

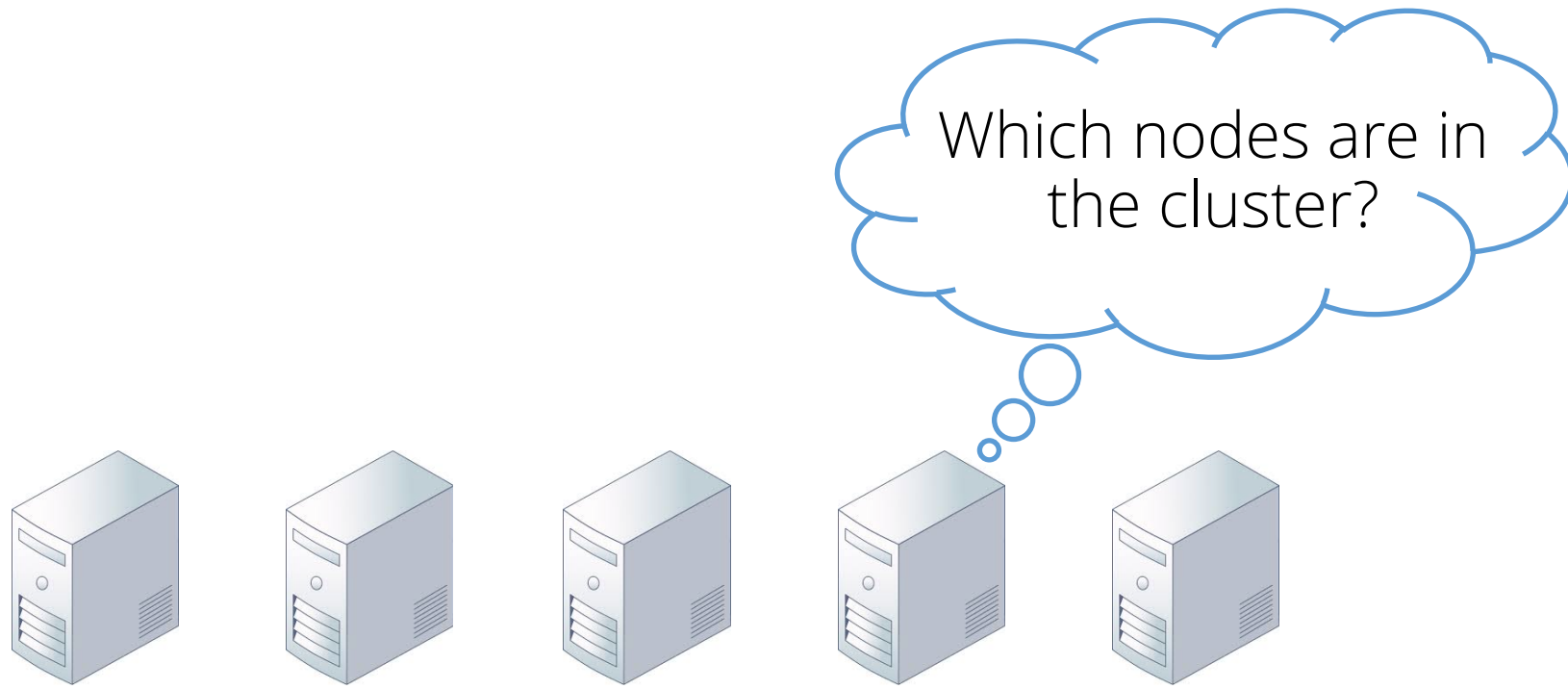
# Stable and consistent membership at scale with Rapid

Lalith Suresh, Dahlia Malkhi, Parikshit Gopalan  
Ivan Porto Carreiro<sup>1</sup>, Zeeshan Lokhandwala<sup>2</sup>

VMware Research

<sup>1</sup>VMware

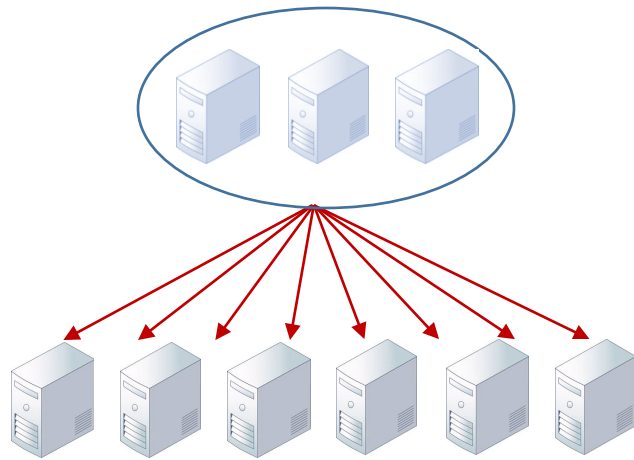
<sup>2</sup>One Concern



# Membership management and failure detection

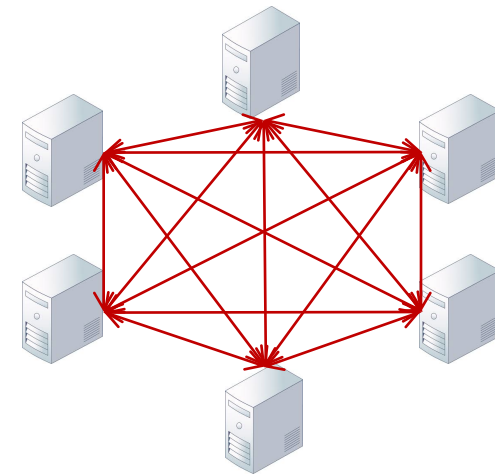
# Types of membership services

Centralized



{ Spark, HDFS,  
Zookeeper-based systems }

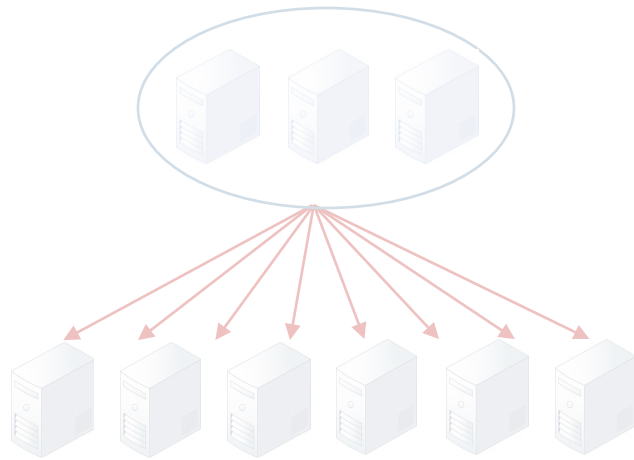
Gossip-based



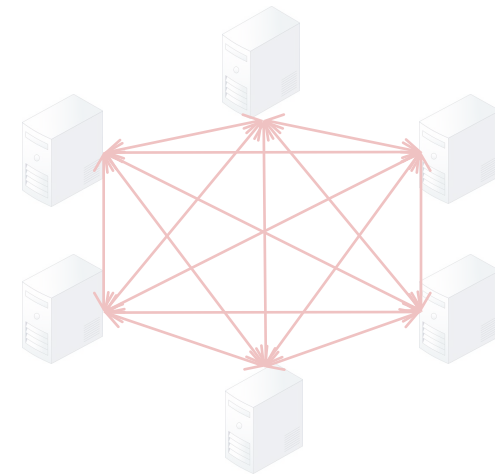
{ Cassandra, Akka,  
Redis Cluster,  
Dynamo }

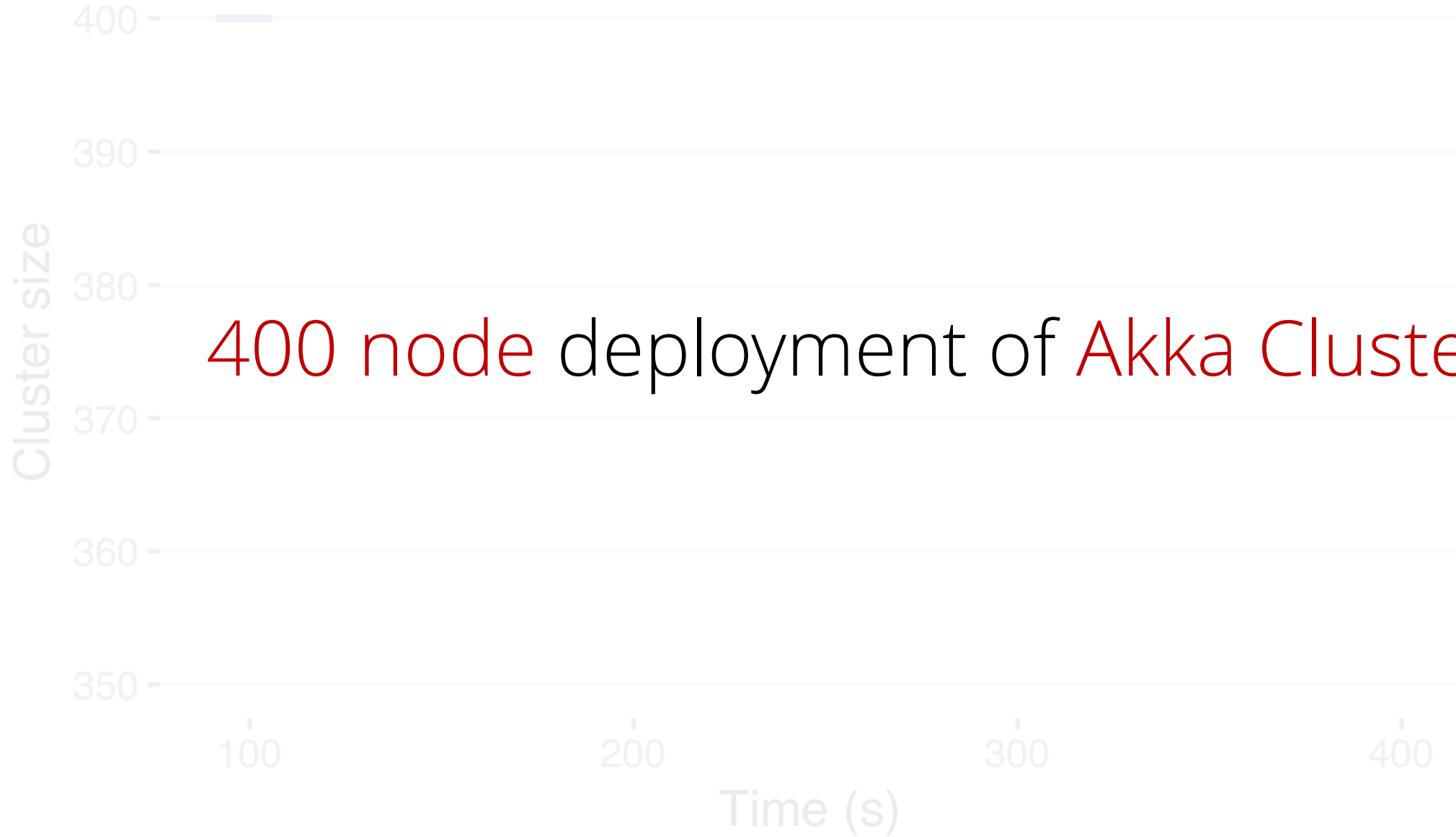
Existing solutions do not provide  
**stability** and **consistency** at **scale**

