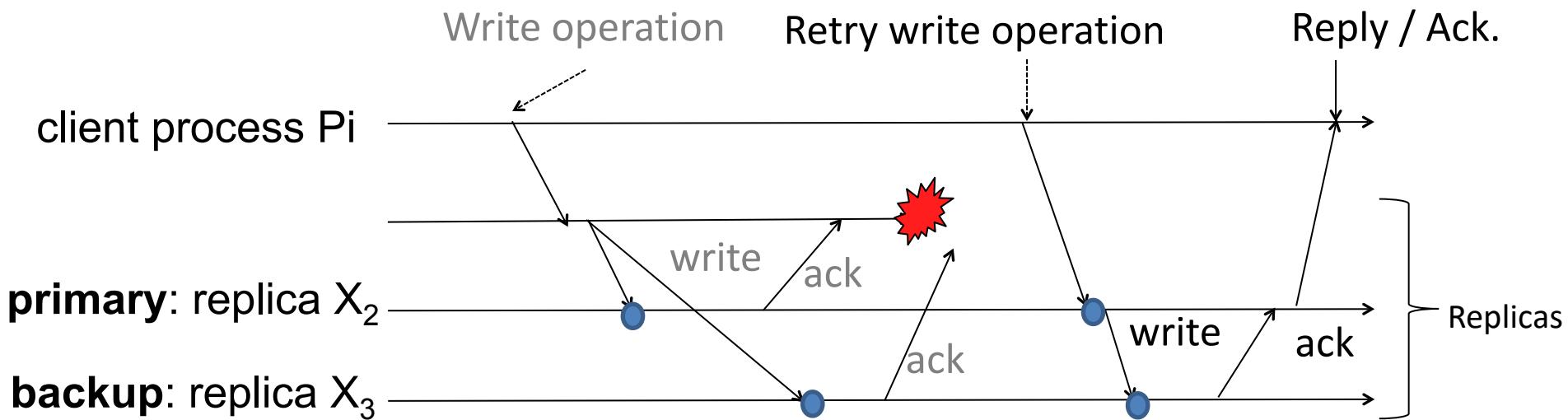
- Timeout mechanism at client triggers retry
- Retry could fail, if against old primary
- Check configuration service for new leader, retry

# Scenario 3

**Primary fails before receiving all write acknowledgements**



# Scenario 3

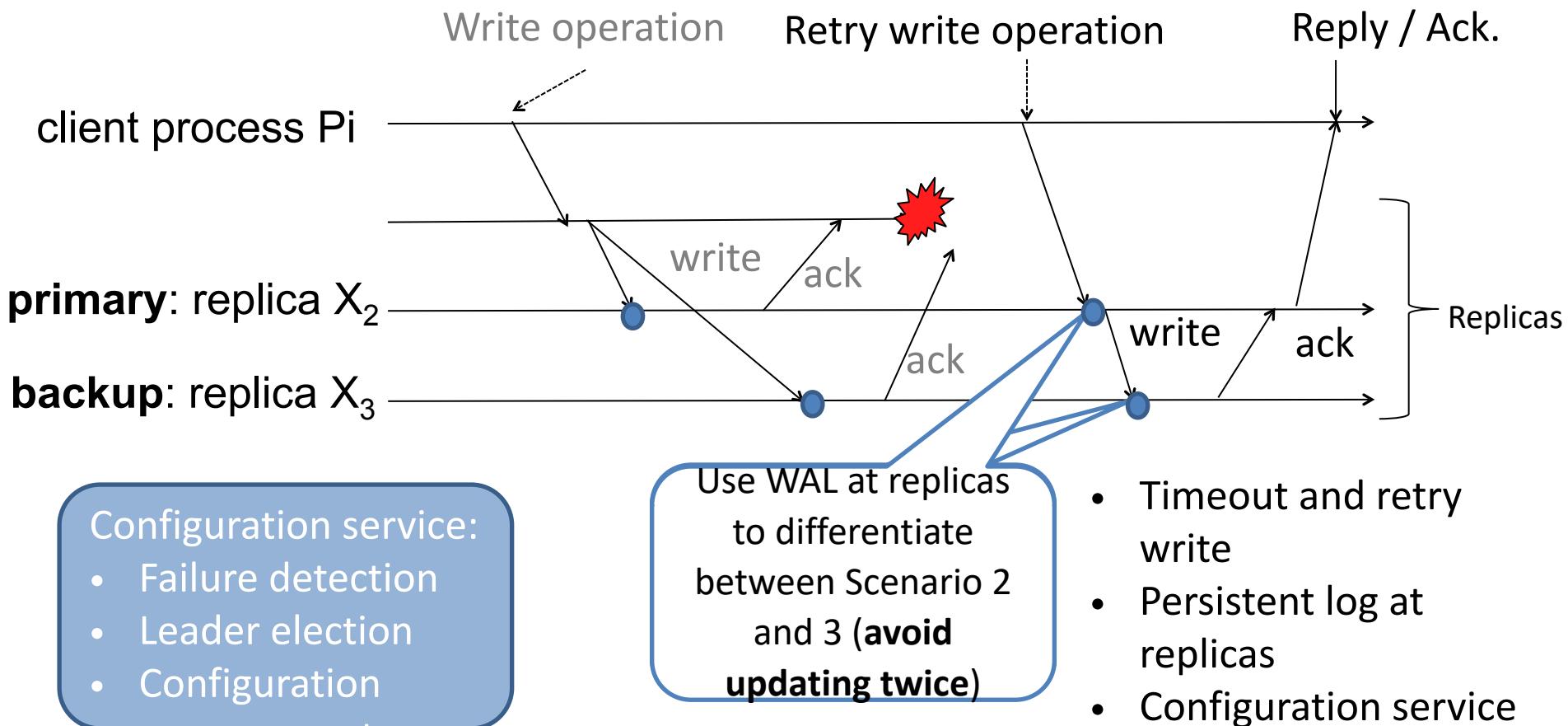


## Configuration service:

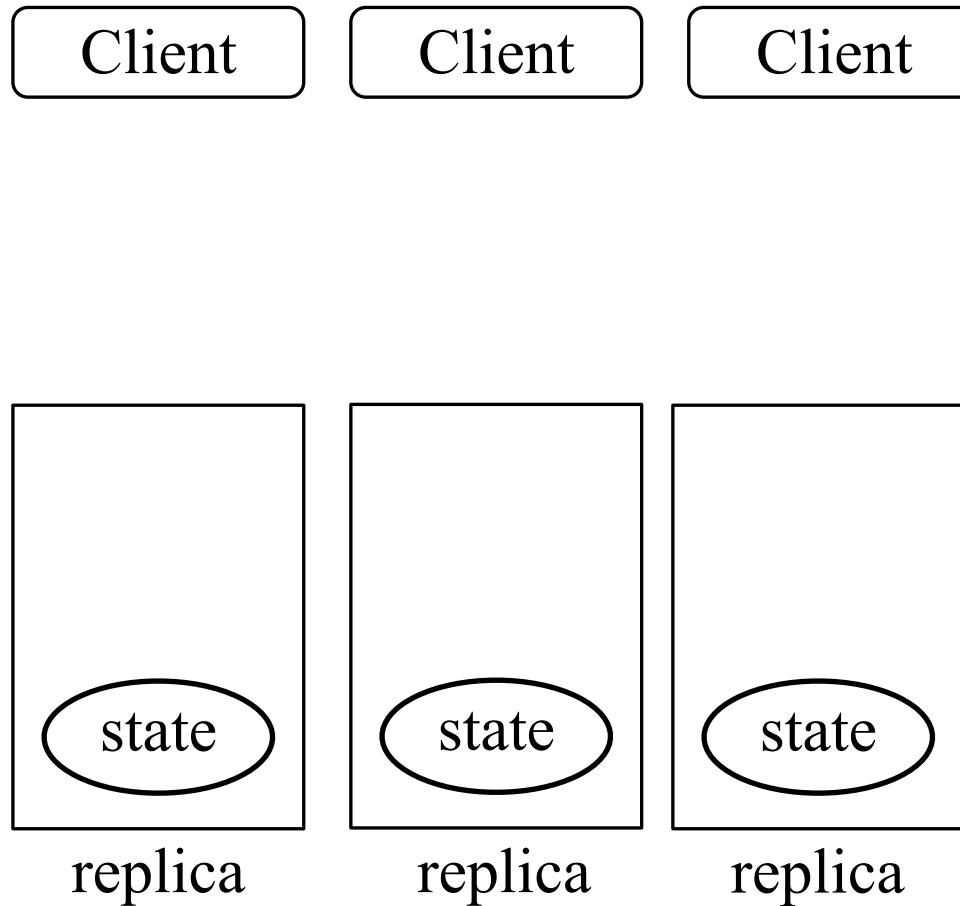
- Failure detection
- Leader election
- Configuration

- Timeout and retry write
- Persistent log at replicas
- Configuration service

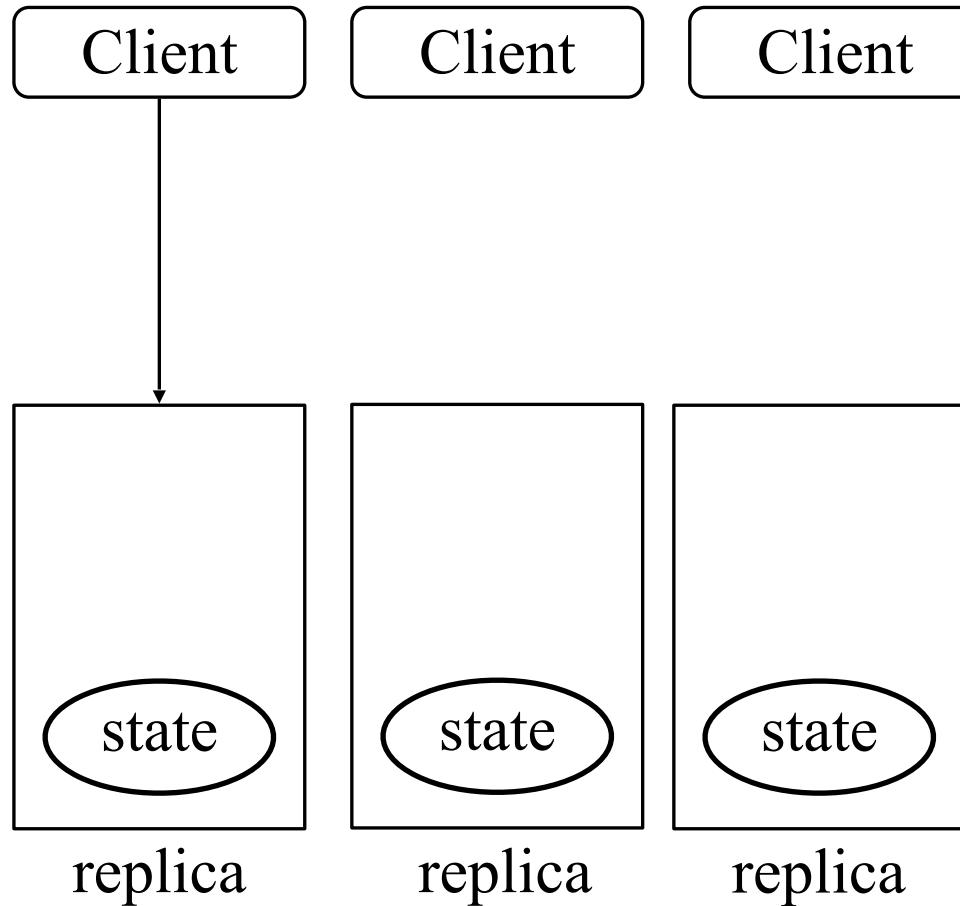
# Scenario 3



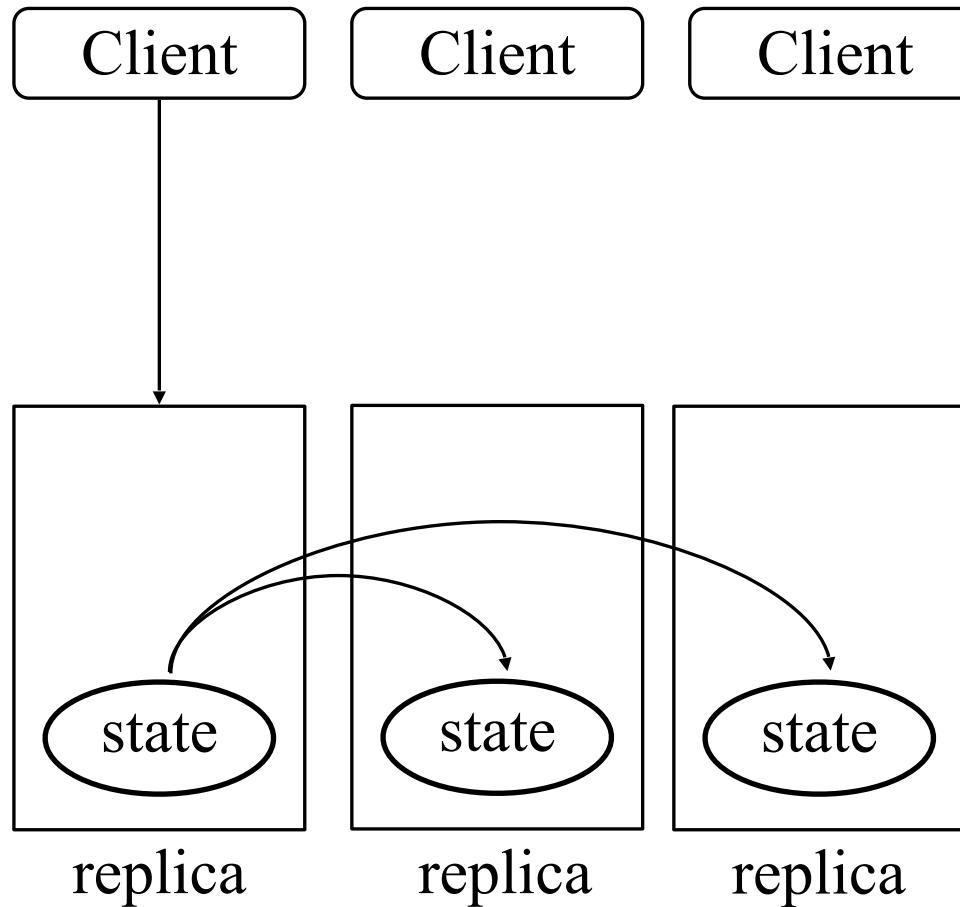
# Multi-primary Replication (MPR)



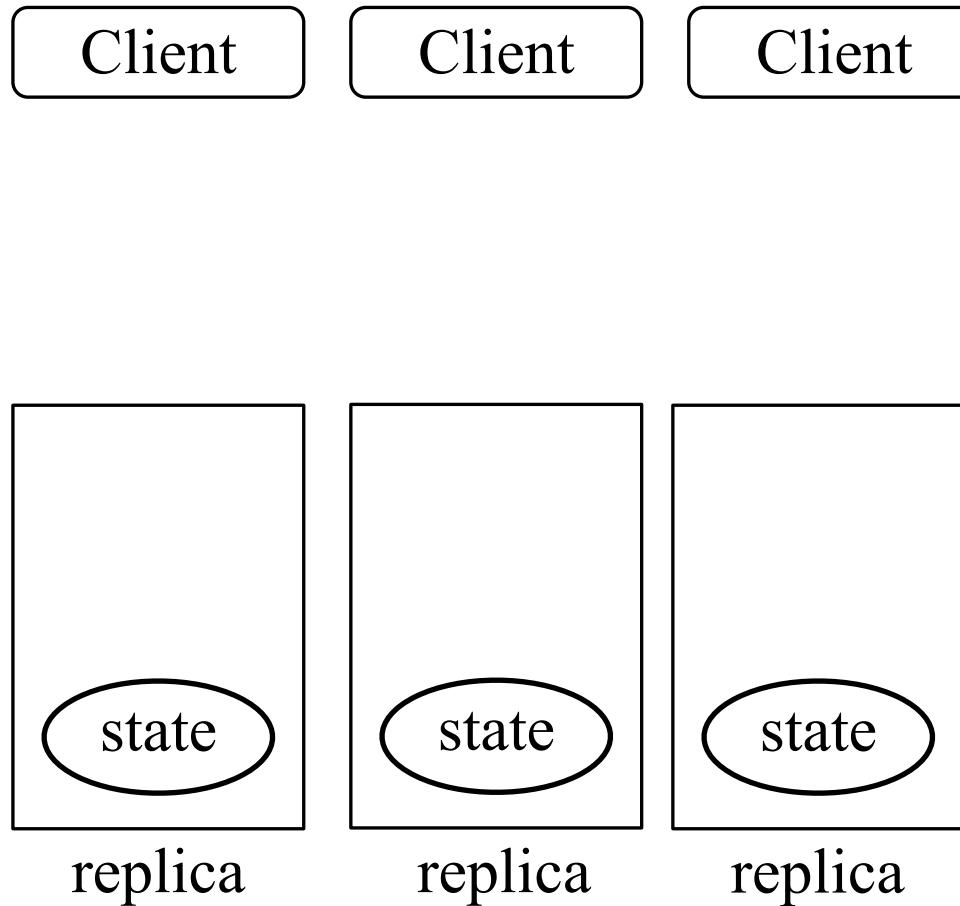
# Multi-primary Replication (MPR)



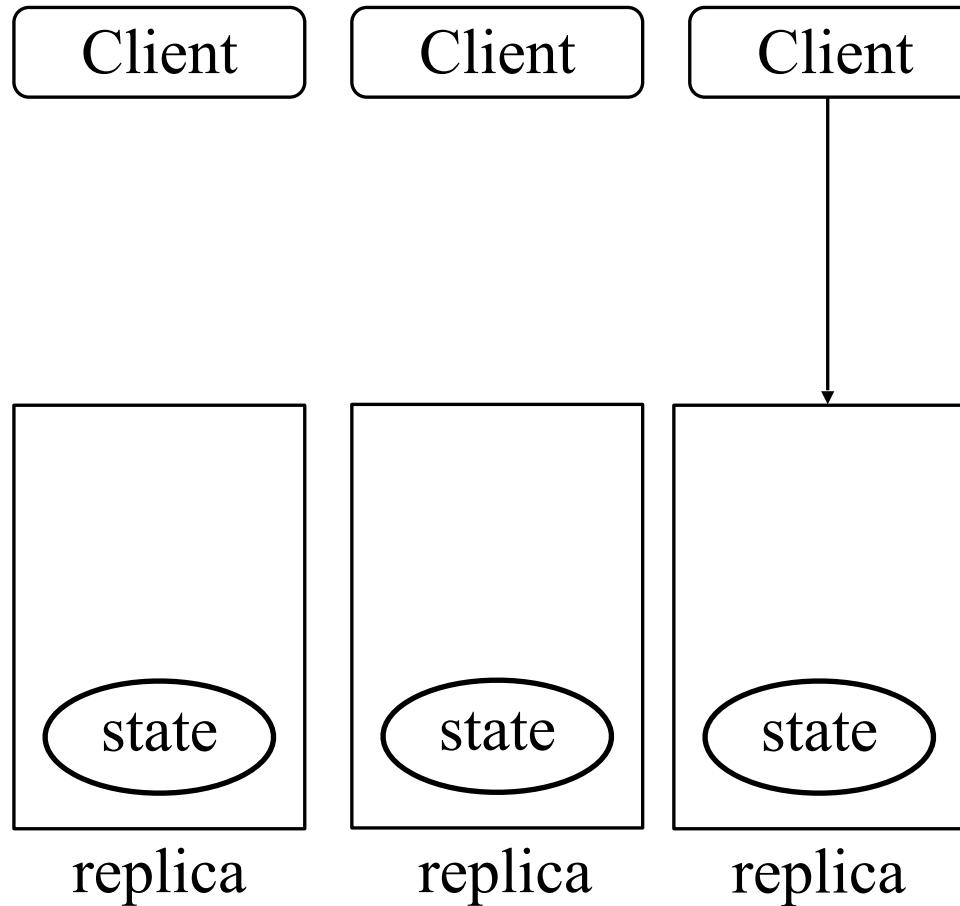
# Multi-primary Replication (MPR)



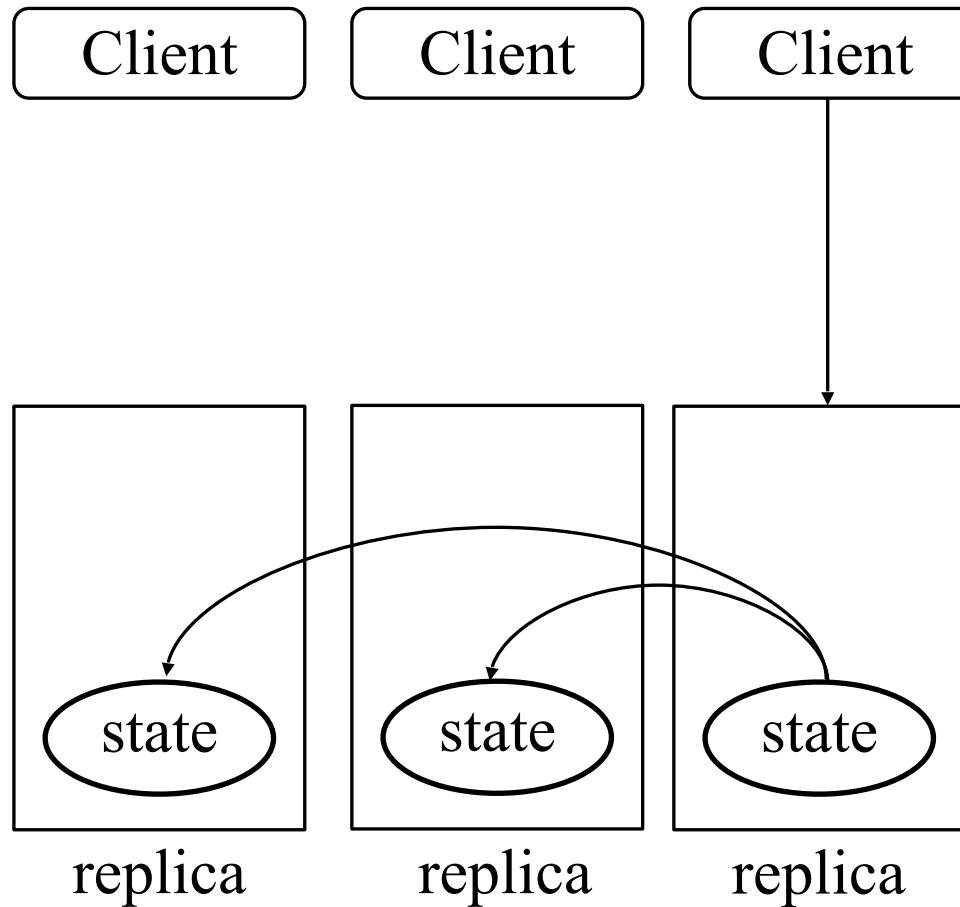
# Multi-primary Replication (MPR)



# Multi-primary Replication (MPR)



# Multi-primary Replication (MPR)



# Multi-primary Replication (MPR)

- Primary-backup approach is not scalable since a single replica handles client requests
  - Inefficient use of replica resources
- Multi-primary approach solves this issue by allowing every replica to handle client requests
  - Replicas have to figure out how to order requests (e.g., **using consensus**)
- If replication is **eager**, replicas have to **agree on order of operations** before they execute any command and respond to clients
  - Can be slow since replica must be locked

# Optimistic Lazy MPR

- To improve response times, **replication** is often **done lazily**
  - Replica first executes locally and returns a response to client right away
  - Replicas asynchronously propagate updates they made
- Also called **optimistic replication**
  - Replicas may diverge, which can introduce inconsistencies, aborts, and rollbacks

# Self-study Questions

- Specify as pseudo code how a configuration service is used in each replication failure scenario.
- Further, specify as pseudo code each the client and replica, including the additional mechanisms needed in each failure scenario.
- Analytically compare all replication patterns in terms of number of messages exchanged for reading and writing, given  $n$  replicas.
- How do the replication patterns compare in terms of potential for concurrency of operations?





