

L20: Replicated state machines with Paxos

Sam Madden
6.033 Spring 2014

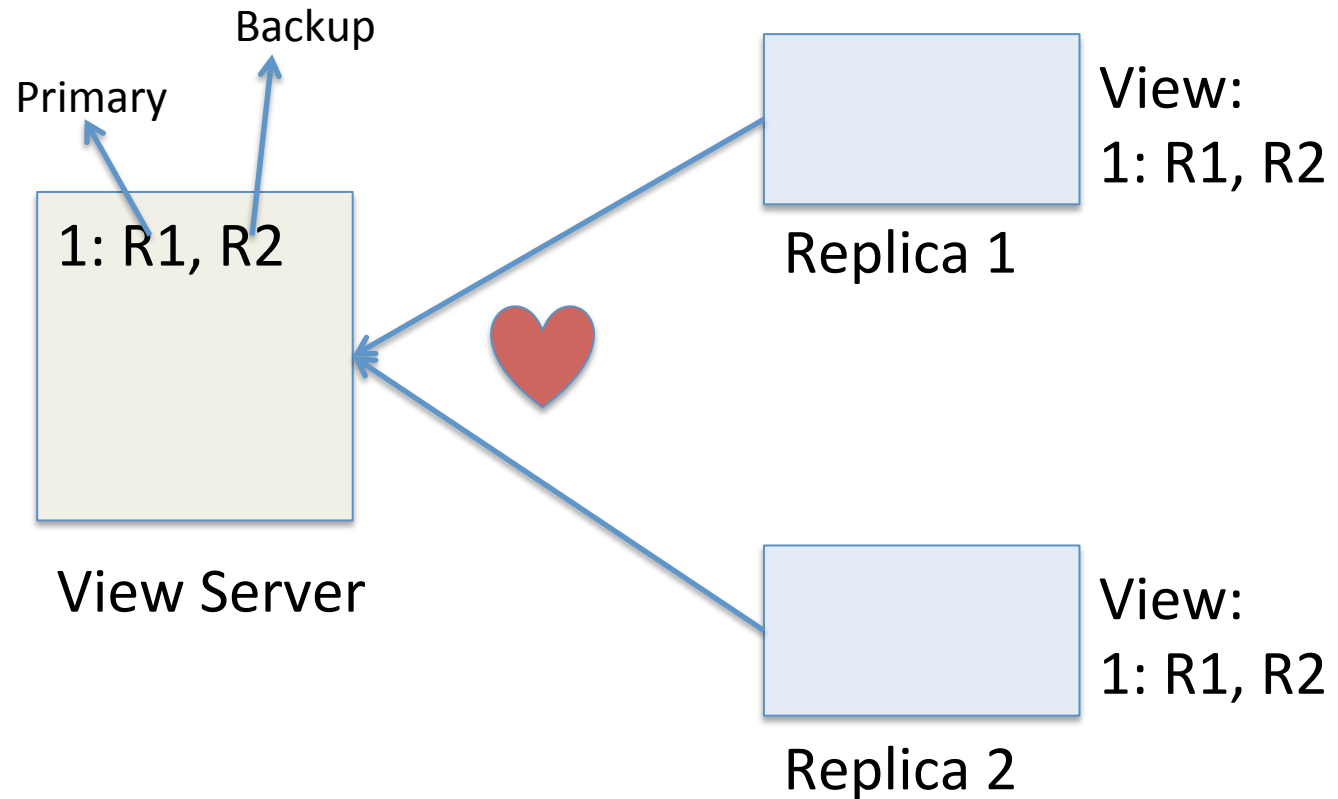
Overall Goal

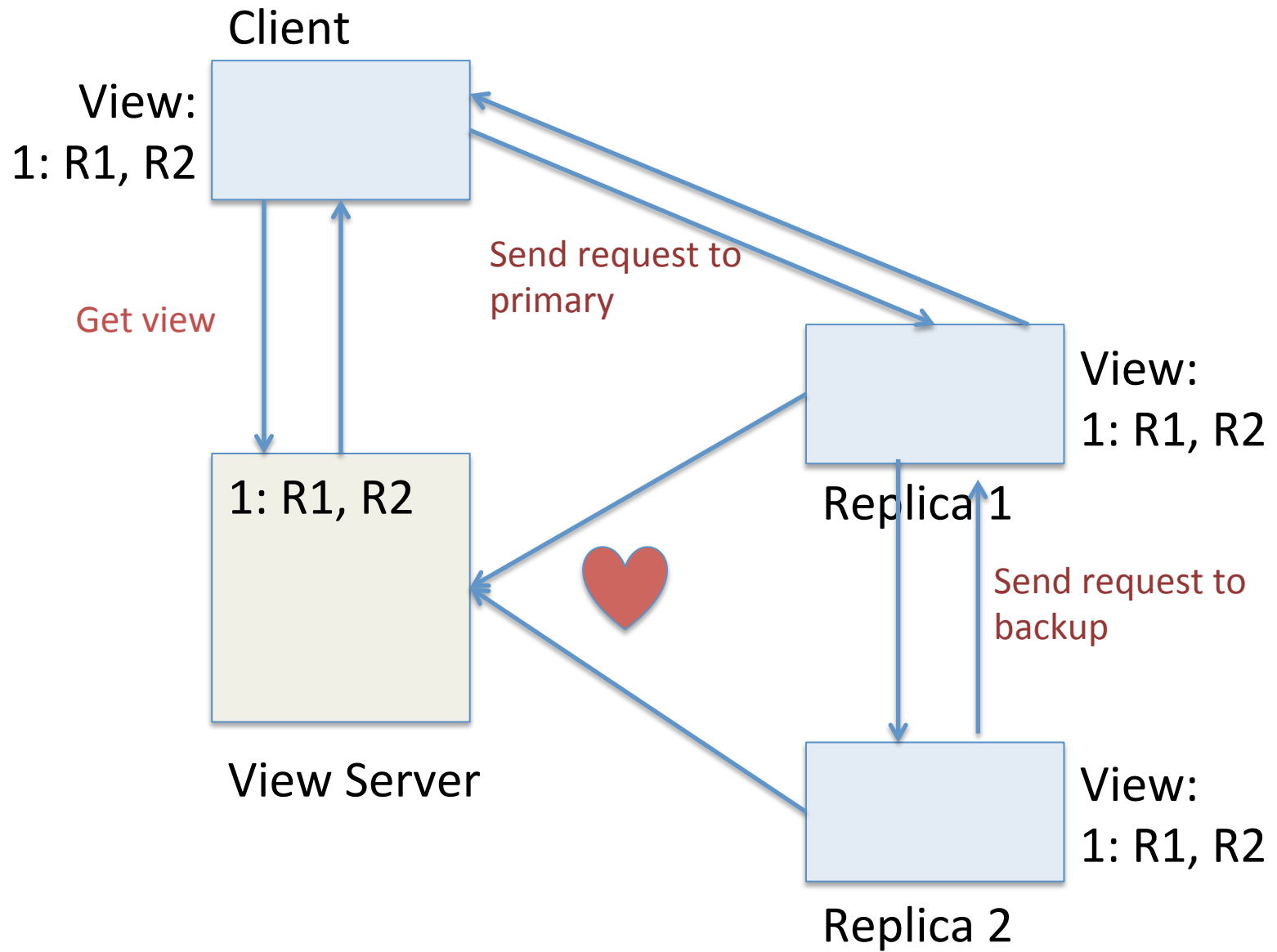
Building a stateful server (e.g., database) that remains *available* in the presence of node failures

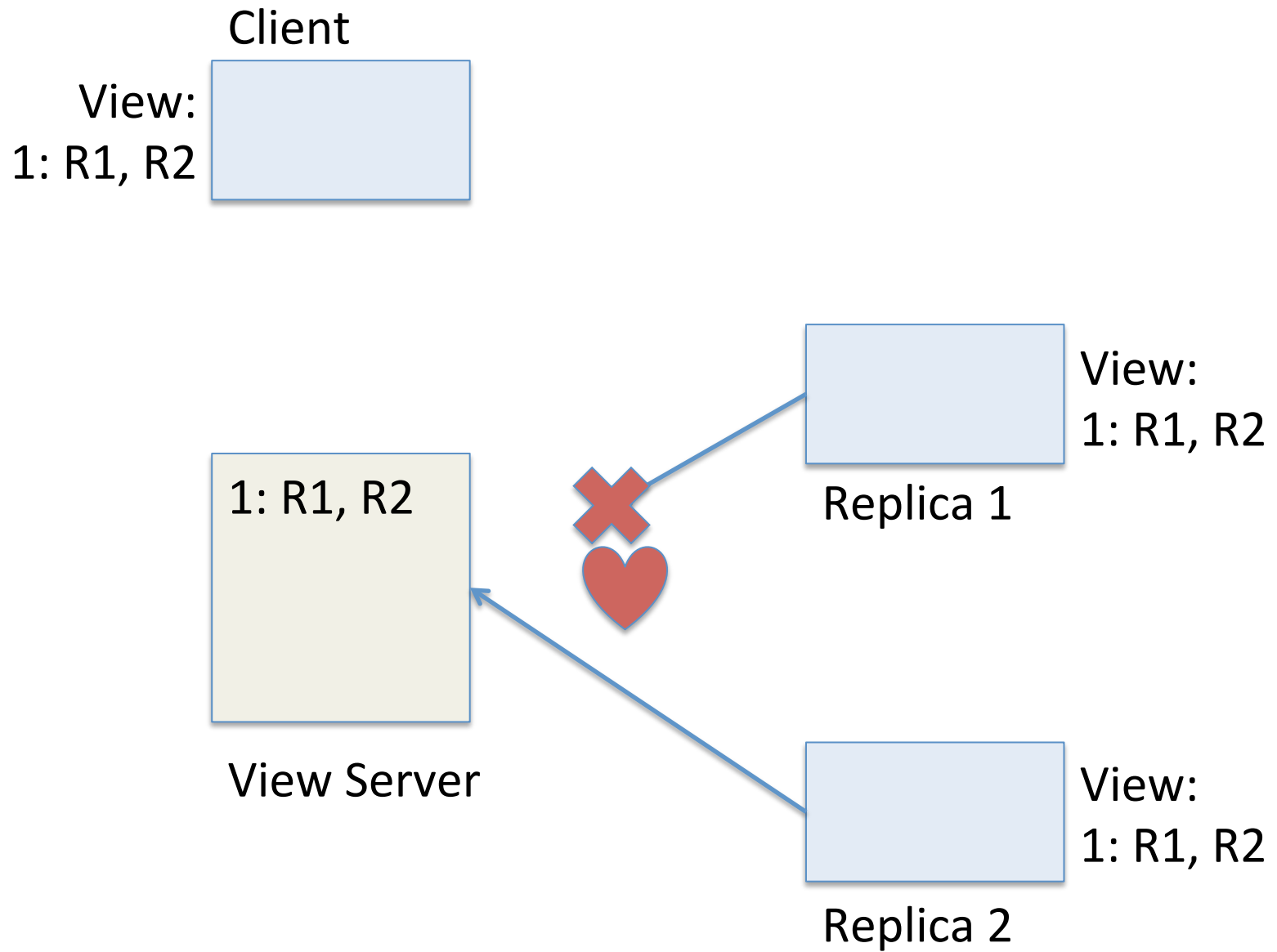
Last time: Replicated State Machine (RSM)