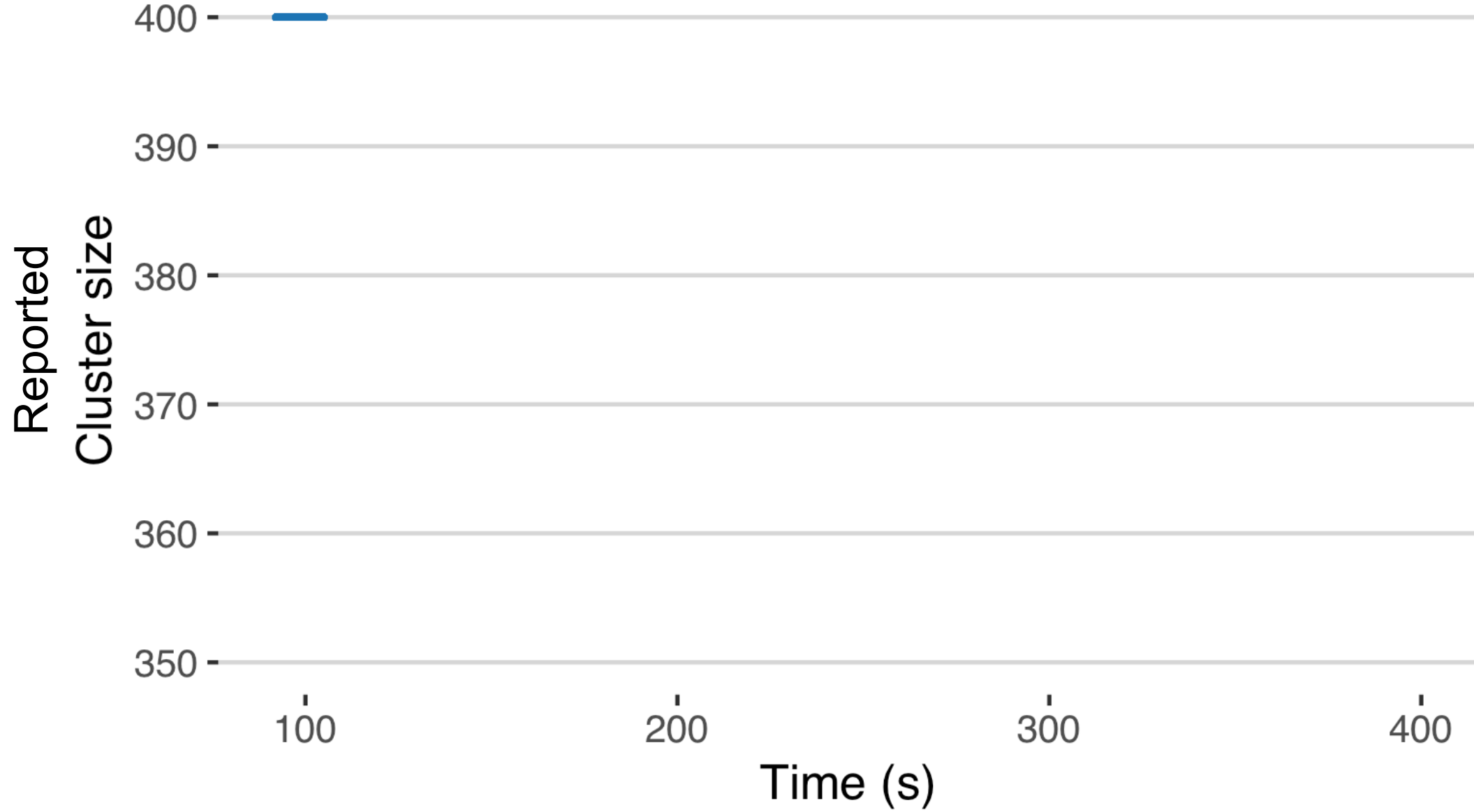
Centralized

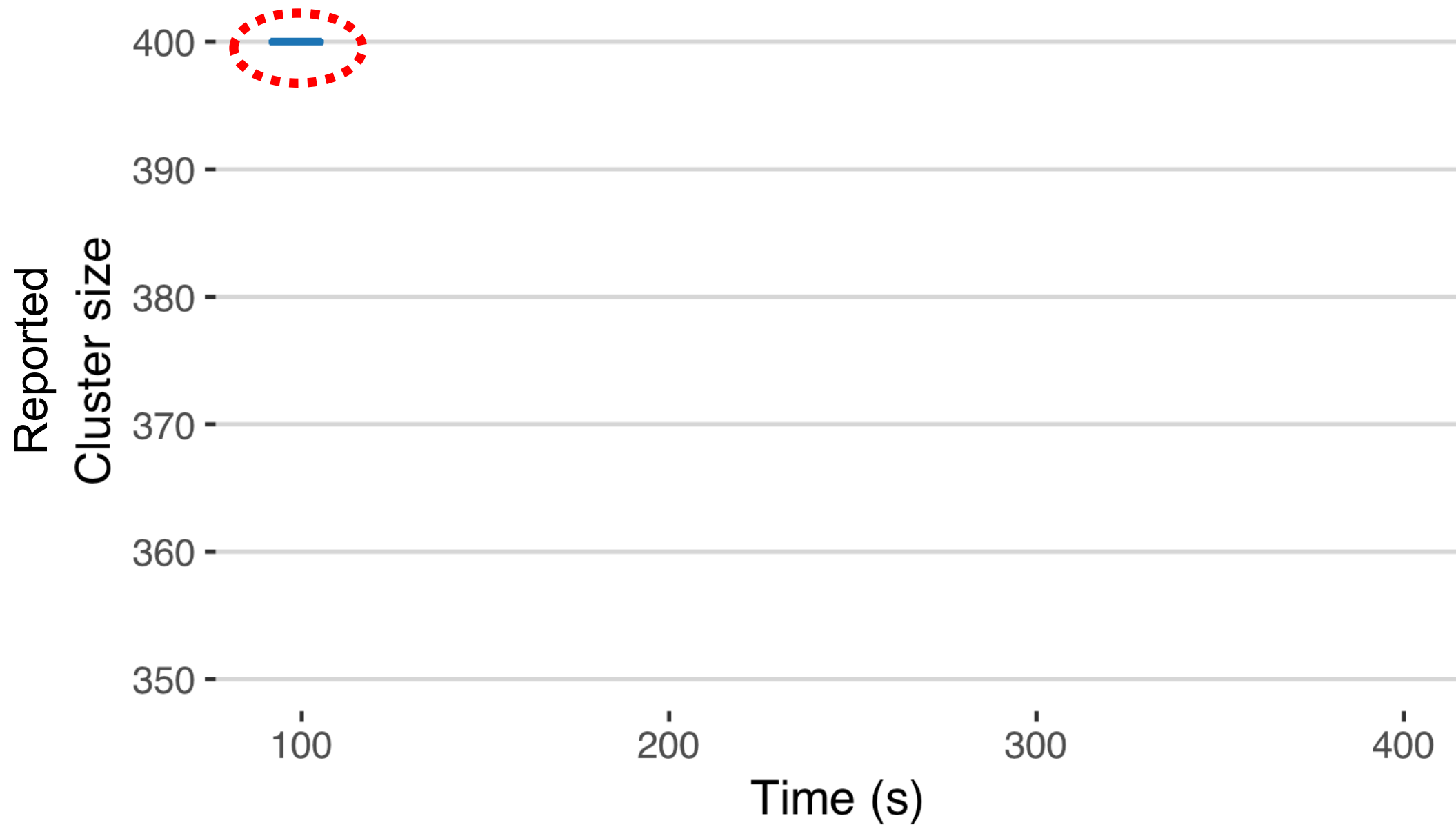


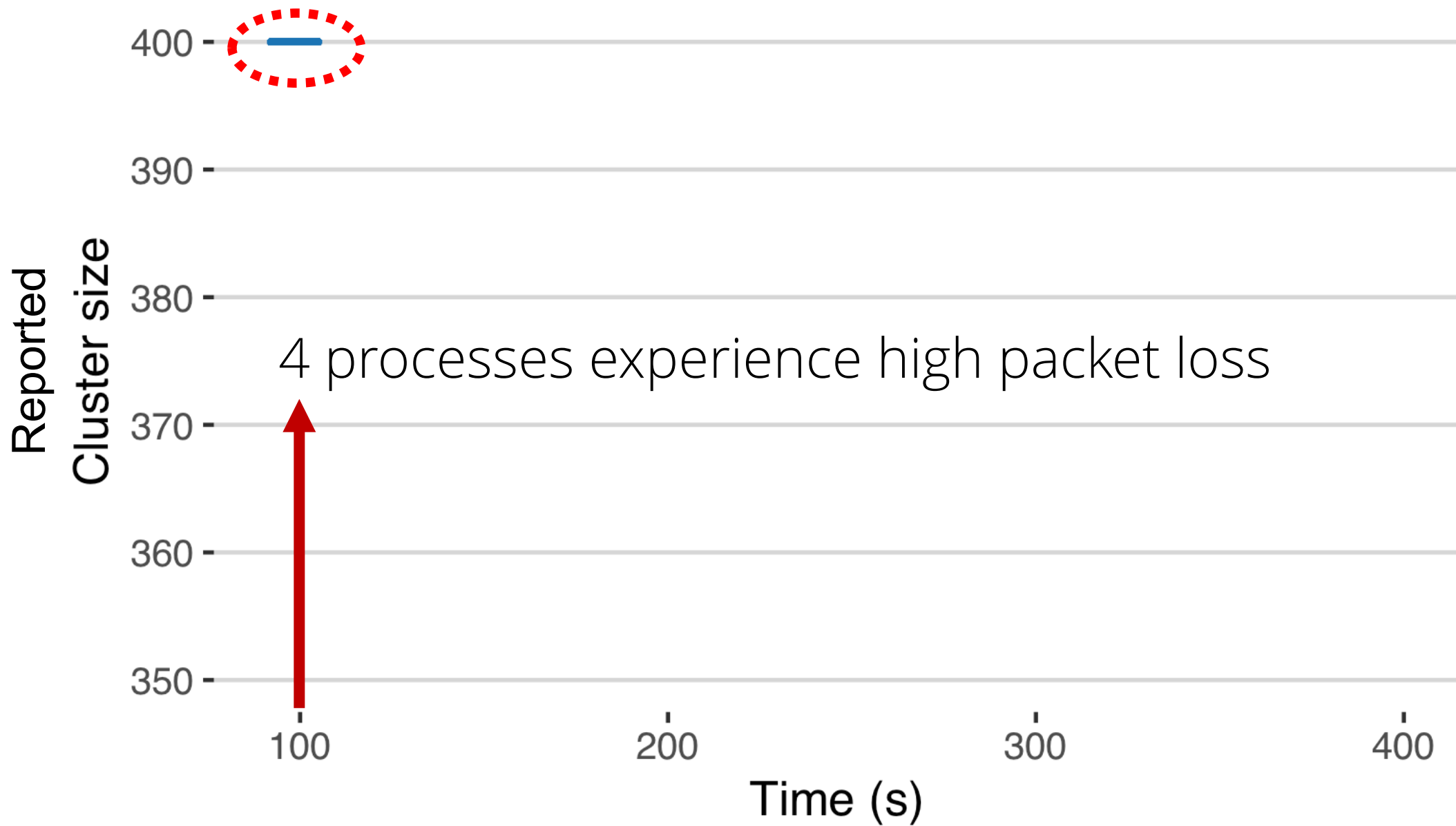
Gossip-based



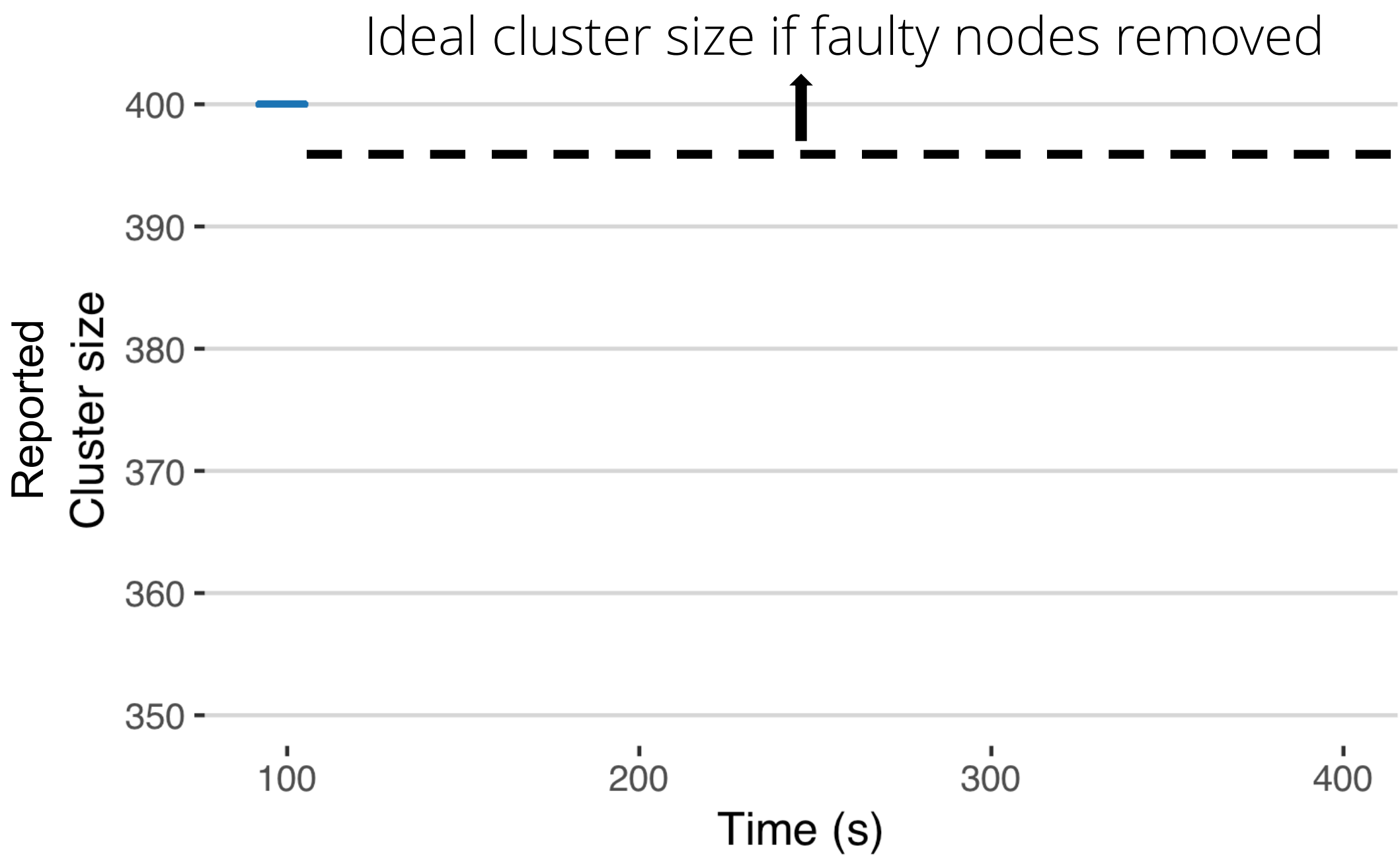


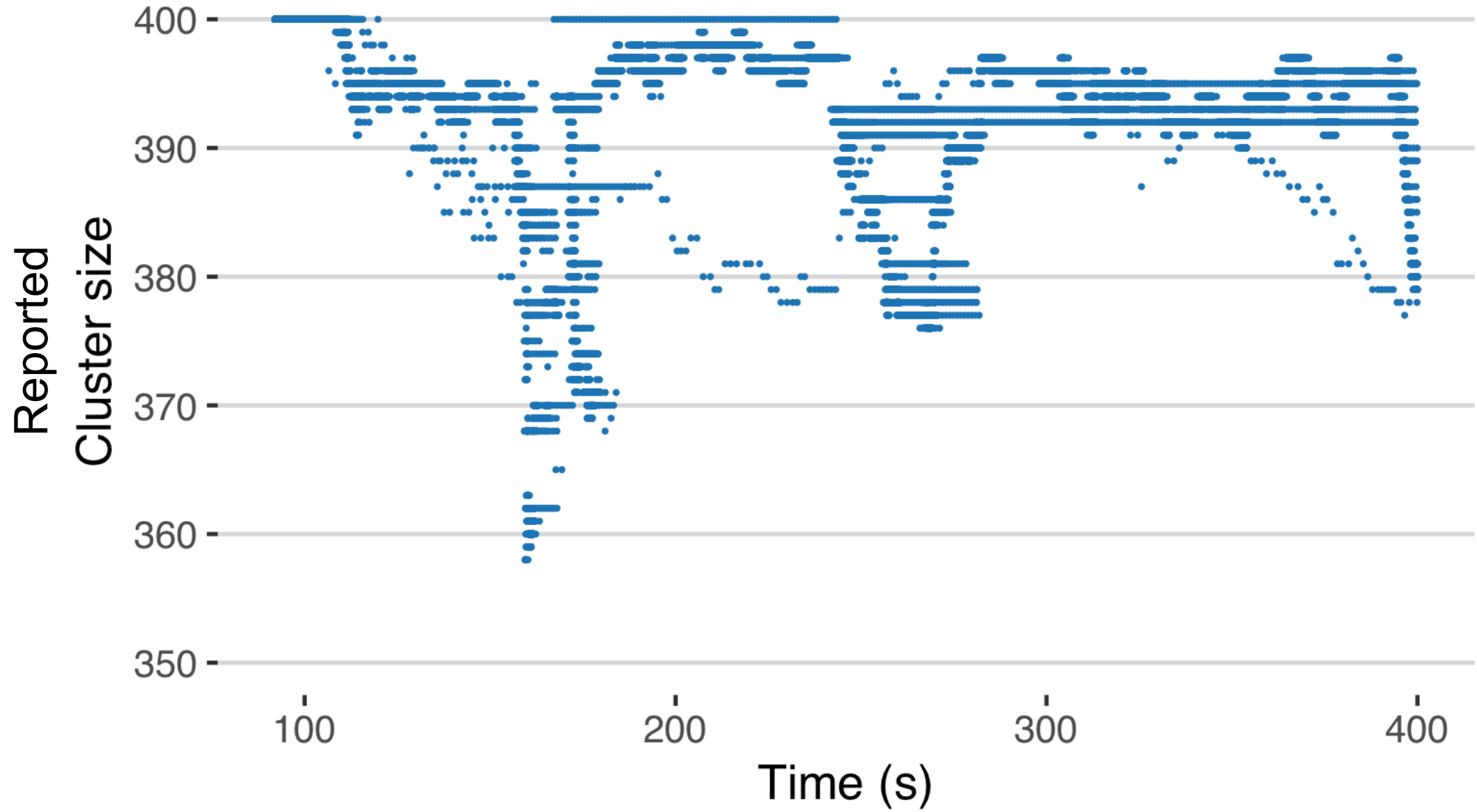


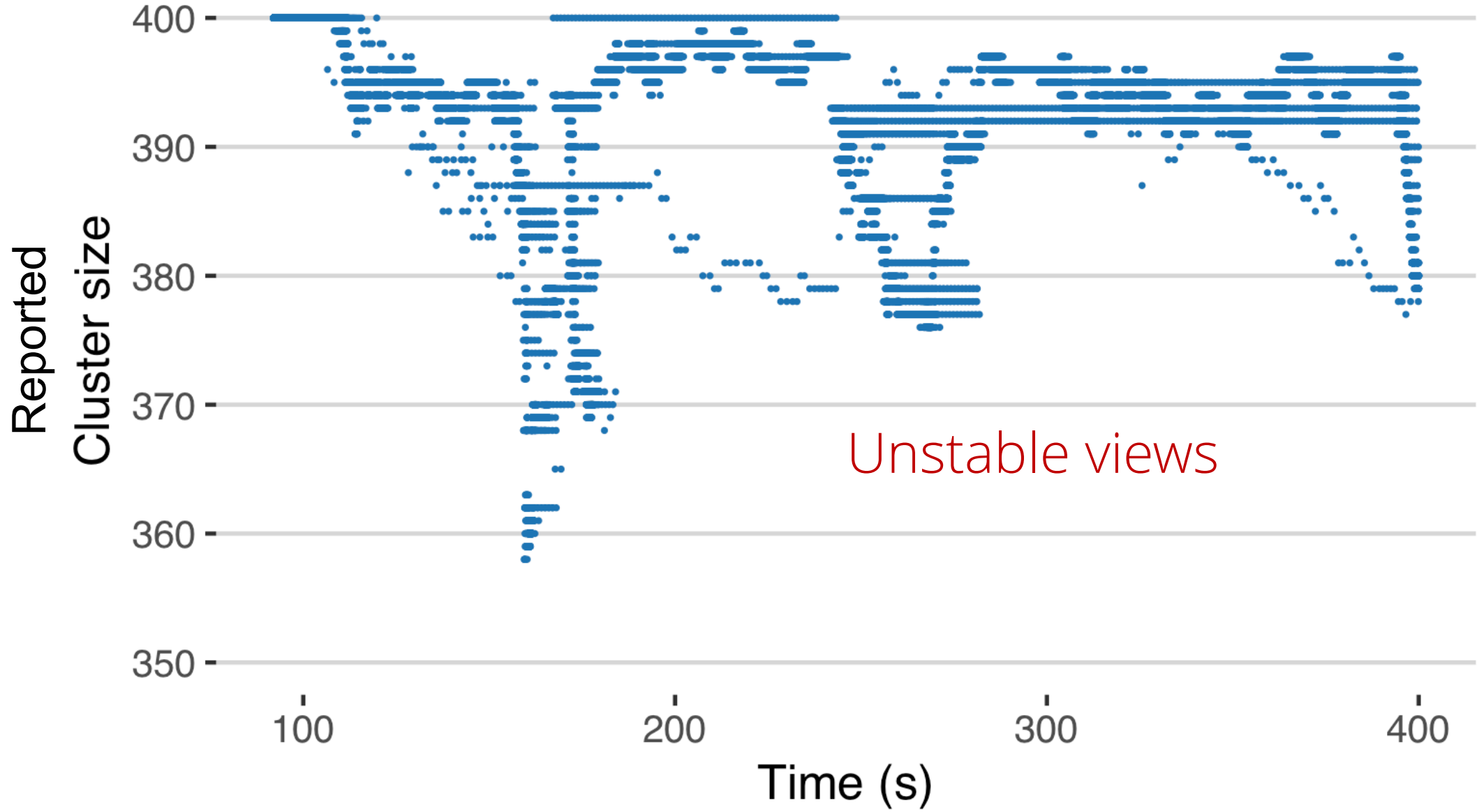


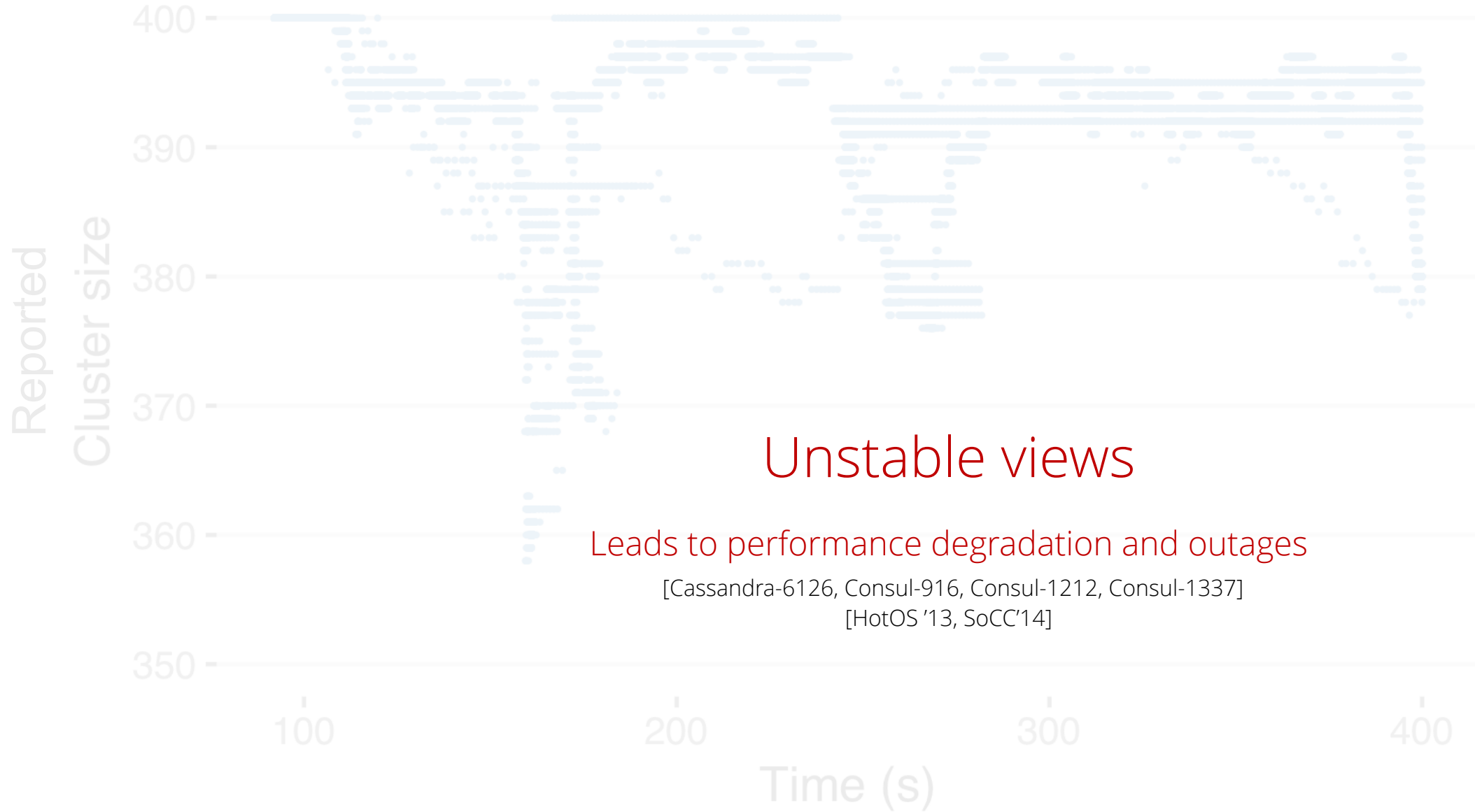


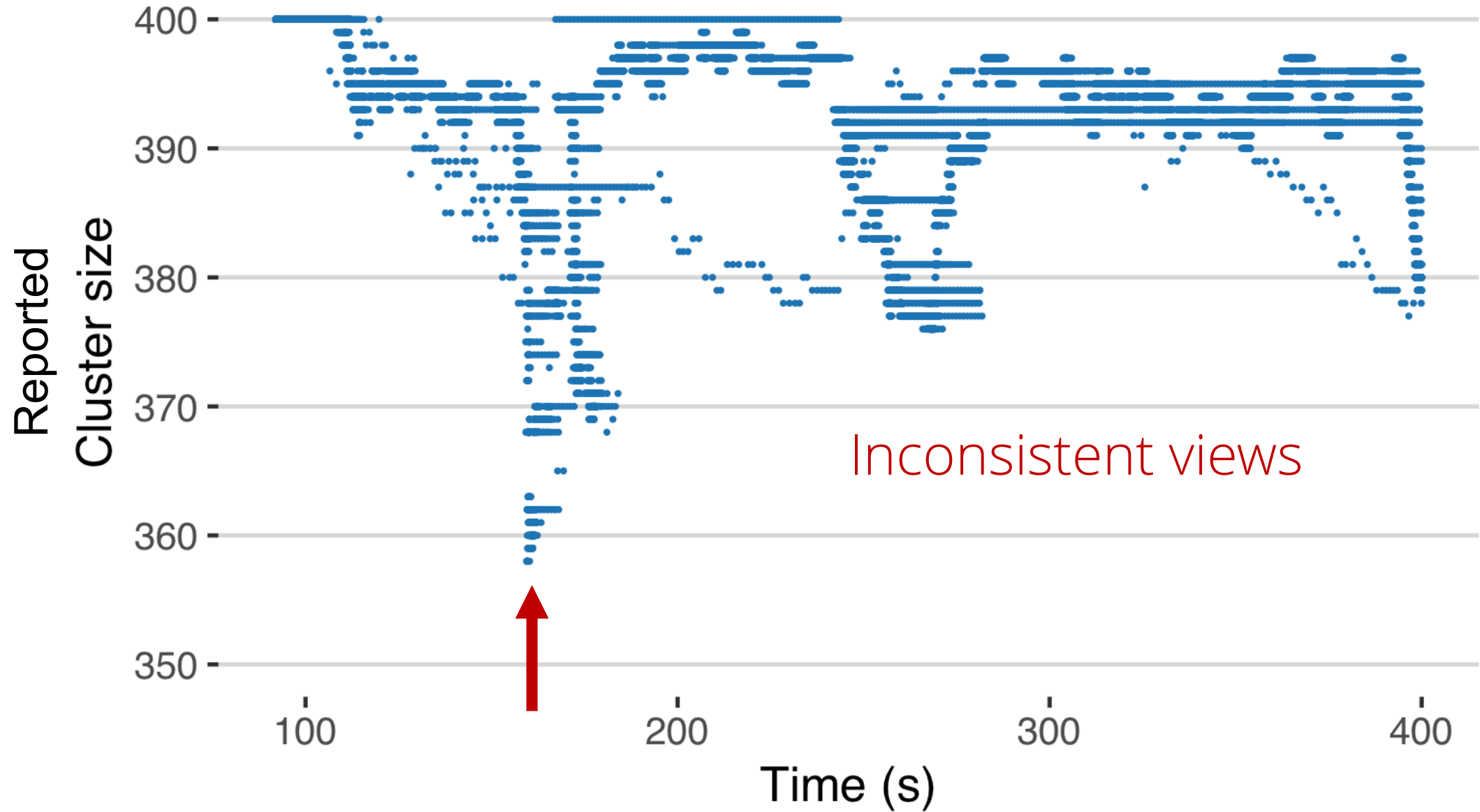