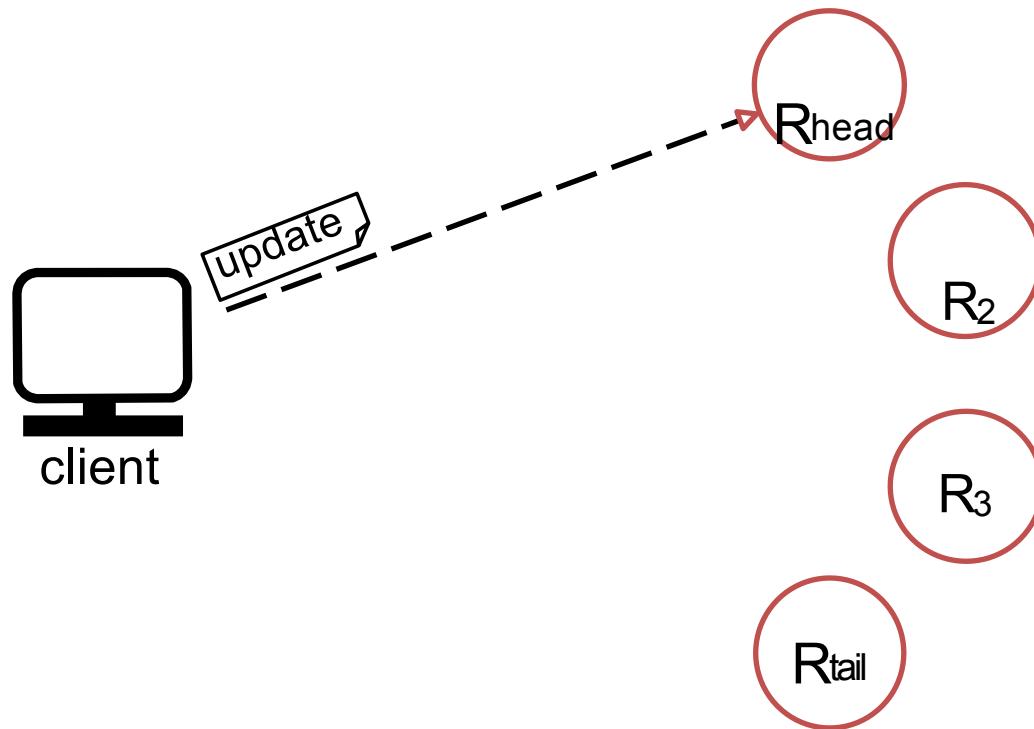
Pixabay.com

# CHAIN REPLICATION

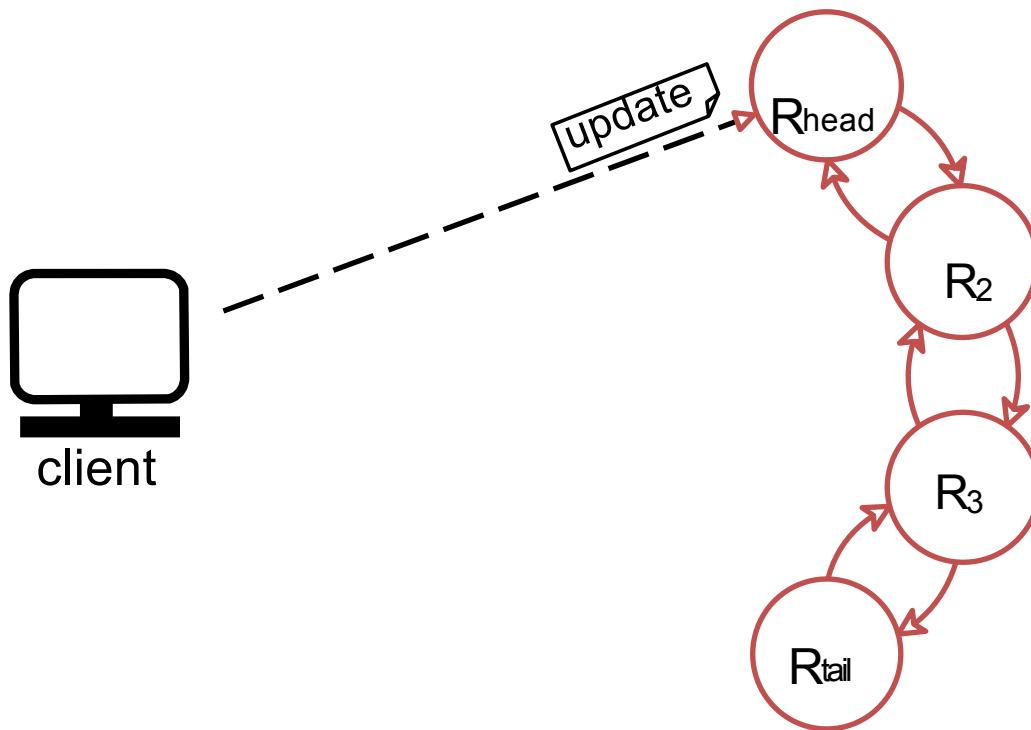
## OVERVIEW

**Inspiration for this lecture taken from a talk given by Deniz Altınbüken.**

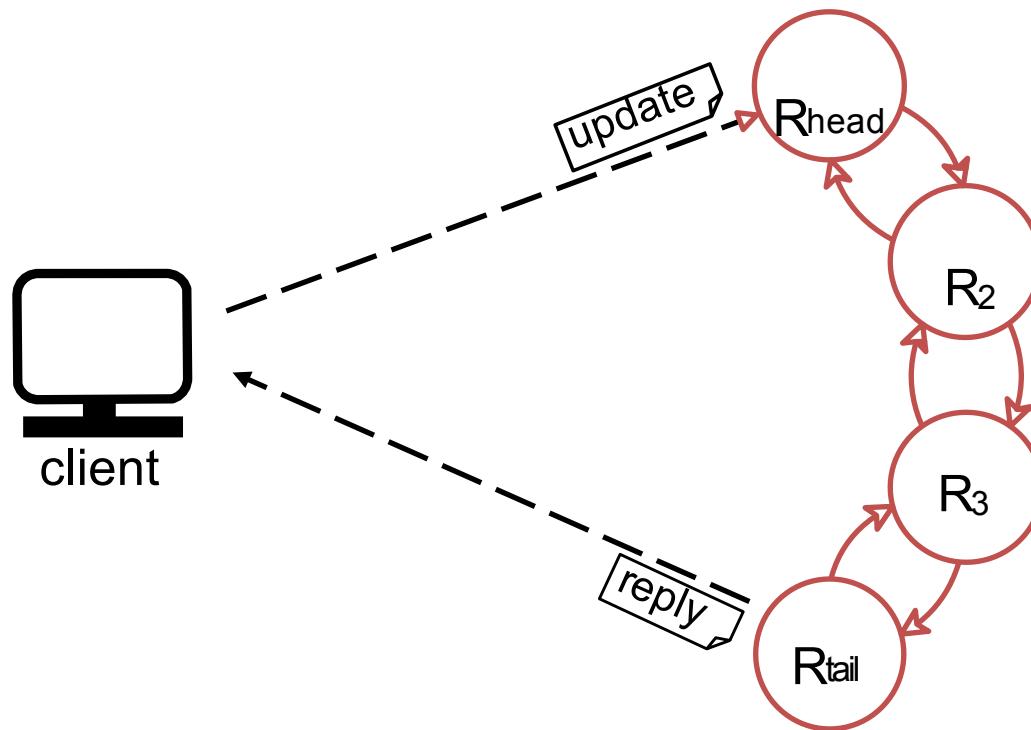
# Chain Replication



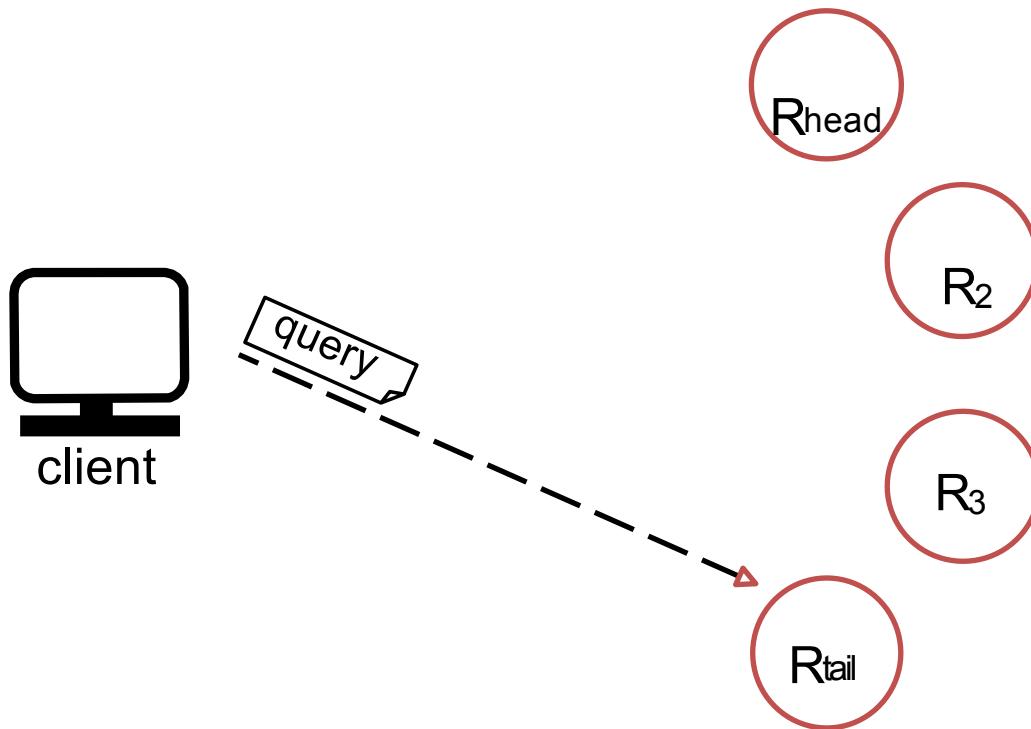
# Chain Replication



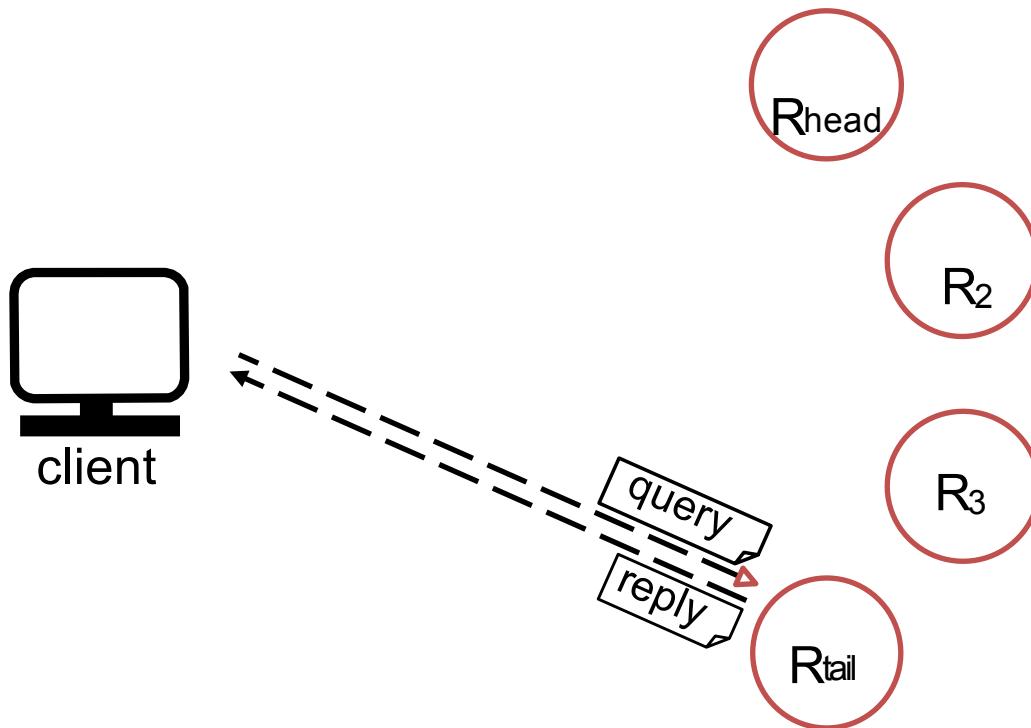
# Chain Replication



# Chain Replication



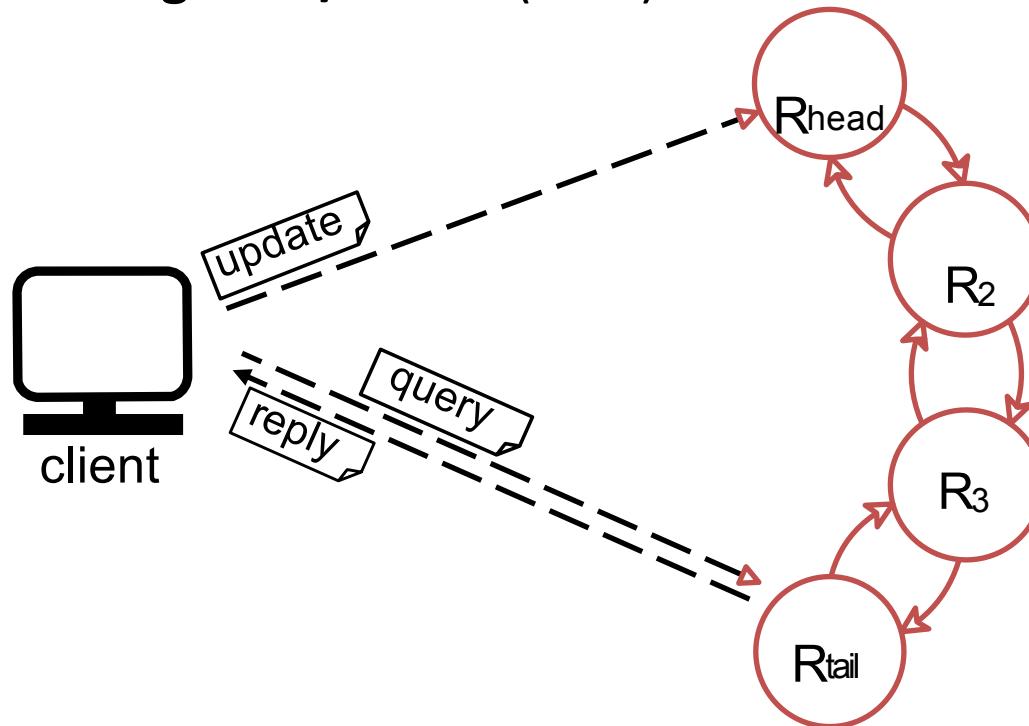
# Chain Replication



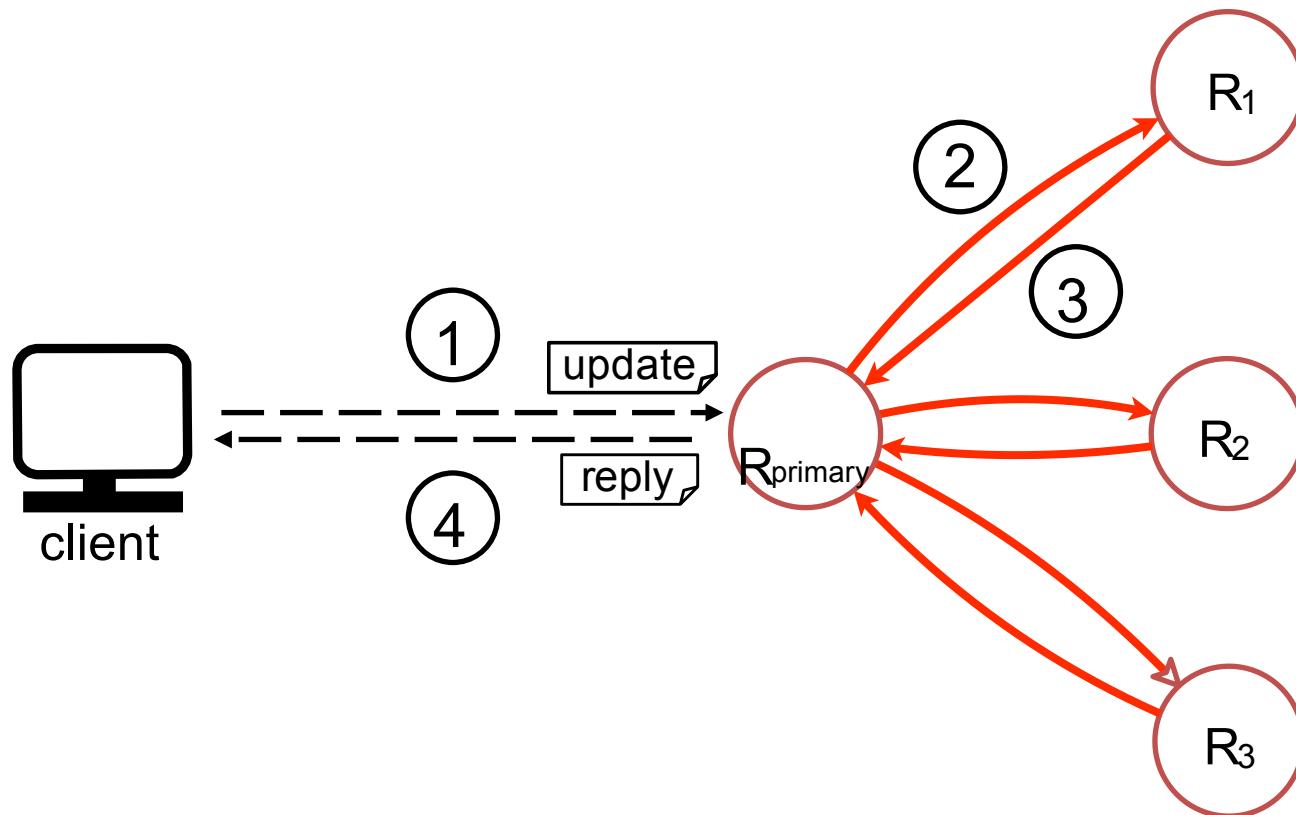
# Chain Replication

**Separate concerns across replicas:**

- Update processing request to head (write)
- Update processing reply from tail (write)
- Query processing from/to tail (read)



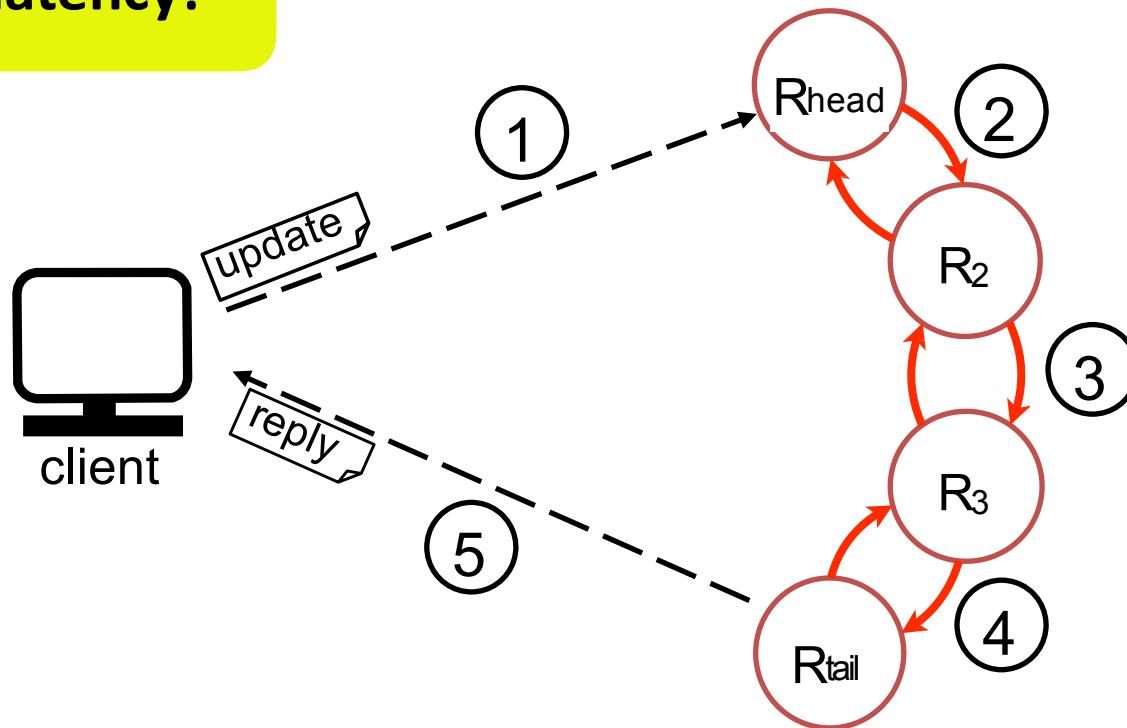
# Primary-Backup Replication



**Four replicas – latency of four messages**

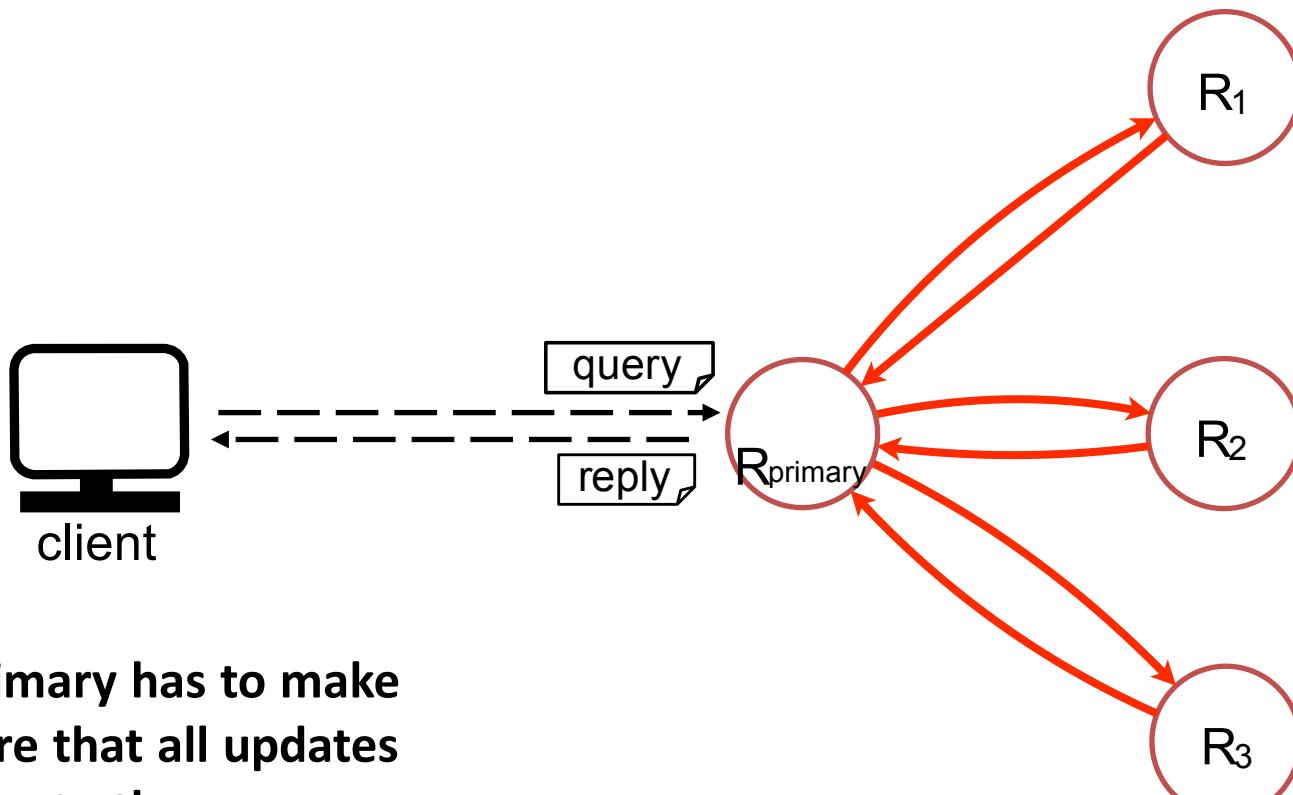
# Chain Replication

Higher latency!