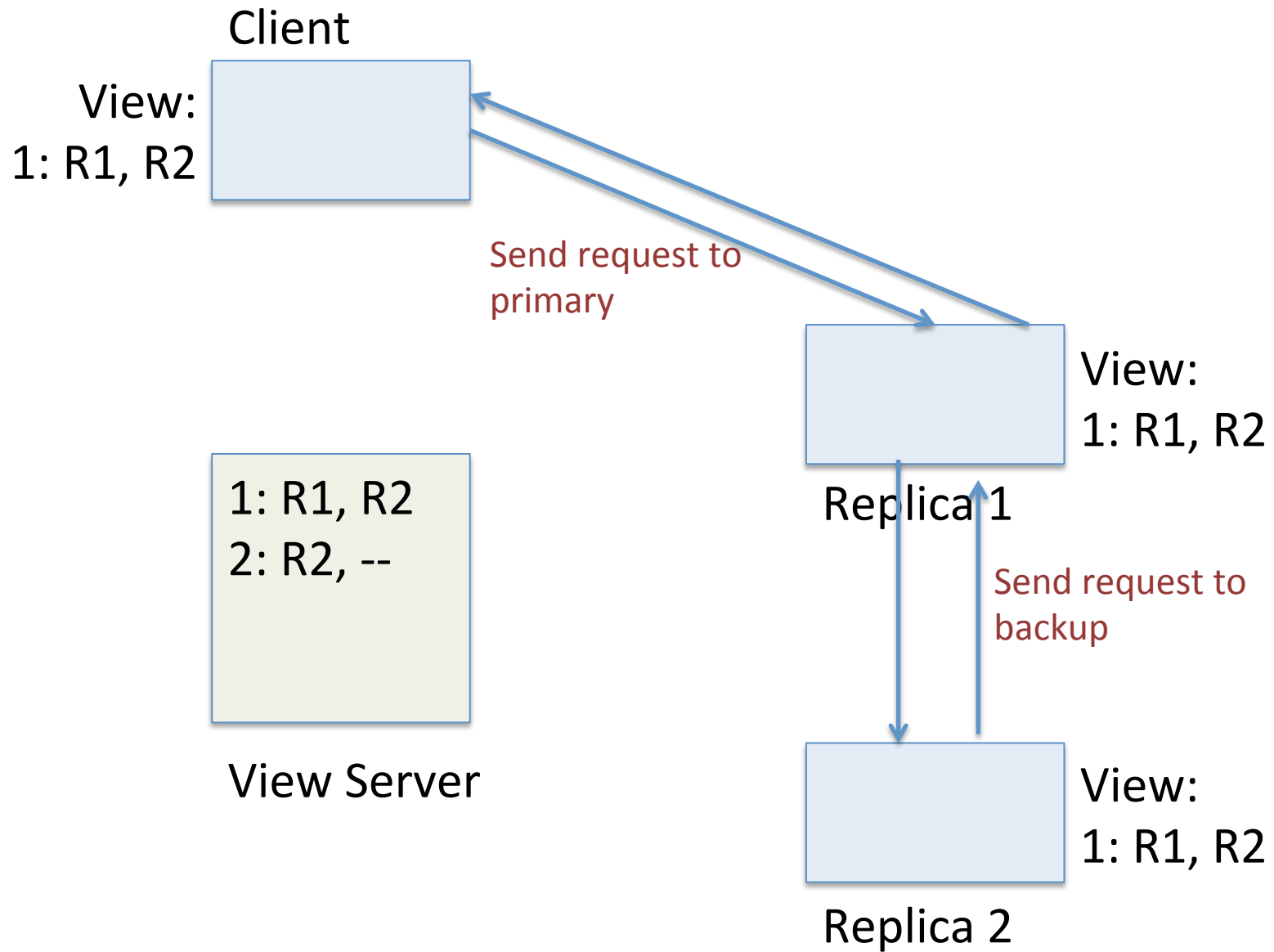
- Key idea: send all replicas same commands in same order
 - This will keep them in sync
- Idea: make one node a primary (coordinator) to order all transactions and send to others
- What if primary fails? Introduce a *view server* that keeps track of current primary / replica

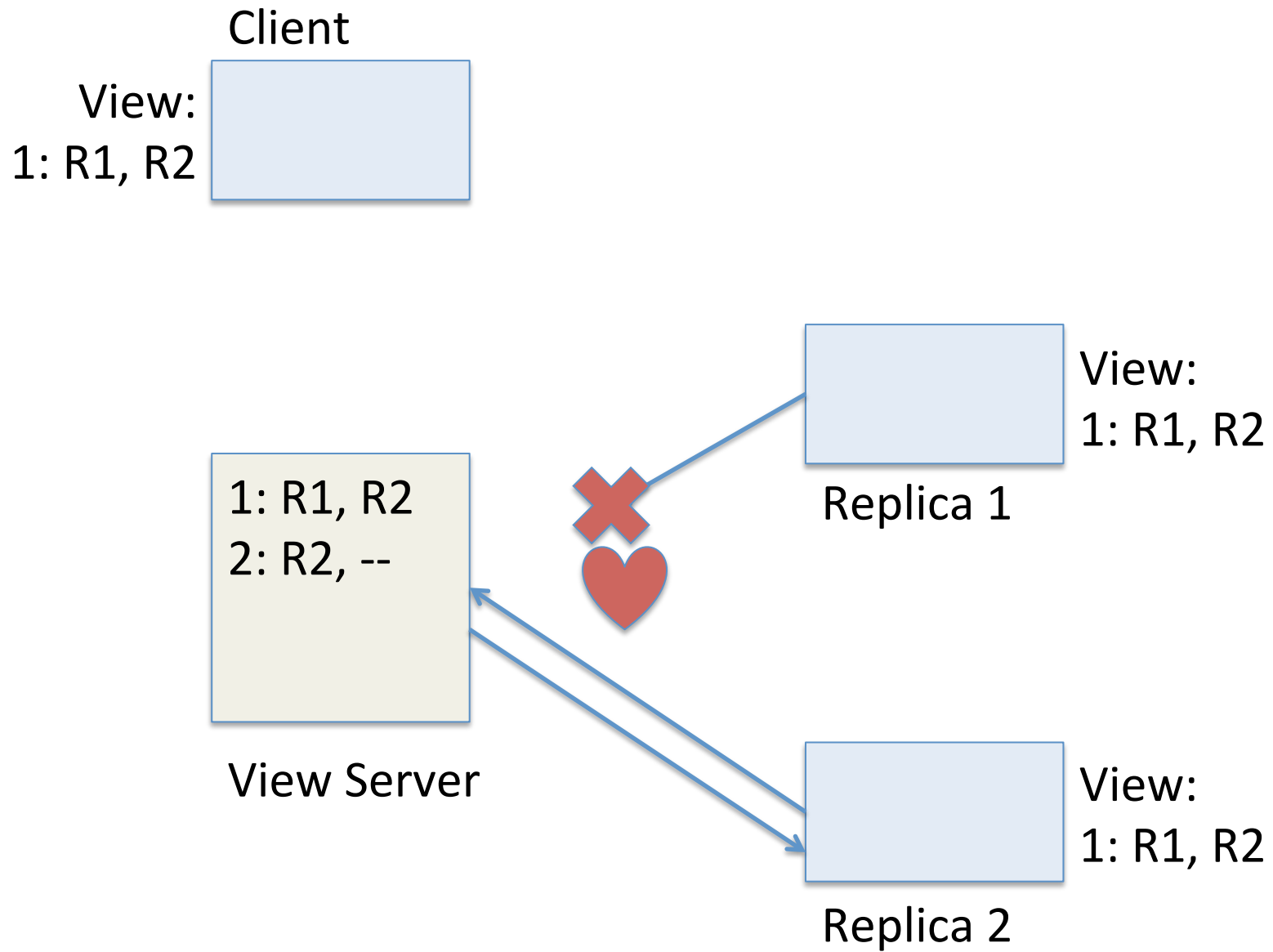
RSM w/ View Server Example

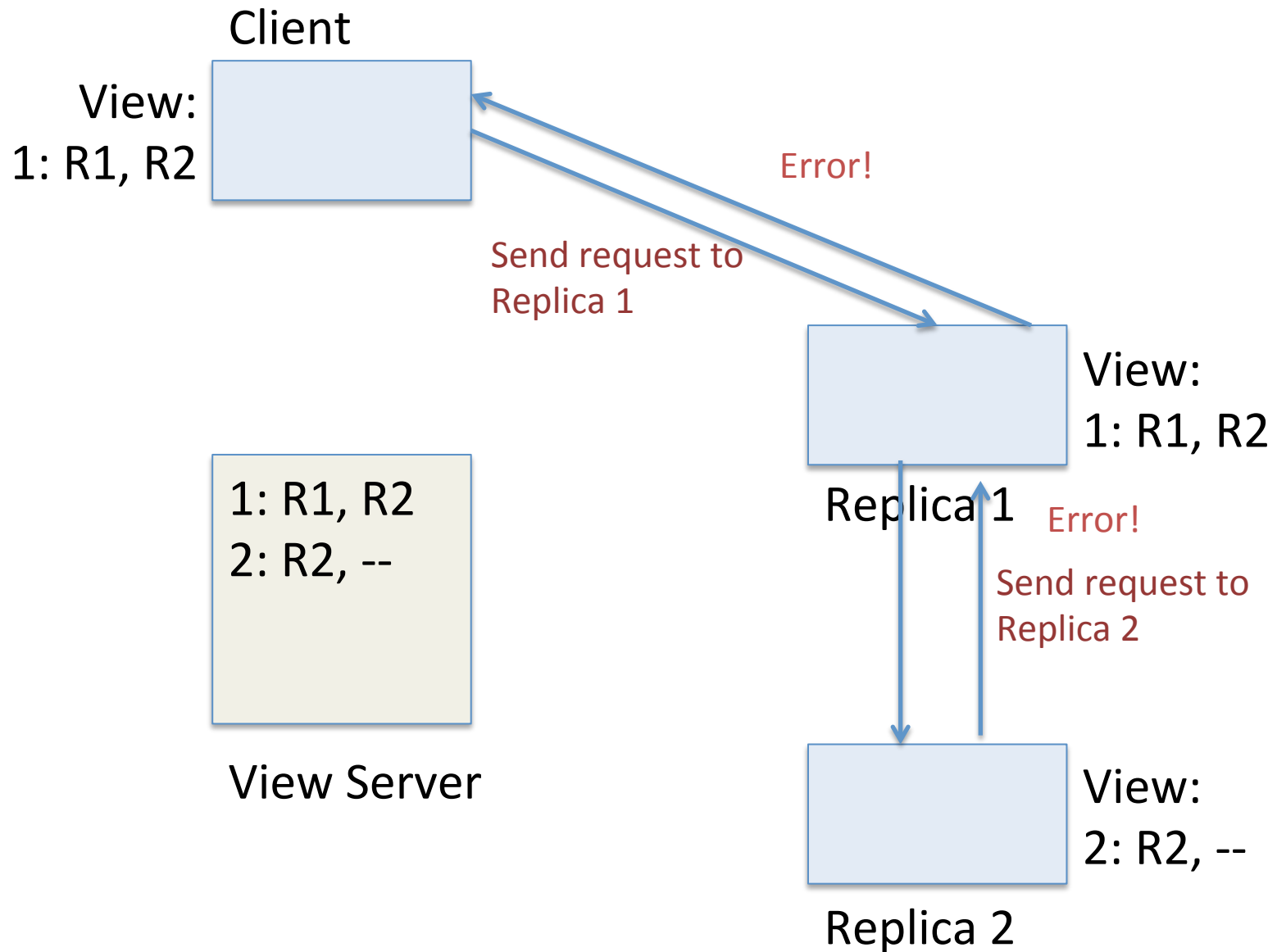






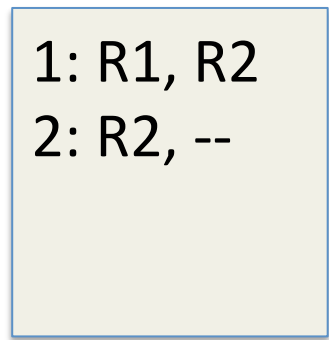
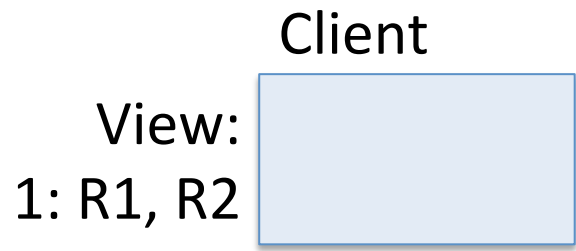




