# Byzantine Fault Tolerance

## CS 425: Distributed Systems

## Fall 2011

**Material drived from slides by I. Gupta and N.Vaidya**

# Reading List

- L. Lamport, R. Shostak, M. Pease, "The Byzantine Generals Problem," ACM ToPLaS 1982.

- M. Castro and B. Liskov, "Practical Byzantine Fault Tolerance," OSDI 1999.

# Byzantine Generals Problem

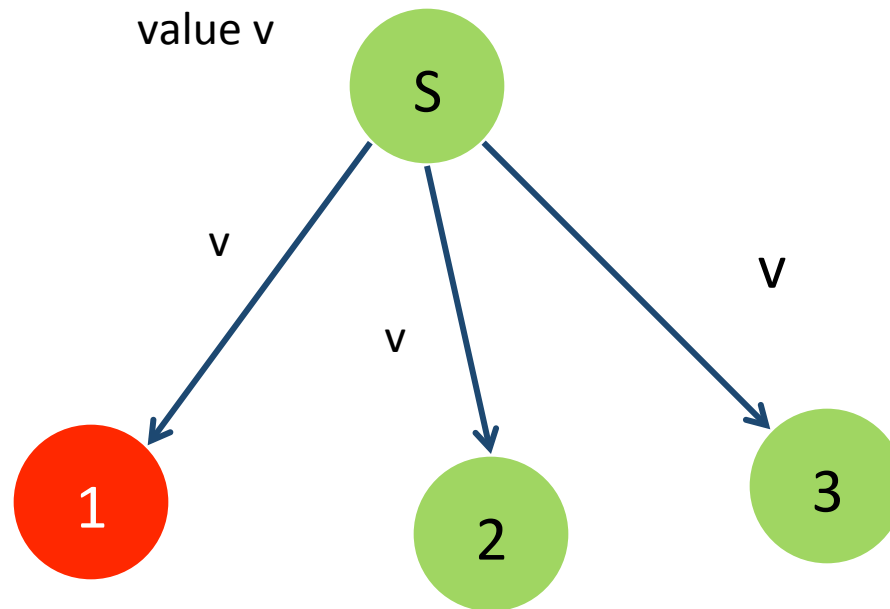A sender wants to send message to $n-1$ other peers

- Fault-free nodes must agree

- Sender fault-free  ➡  agree on its message

- Up to $f$ failures

# Byzantine Generals Problem

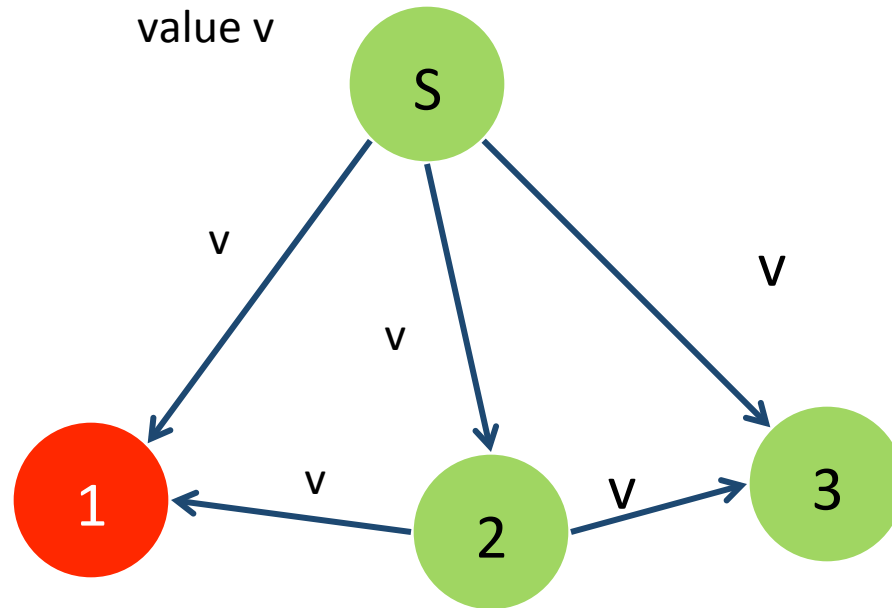A sender wants to send message to $n-1$ other peers