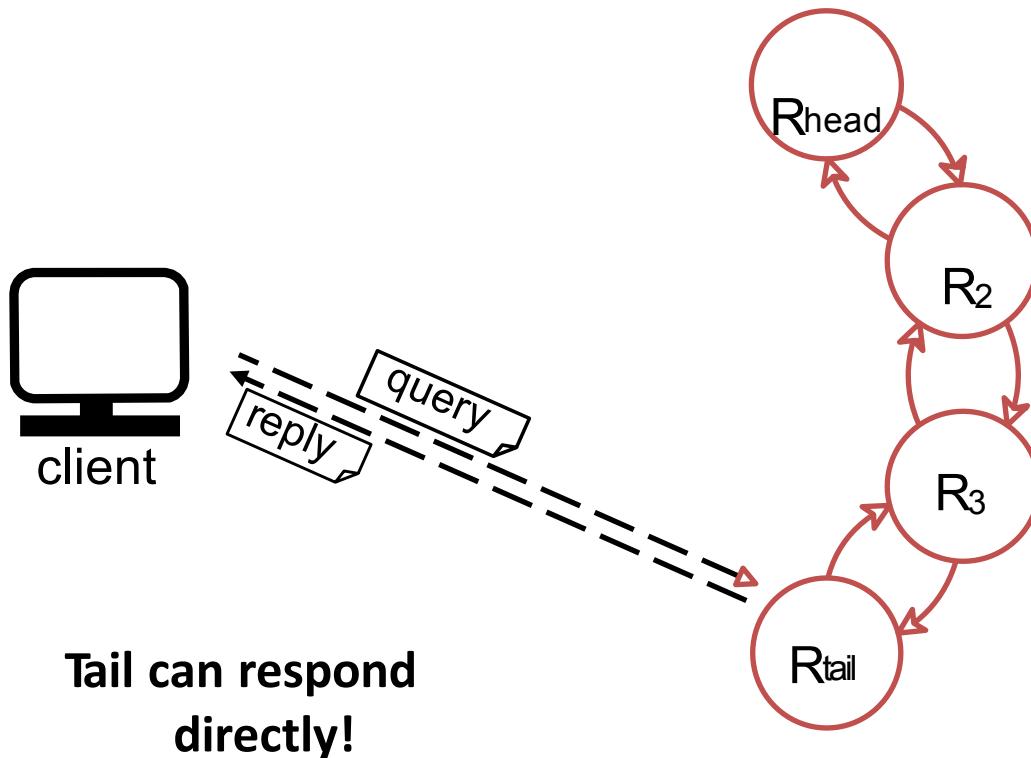
Four replicas – latency of five messages

# Primary-Backup Replication



**Primary has to make  
sure that all updates  
prior to the query are  
completed!**

# Chain Replication

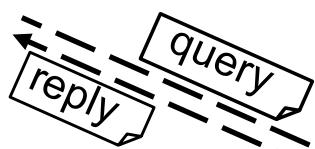


# Chain Replication

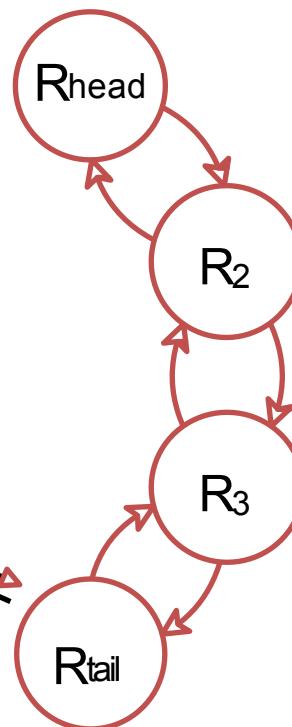
Higher throughput!



client



Tail can respond directly!



# Fault-tolerance in Chain Replication

- Need  $f + 1$  nodes to tolerate  $f$  failures



Pixabay.com

# Self-study Questions

- How many node failures can the primary-backup approach sustain without disrupting service?
- Discuss pros and cons of chain, primary-backup, active, and multi-primary replication?
- Quantify the number of messages exchanged by drawing on  $n$  replicas for each of the replication schemes for  $m$  updates.
- Discuss the replica balance across all approaches?





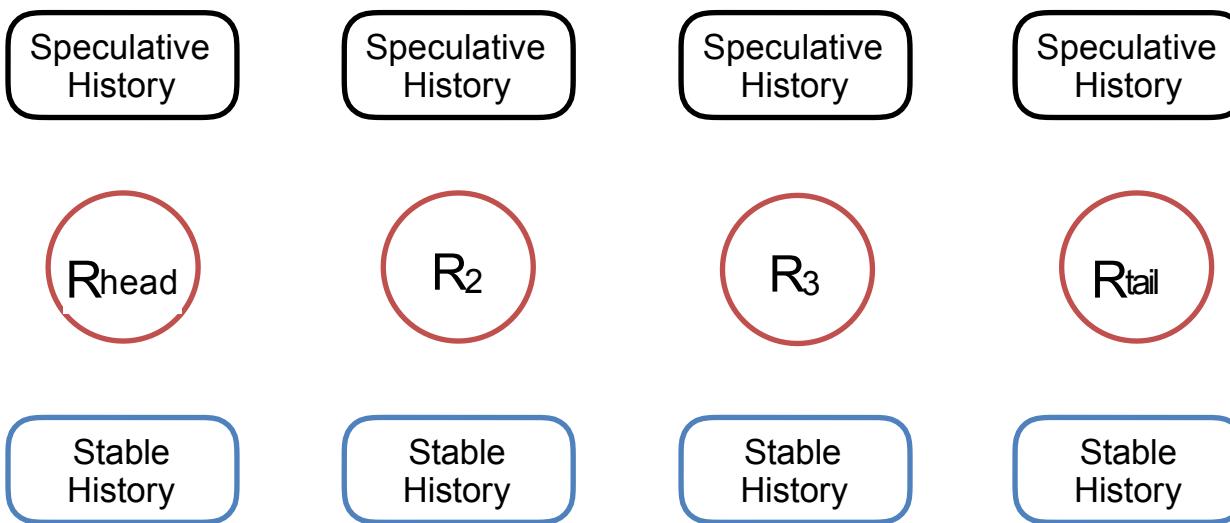
Pixabay.com

# CHAIN REPLICATION

## UPDATE & QUERY OPERATIONS

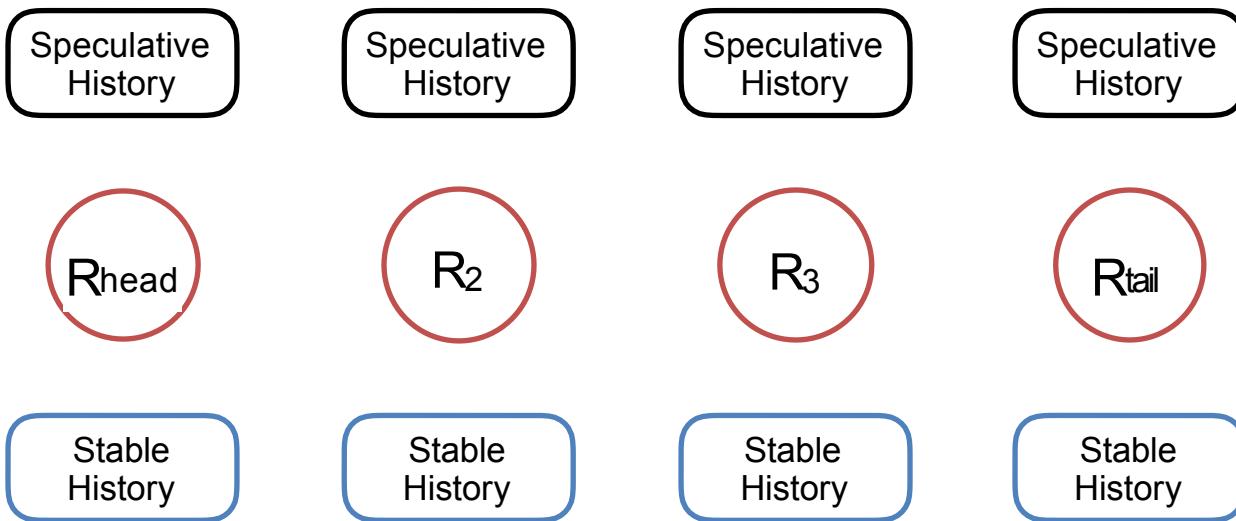
Inspiration for this lecture taken from a talk given by Deniz Altınbüken.