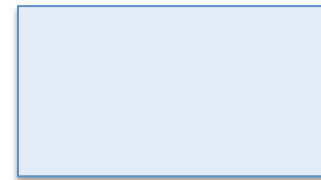


Moment R2 learns about new view it becomes active

Replica 2 refuses to process request

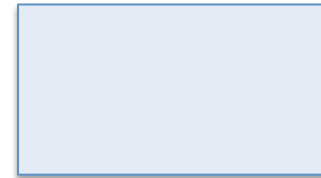


View Server



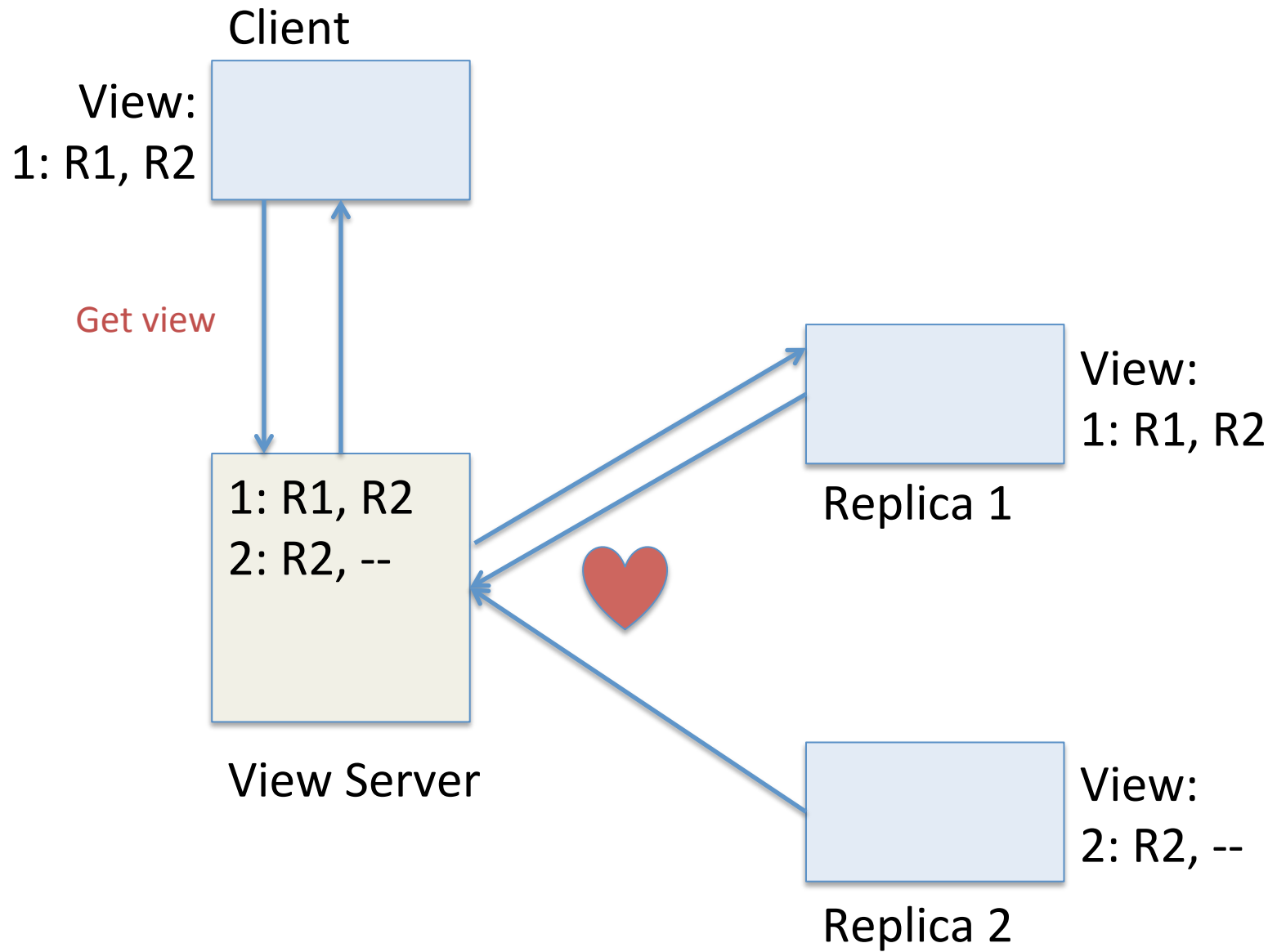
View:
1: R1, R2

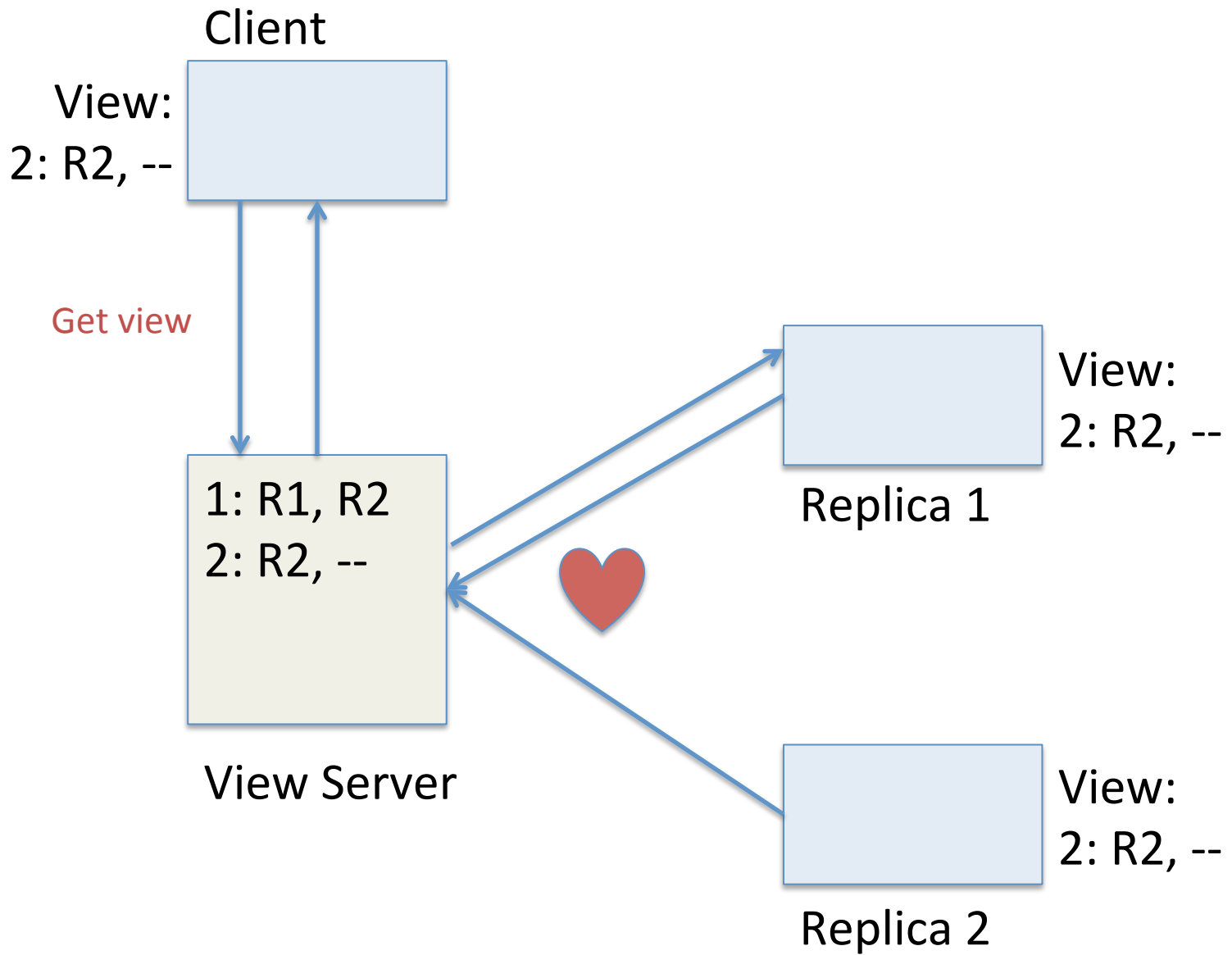
Replica 1

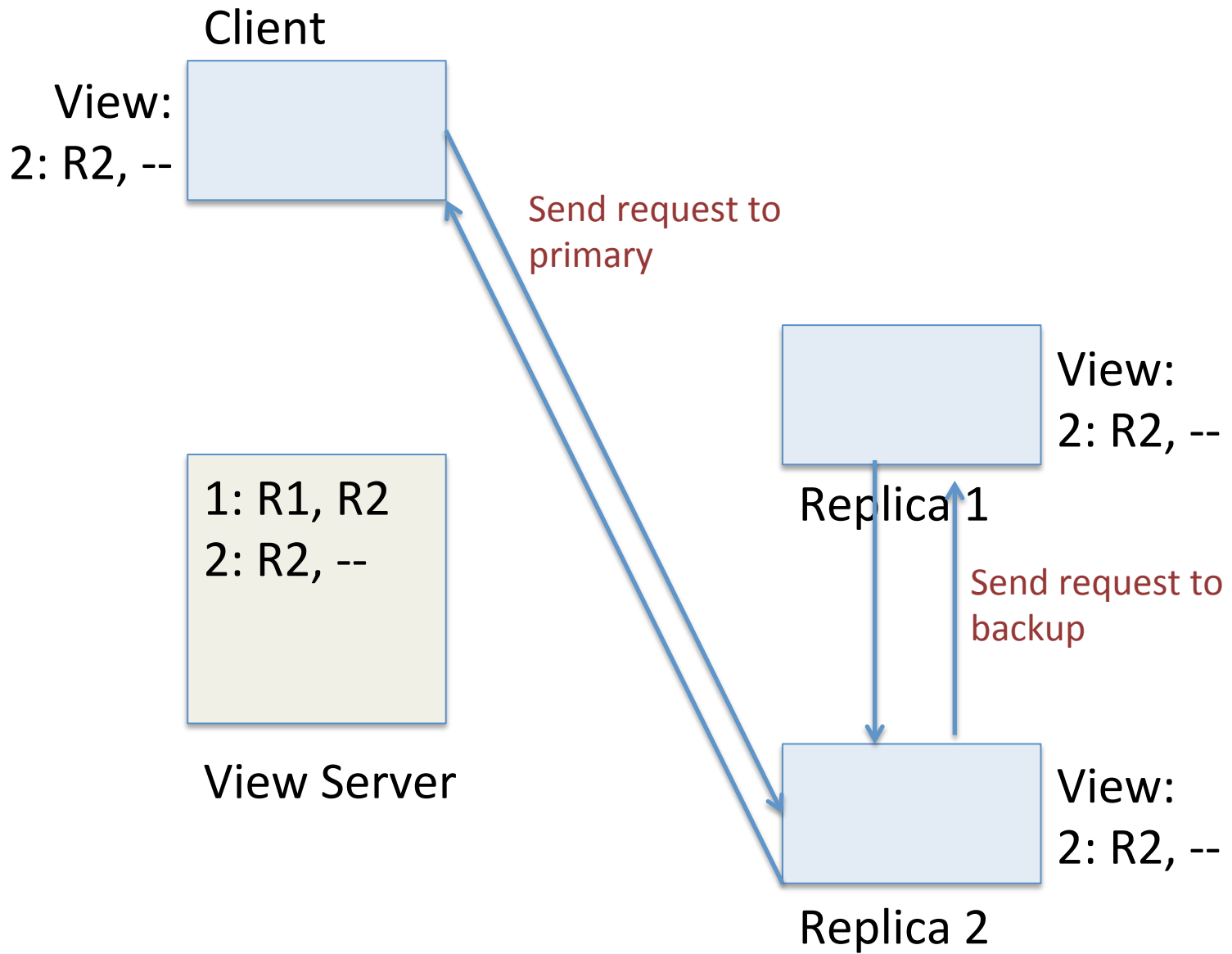


View:
2: R2, --

Replica 2







Problem: How to make view server
fault tolerant

Paxos Goal

$N \geq 3$ nodes, trying to *agree* about some value
(e.g., the next view server for the RSM protocol)

Paxos properties

- All nodes agree on a value, despite node failures, network failures, delays
 - E.g., X is the next operation to execute
 - E.g., Y is the next primary
- Fault tolerant: succeeds if fewer than $N/2$ nodes fail
 - Liveness not guaranteed
- Assumption: nodes are fail-stop

Setup

Servers each run 3 processes:

Proposer – proposes new values

Acceptor – accepts (or rejects) proposals

Learner – client waiting for new value

Paxos rule

- If a majority of nodes in an earlier proposal number accepted a value, later proposals must accept the same value
- State maintained by acceptor:
 - N_p : largest proposal seen in prepare
 - N_a : largest proposal seen in accept
 - V_a : value accepted for proposal N_a
- State must be persistent across reboot

Propose(V):

choose unique N , $> N_p$

send **Prepare**(N) to acceptors

if **Prepare_OK**(N_a , V_a) from majority:

$V' = V_a$ with highest N_a , or V if none

send **Accept**(N , V') to acceptors

if **Accept_OK**(N) from majority:

send **Decided**(V') to learners