- Fault-free nodes must agree

- Sender fault-free ➜ agree on its message

- Up to $f$ failures
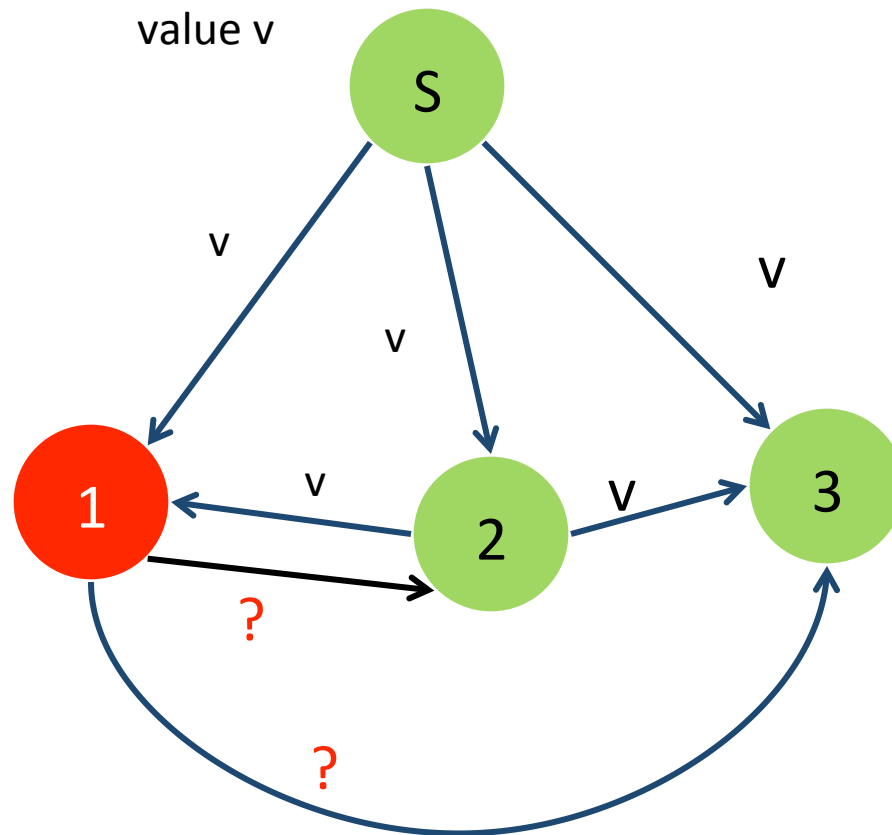
# Byzantine Generals Algorithm