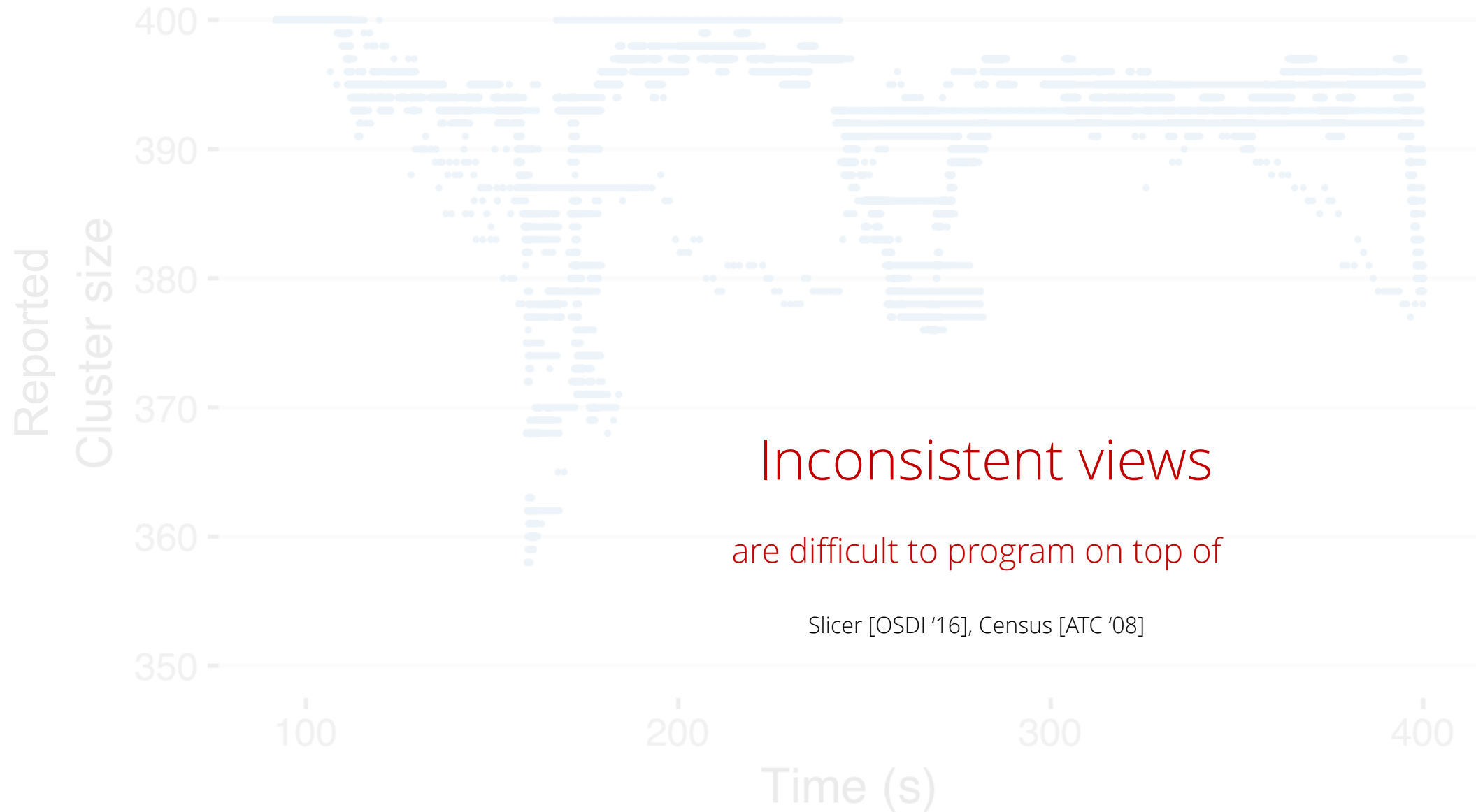












# Rapid

Stable and consistent membership at scale

# Rapid

Stable and consistent membership at scale



Robust against asymmetric network failures, flip-flops, packet loss etc.



# Rapid

Stable and consistent membership at scale



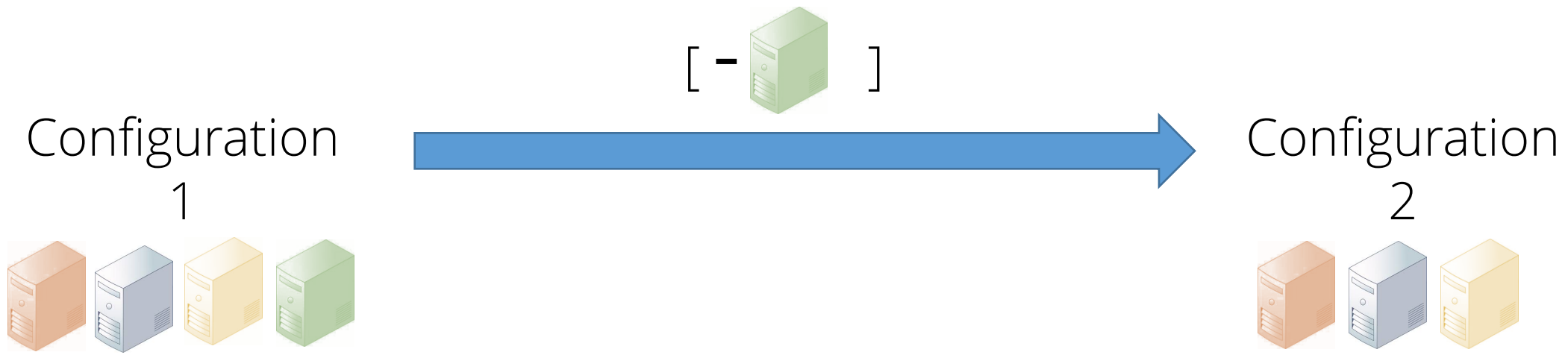
Processes see the same  
sequence of membership  
changes