Paxos

Proposer

Prepare(N):

if $N > N_p$:

log $N_p = N$

reply **Prepare_OK**(N_a , V_a)

- **N_p** : largest proposal seen in prepare
- **N_a** : largest proposal seen in accept
- **V_a** : value accepted for proposal N_a

Acceptor

Accept(N, V):

if $N \geq N_p$:

log $N_a = N$, log $V_a = V$

reply **Accept_OK**(N_a , V_a)

Propose(V):

choose unique N , $> N_p$

send **Prepare**(N) to acceptors

if **Prepare_OK**(N_a , V_a) from majority:

$V' = V_a$ with highest N_a , or V if none

send **Accept**(N , V') to acceptors

if **Accept_OK**(N) from majority:

send **Decided**(V') to learners

Paxos

Proposer

Prepare(N):

if $N > N_p$:

log $N_p = N$

reply **Prepare_OK**(N_a , V_a)

- **N_p** : largest proposal seen in prepare
- **N_a** : largest proposal seen in accept
- **V_a** : value accepted for proposal N_a

Acceptor

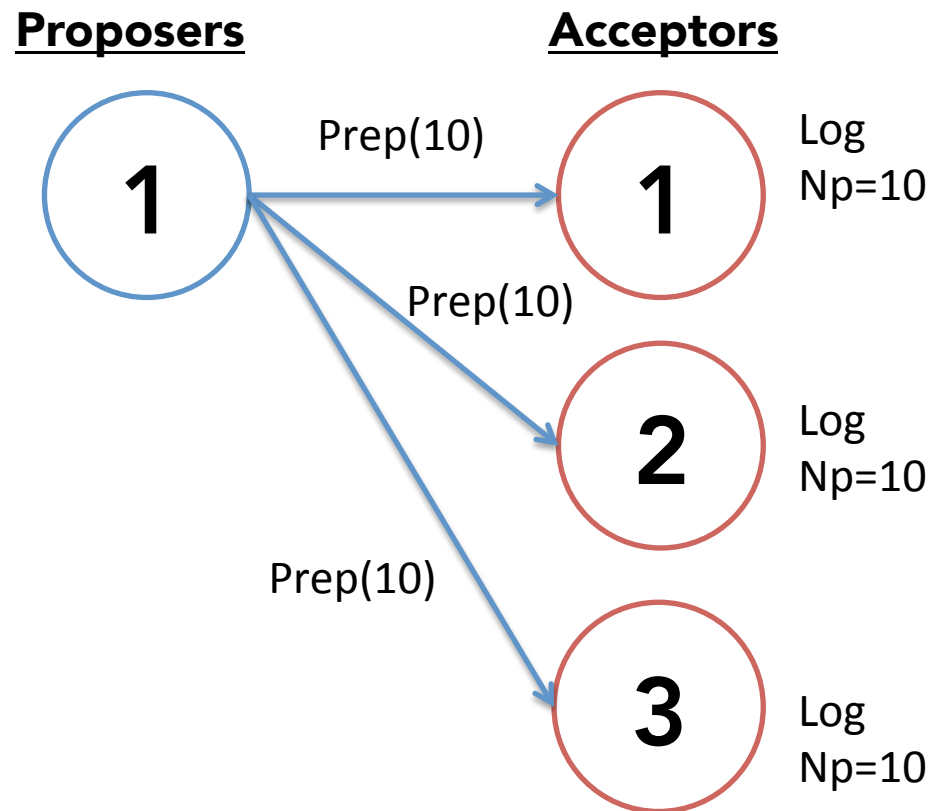
Accept(N, V):

if $N \geq N_p$:

log $N_a = N$, log $V_a = V$

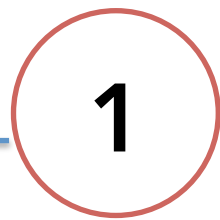
reply **Accept_OK**(N_a , V_a)

Example 1



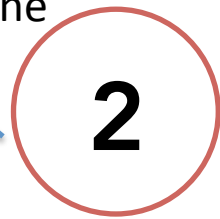
Proposers

Acceptors



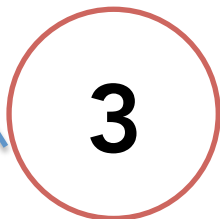
Log
Np=10

Ok, None



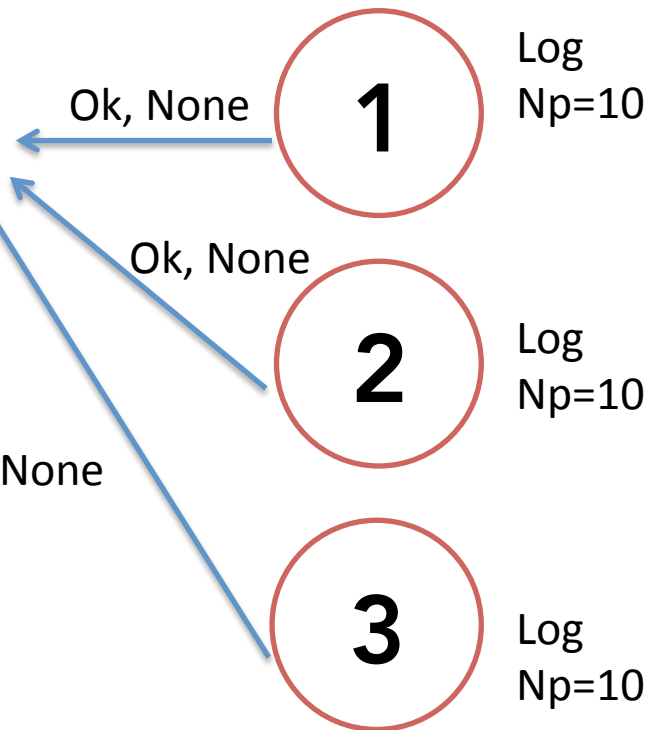
Log
Np=10

Ok, None



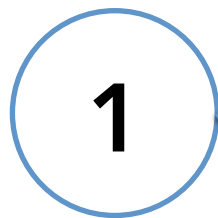
Log
Np=10

Ok, None

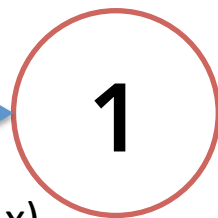


Proposers

Acceptors



Acc(10,x)

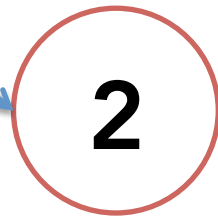


Log
Np=10

Log
Na=10

Log
Va=x

Acc(10,x)

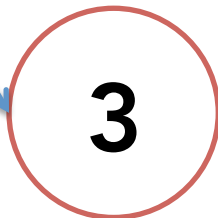


Log
Np=10

Log
Na=10

Log
Va=x

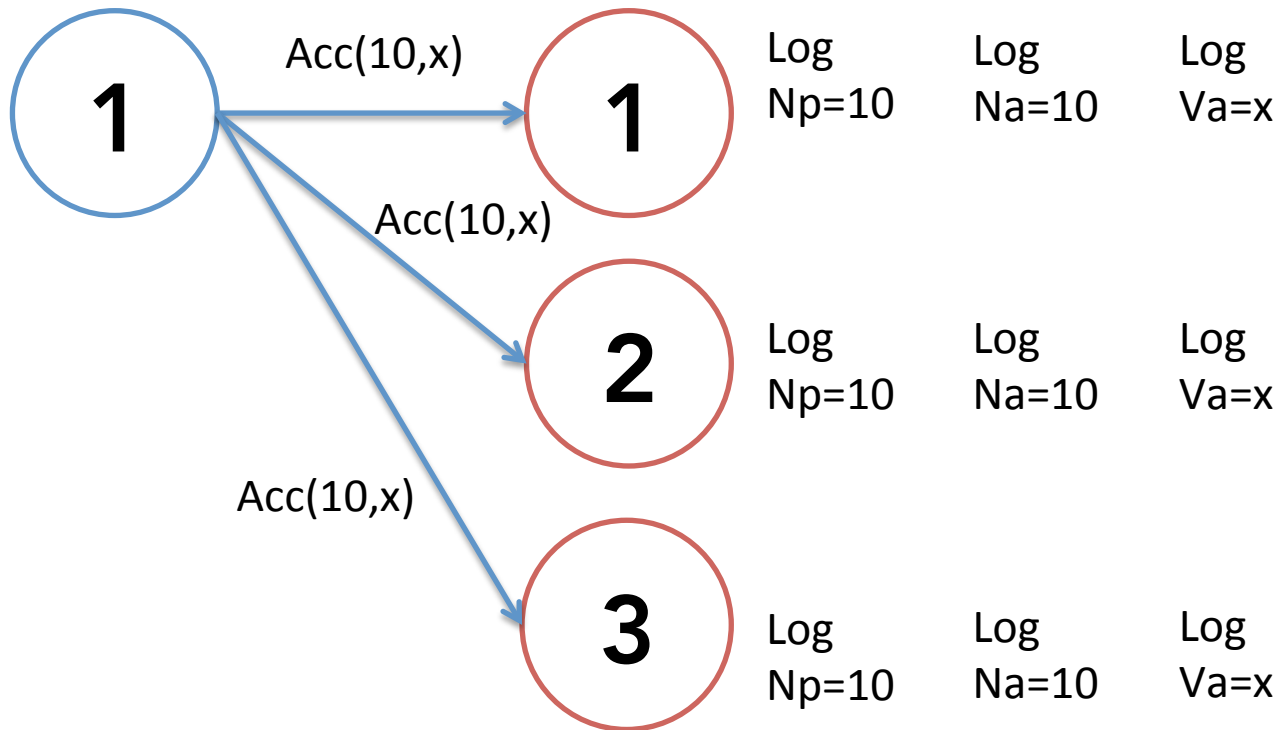
Acc(10,x)