# Updates



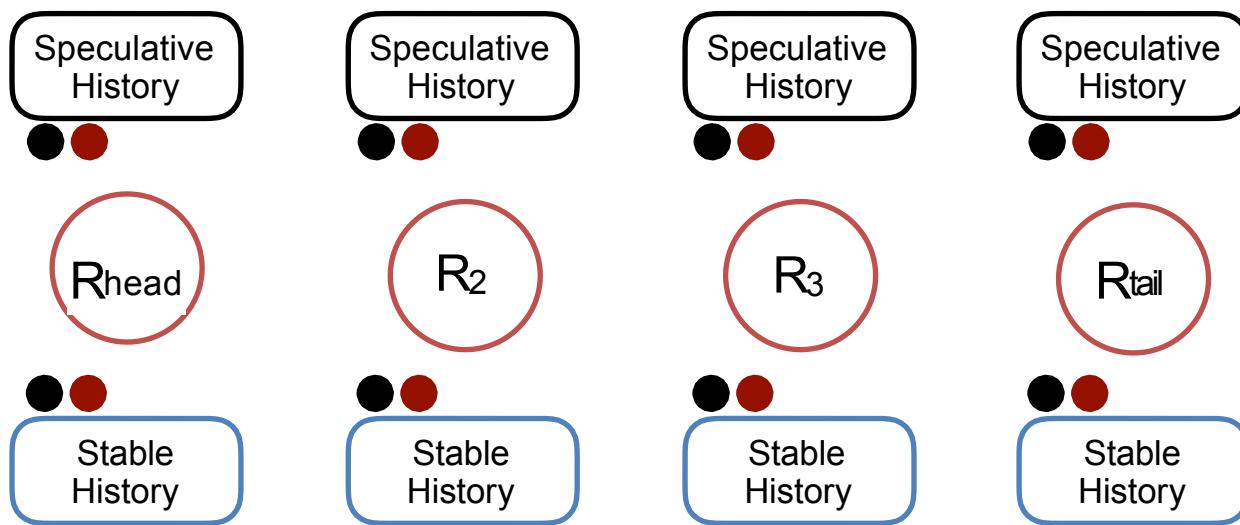
# Updates

$R_2$  is the predecessor of  $R_3$

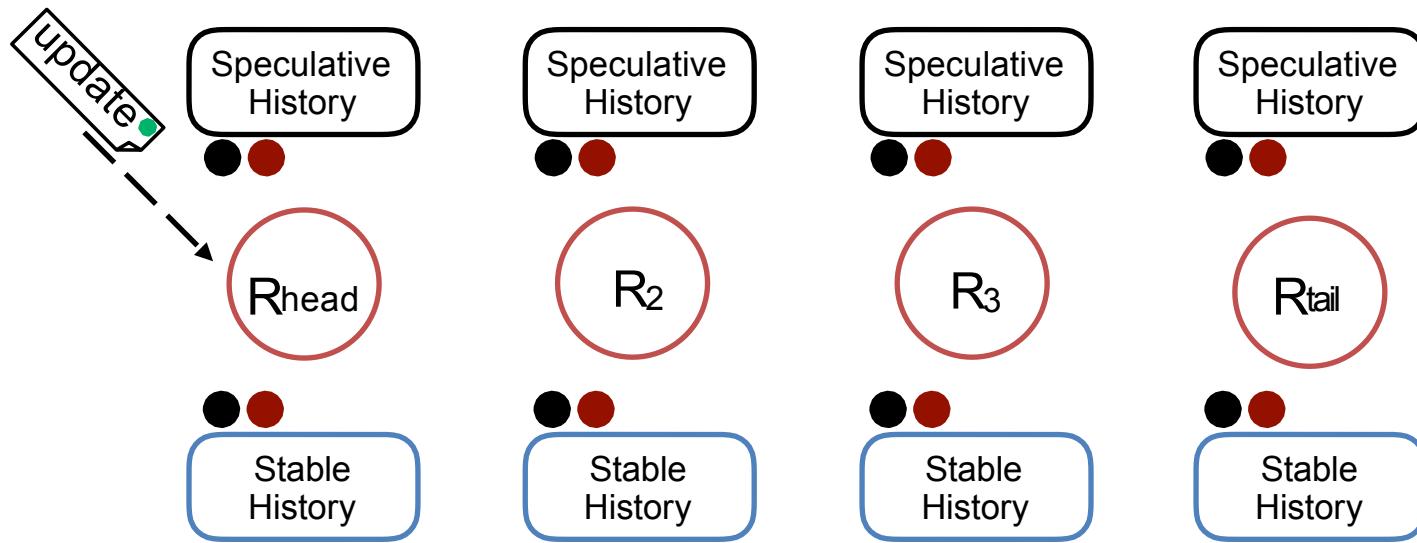


$R_3$  is the successor of  $R_2$

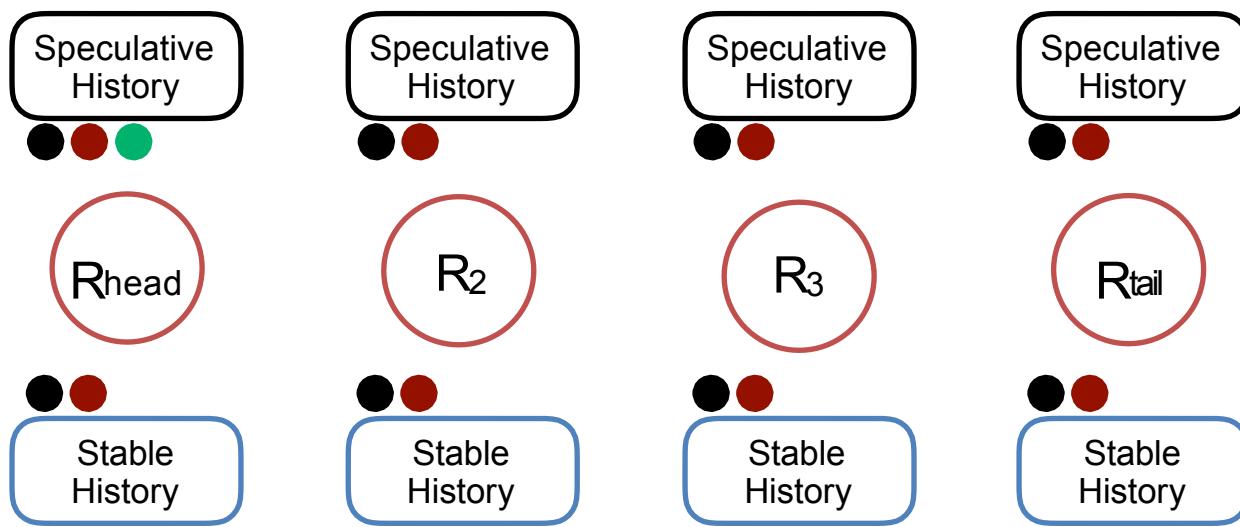
# Updates



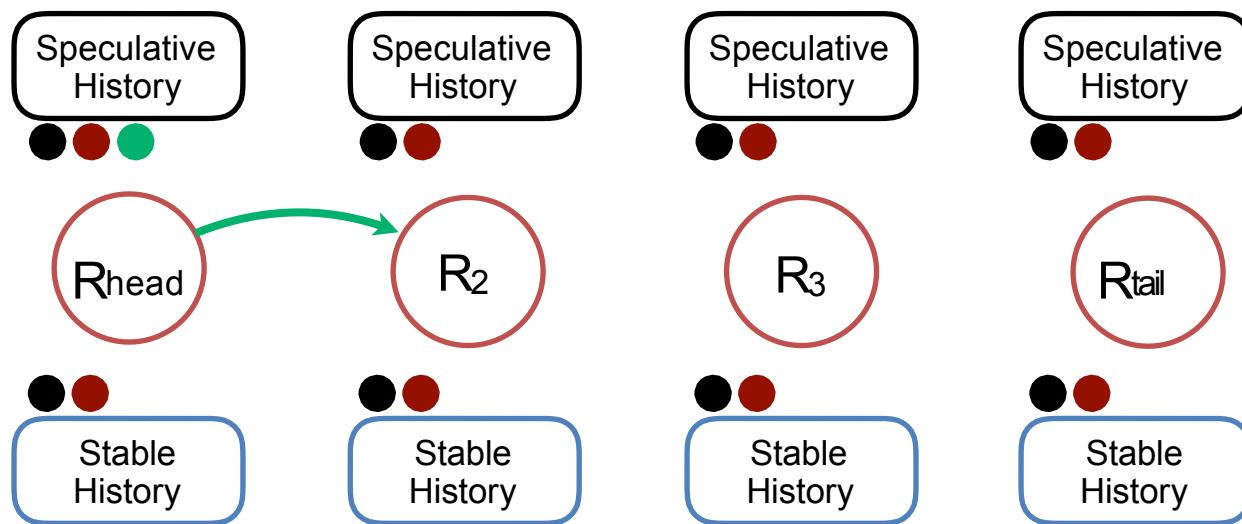
# Updates



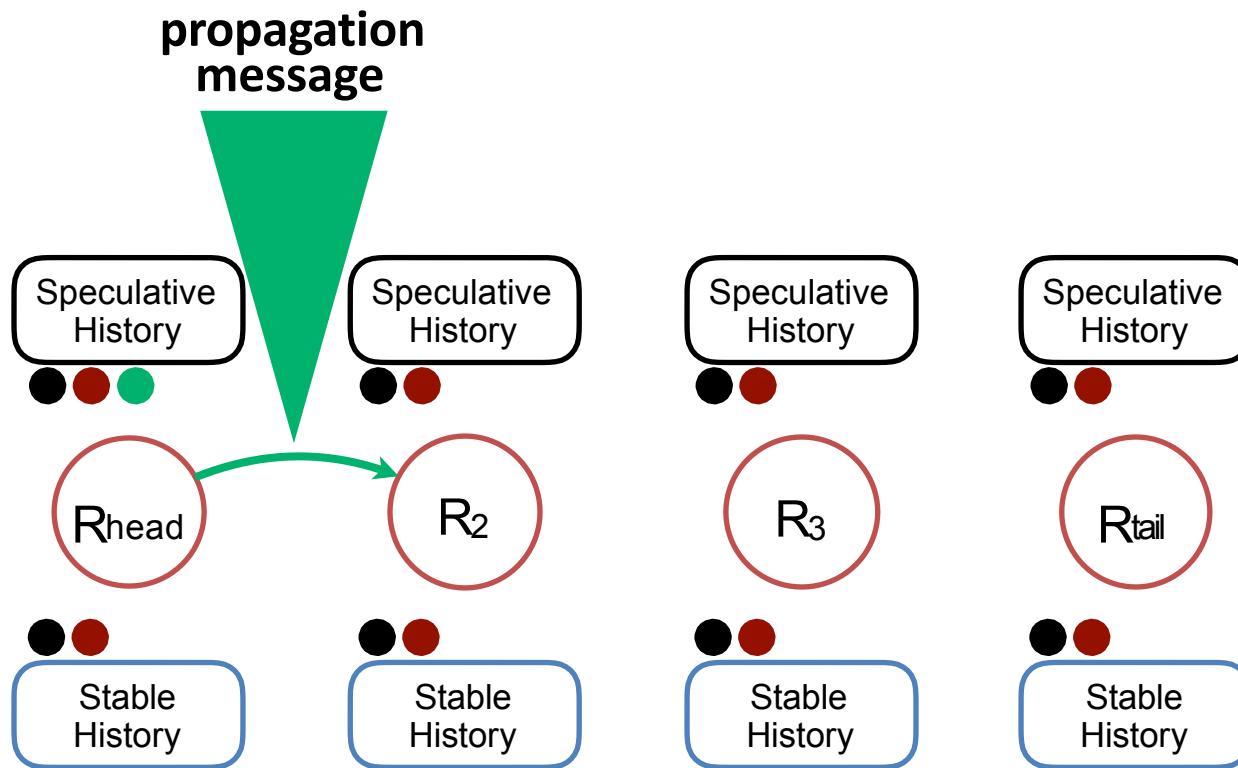
# Updates



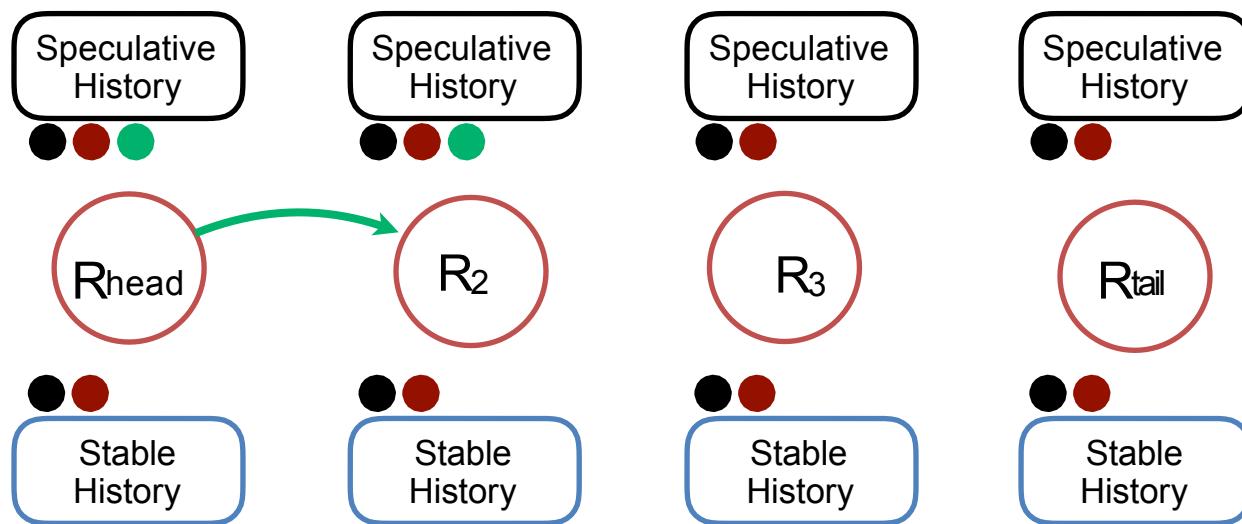
# Updates



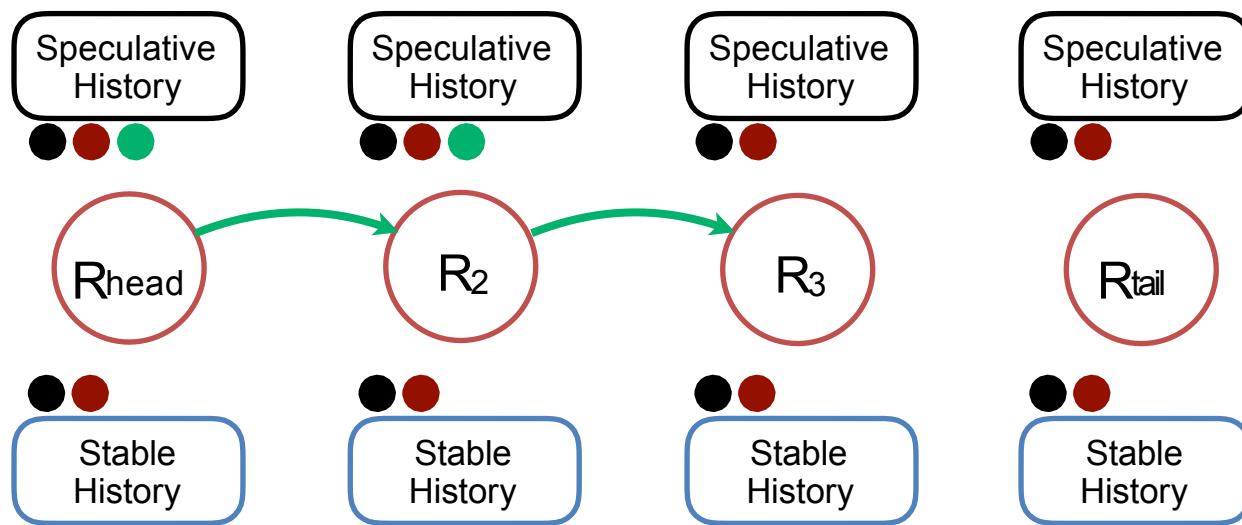
# Updates



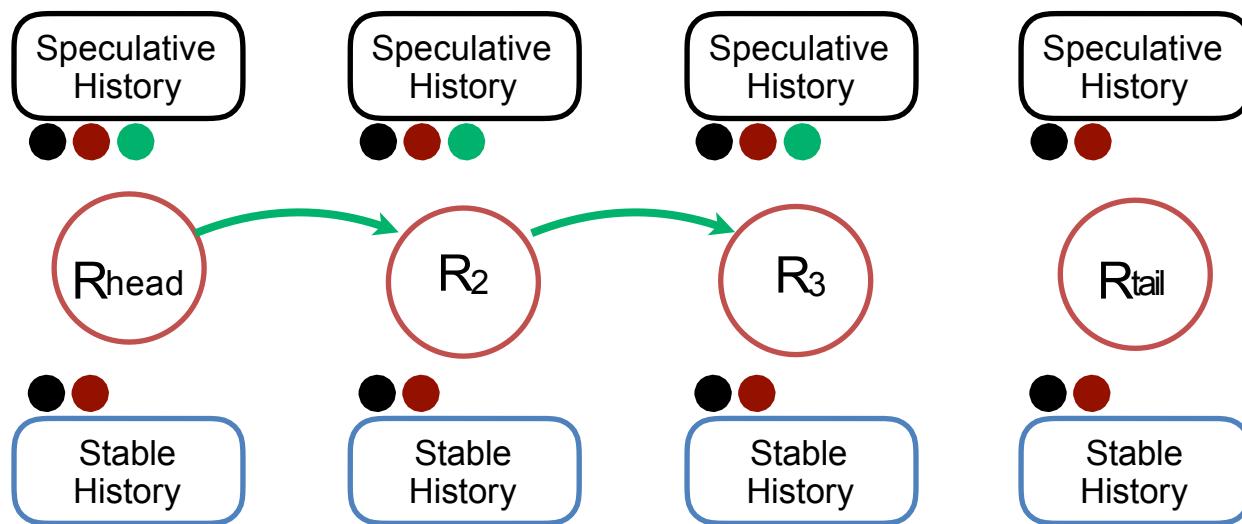
# Updates



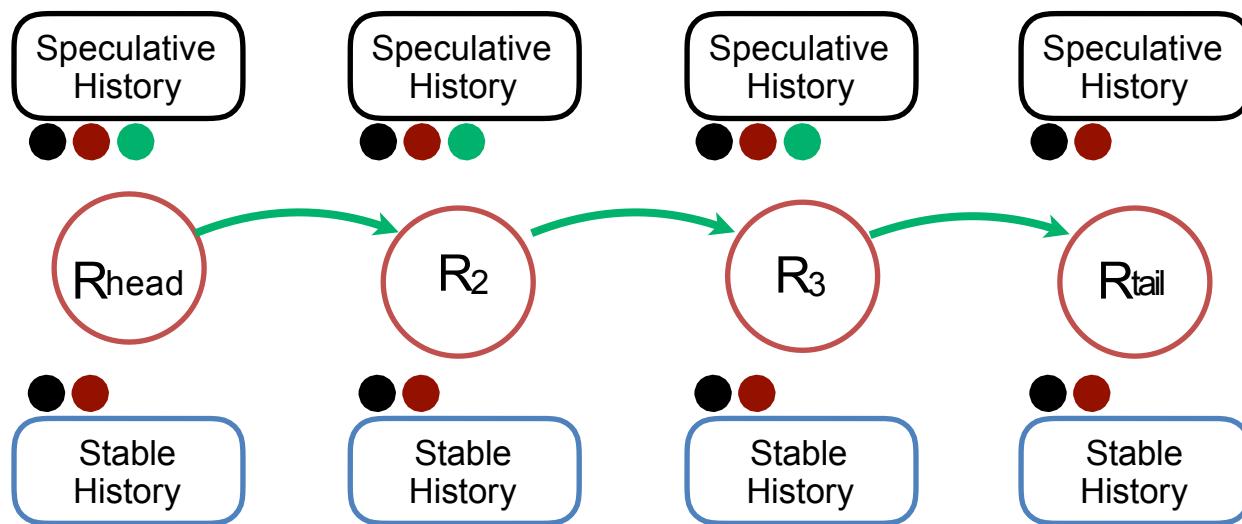
# Updates



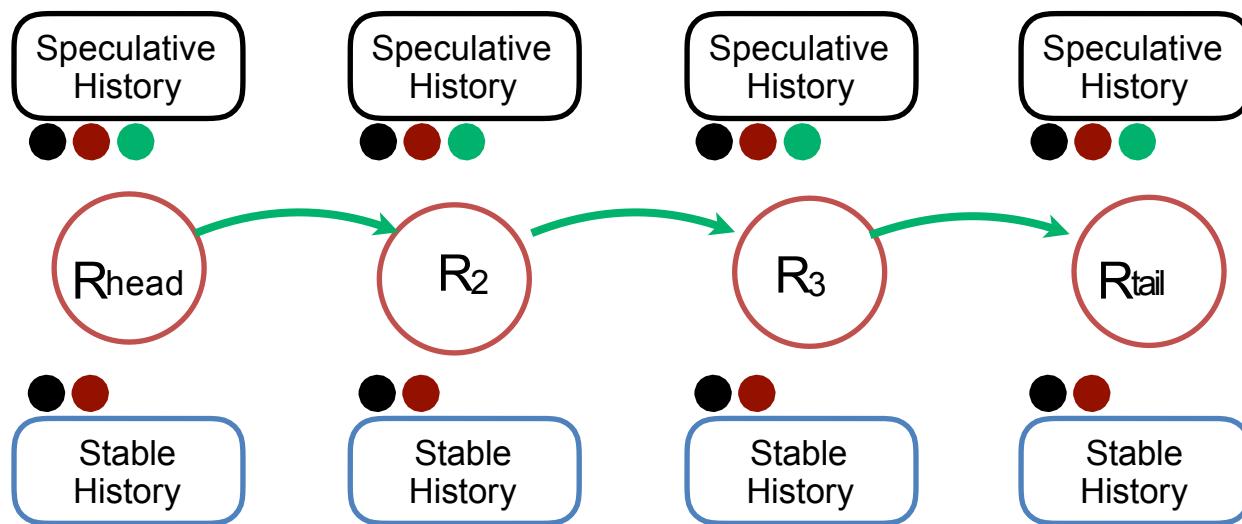
# Updates



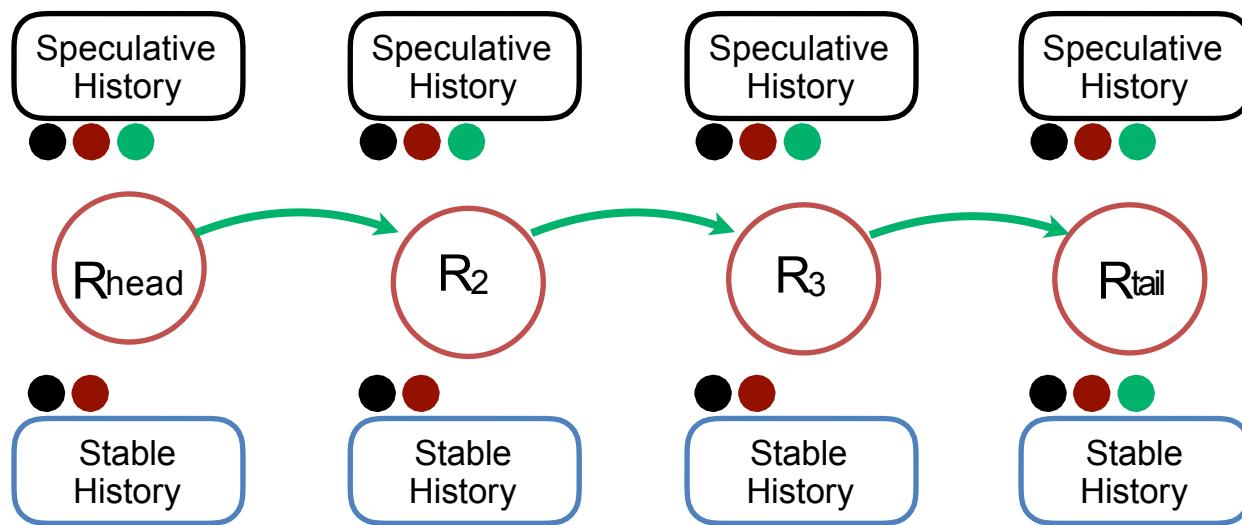
# Updates



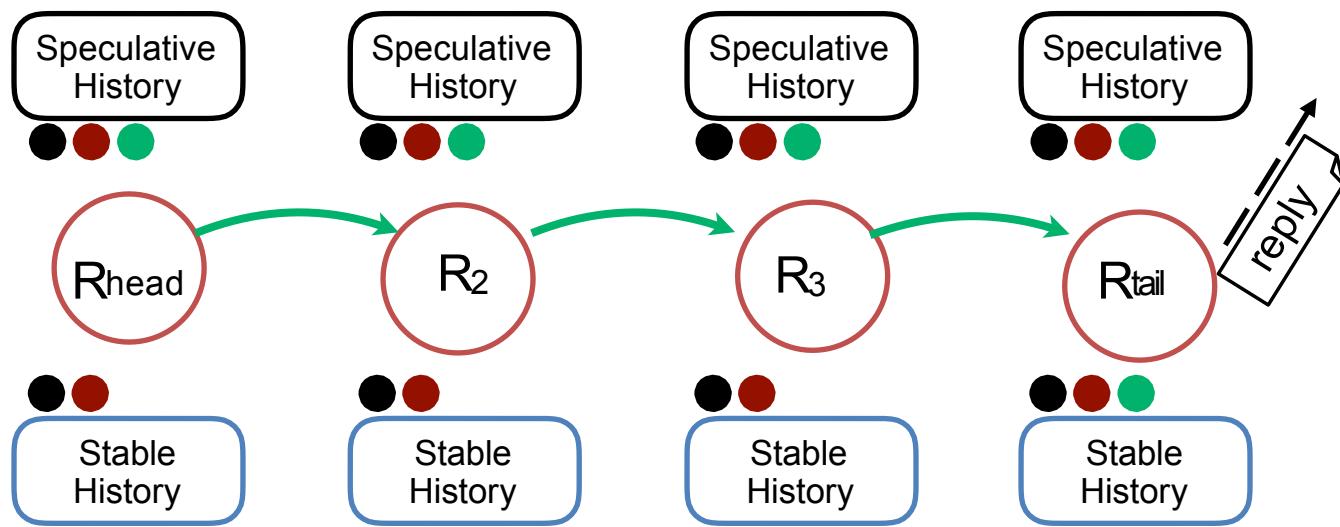
# Updates



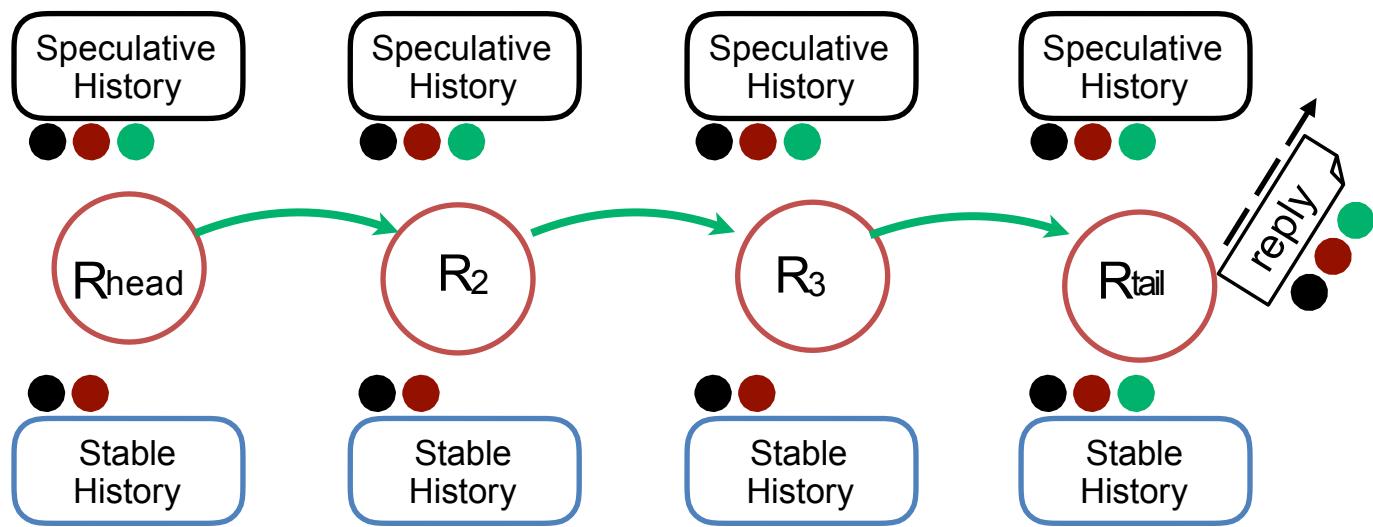
# Updates



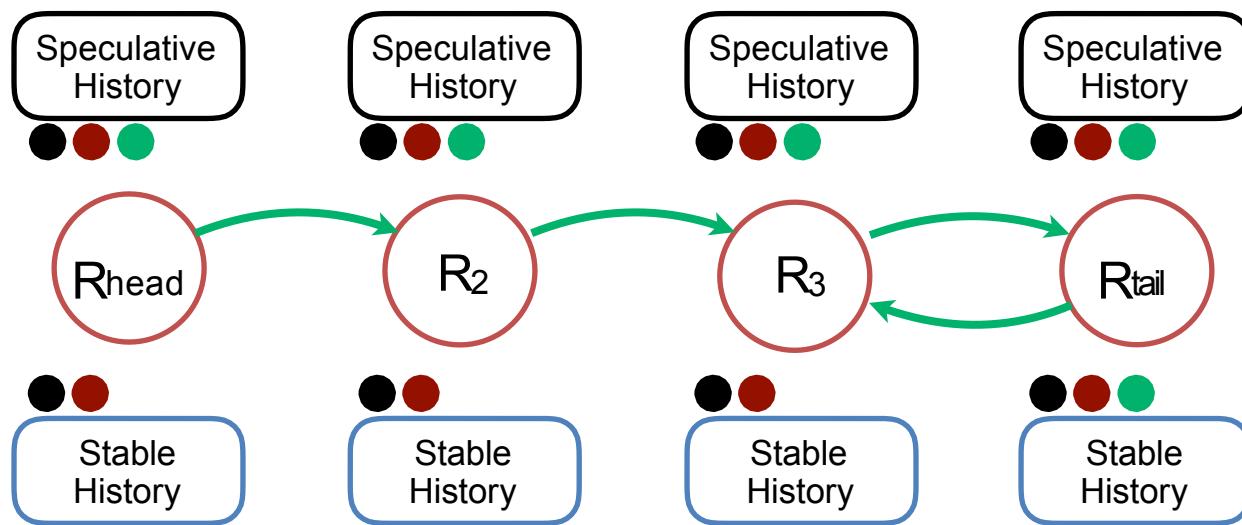
# Updates



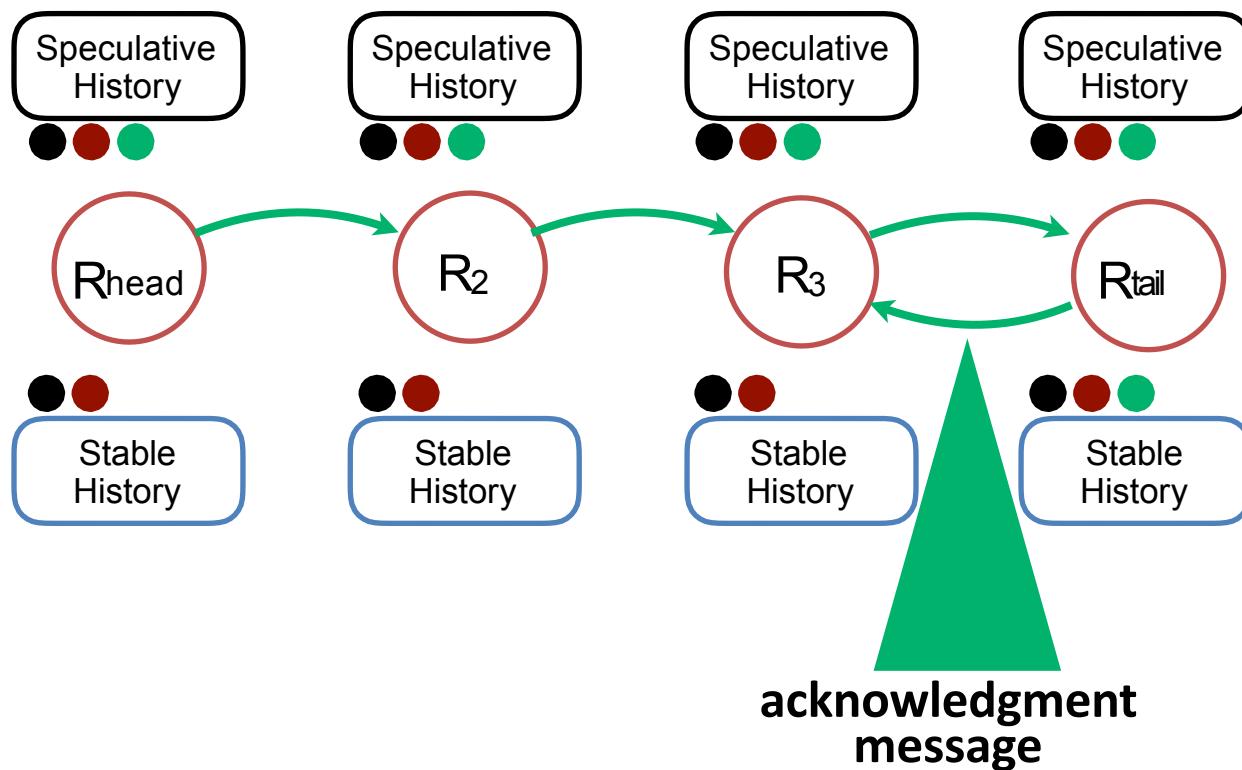
# Updates



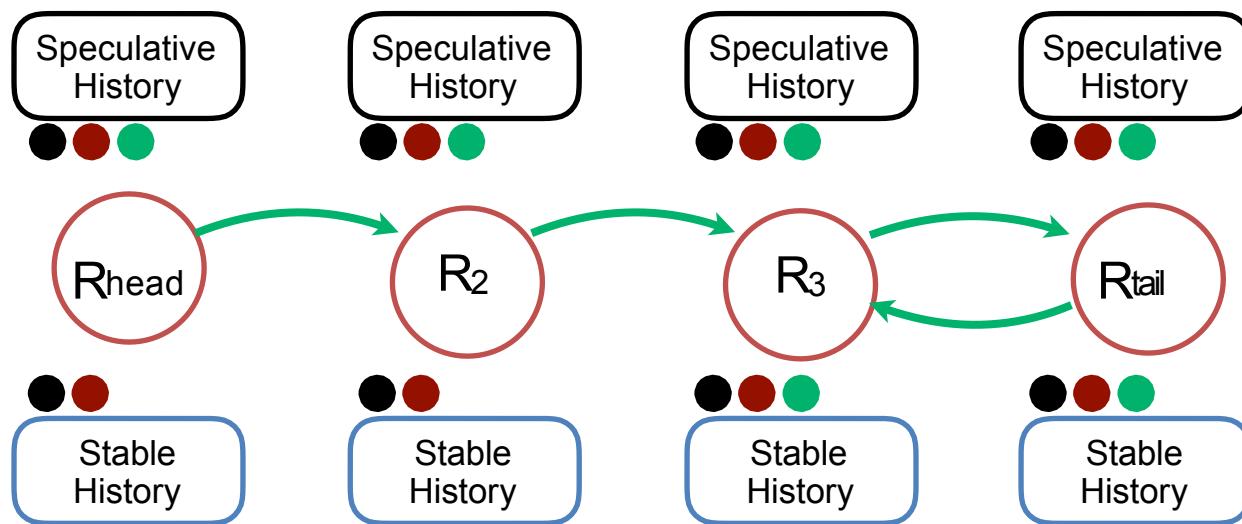
# Updates



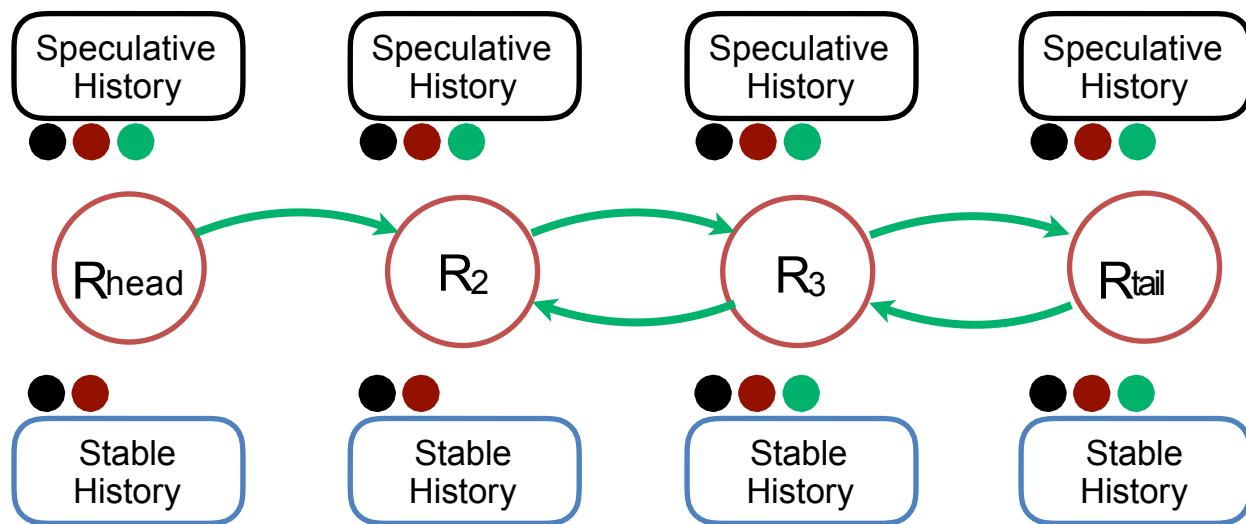
# Updates



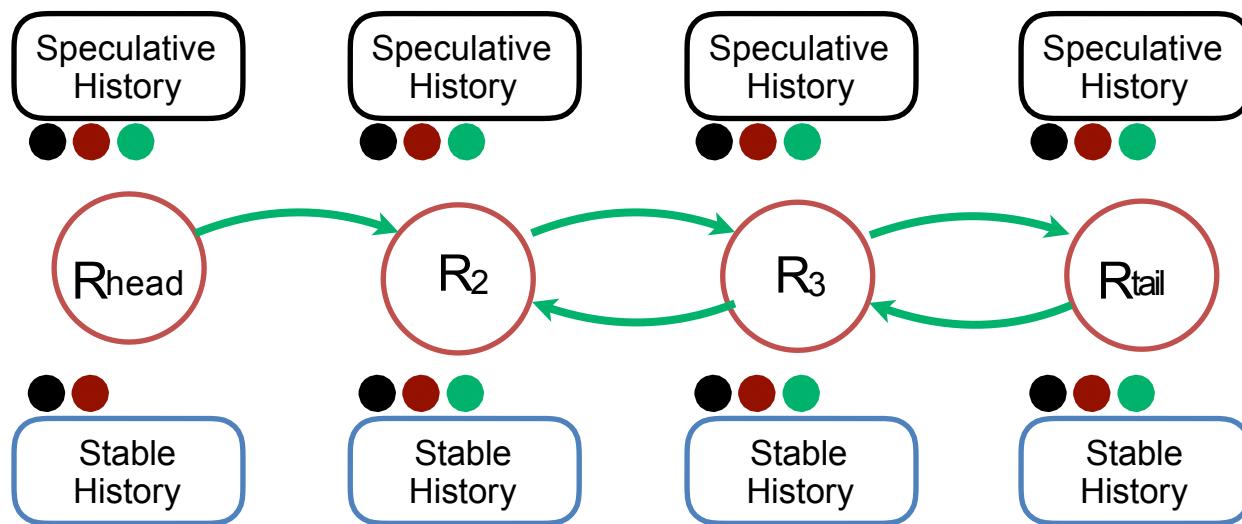
# Updates



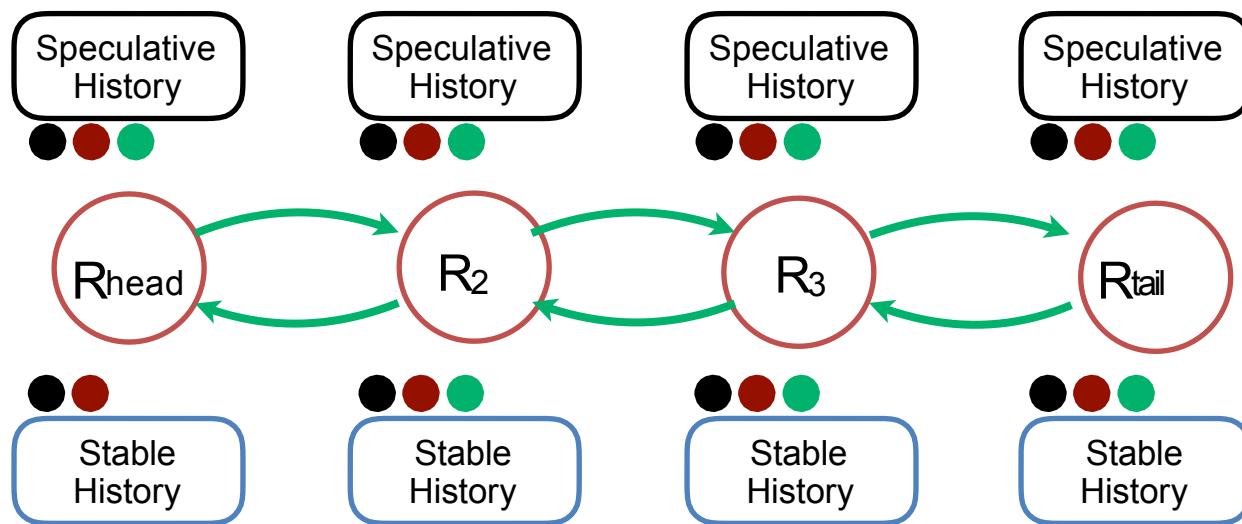
# Updates



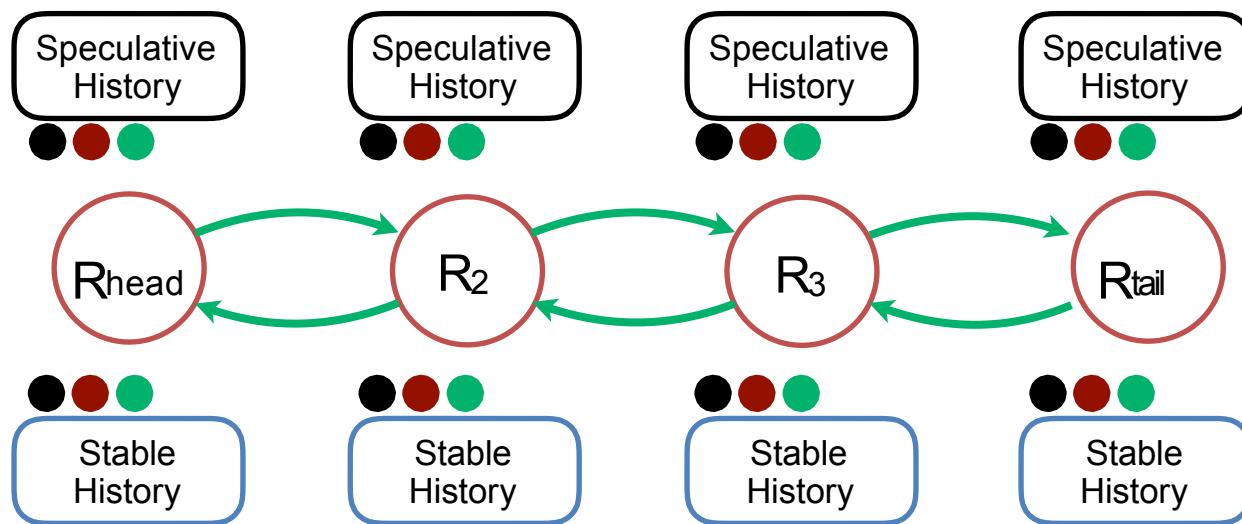
# Updates



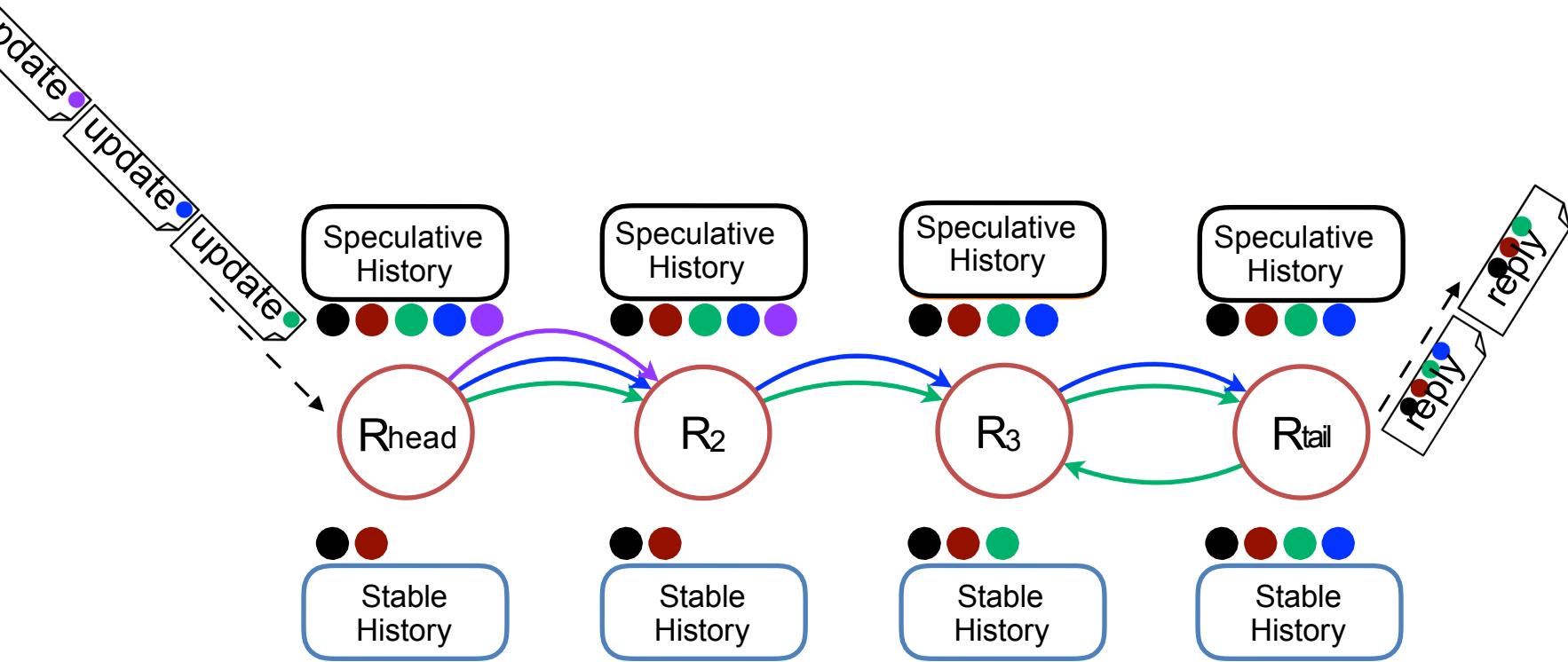
# Updates



# Updates

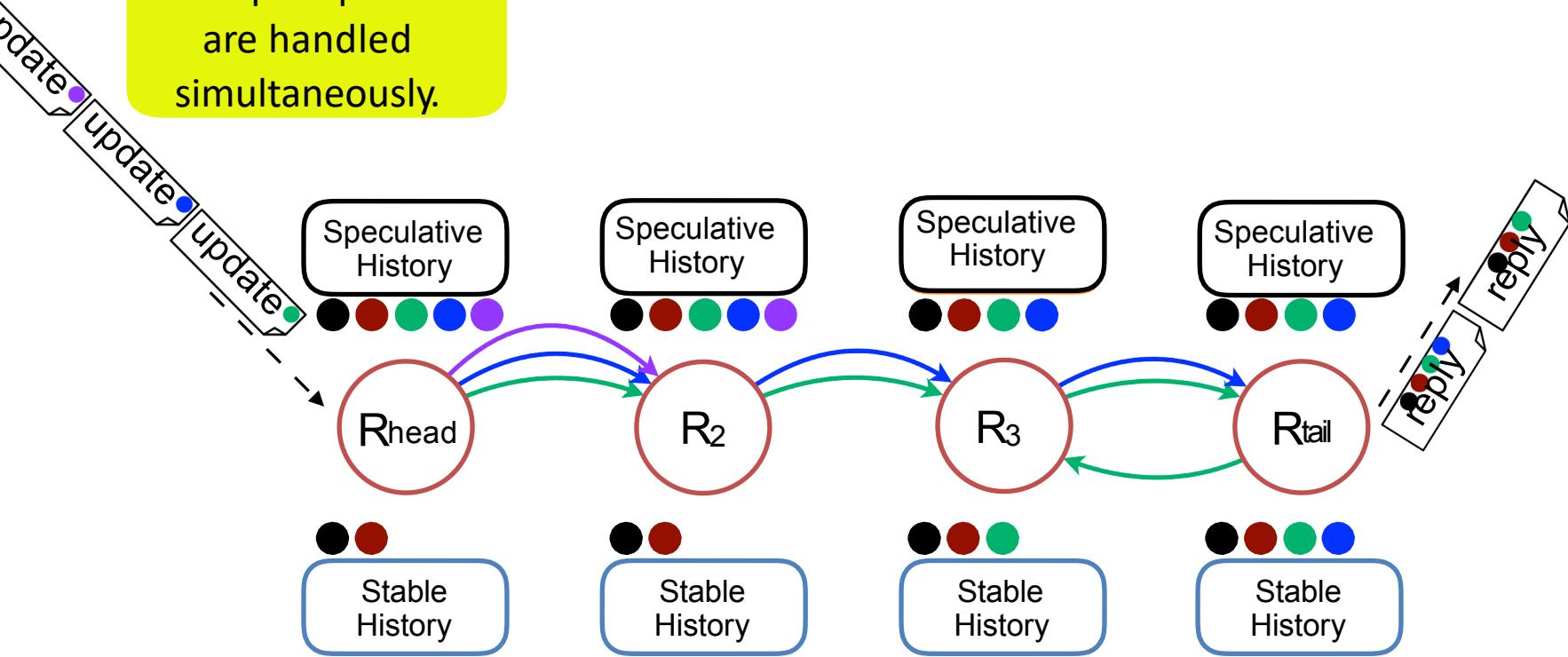