# Rapid

Stable and consistent membership at scale



Bootstraps 2000 nodes 2-5x  
faster than Zookeeper and  
Memberlist



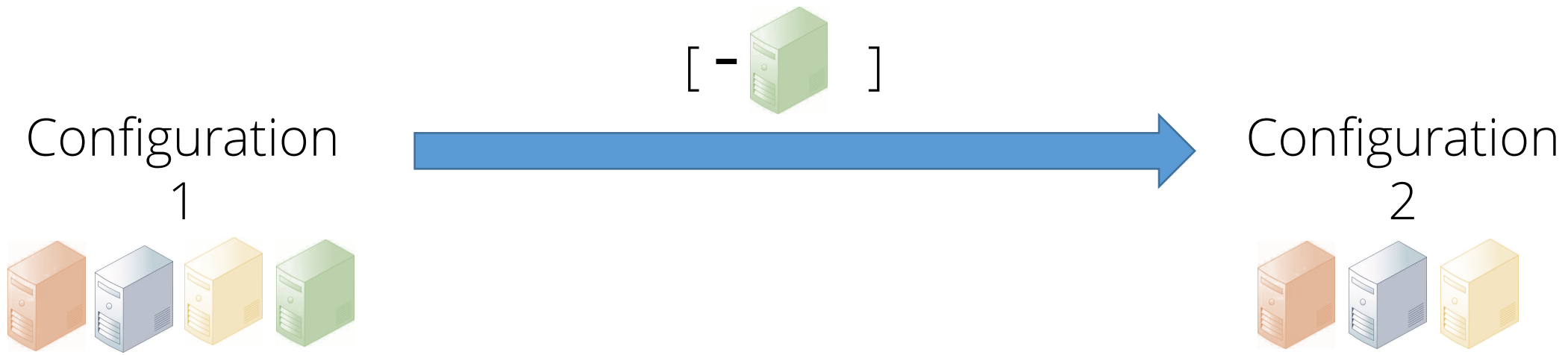
Expander graph-based monitoring

Multi-process cut detection

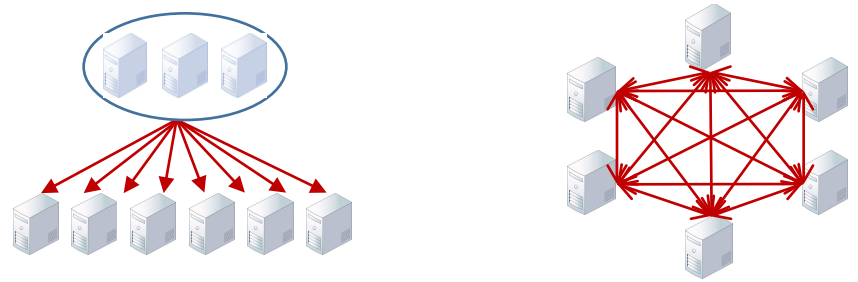
Fast path to consensus

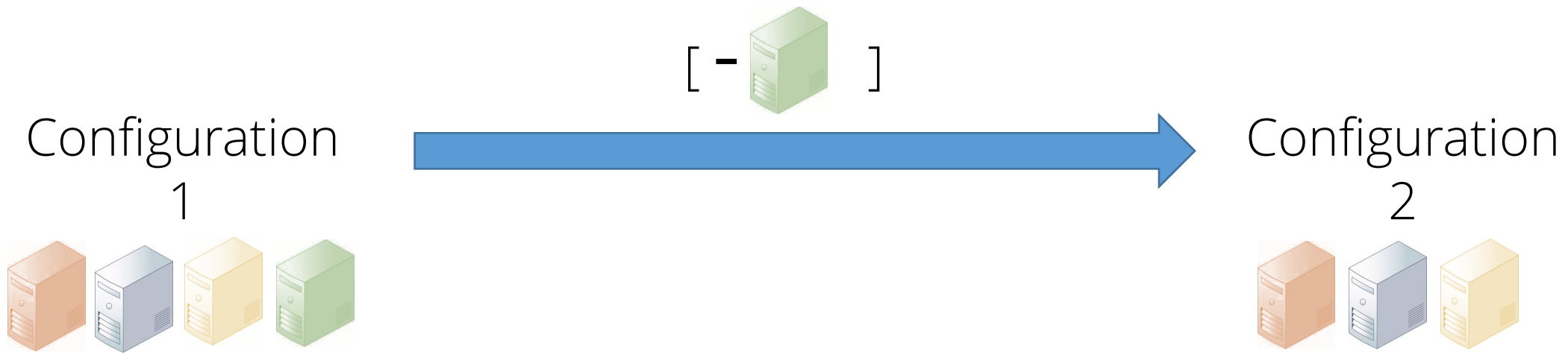
Stability

Consistency

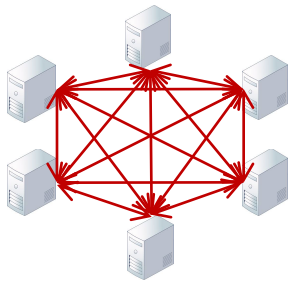


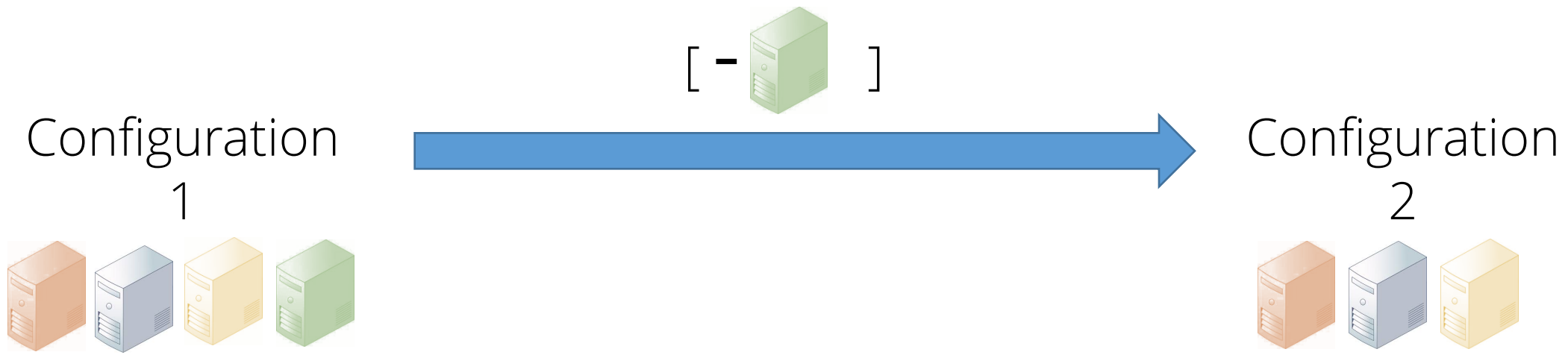
Rapid runs in both centralized and decentralized configurations