Log
Np=10

Log
Na=10

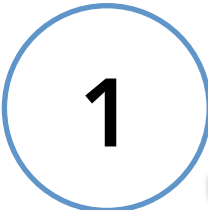
Log
Va=x



Commit point when majority of acceptors log value

Proposers

Acceptors



Log
Np=10

Log
Na=10

Log
Va=x

Ok



Log
Np=10

Log
Na=10

Log
Va=x

Ok

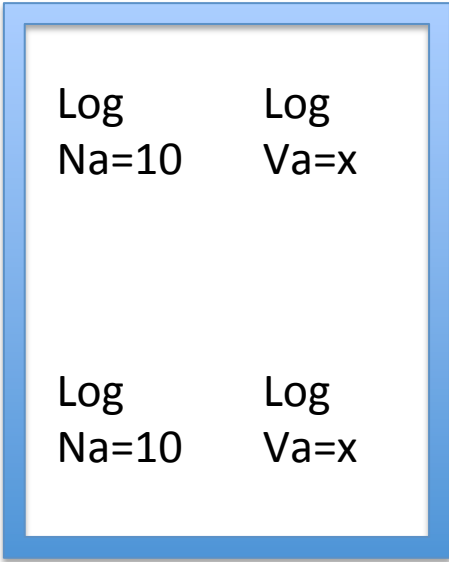


Log
Np=10

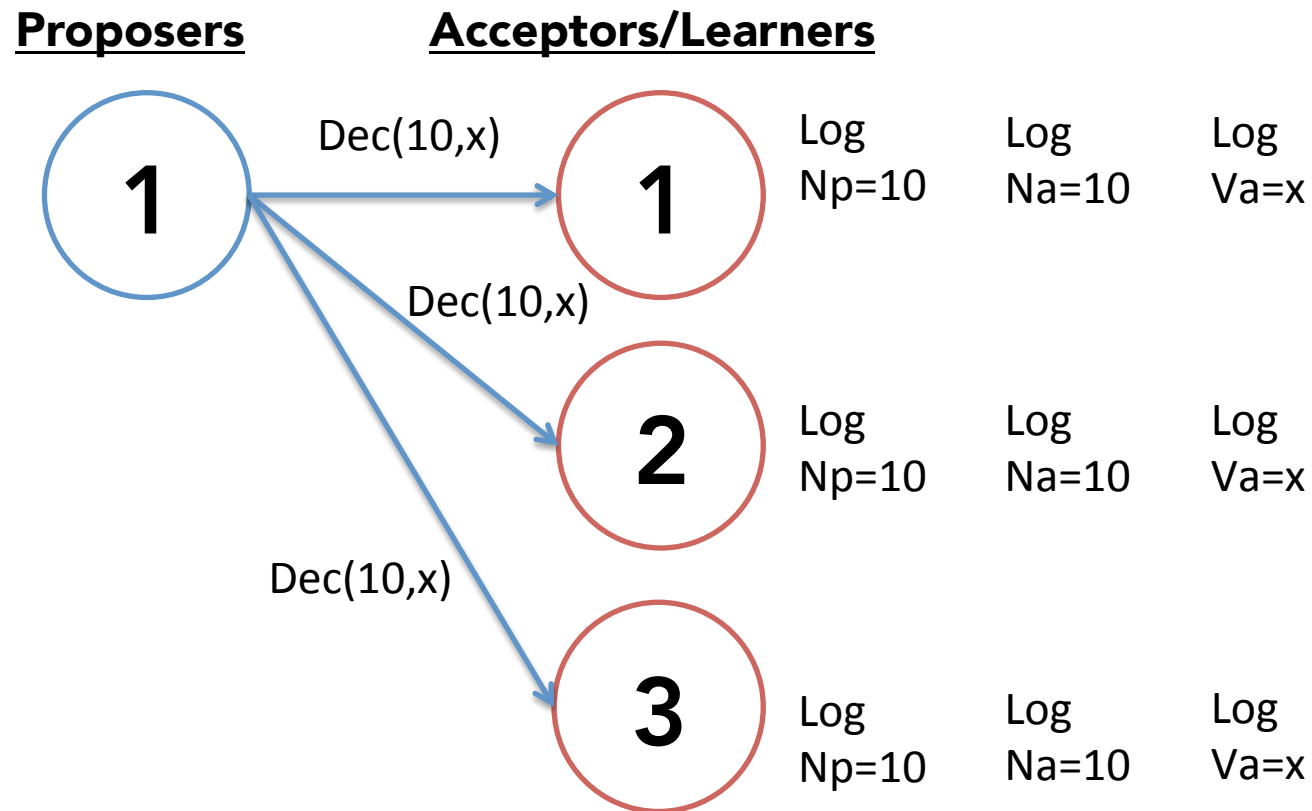
Log
Na=10

Log
Va=x

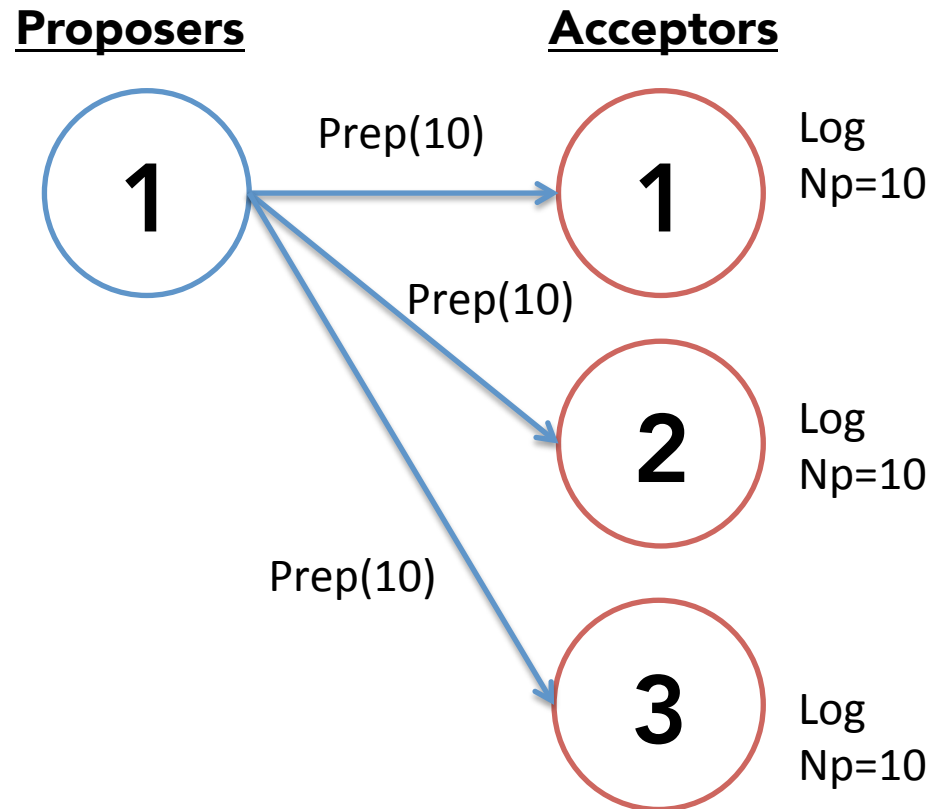
Ok



A1—3 can share with learners

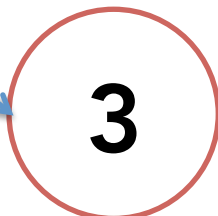
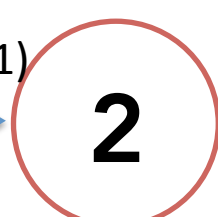
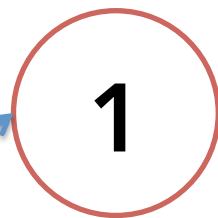
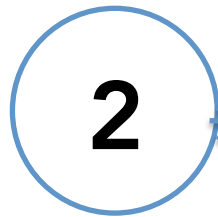
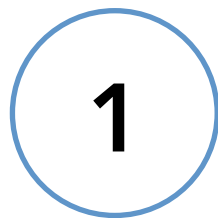


Example 2



Proposers

Acceptors



Prep(11)

Prep(11)

Prep(11)

Log
Np=10

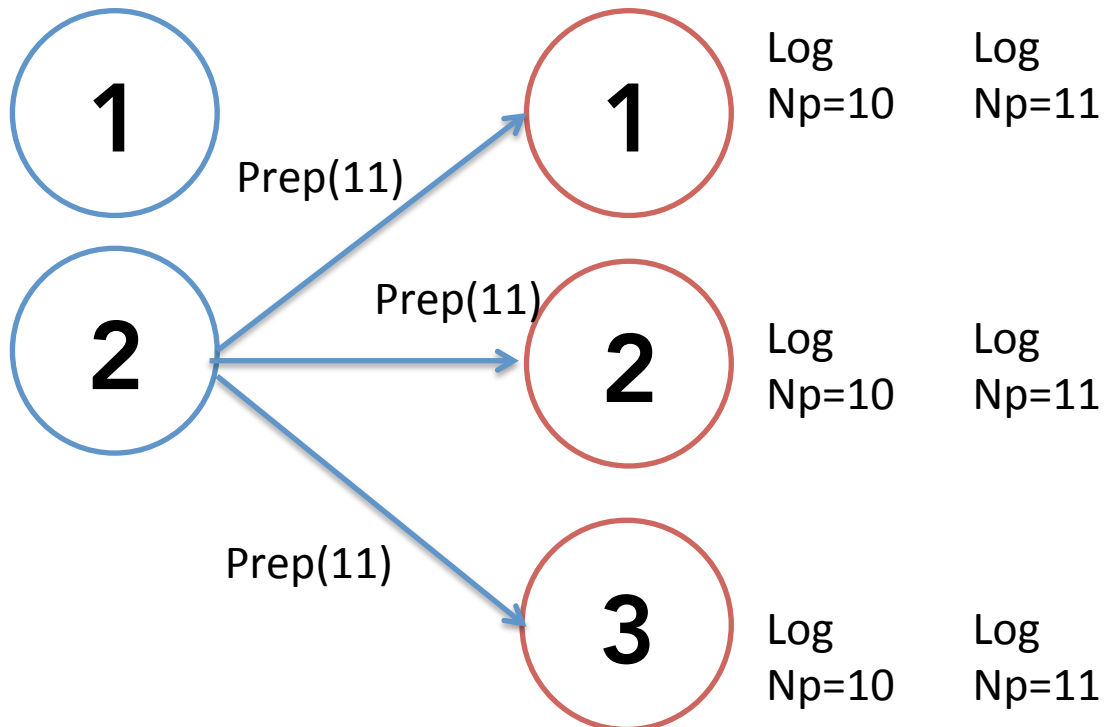
Log
Np=10

Log
Np=10

Log
Np=11

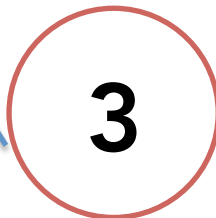
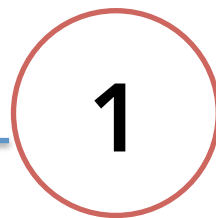
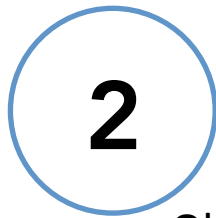
Log
Np=11