# Updates

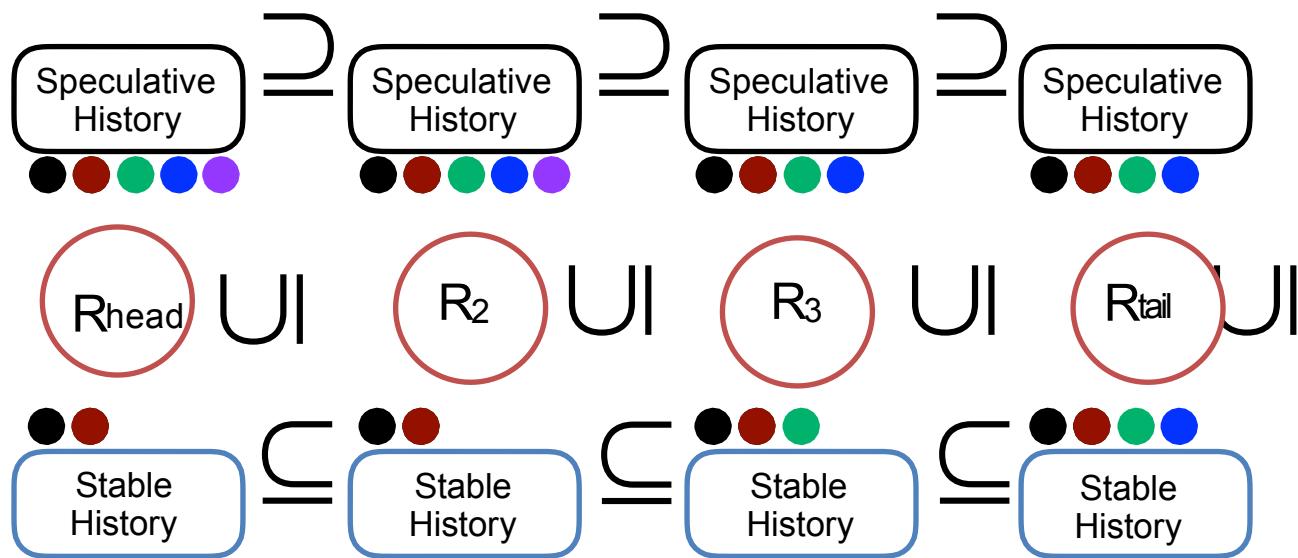


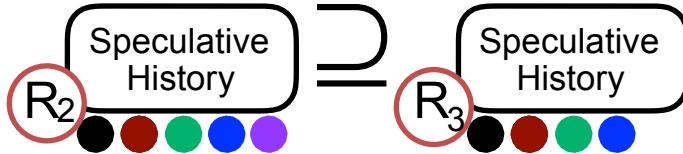
# Updates

Multiple updates  
are handled  
simultaneously.



# Updates



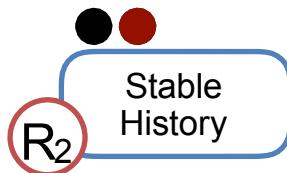


The speculative history of a node's successor is a **subset** of that node's speculative history.

---

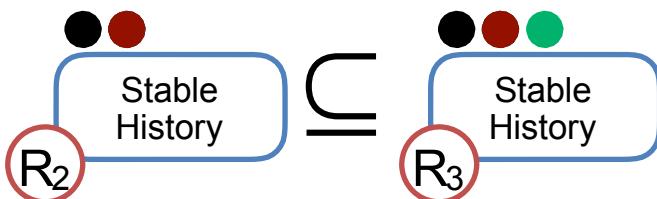


U|



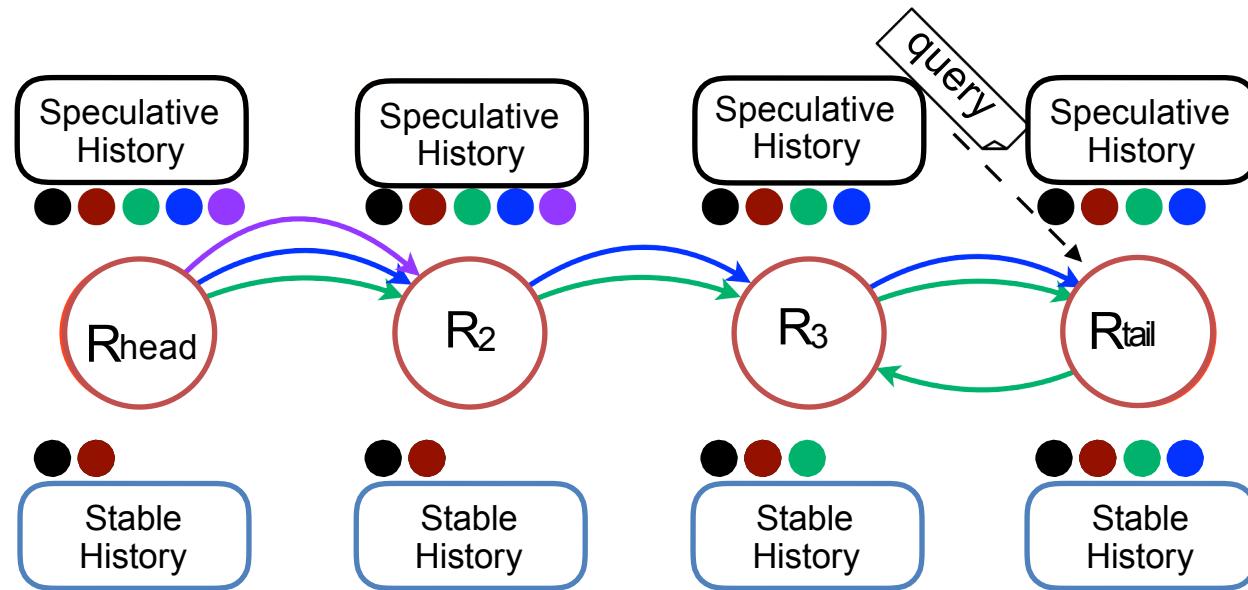
The speculative history of a node is a **superset** of its stable history.

---

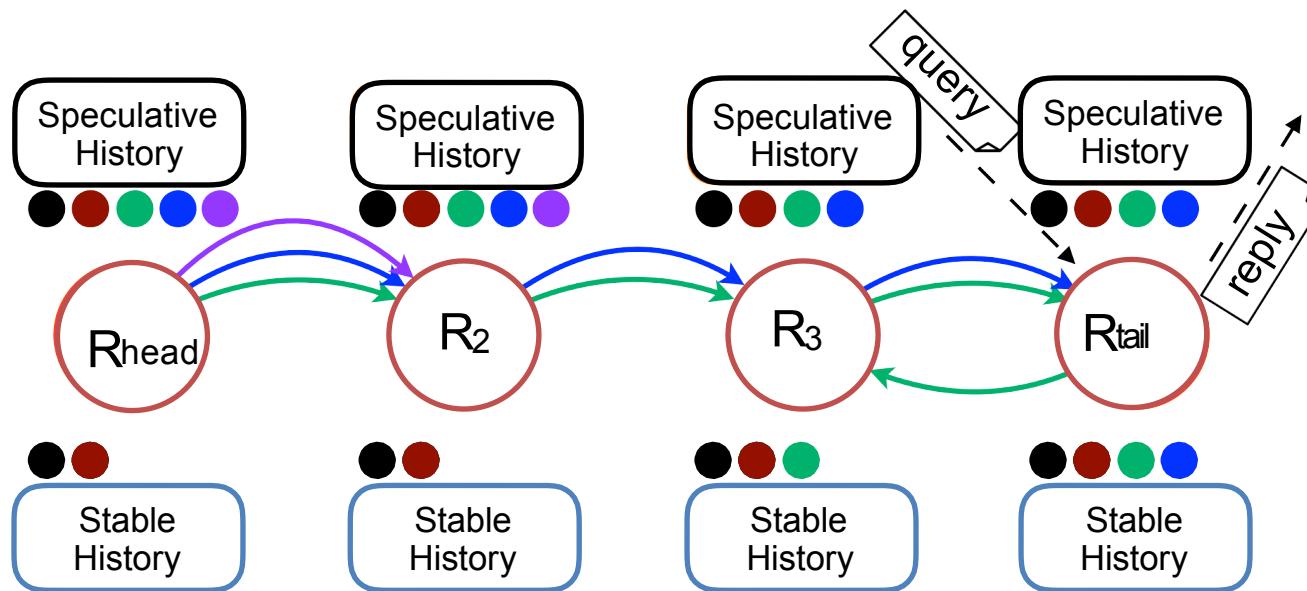


The stable history of a node's successor is a **superset** of that node's stable history.

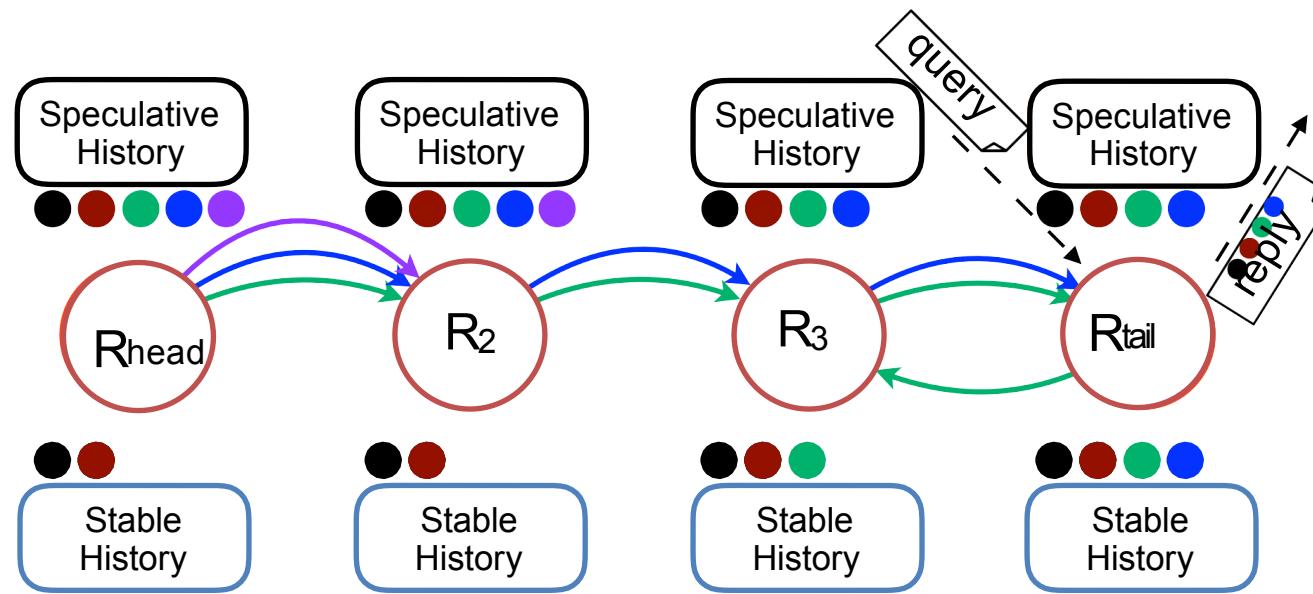
# Queries



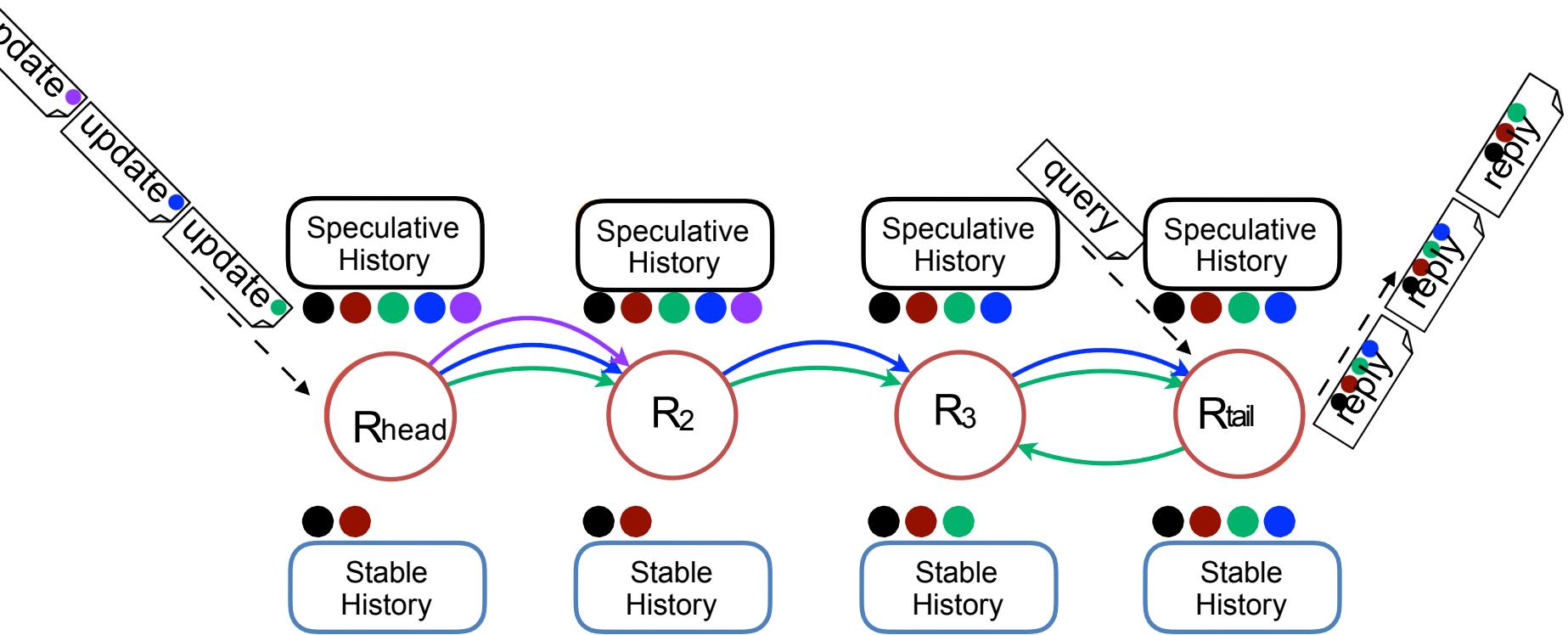
# Queries



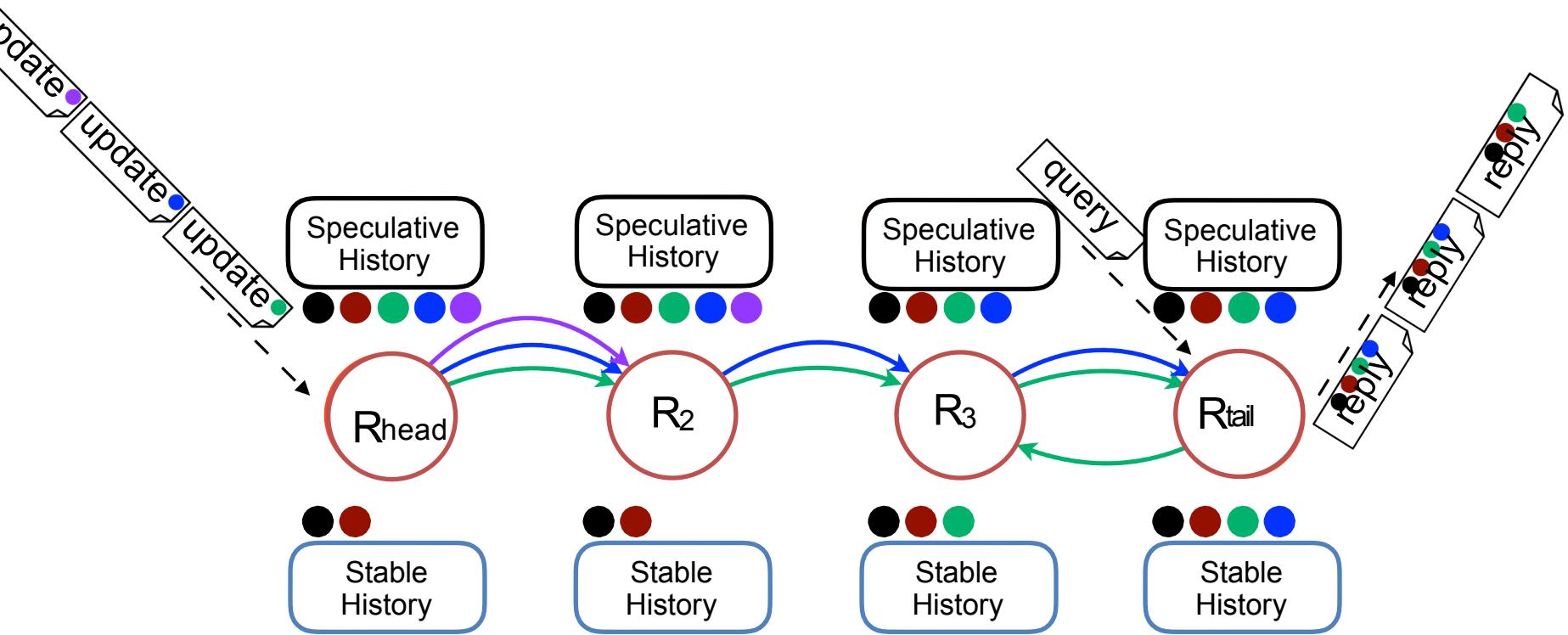
# Queries



# Queries



# Queries



**The tail is the point of linearization!**

# Self-study Questions

- Do speculative histories have to be persisted to disk?
- Do stable histories have to be persisted to disk?
- While the subset relationships are neat, do they serve an ulterior purpose?
- Could a reply to an update be sent before transitioning the update to the stable history of the tail?
- Can chain replication work effectively with a single replica?