**This Talk:** decentralized design and failures





Monitoring Overlay

Expander graph-based monitoring

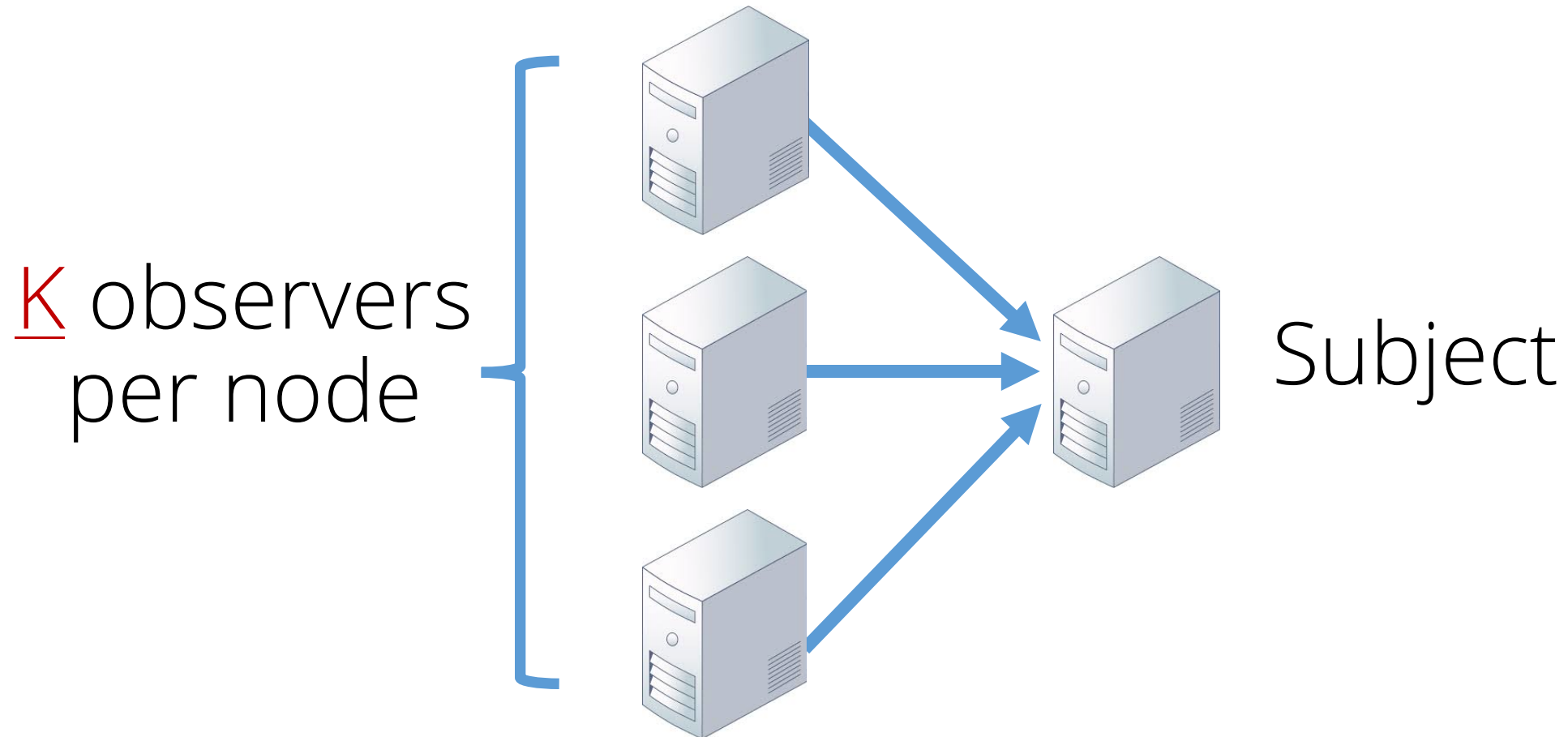
Membership change proposal

Multi-process cut detection

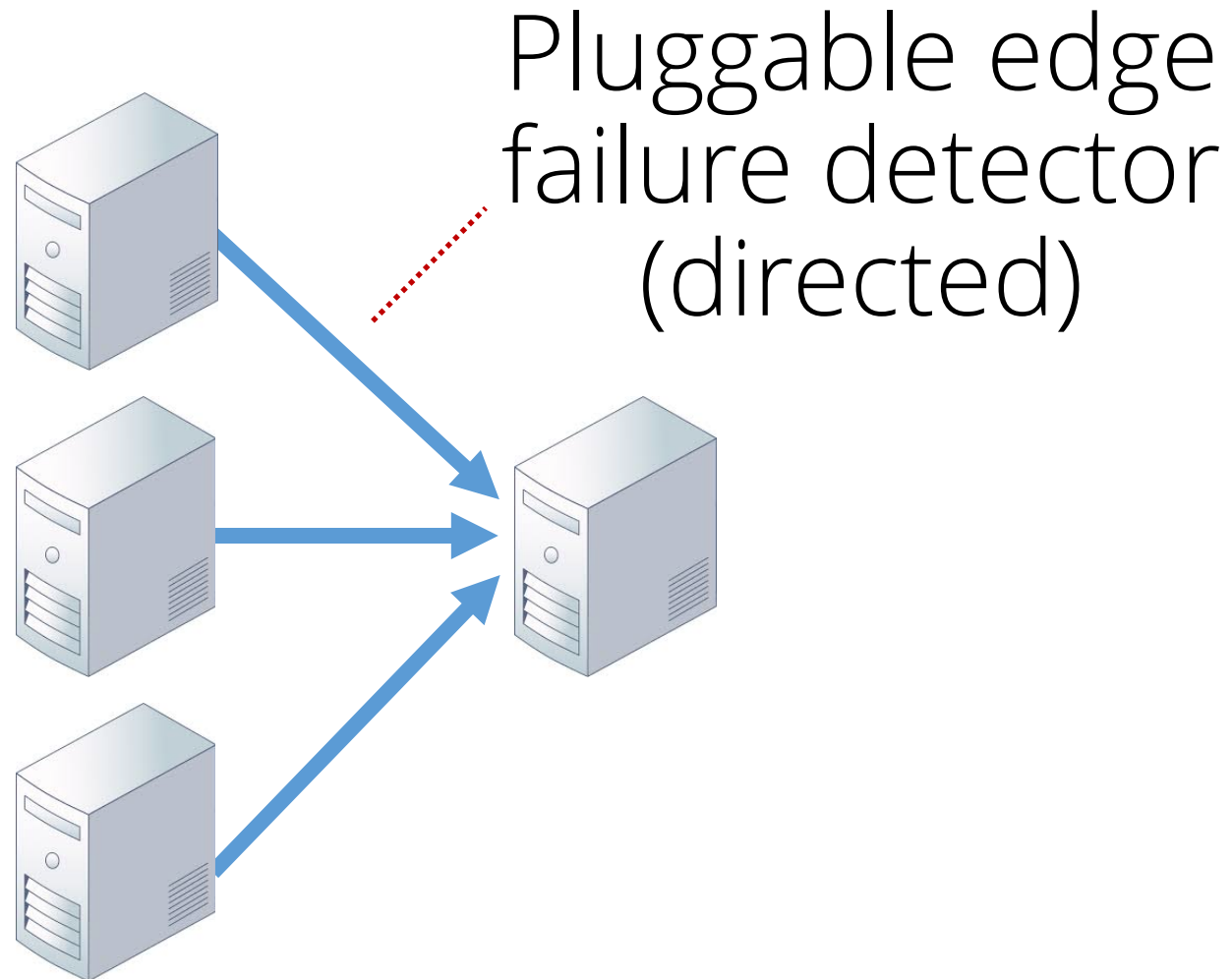
View change with consensus

Fast path to consensus

# Expander-graph based monitoring

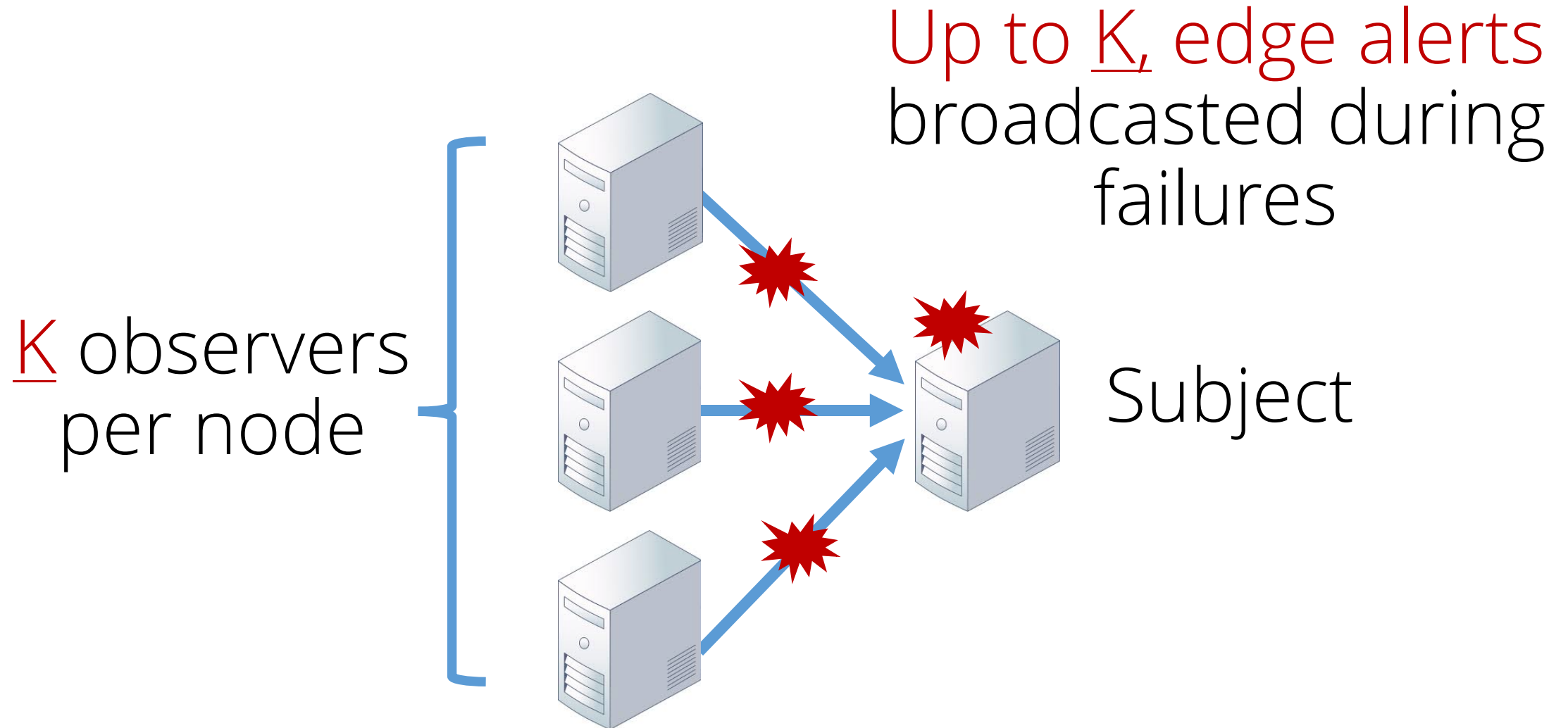


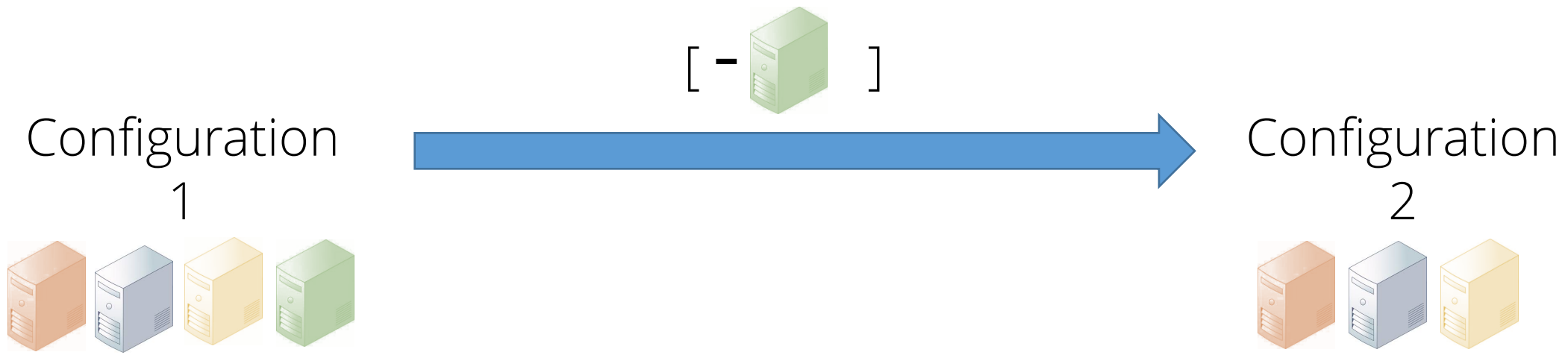
# Expander-graph based monitoring





# Expander-graph based monitoring





Monitoring Overlay

Expander graph-based monitoring



Membership change proposal

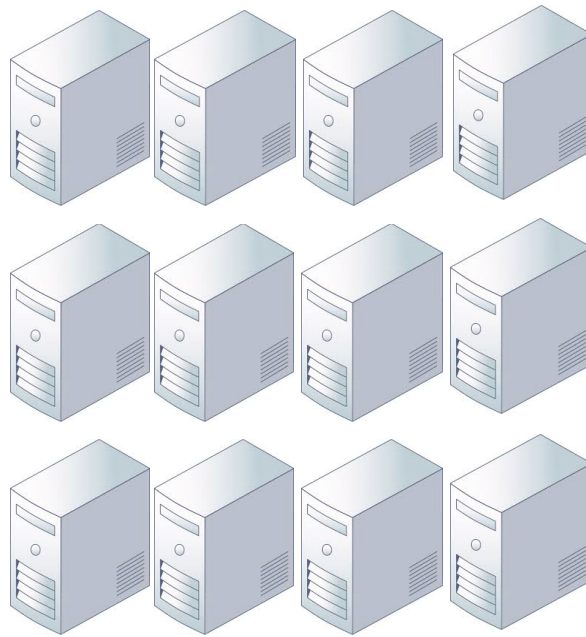
Multi-process cut detection



View change with consensus

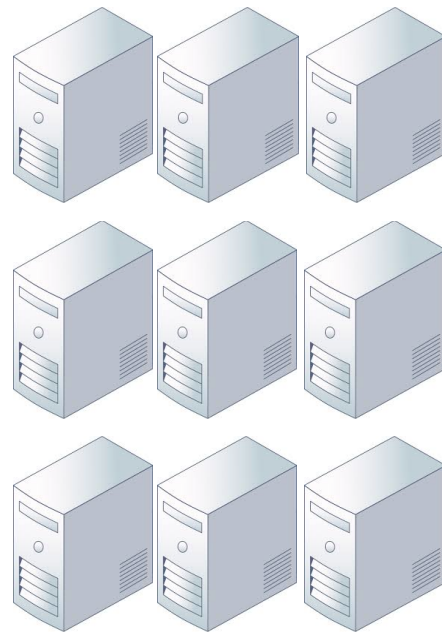
Fast path to consensus

# Multi-process cut detection

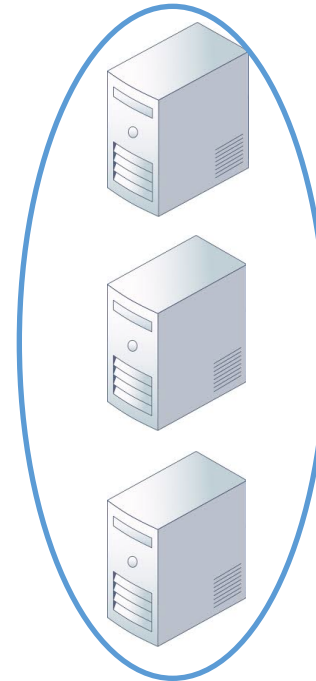


Expander-based  
monitoring overlay

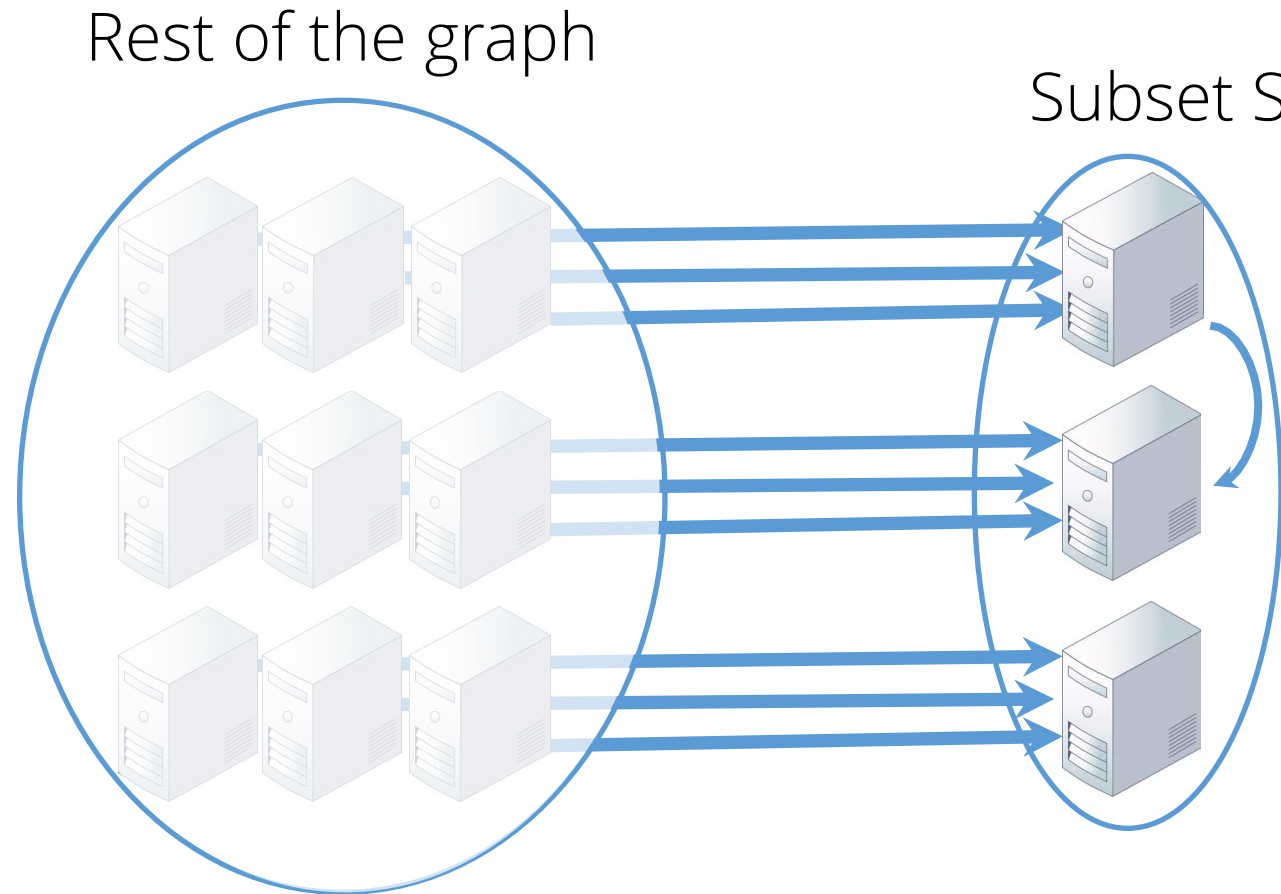
# Multi-process cut detection



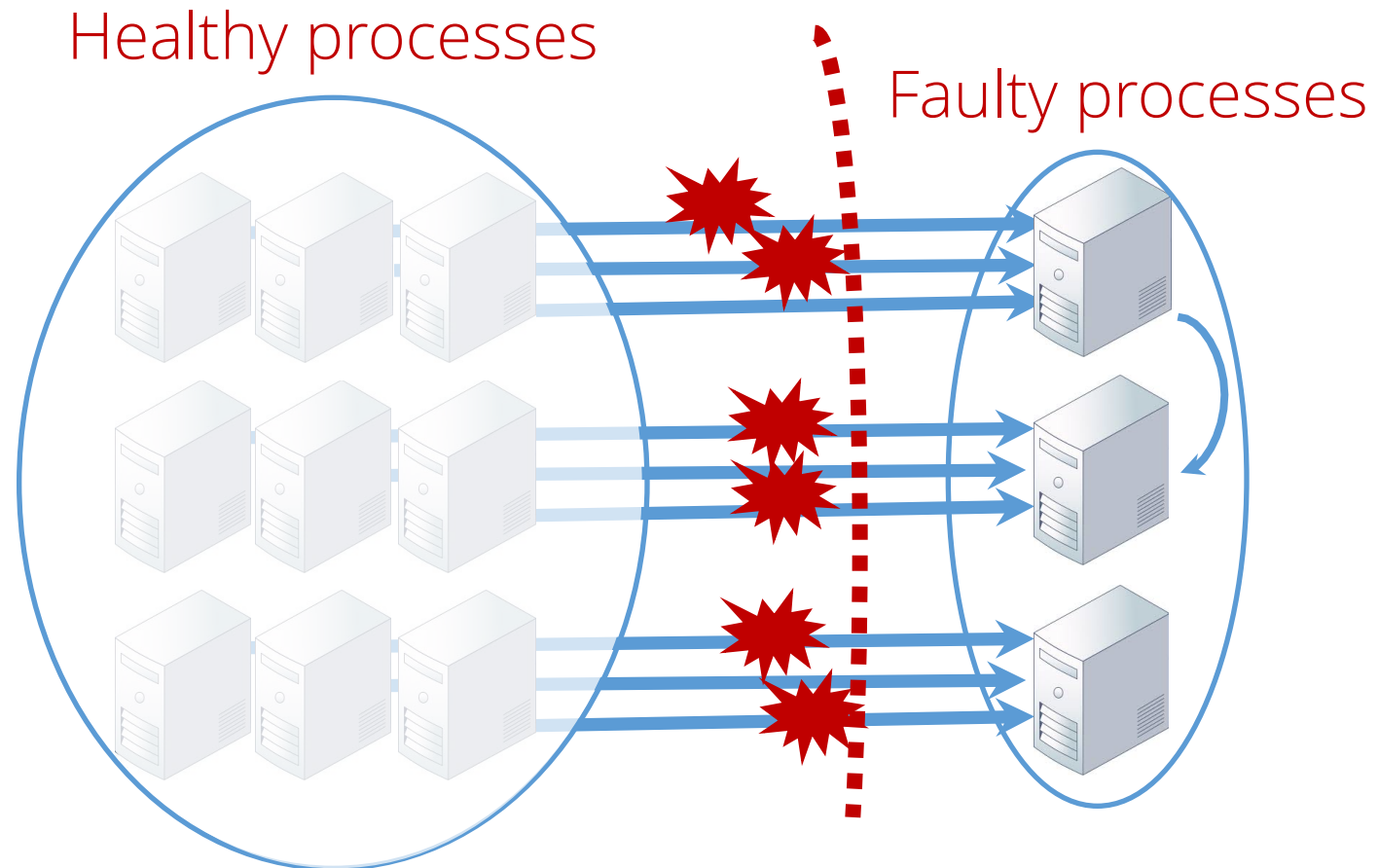
Subset S



# Multi-process cut detection



# Multi-process cut detection

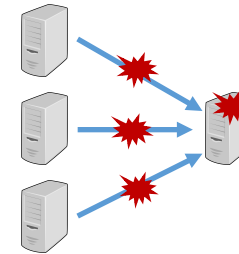


# Multi-process cut detection

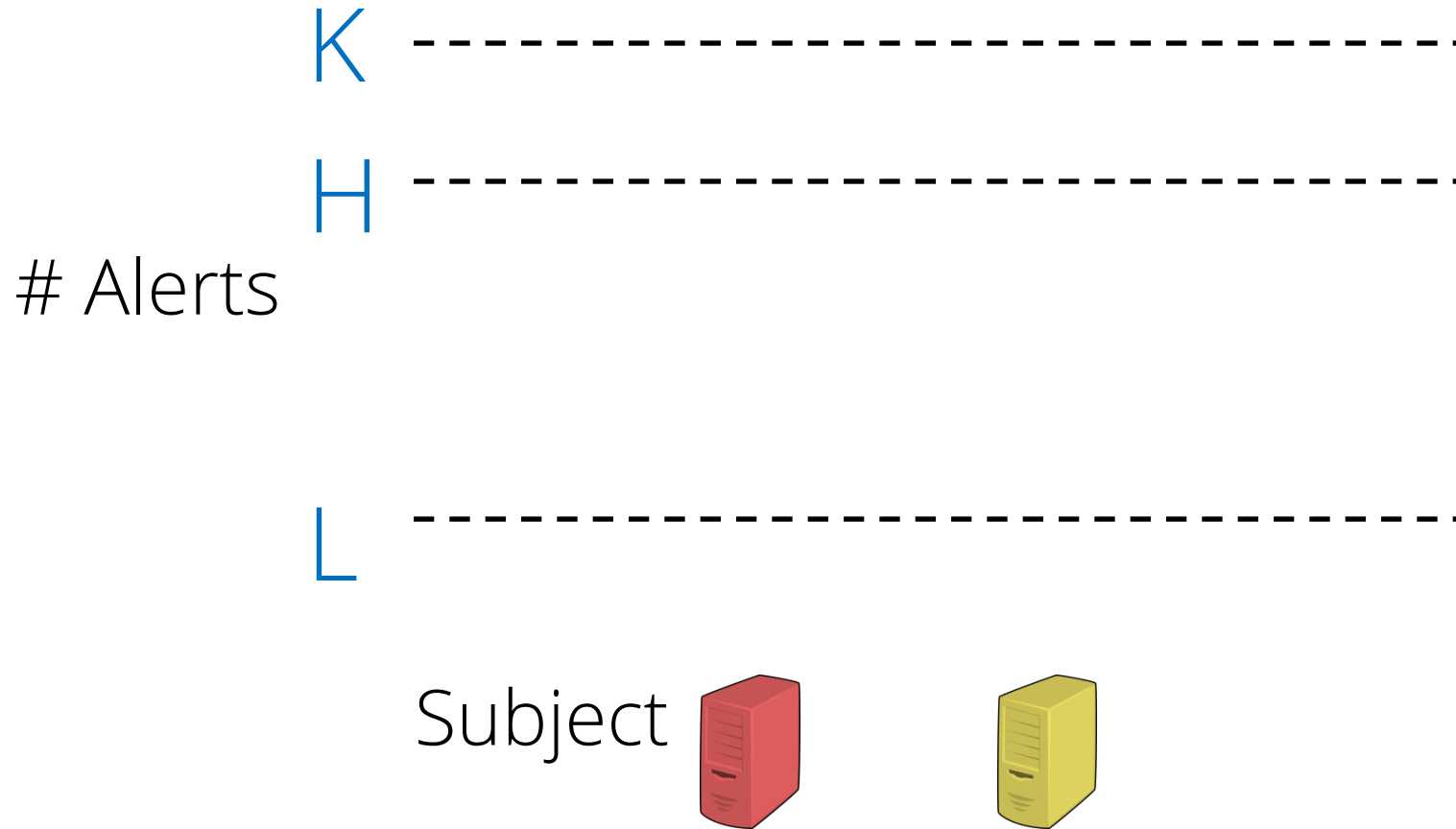