Log
Np=11



Proposers

Acceptors



Log
Np=10

Log
Np=11

Log
Np=10

Log
Np=11

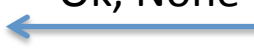
Log
Np=10

Log
Np=11

Ok, None

Ok, None

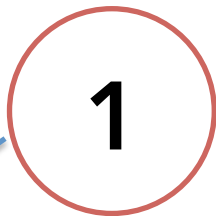
Ok, None



Proposers



Acceptors



Ok, None

Ok, None

Ok, None

Log
Np=10

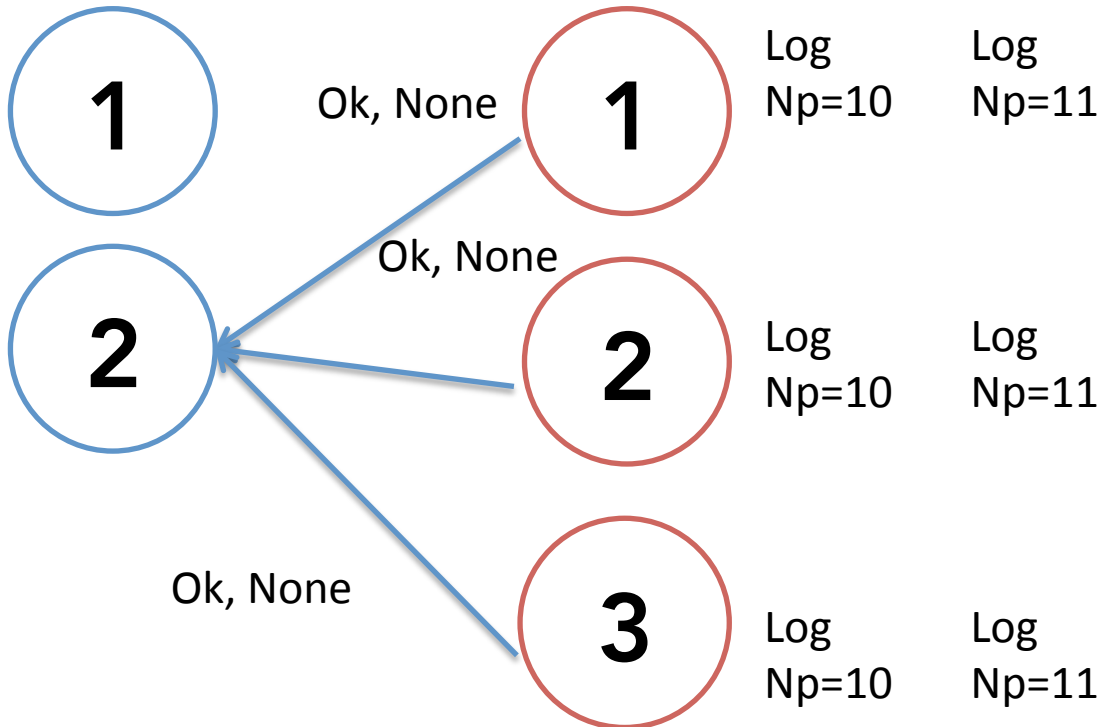
Log
Np=10

Log
Np=10

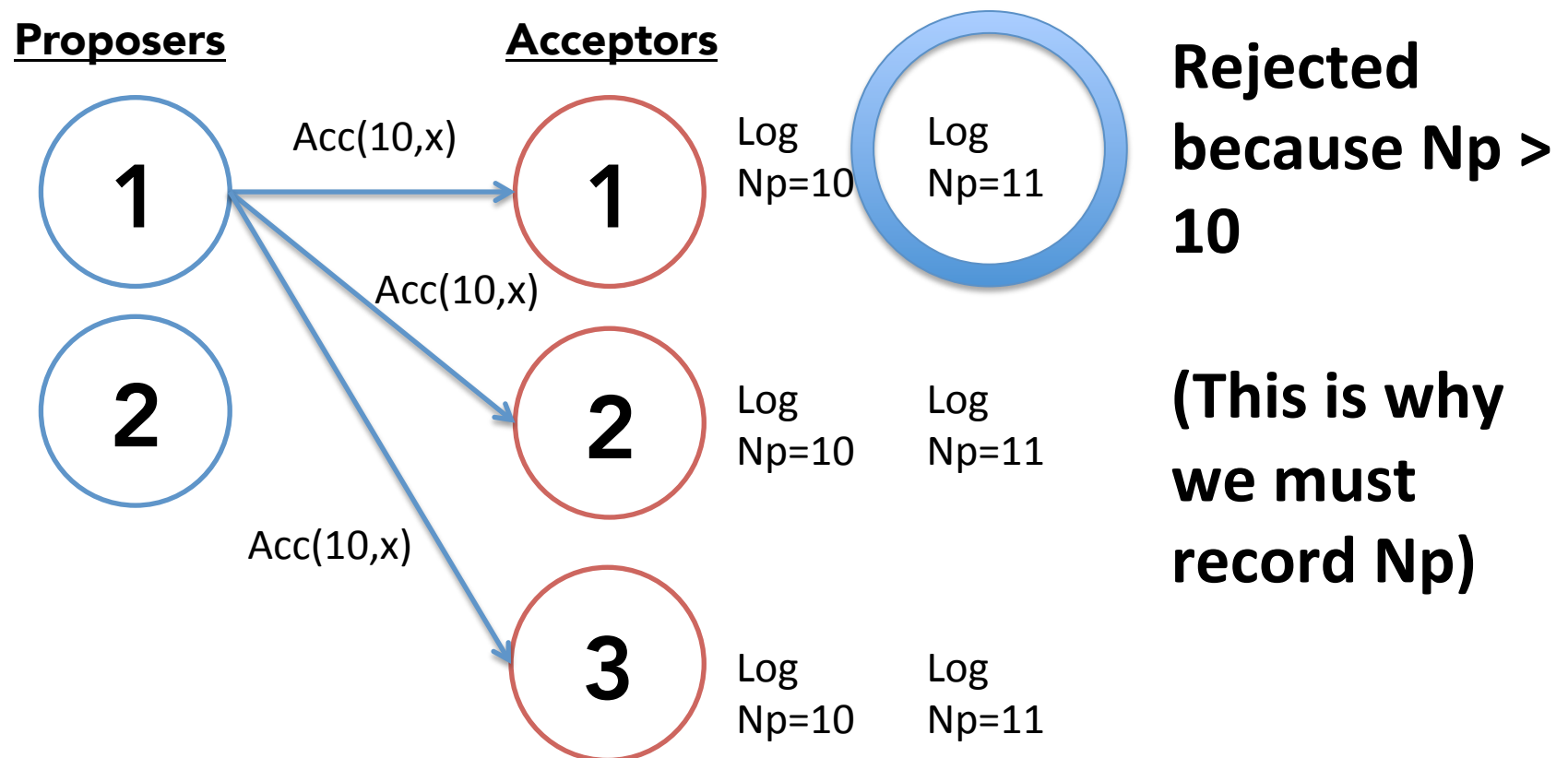
Log
Np=11

Log
Np=11

Log
Np=11

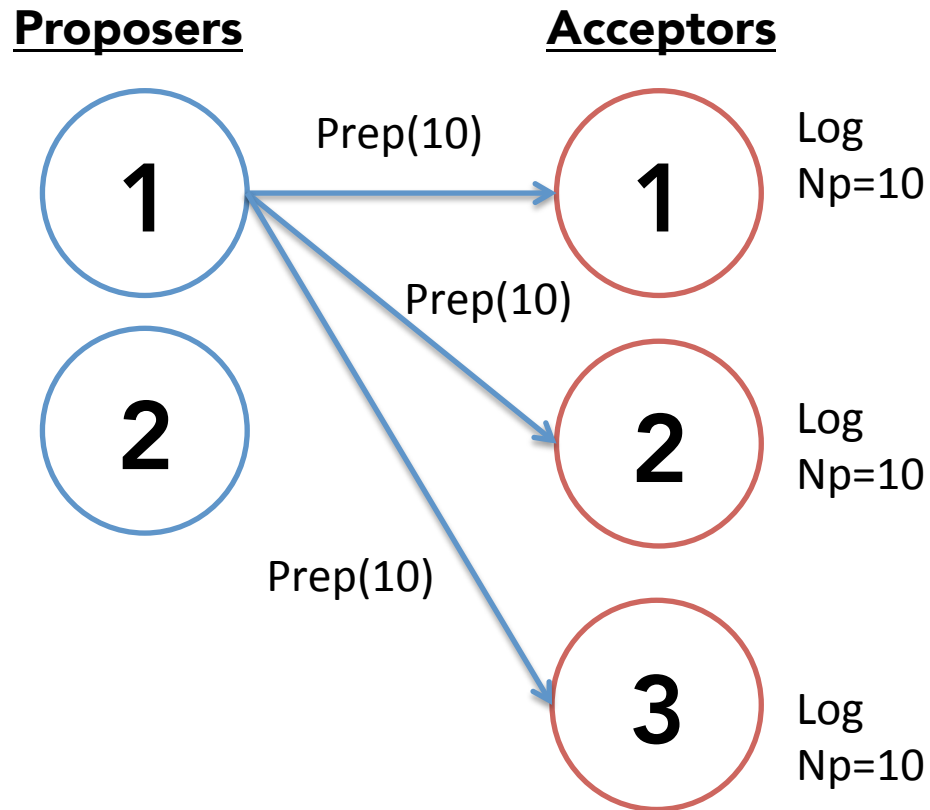


P1 attempts commit first, fails



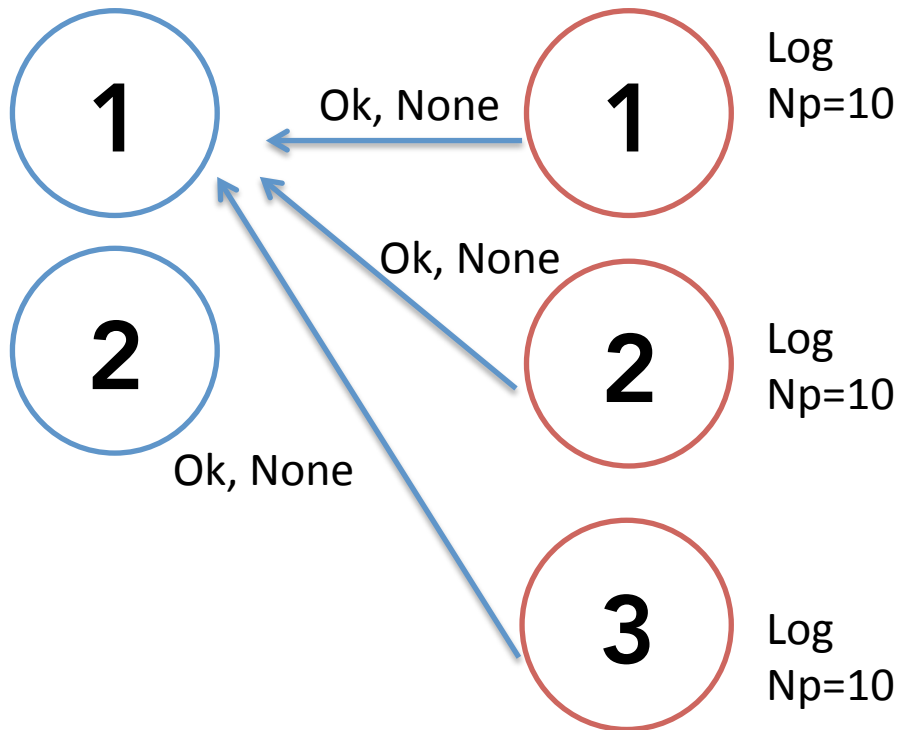
Decided Message Omitted for Brevity

Example 3



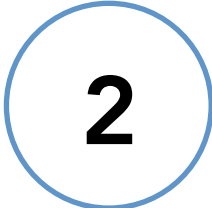
Proposers

Acceptors



Proposers

Acceptors



Acc(10,x)

Acc(10,x)

Acc(10,x)

Log
Np=10

Log
Np=10

Log
Np=10

Log
Na=10

Log
Na=10

Log
Na=10

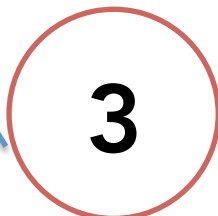
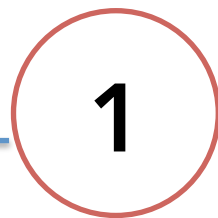