Pixabay.com

# CHAIN REPLICATION

## FAILUR HANDLING & RECONFIGURATIONS

Inspiration for this lecture taken from a talk given by Deniz Altınbüken.

# Failure Handling & Reconfigurations

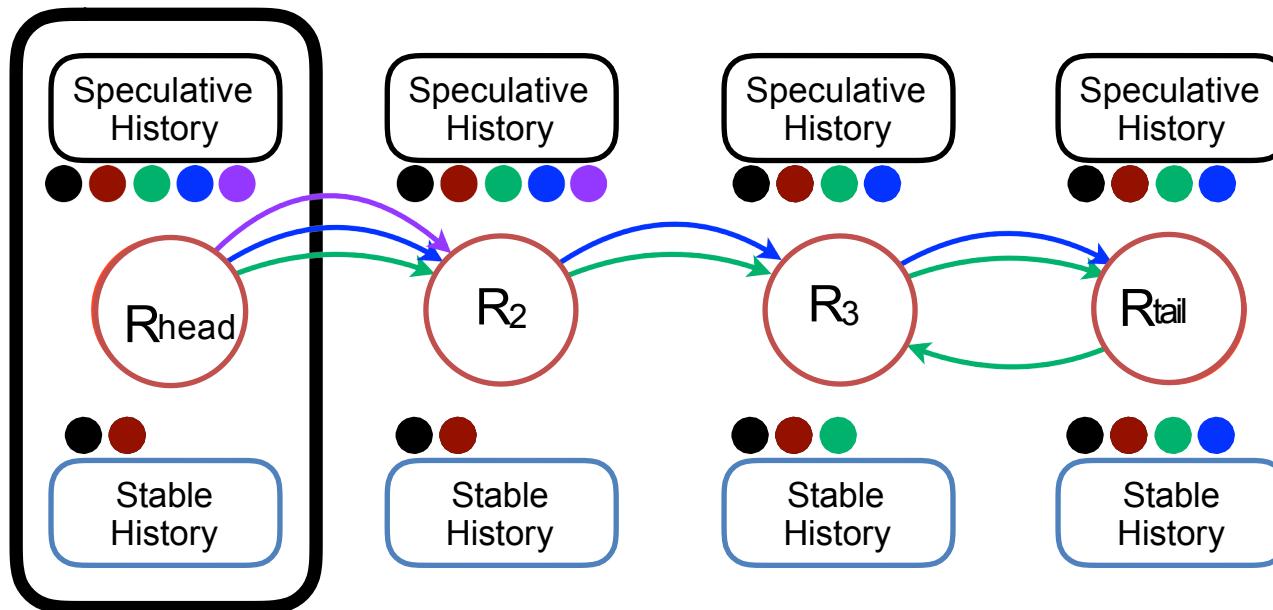
- Chain replication operations:
  - Updates
  - Queries
  - Failures
  - Reconfigurations

# Failures

- Head failure
- Middle node failure
- Tail failure

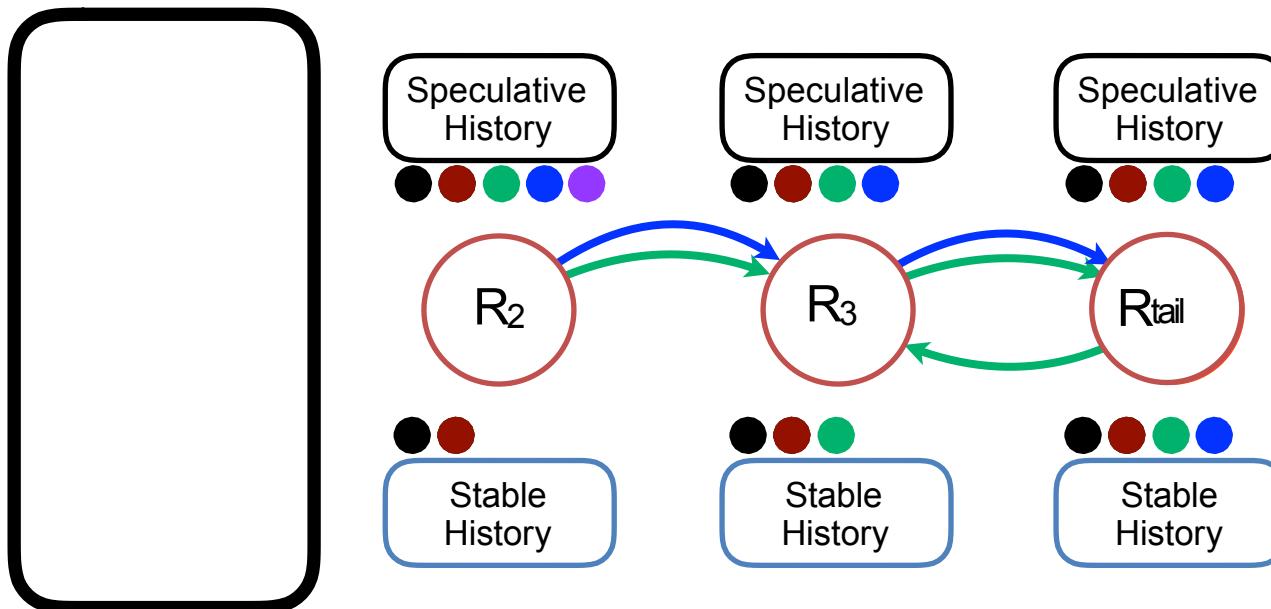


# Head Failure I



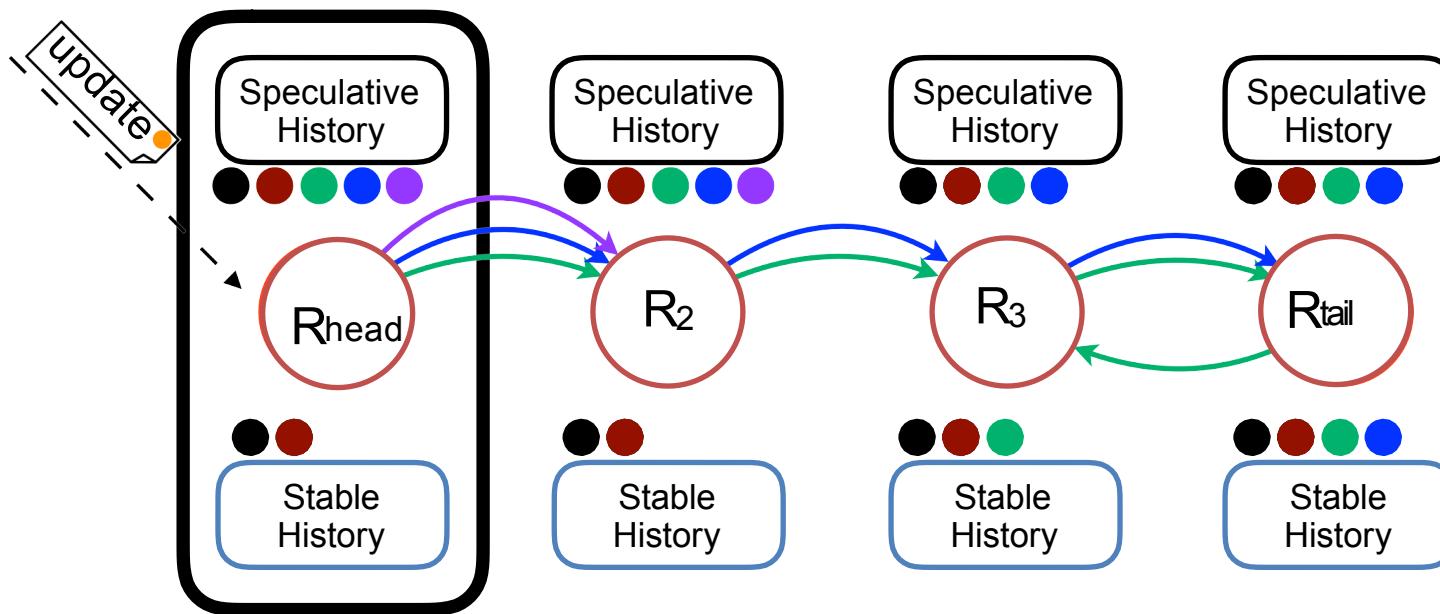
# Head Failure I

$R_2$  becomes new head



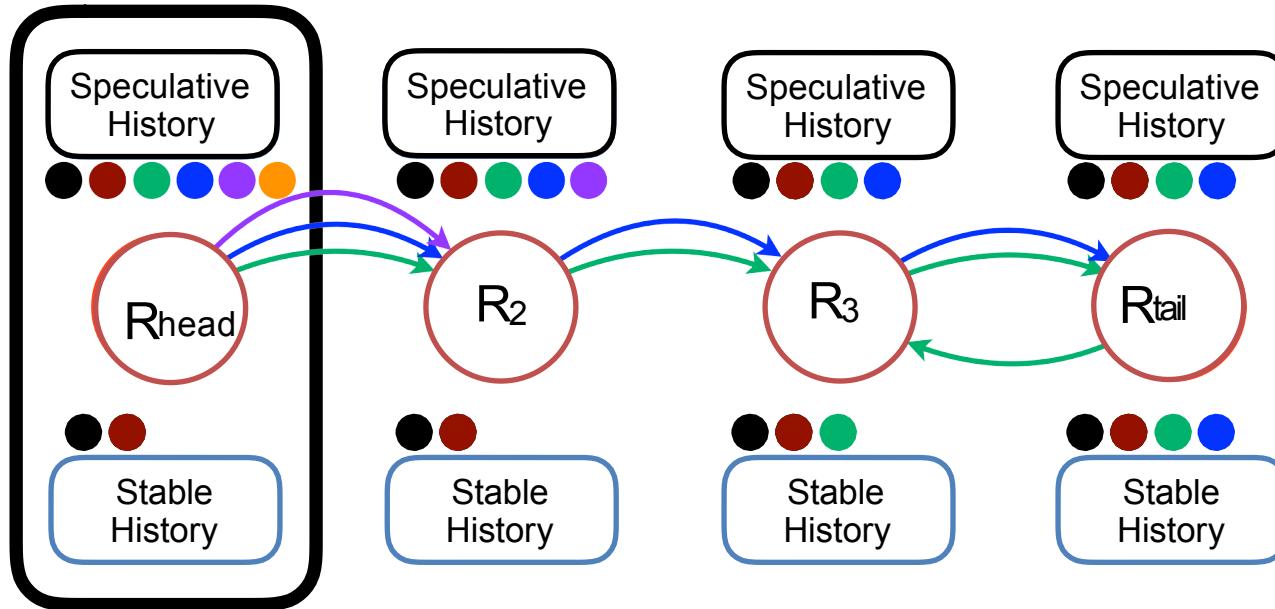
# Head Failure I

In-flight, non-propagated updates



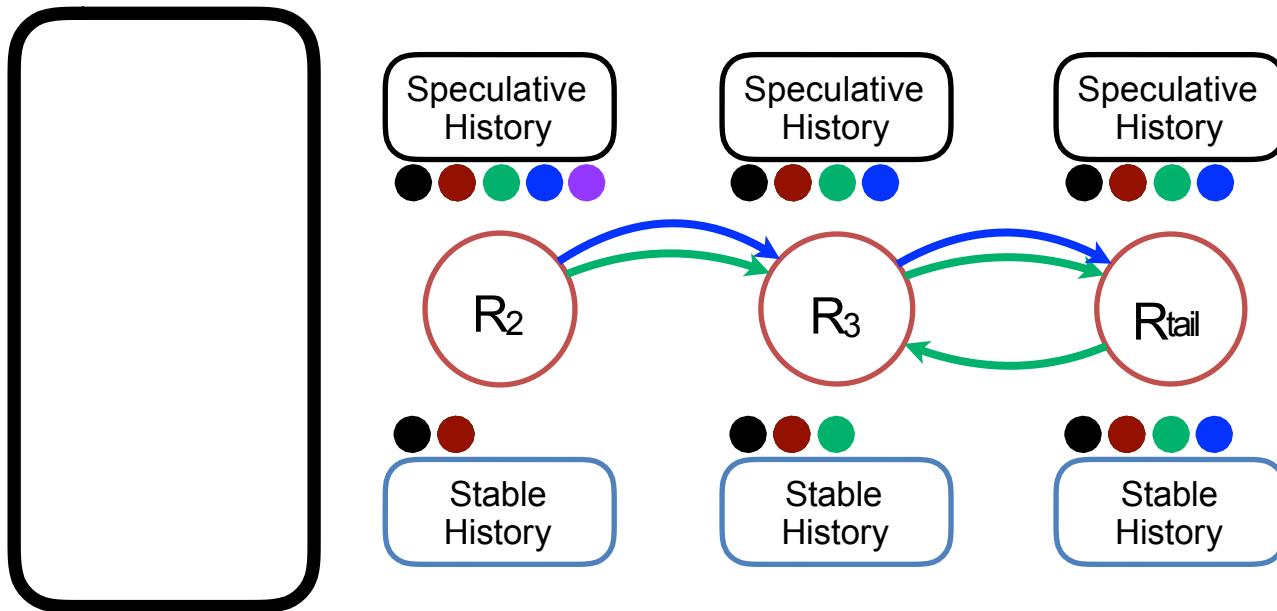
# Head Failure I

In-flight, non-propagated updates



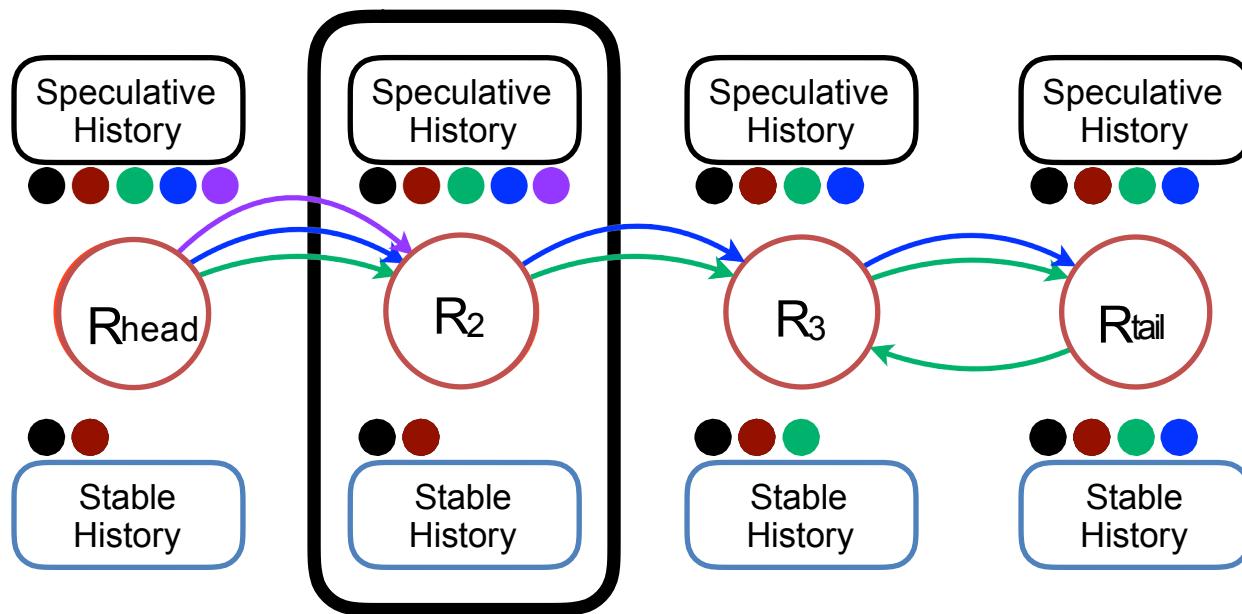
# Head Failure I

In-flight, non-propagated updates

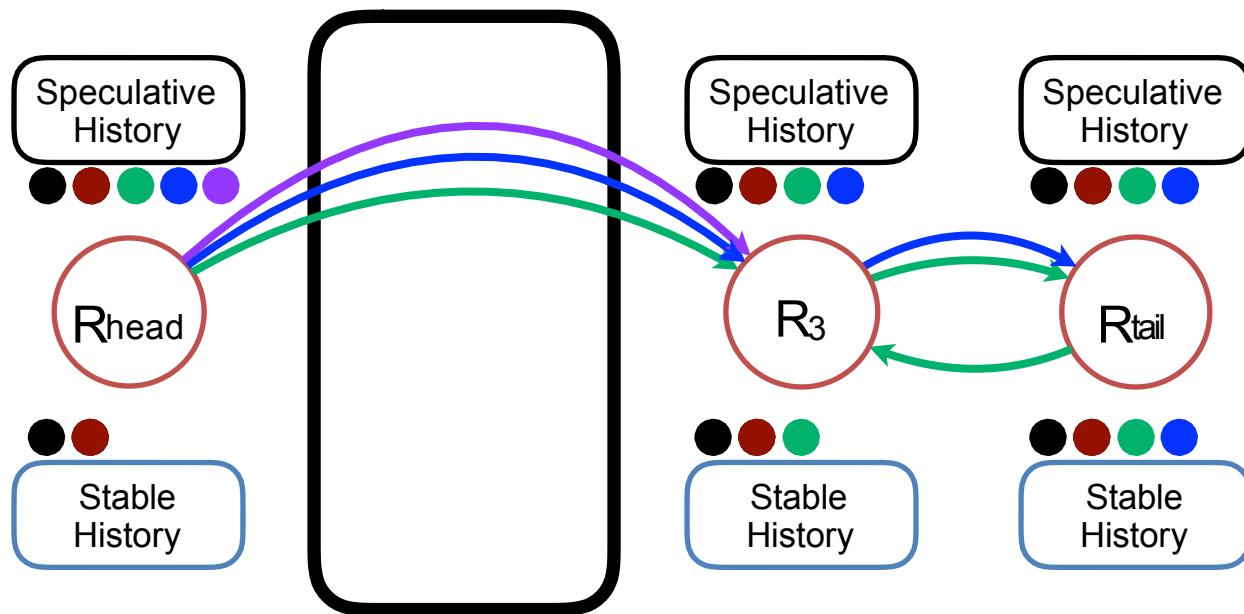


**Client would not receive a reply, timeout, and retry**

# Middle Node Failure II

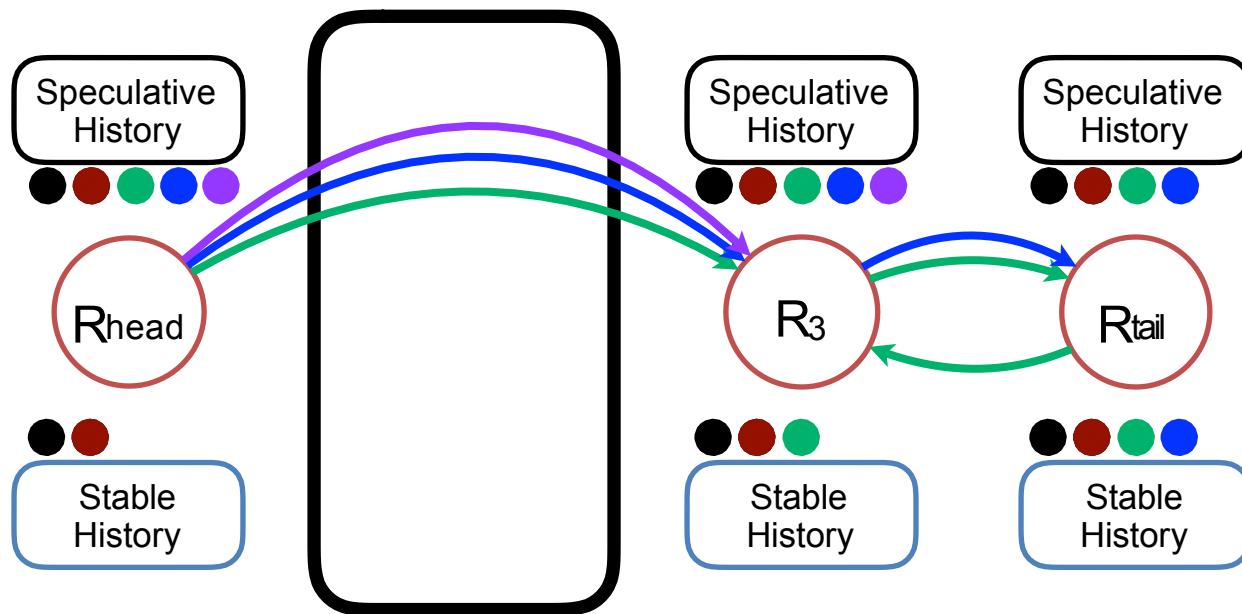


# Middle Node Failure II



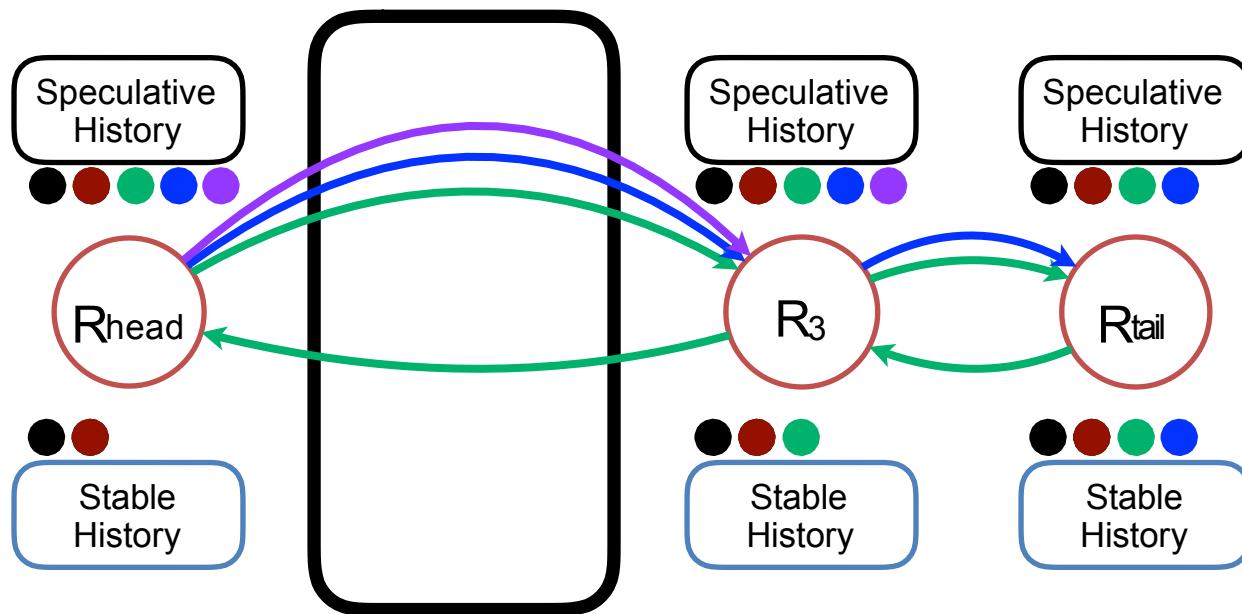
Predecessor needs to talk to failed node's successor