K -----

H -----

L -----

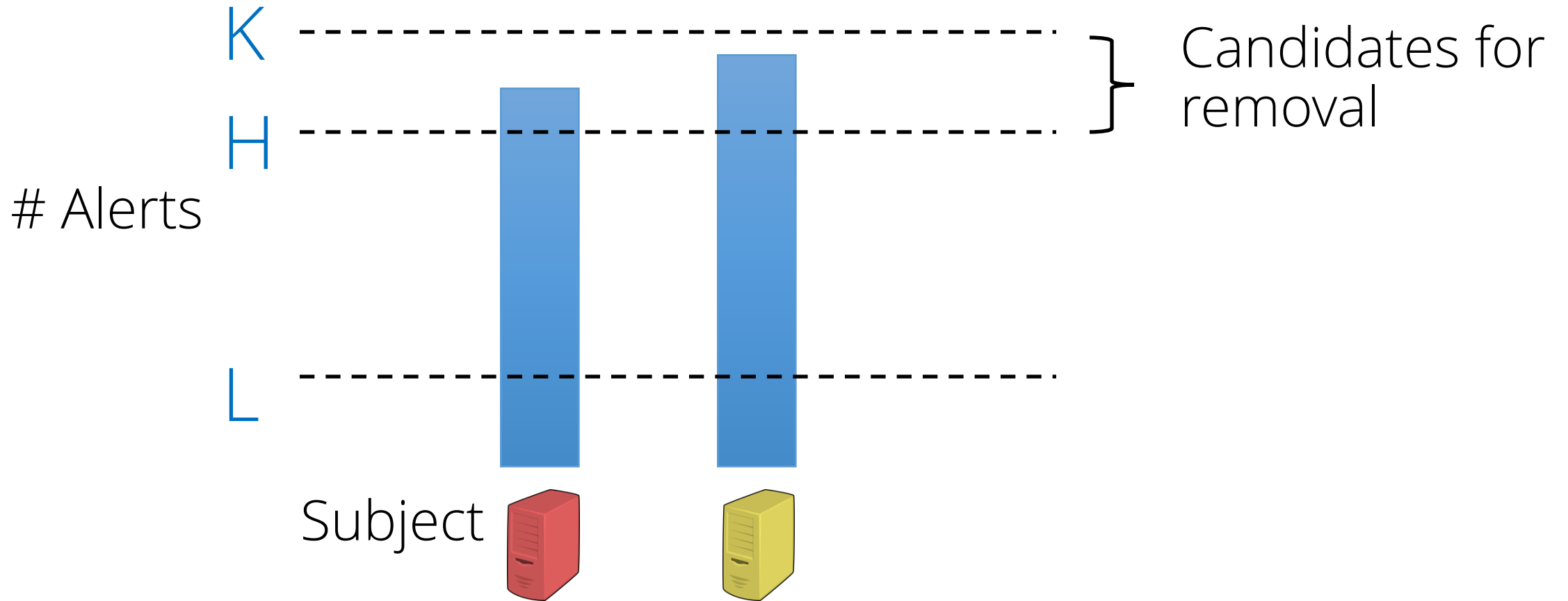


# Multi-process cut detection

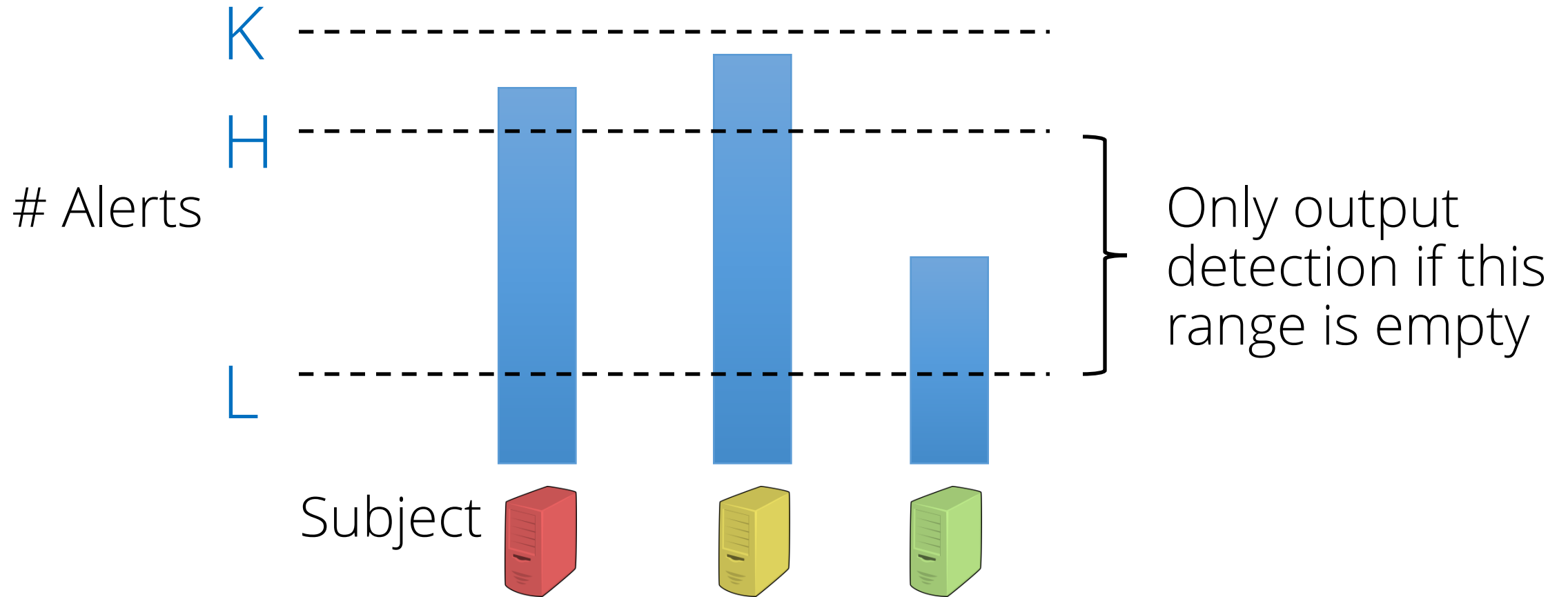




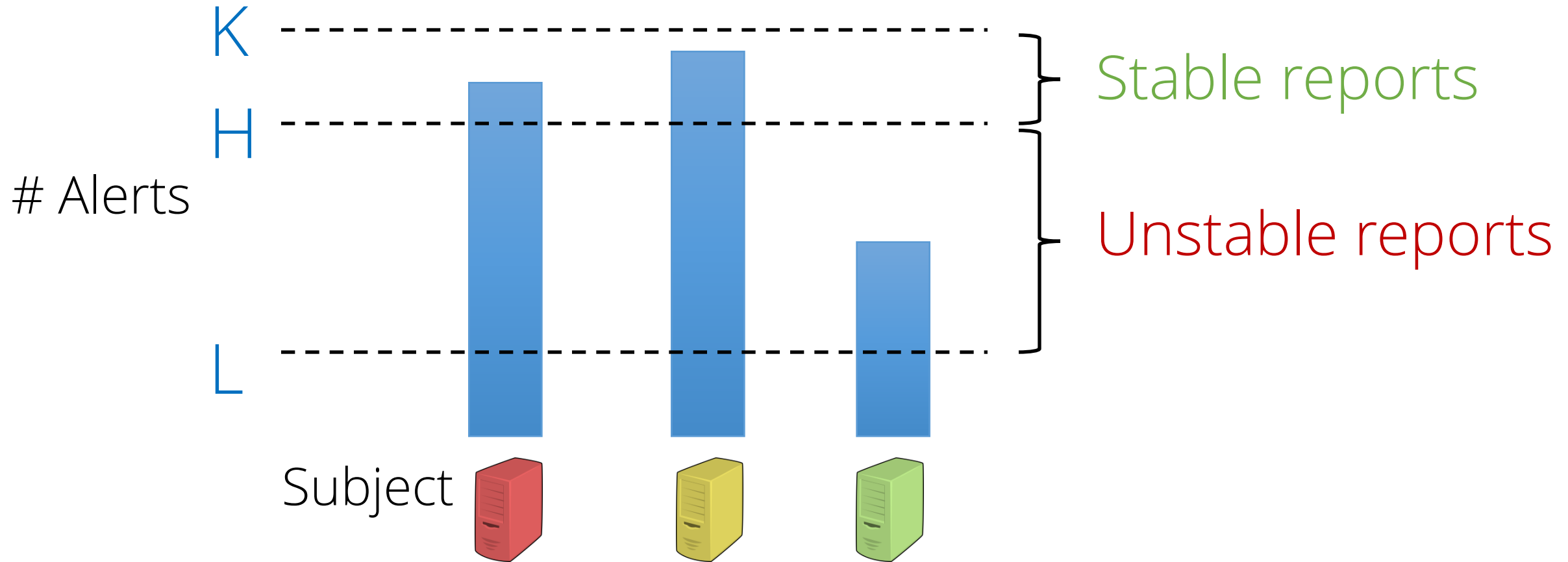
# Multi-process cut detection



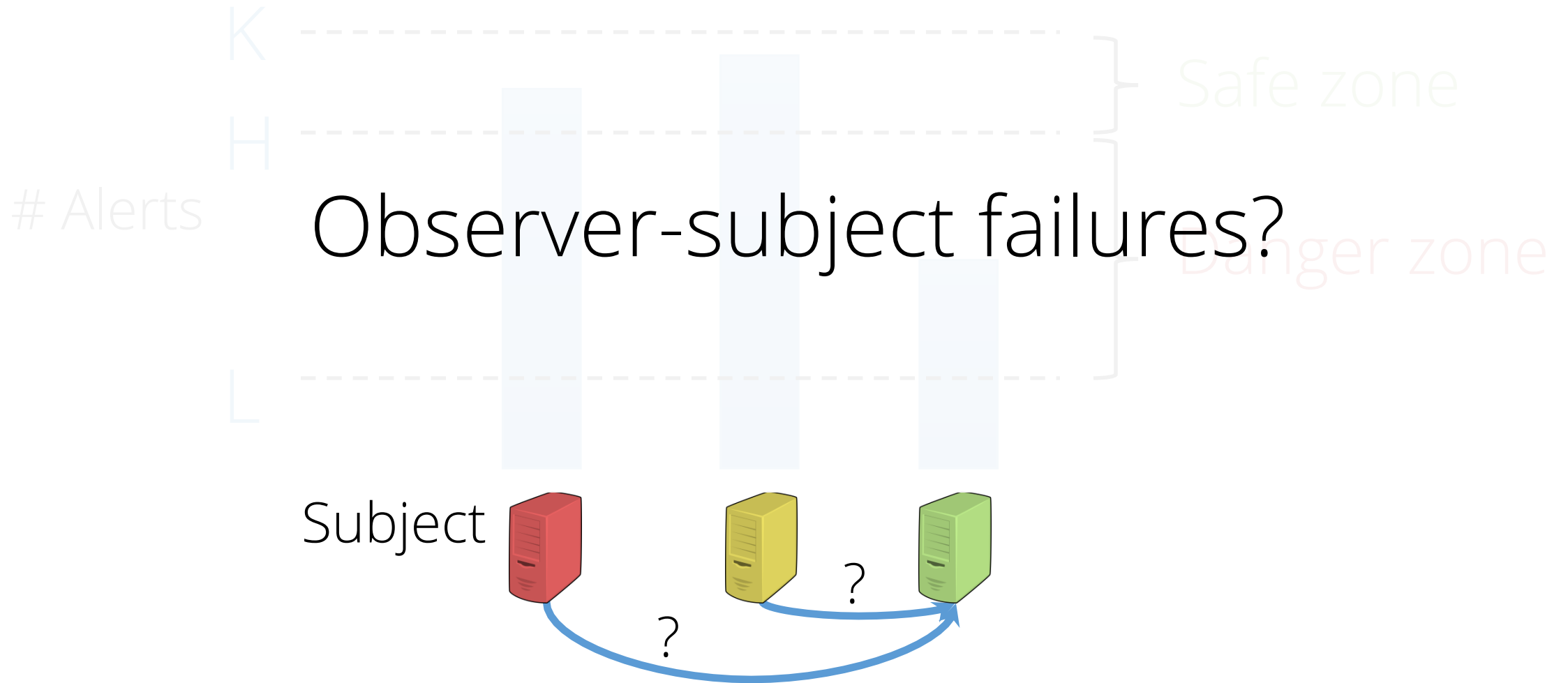
# Multi-process cut detection



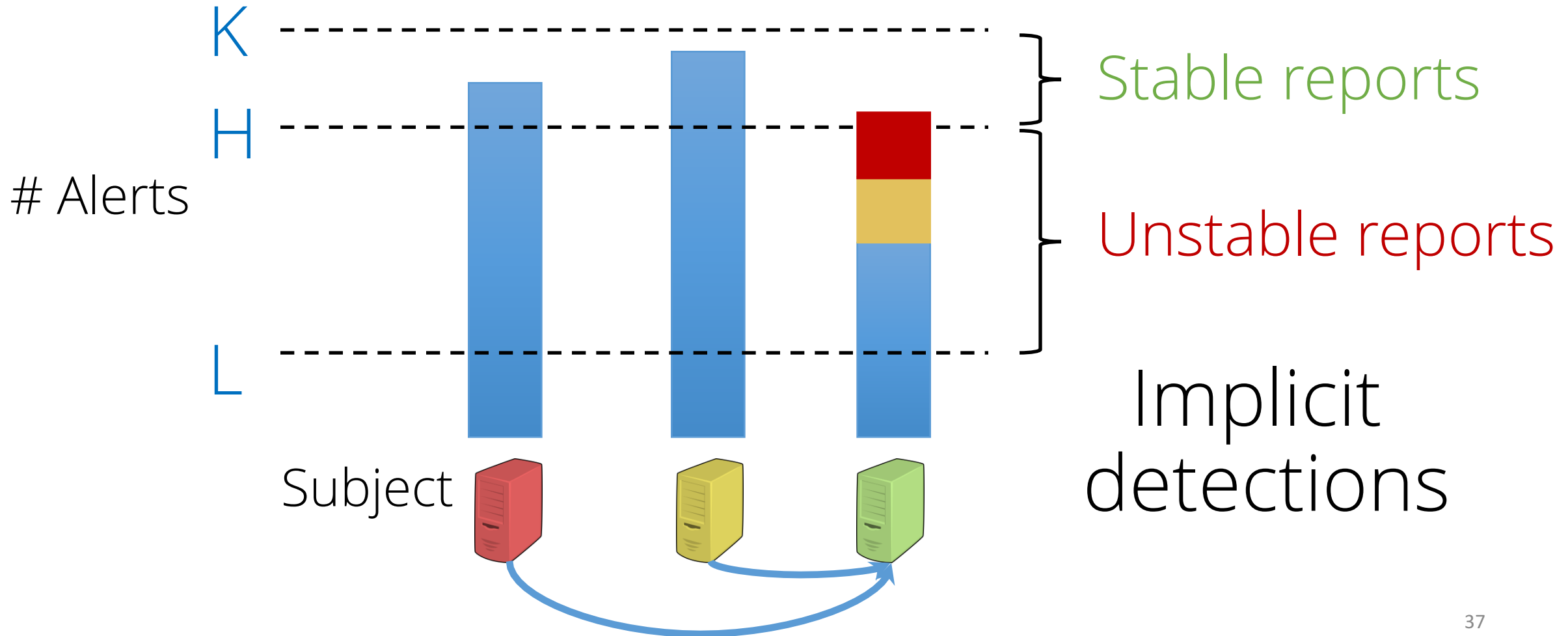
# Multi-process cut detection



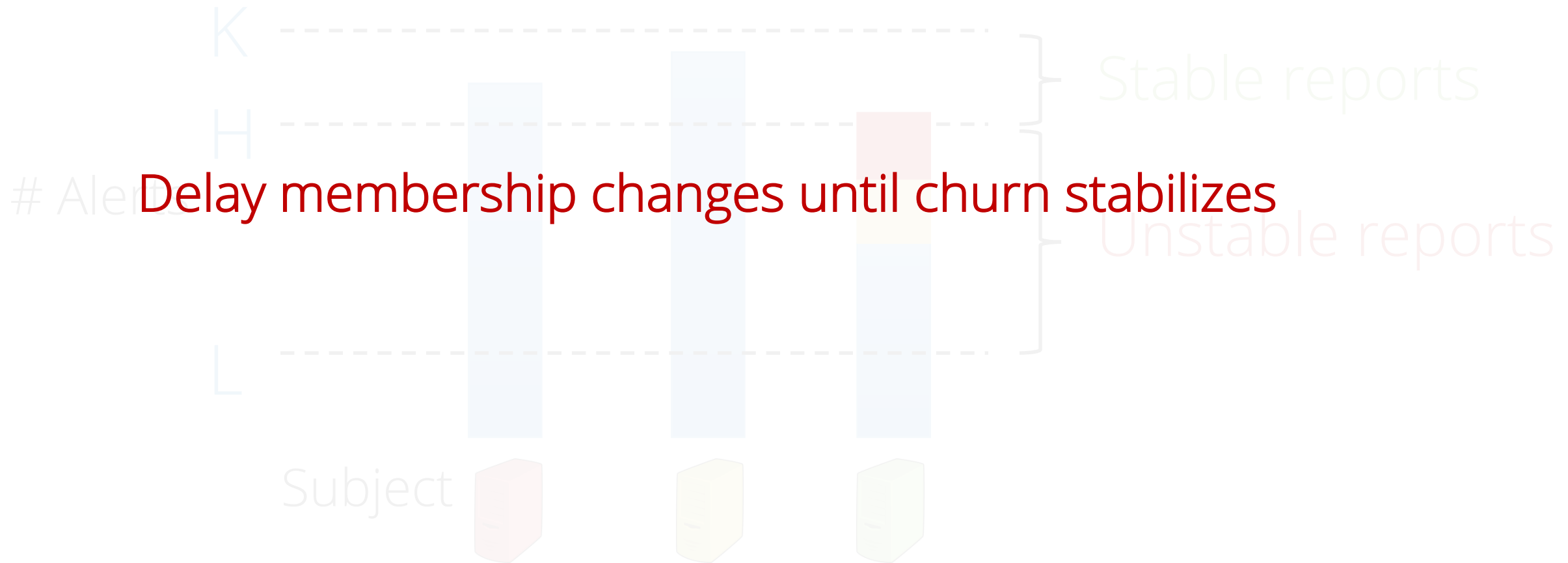
# Multi-process cut detection



# Multi-process cut detection

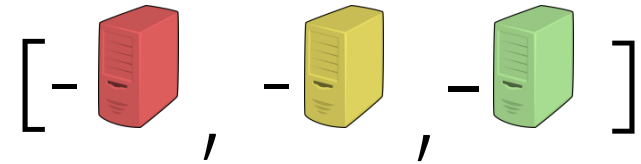


# Multi-process cut detection



# Almost-everywhere agreement

All processes output the same cut



with high probability

# Almost-everywhere agreement

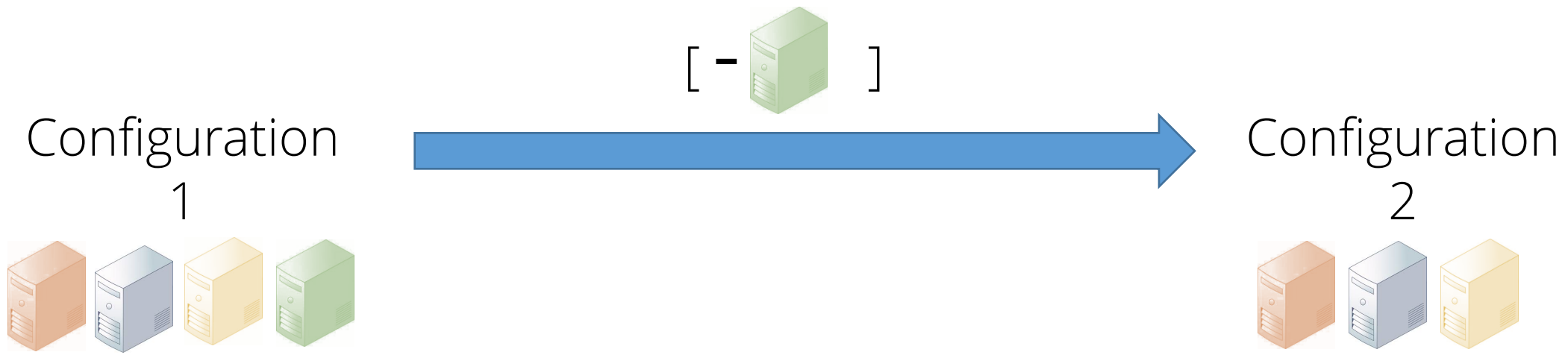
1000 processes, 8 failures,  $K=10$

$H=9$   
|  
0%  
conflicts  
|  
 $L=1$

$H=9$   
|  
0%  
conflicts  
|  
 $L=2$

$H=9$   
|  
0.06%  
conflicts  
|  
 $L=4$





Monitoring Overlay

Expander graph-based monitoring



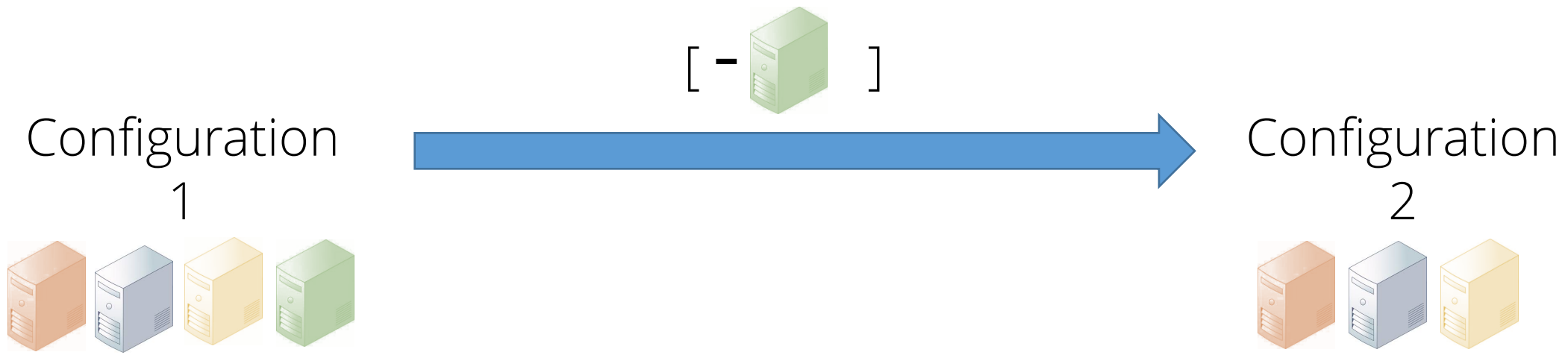
Membership change proposal

Multi-process cut detection



View change with consensus

Fast path to consensus



Monitoring Overlay