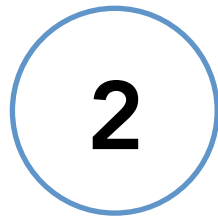
Log
Va=x

Log
Va=x

Log
Va=x

Proposers

Acceptors



Ok

Ok

Ok

Log
Np=10

Log
Na=10

Log
Va=x

Log
Np=10

Log
Na=10

Log
Va=x

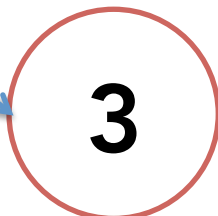
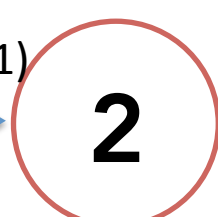
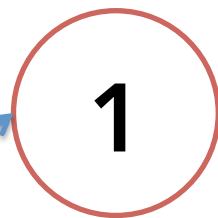
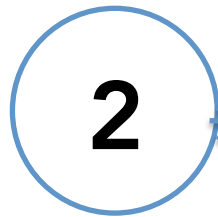
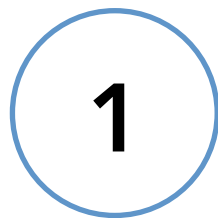
Log
Np=10

Log
Na=10

Log
Va=x

Proposers

Acceptors



Prep(11)

Prep(11)

Prep(11)

Log
Np=10

Log
Na=10

Log
Va=x

Log
Np=11

Log
Np=10

Log
Na=10

Log
Va=x

Log
Np=11

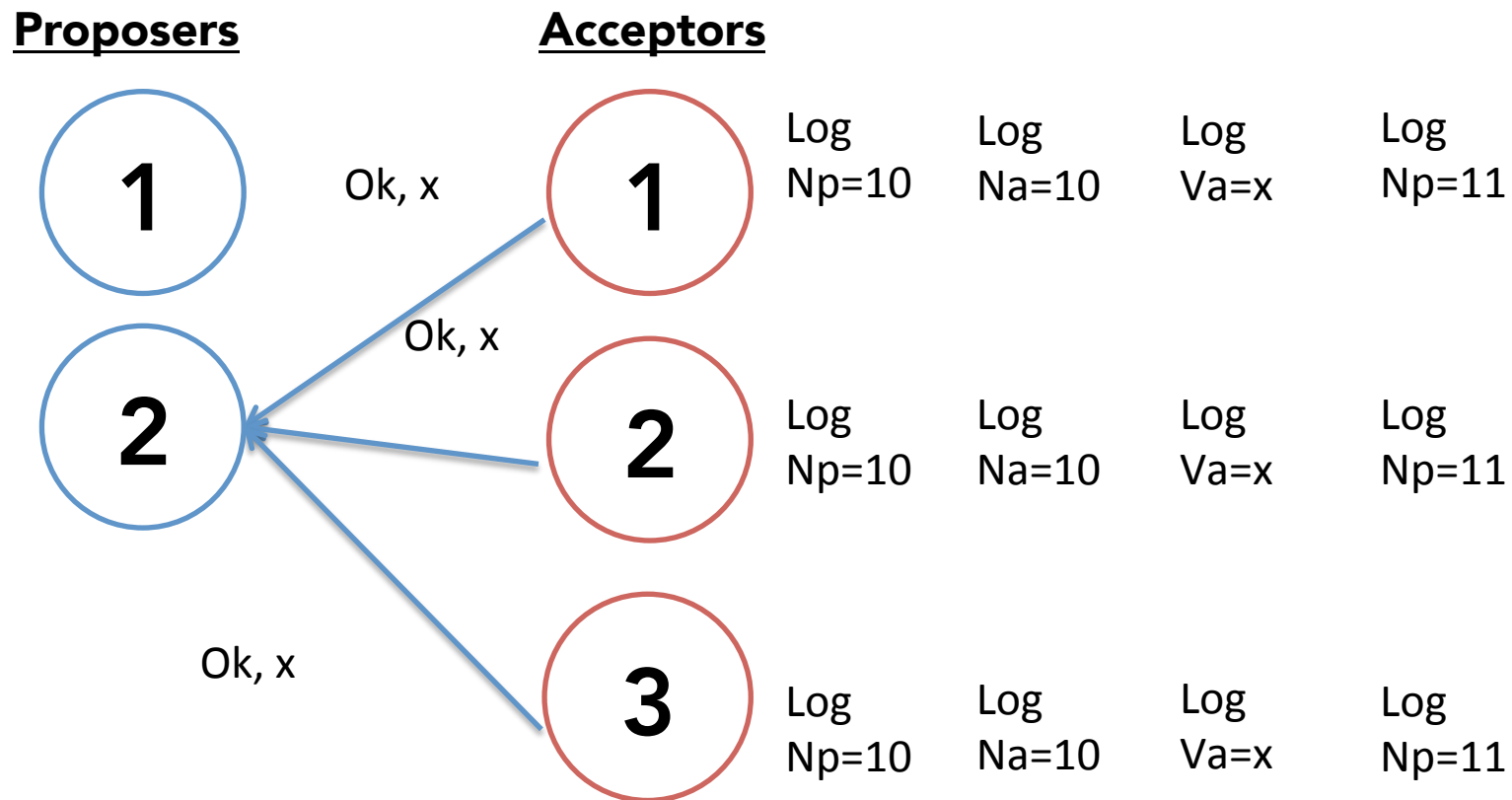
Log
Np=10

Log
Na=10

Log
Va=x

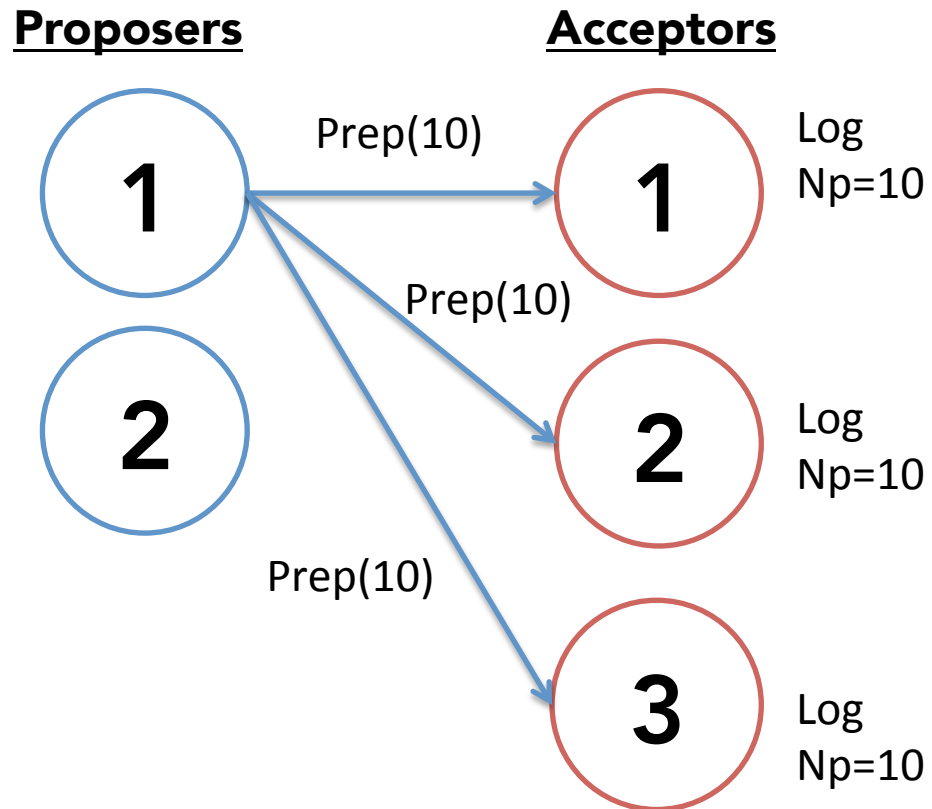
Log
Np=11

New proposer learns previously committed value



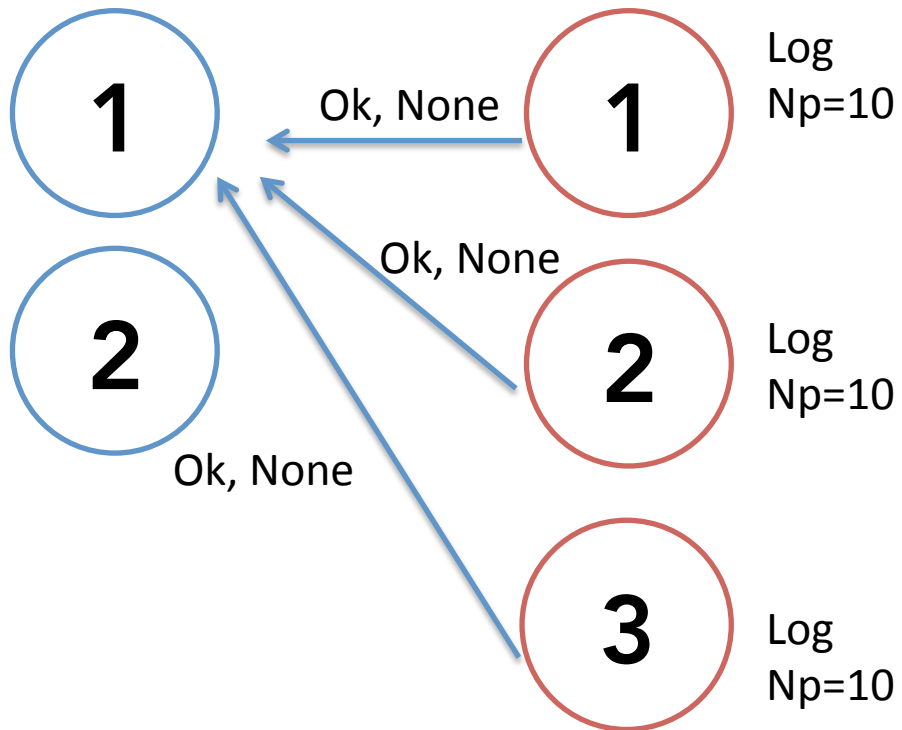
Decided Message Omitted for Brevity

Example 4



Proposers

Acceptors



Proposers

Acceptors



Acc(10,x)