# Middle Node Failure II

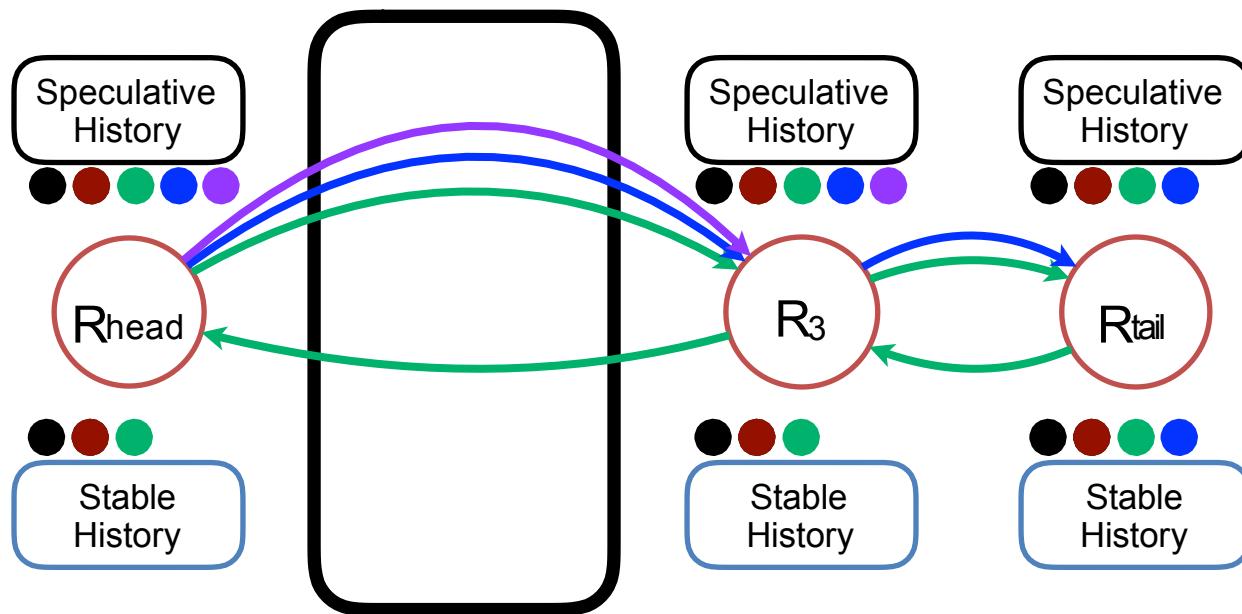


**Predecessor propagates update to new successor**

# Middle Node Failure II

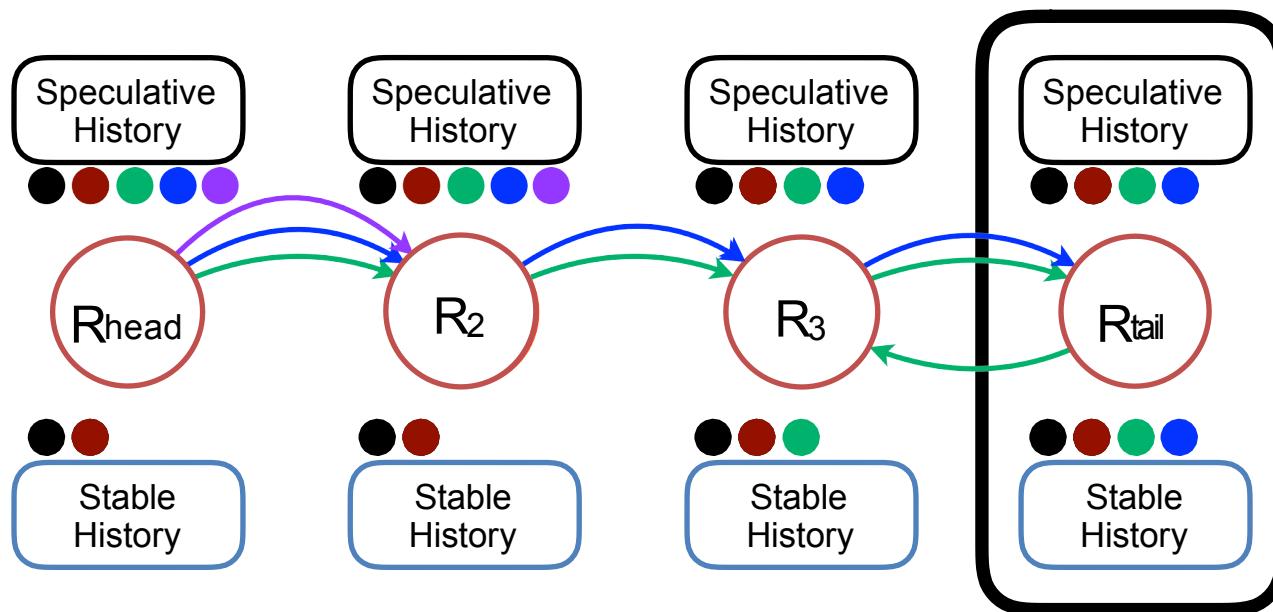


# Middle Node Failure II

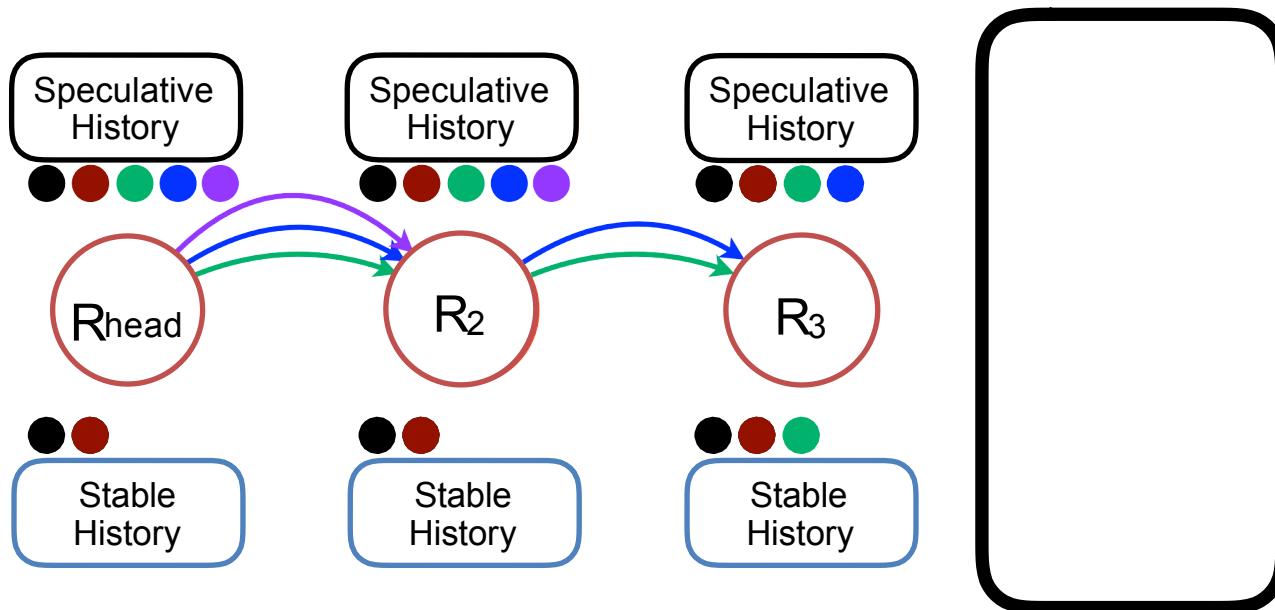


**Successor propagates in-flight acknowledgements to new predecessor**

# Tail Failure III

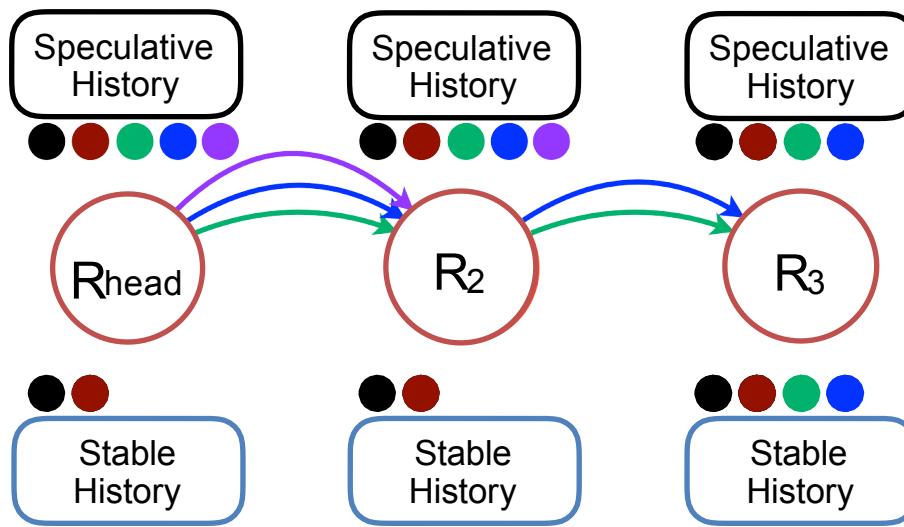


# Tail Failure III



**R<sub>3</sub> becomes new tail**

# Tail Failure III



**R<sub>3</sub> flushes its speculative history s.t. stable  
equals speculative history again**

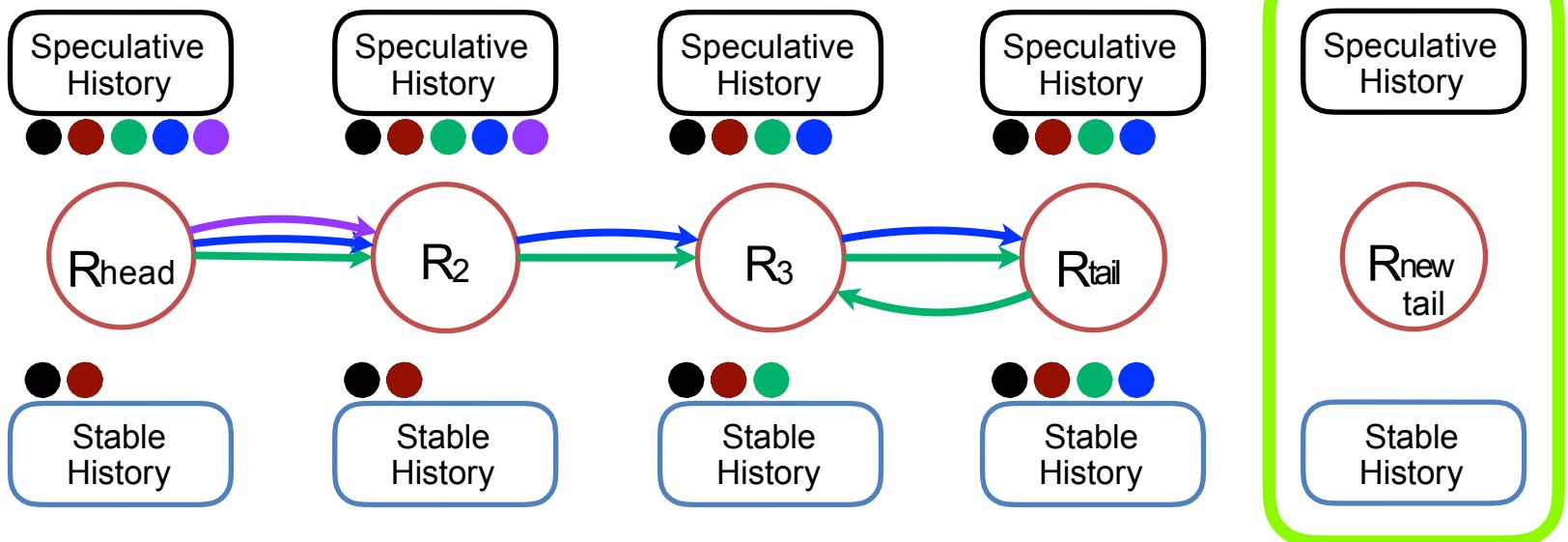
# Reconfigurations

- Adding a new node for failure recovery
- Adding a new node to extend topology
- Setting up a new chain of replicas

# Adding a New Node I

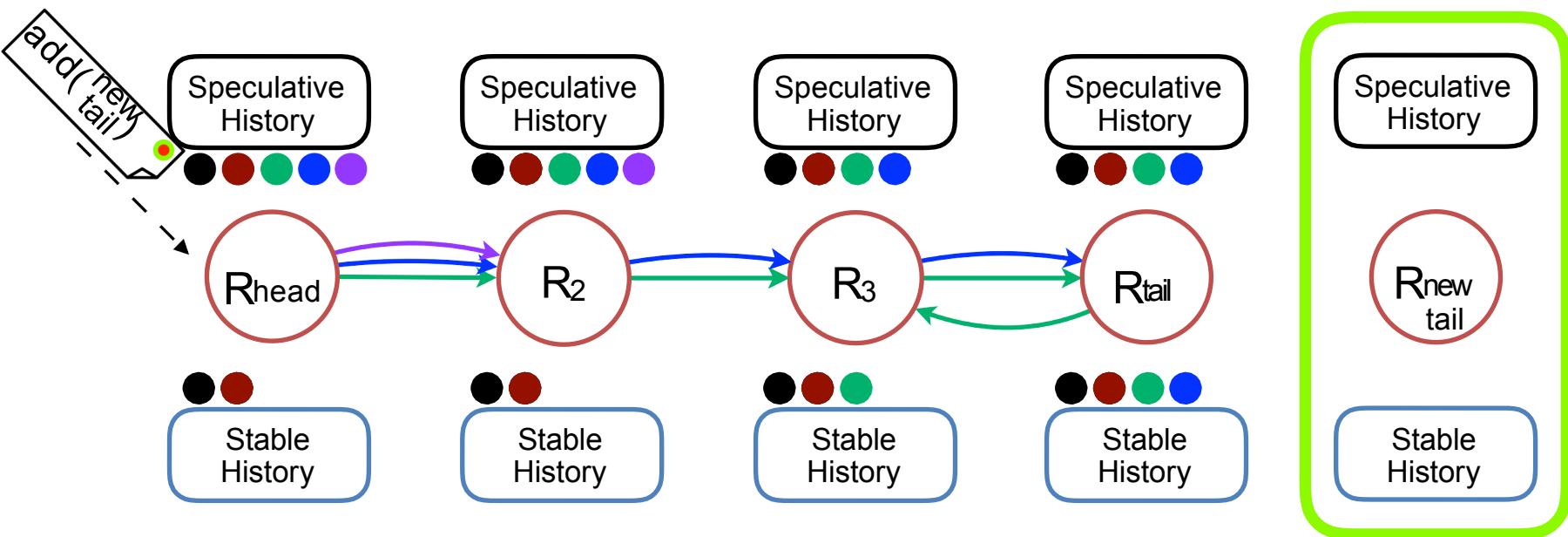
## A Configuration Change

Adding an initially empty node



# Adding a New Node II

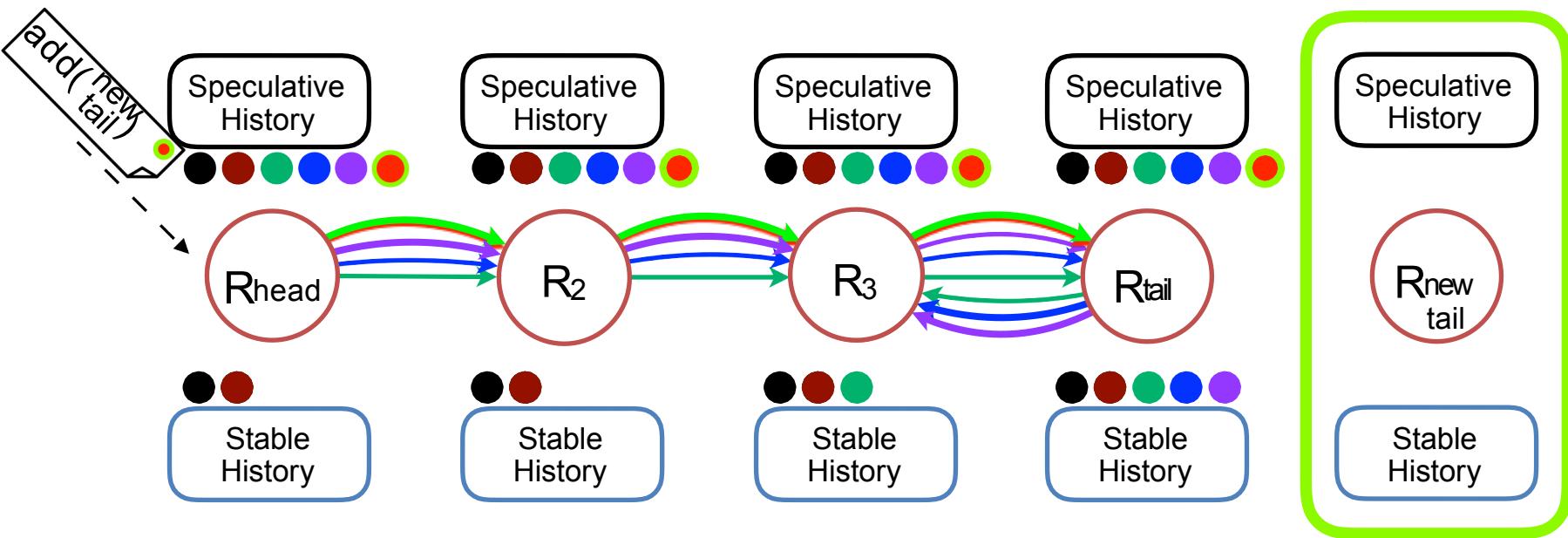
## A Configuration Change



- New nodes are added to chain with special **configuration updates**, added to histories: **add(nodeid)**
- Entire chain is build in this manner

# Adding a New Node

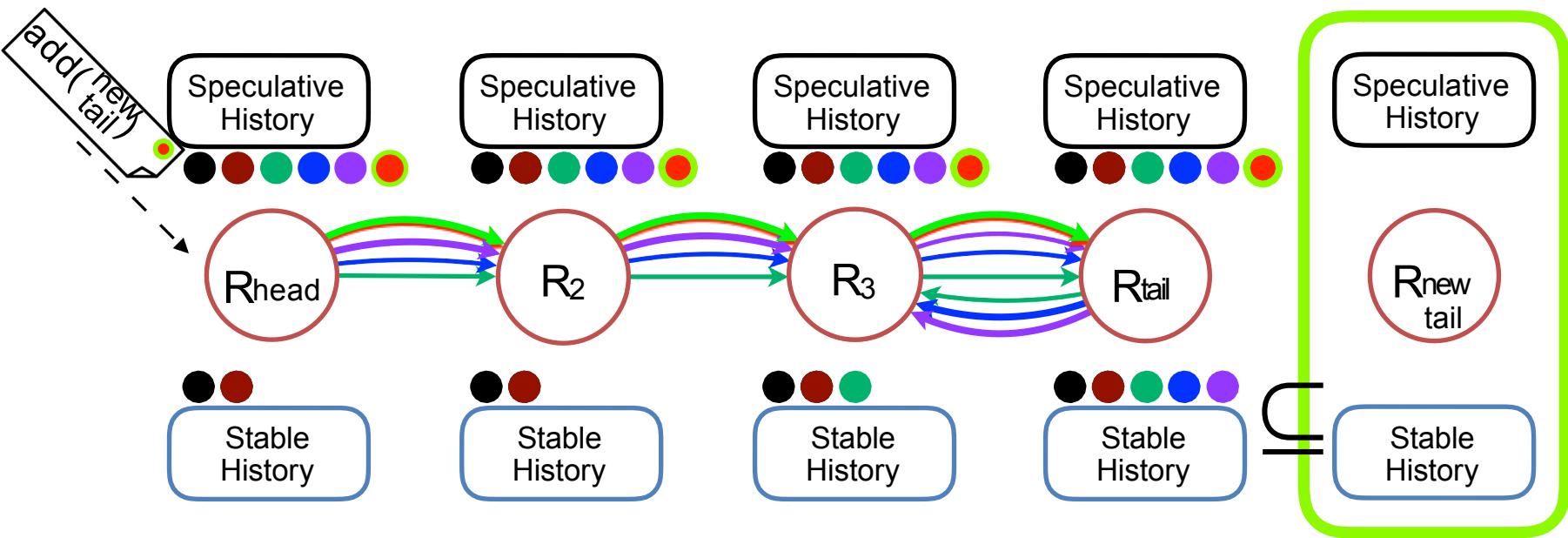
## Inferring Configuration



- By looking at **order of these updates**, a node can **determine configuration of chain**
- Old tail discovers it no longer is the tail (via receipt of )

# Adding a New Node I

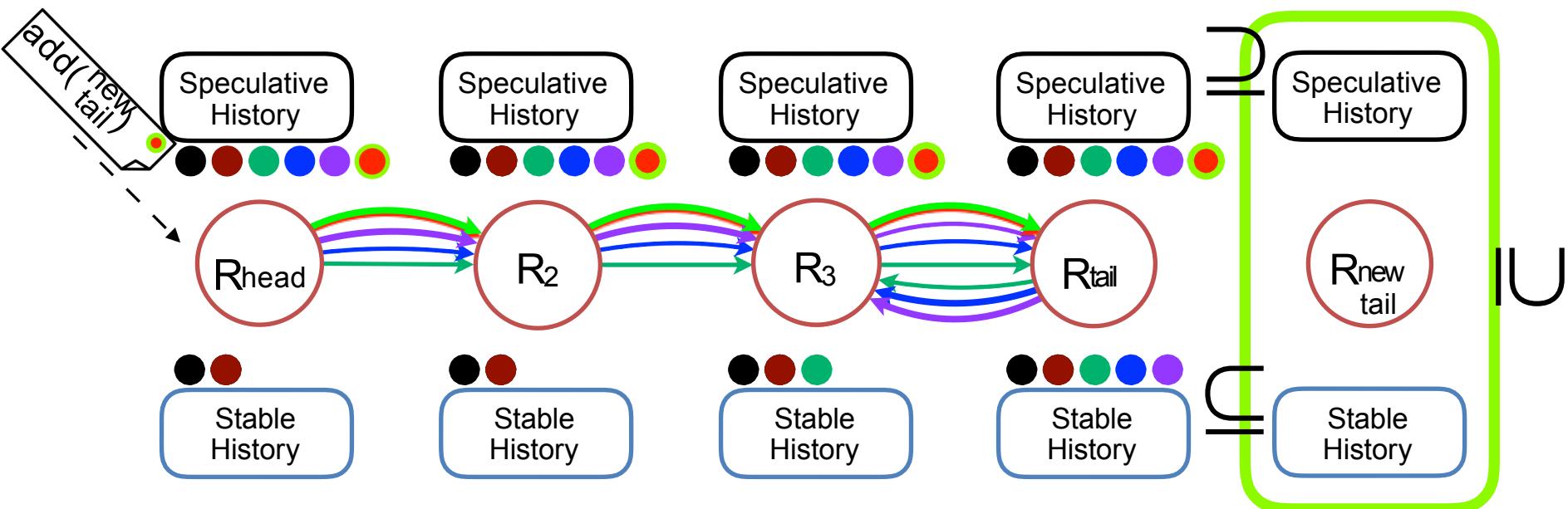
## Relationship Among Histories



- **Stable history of new tail should be superset of stable history of old tail**