Expander graph-based monitoring



Membership change proposal

Multi-process cut detection



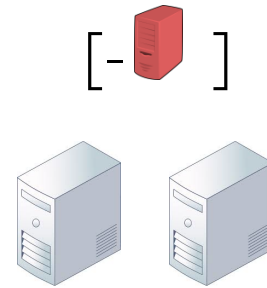
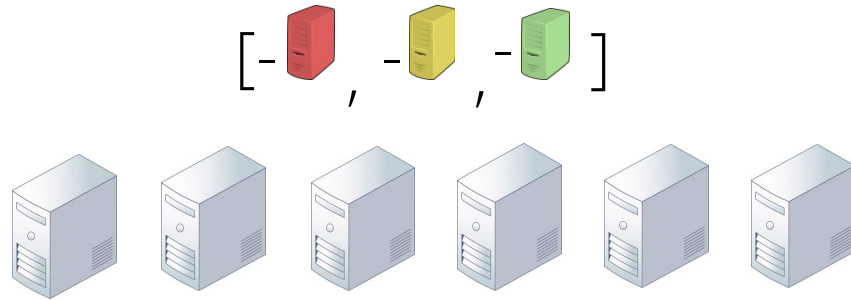
View change with consensus

Fast path to consensus

Almost-everywhere  
agreement



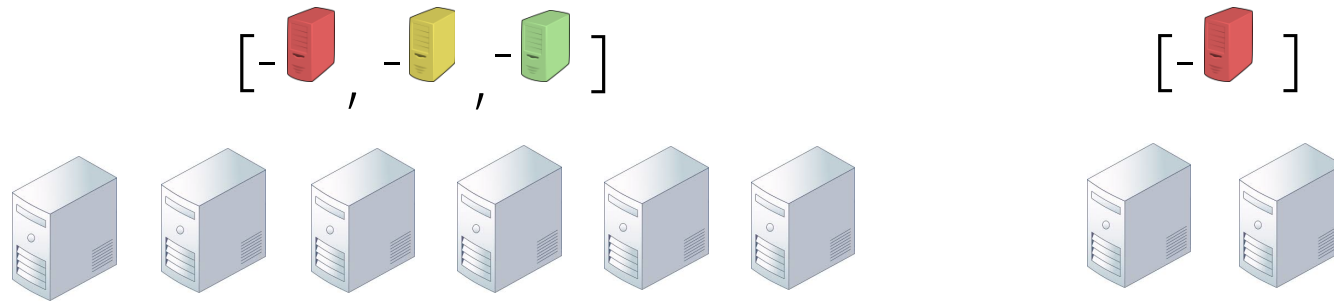
Full  
agreement



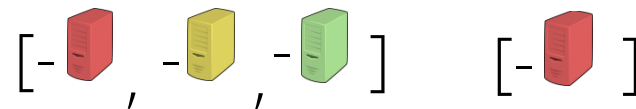
Almost-everywhere agreement



Full agreement



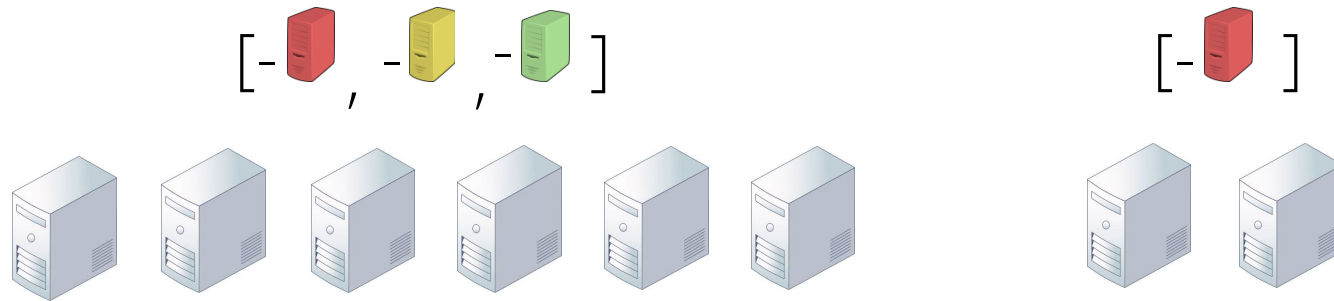
Every node counts #votes per-proposal



Almost-everywhere  
agreement



Full  
agreement



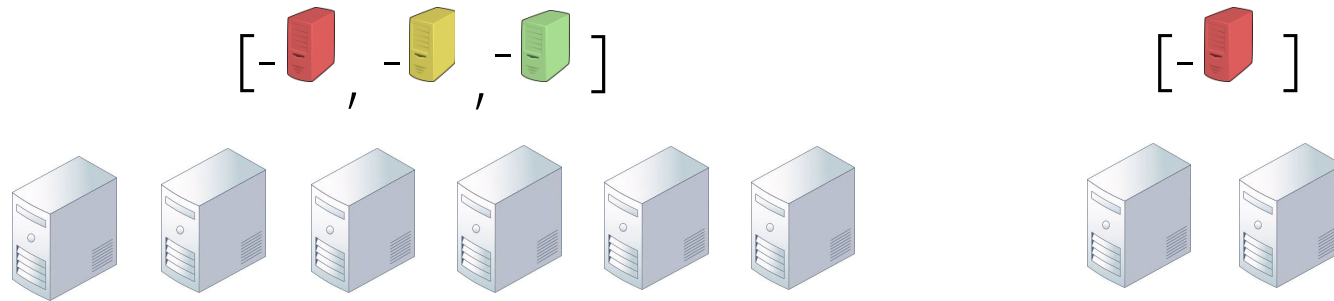
Gossip-based  
Counting protocol

Decide if **Fast Paxos quorum** ( $> \frac{3n}{4}$  nodes) of identical votes

# Almost-everywhere agreement



# Full agreement



Gossip-based  
Counting protocol

1000 processes, 10 node membership change  
~11 KB bandwidth usage per node for 1 second  
(Memberlist uses ~8 KB/s)

# Evaluation

Implementation: ~2700 LOC in Java (~2600 LOC of tests)

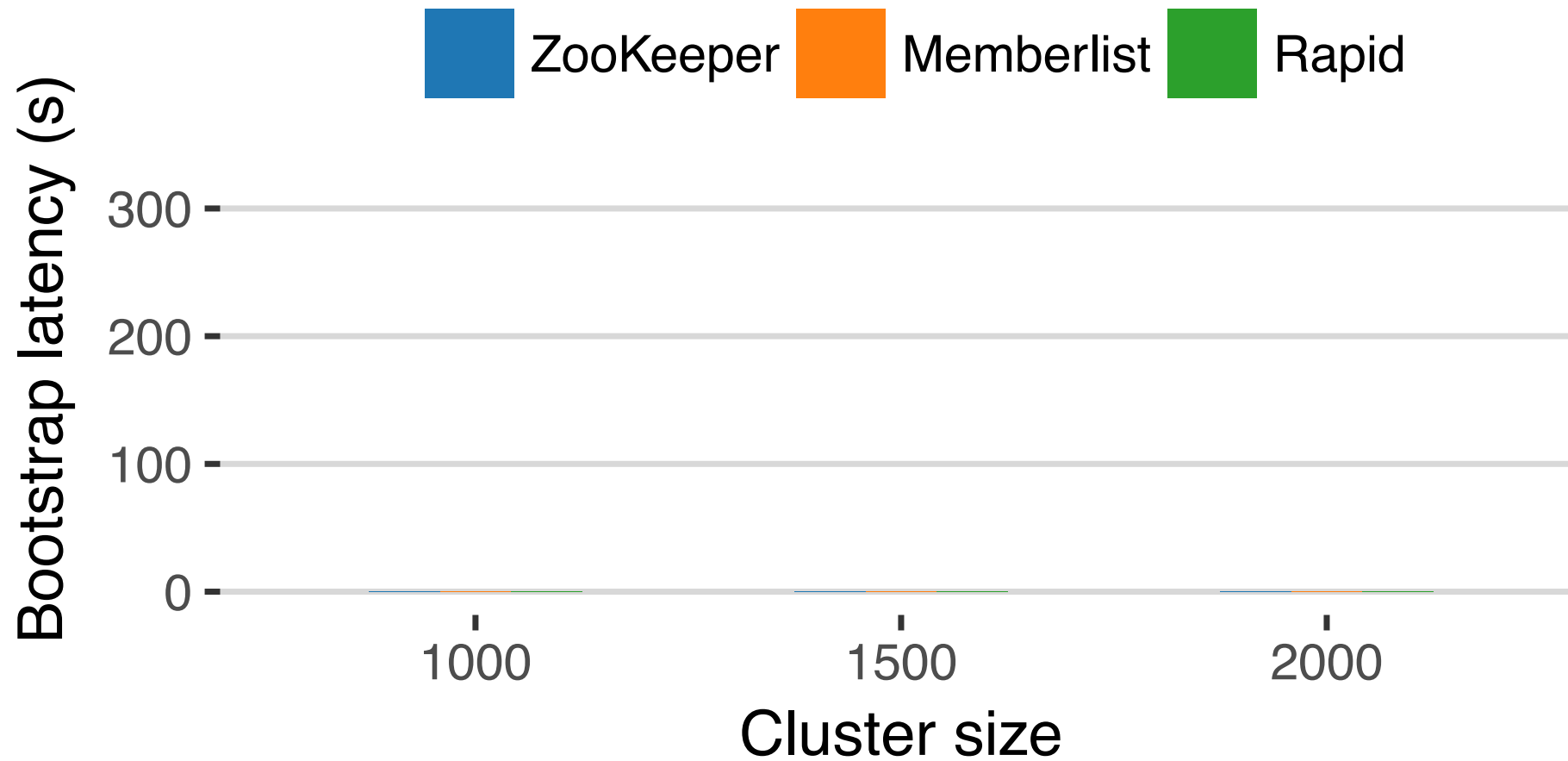
[github.com/lalithsuresh/rapid](https://github.com/lalithsuresh/rapid)

Compared against 3-node Zookeeper cluster and Memberlist.