Acc(10,x)

Acc(10,x)



Log
Np=10

Log
Na=10

Log
Va=x

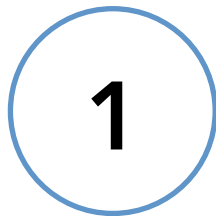
Log
Np=10

Log
Na=10

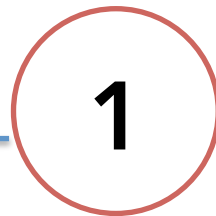
Log
Va=x

Log
Np=10

Proposers



Acceptors



Log
Np=10

Log
Na=10

Log
Va=x

← Ok

← Ok

Log
Np=10

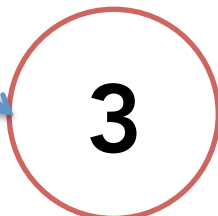
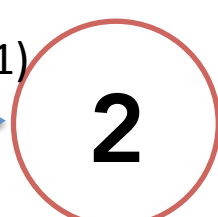
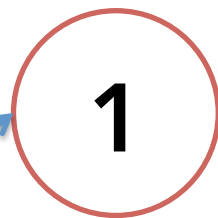
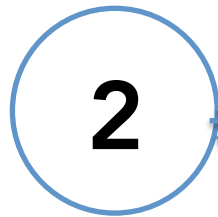
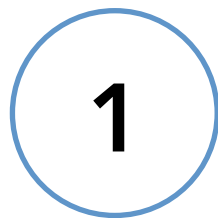
Log
Na=10

Log
Va=x

Log
Np=10

Proposers

Acceptors



Prep(11)

Prep(11)

Prep(11)

Log
Np=10

Log
Na=10

Log
Va=x

Log
Np=11

Log
Np=10

Log
Na=10

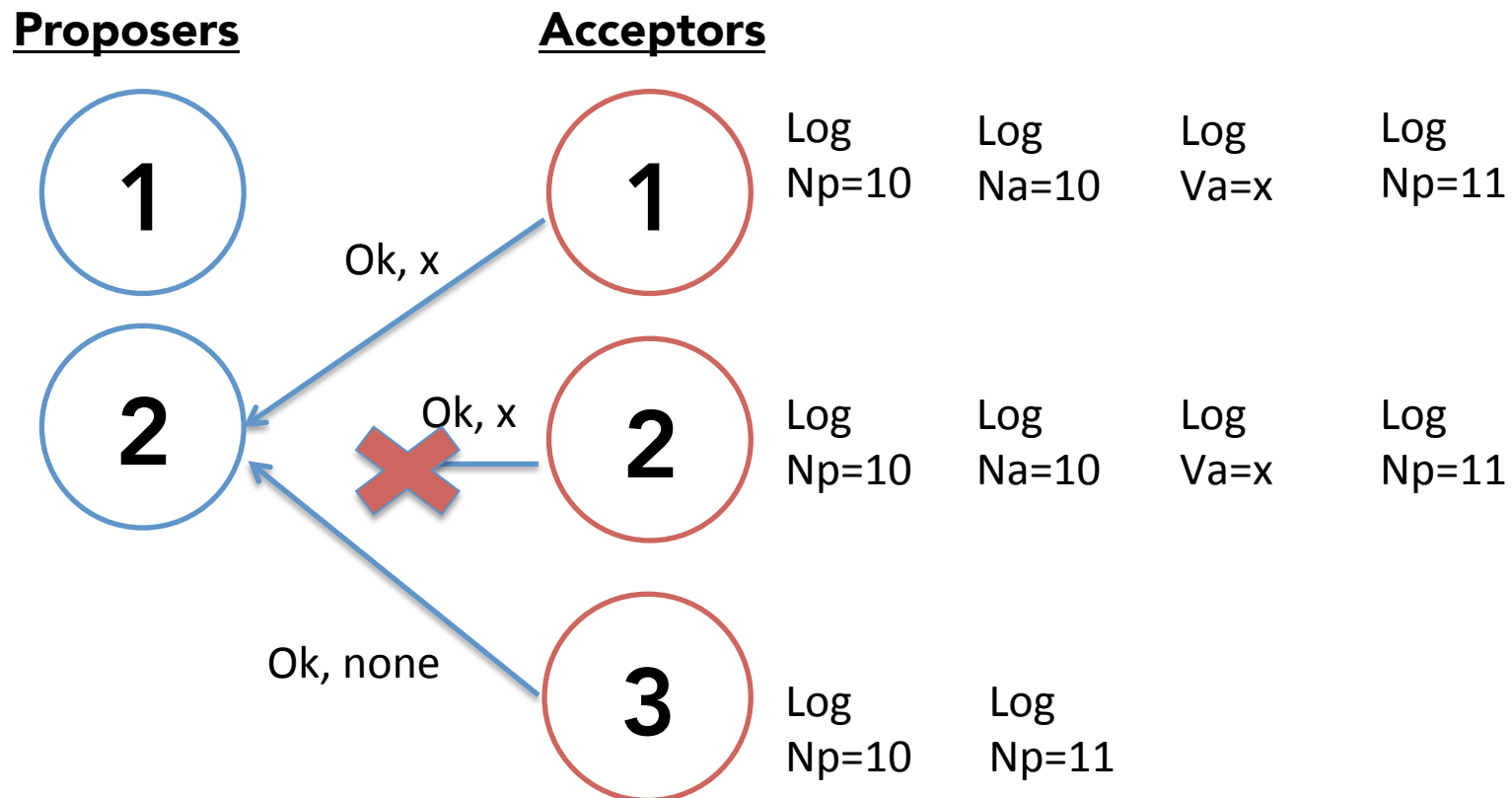
Log
Va=x

Log
Np=11

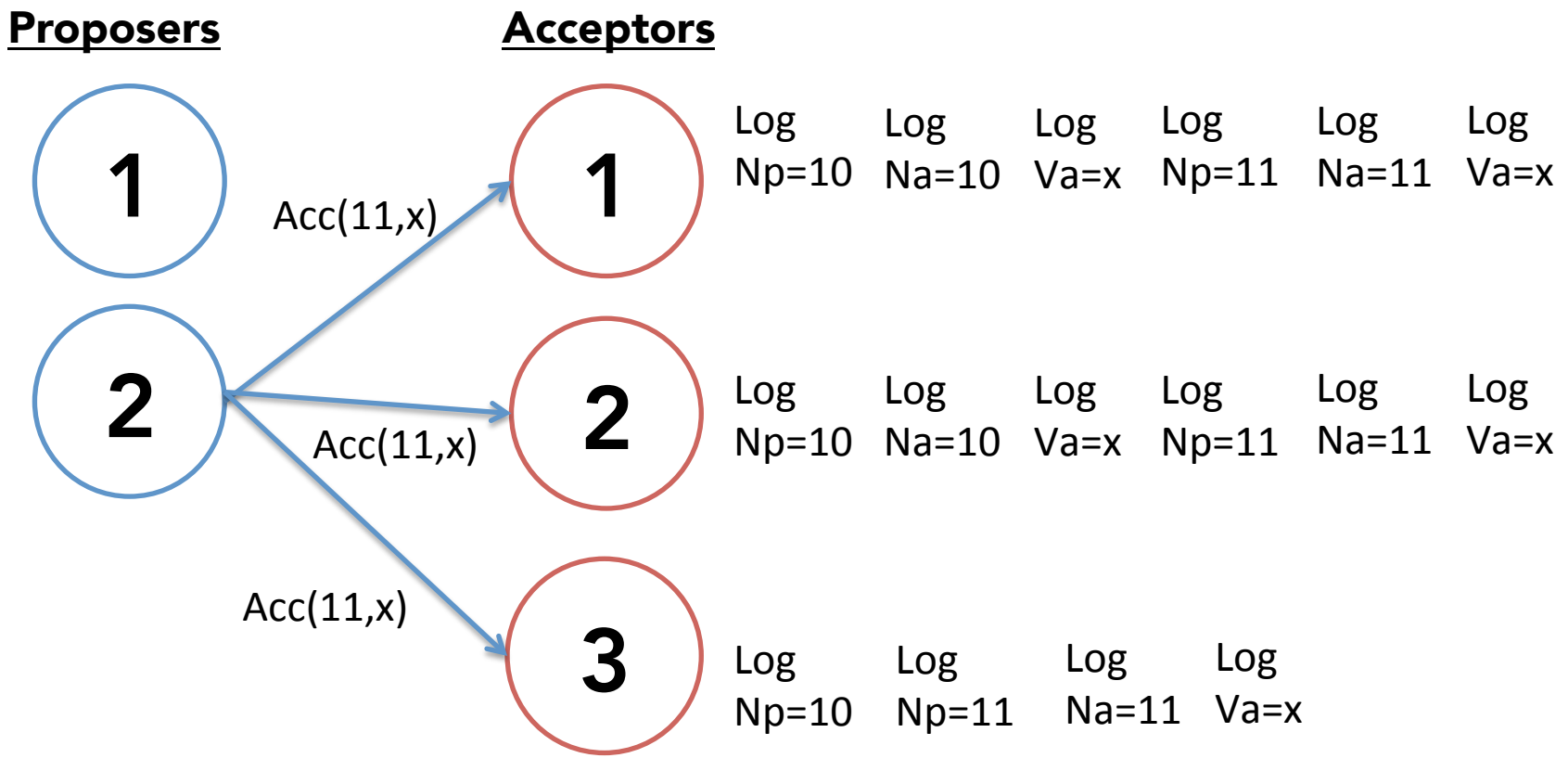
Log
Np=10

Log
Np=11

P2 learns value is x; has to propagate it even if hears from non-majority



P2 learns committed value is x, tells A3



Decided Message Omitted for Brevity

Summary

- Consistency: single-copy semantics
- Replicated state machines provide single-copy
 - Key issue: agreeing on order of operations
 - Hard case: network partition
- Paxos allows replicas to reach consensus, in presence of machine and network failures
 - Widely used in practice [Chubby, ZooKeeper, etc.]