# Adding a New Node II

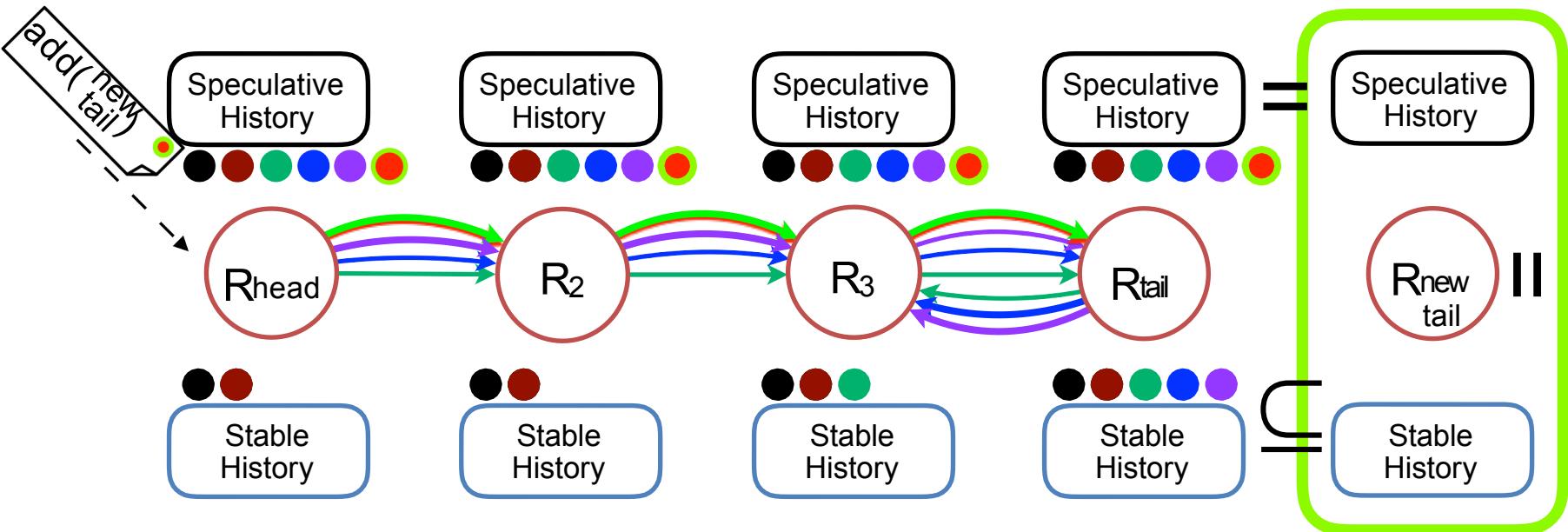
## Relationship Among Histories



- **Speculative history** of new tail should be a **superset** of its **stable history**

# Adding a New Node III

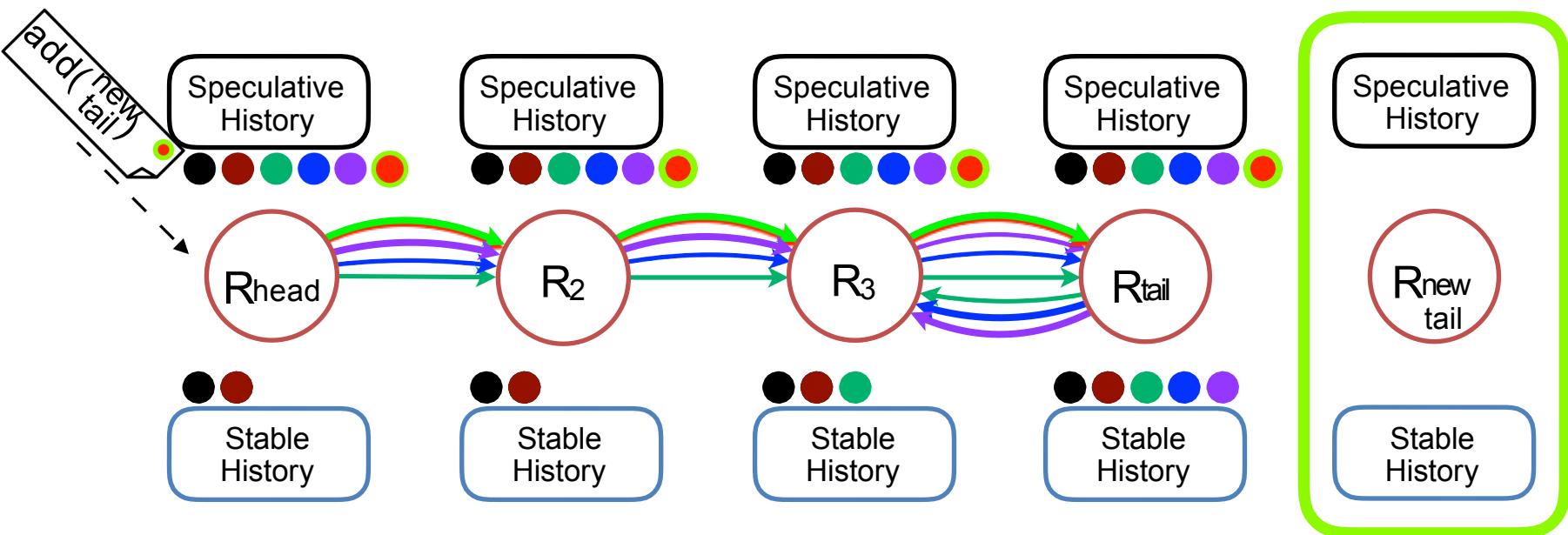
## Relationship Among Histories



- **Speculative and stable histories** of new tail should become **equal** to the **speculative history of old tail**

# Adding a New Node IV

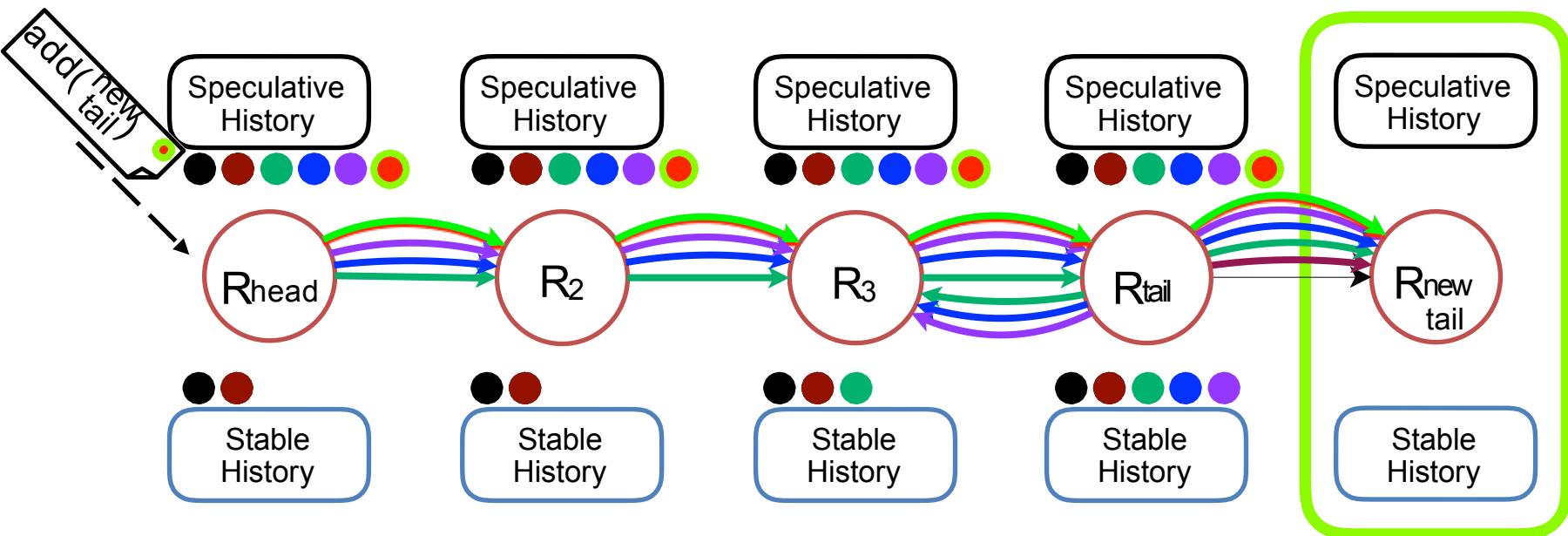
## Relationship Among Histories



- **Old tail** should not answer to queries when the new tail does

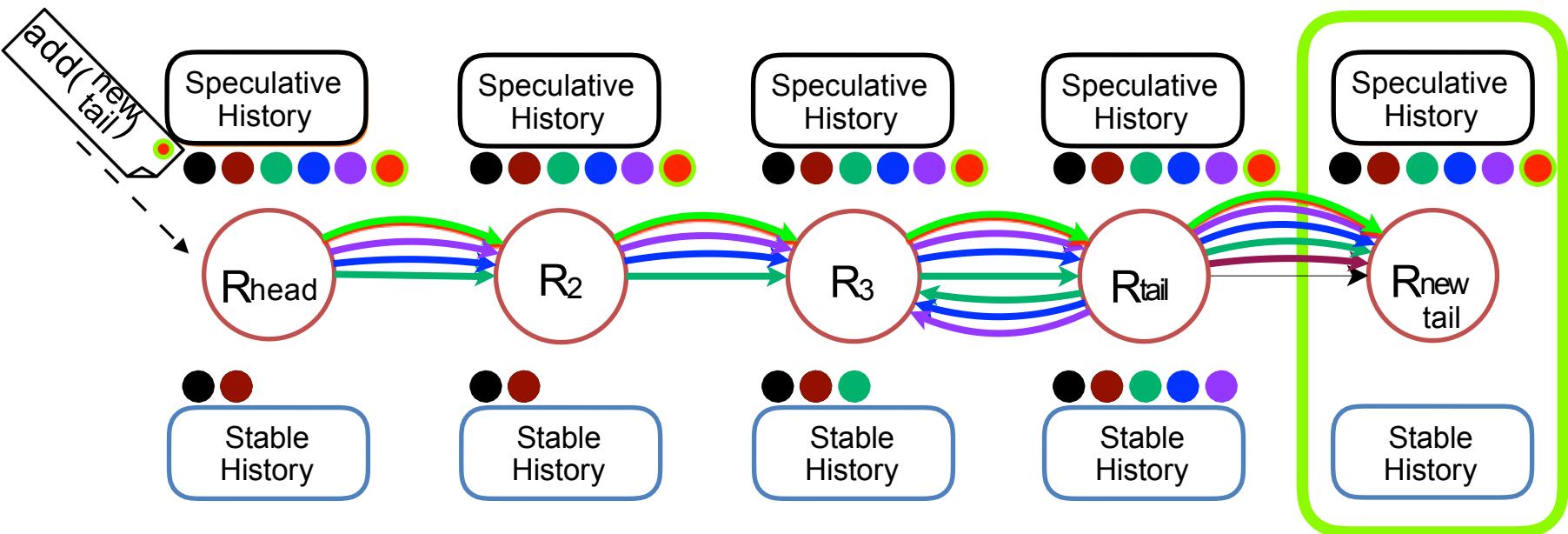
# Adding a New Node

## Flush History to New Tail



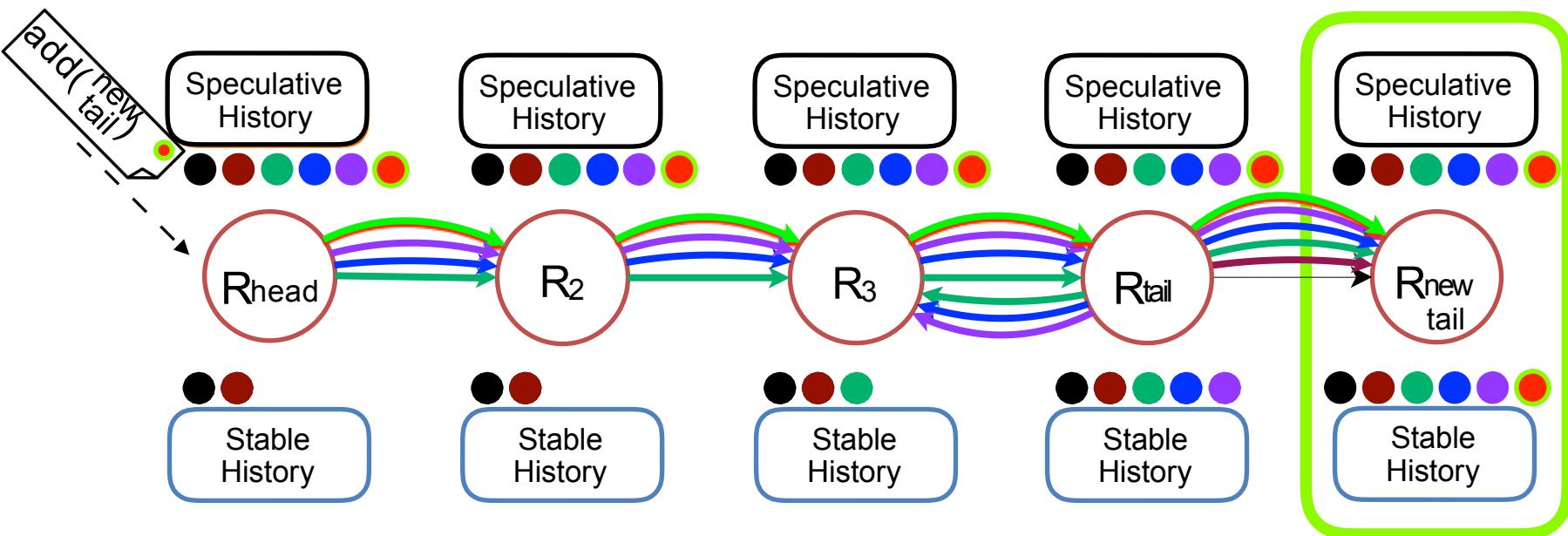
# Adding a New Node

## Flush History to New Tail



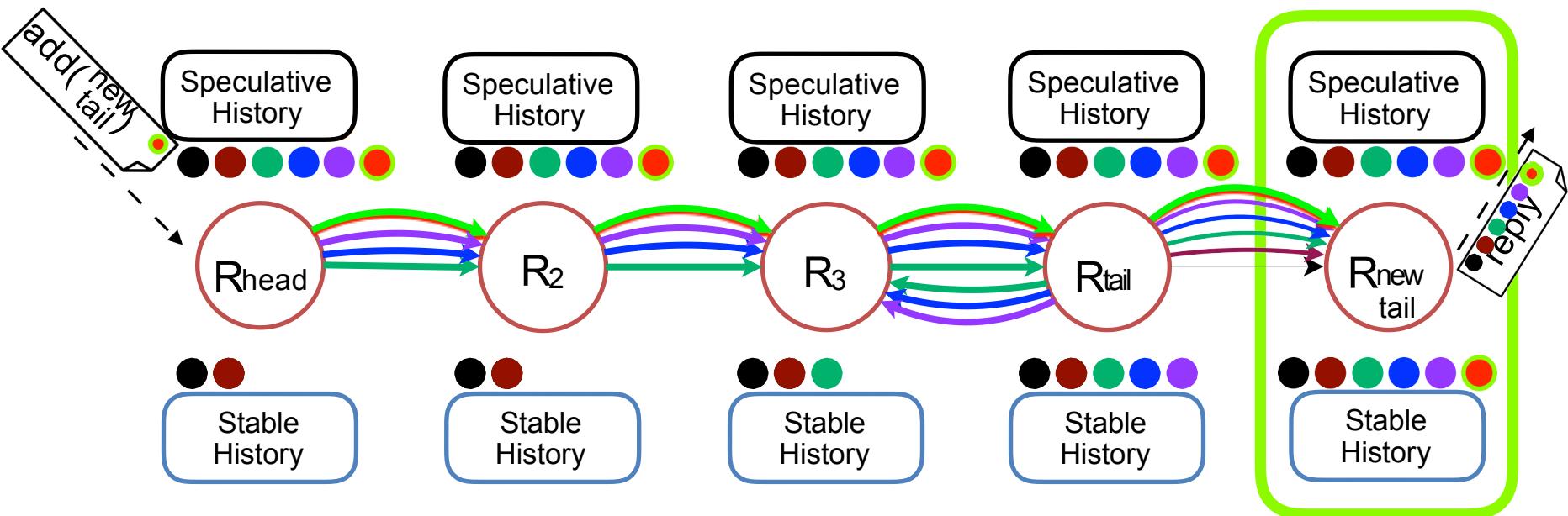
# Adding a New Node

## Copy Speculative onto Stable History



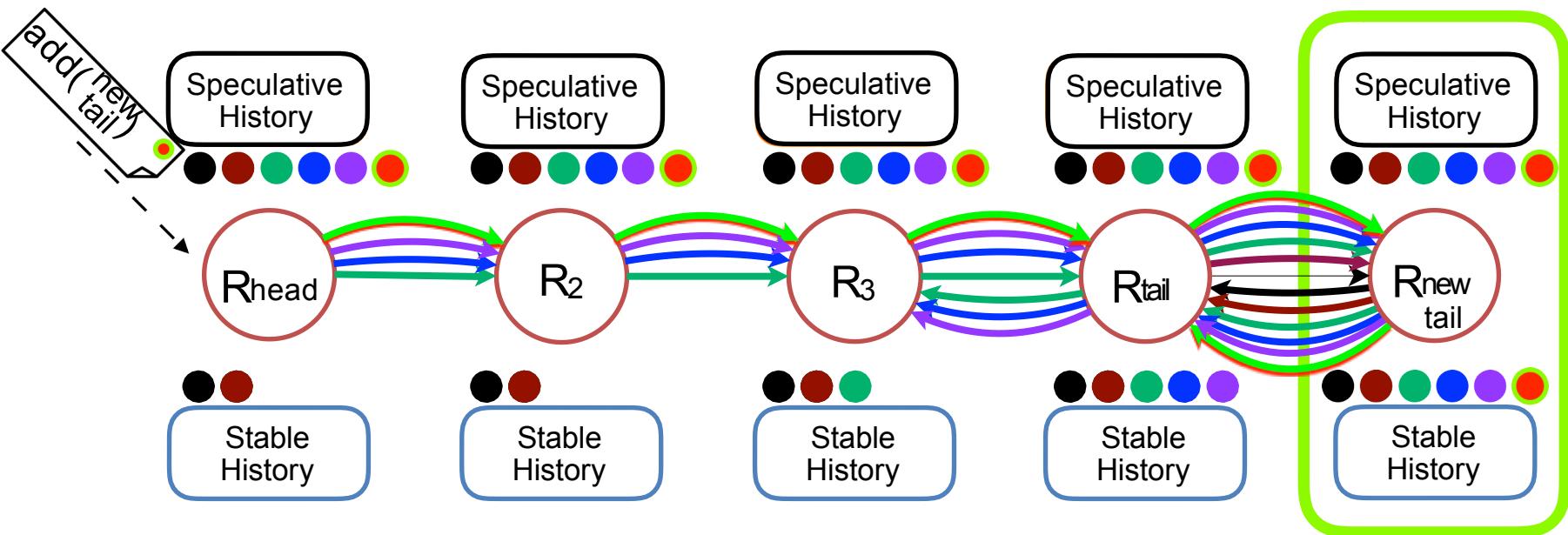
# Adding a New Node

Respond to Queries and Acknowledge Updates



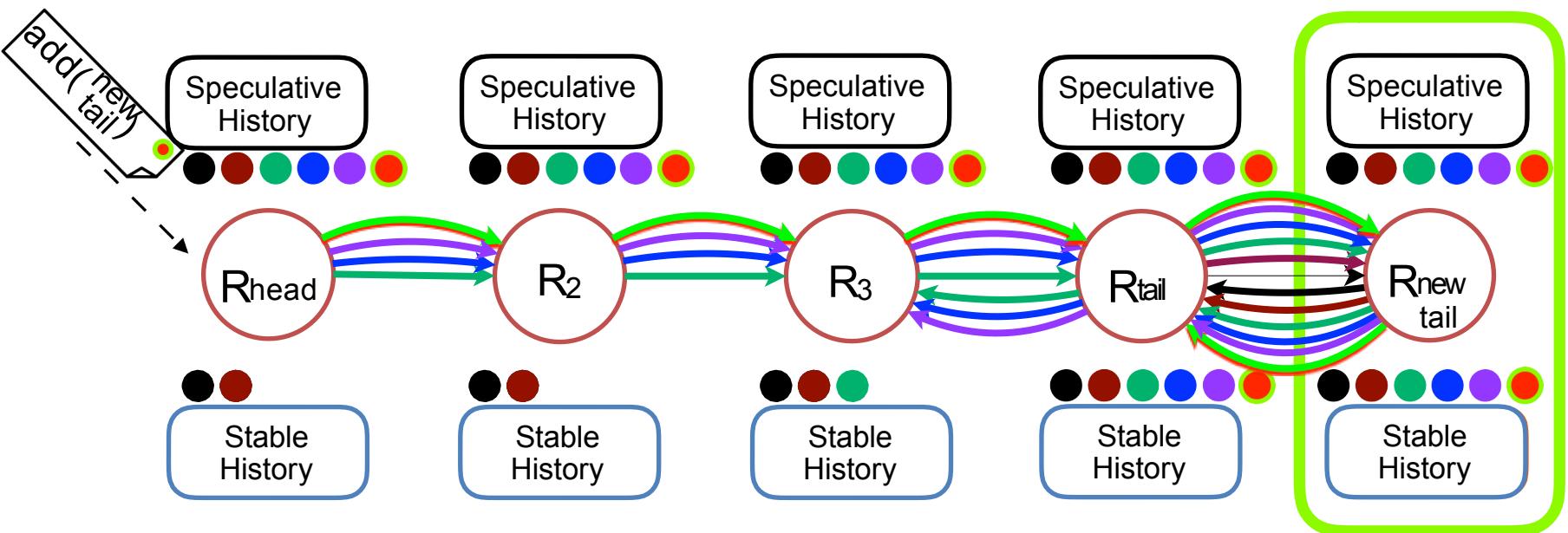
# Adding a New Node I

## Propagate Acknowledgements



# Adding a New Node II

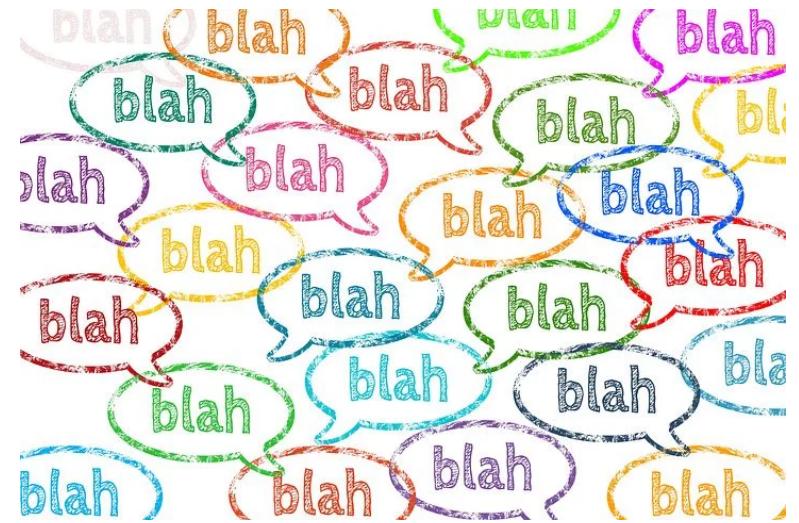
## Propagate Acknowledgements



# Self-study Questions

- Given a topology  $n$  replicas,  $n>2$ , how does replica  $i$  know about replica  $i+2$ , assuming replica  $i+1$  failed?
- Go through the motions of constructing a topology with `add(nodeId)` messages.
- When does a replica know for sure a node has been permanently added?

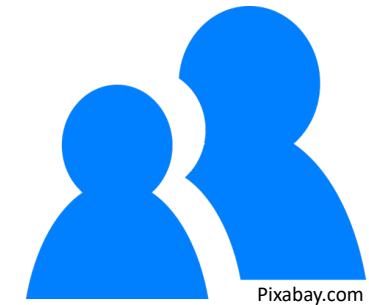




Pixabay.com

# GOSSIPING

# Gossiping protocols



- Disseminate information in **incremental manner**
  - Avoid overloading nodes with heavy broadcast messages
  - Drawback is longer propagation time for information
- Each node maintains a **partial view** of other nodes
- Each node chooses **random** nodes from its view to exchange information with
  - Application data (e.g., current state)
  - Its partial view (e.g., topology information)
- Nodes update their state and partial view based on the information received
- Gossiping happens **periodically** and **non-deterministically**
- Used in Cassandra for propagating status of each node, failure information and metadata

# Lazy Replication Using Gossiping

- Replicas gossip about operations processed
- Replicas **reconcile** (compare) their operation logs and each apply any operations not yet seen
- Former step is highly application dependent
- Assumes updates can be applied in any order
- If system processes no more operations. then each replica **eventually** converges to the same state by gossiping enough times



Pixabay.com

# Self-study Questions

- Draw out any topology of nodes and inject a message at a randomly chosen node, compare a broadcast (send to all neighbours versus a gossip (send to some neighbours):
  - How many messages are required?
  - How long does it take for all nodes to be up-to-date?
- Have each node in your topology maintain a data structure (e.g., counter, list, array, set, etc.), inject data structure updates at random nodes and propagate these updates via gossip:
  - What is the net result?
  - Does your replicated system converge at each replica to one and the same state (data structure)?