Experiments run on 100 VMs  
(2 cores, 4GB RAM each)

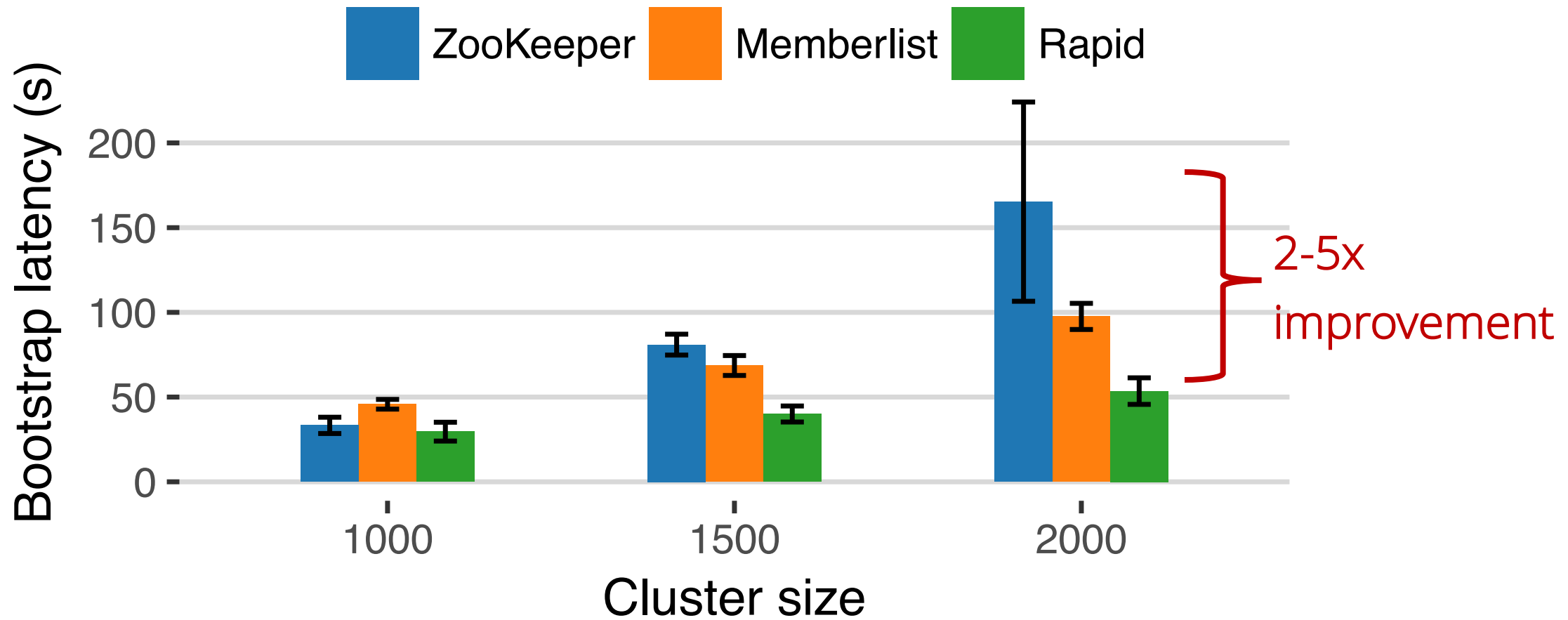
Not showing Akka Cluster because it did not scale past 500 nodes.

# Bootstrap times

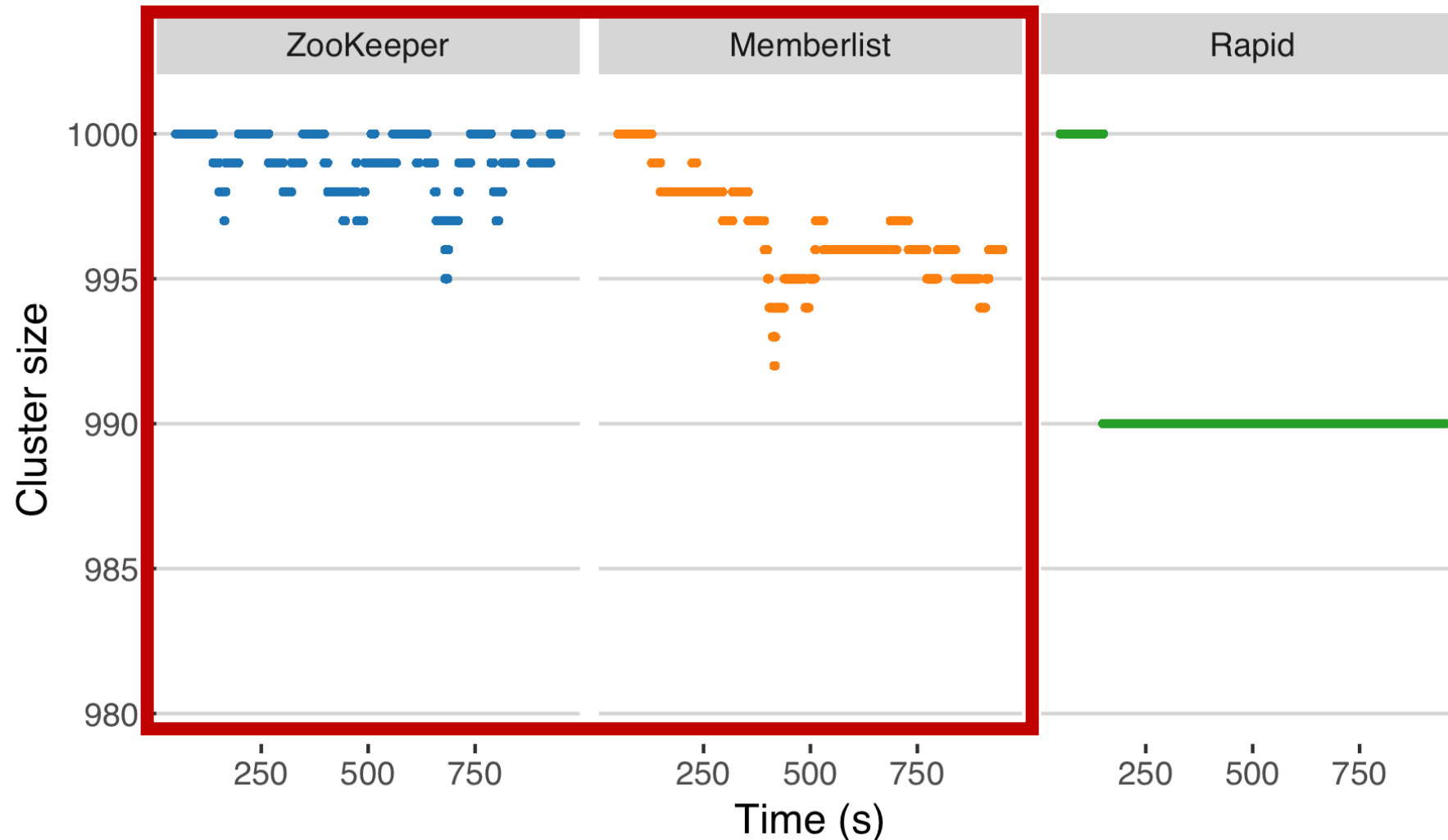




# Bootstrap times

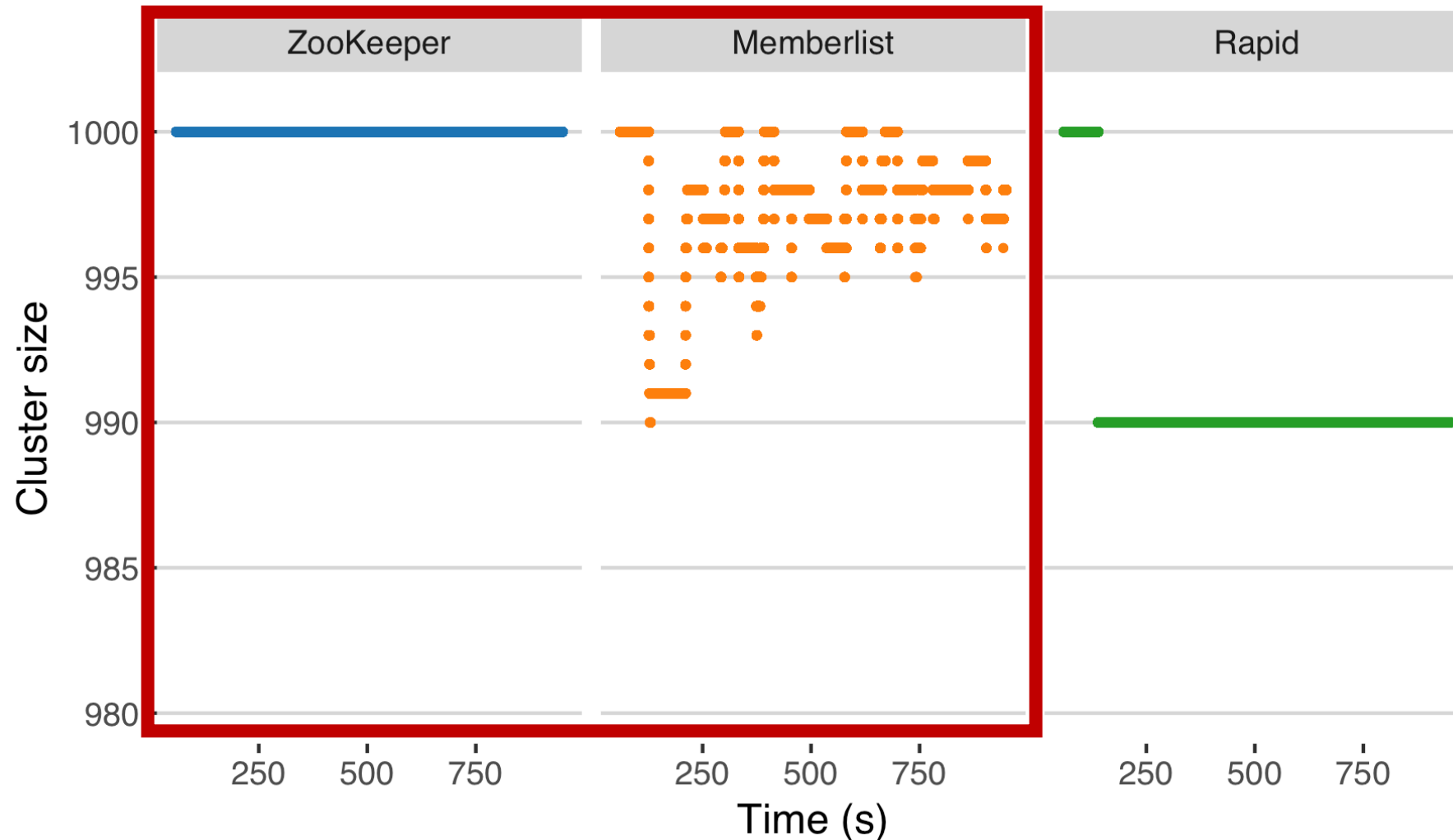


1% of processes experience  
high packet loss

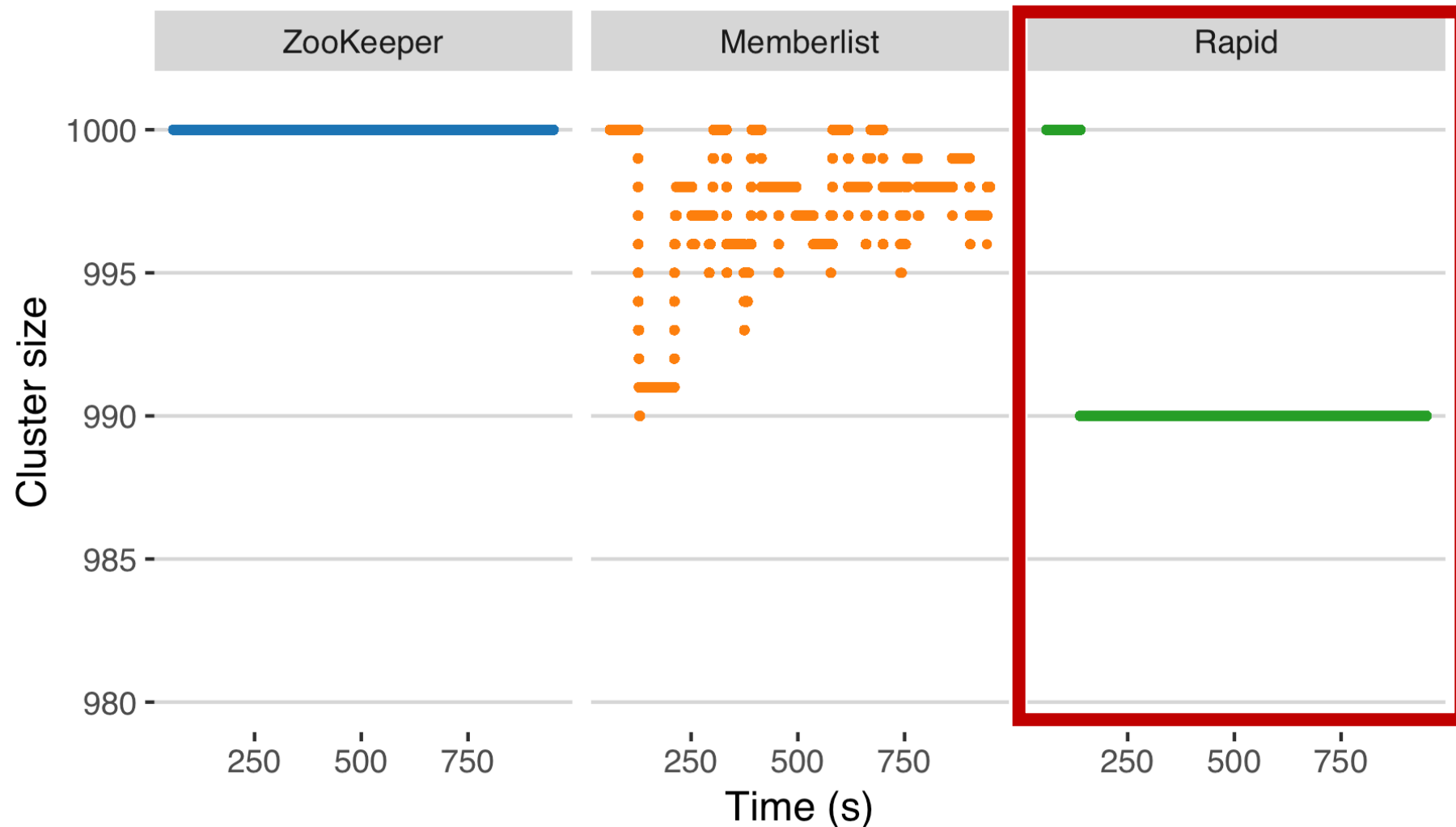




# 1% of processes experience one way network partition



# 1% of processes experience one way network partition



# Rapid

Stable and consistent membership at scale