value v
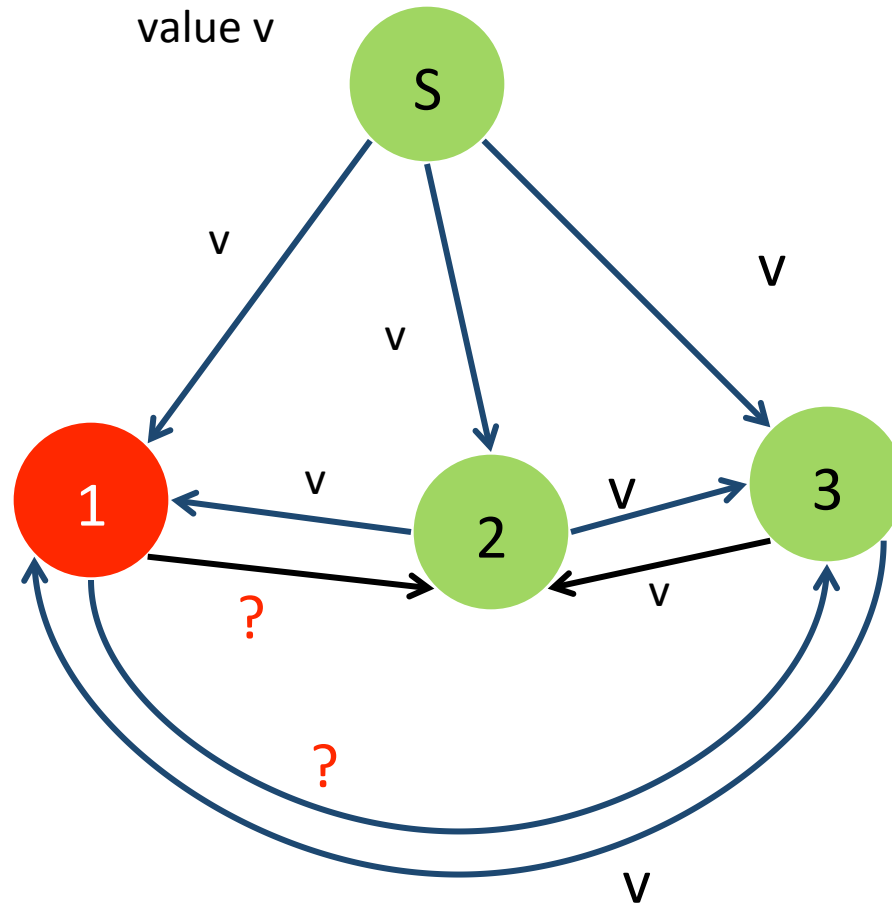
S

v

v

v

Faulty peer

1

2

3

# Byzantine Generals Algorithm

# Byzantine Generals Algorithm

# Byzantine Generals Algorithm

# Byzantine Generals Algorithm

# Byzantine Generals Algorithm

value v

S

Majority
vote results
in correct
result at
good peers

1

2

3

v

v

v

v

v

v

v

?

?

x

v

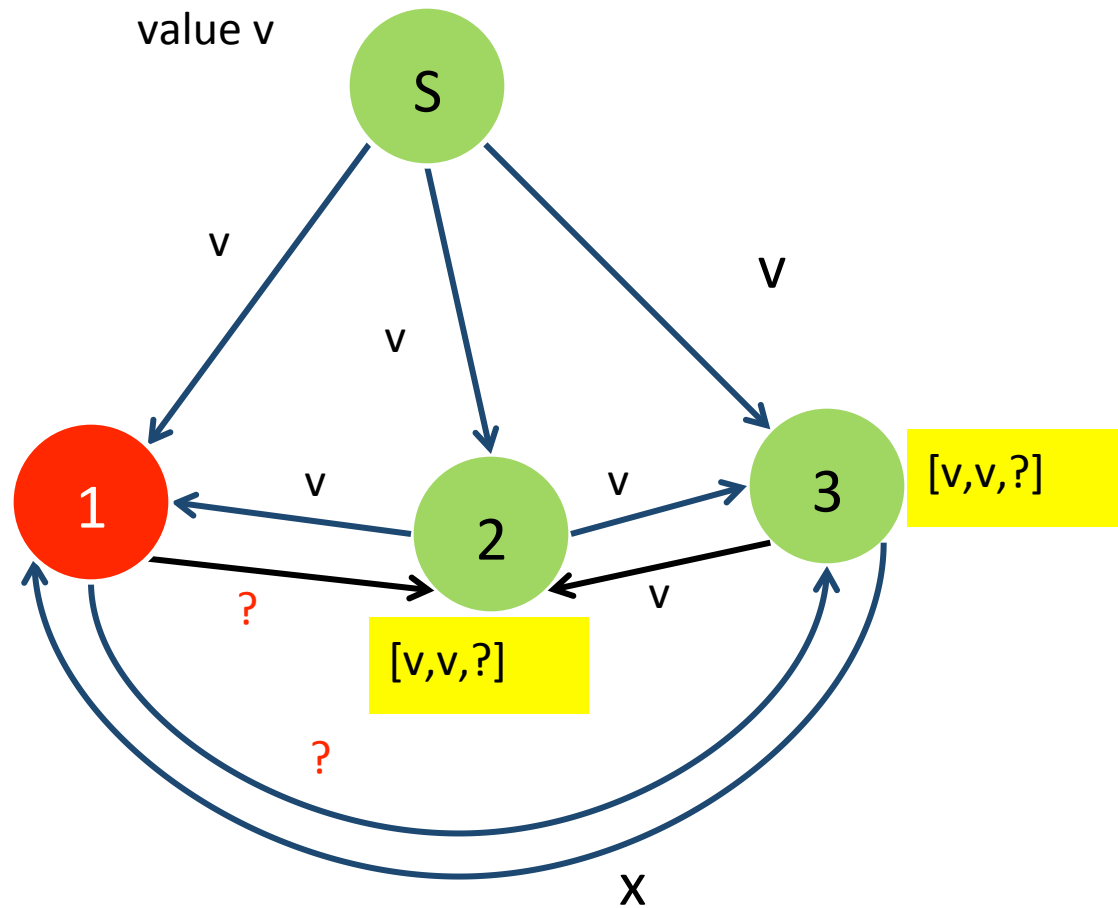v

# Byzantine Generals Algorithm

Faulty source

# Byzantine Generals Algorithm

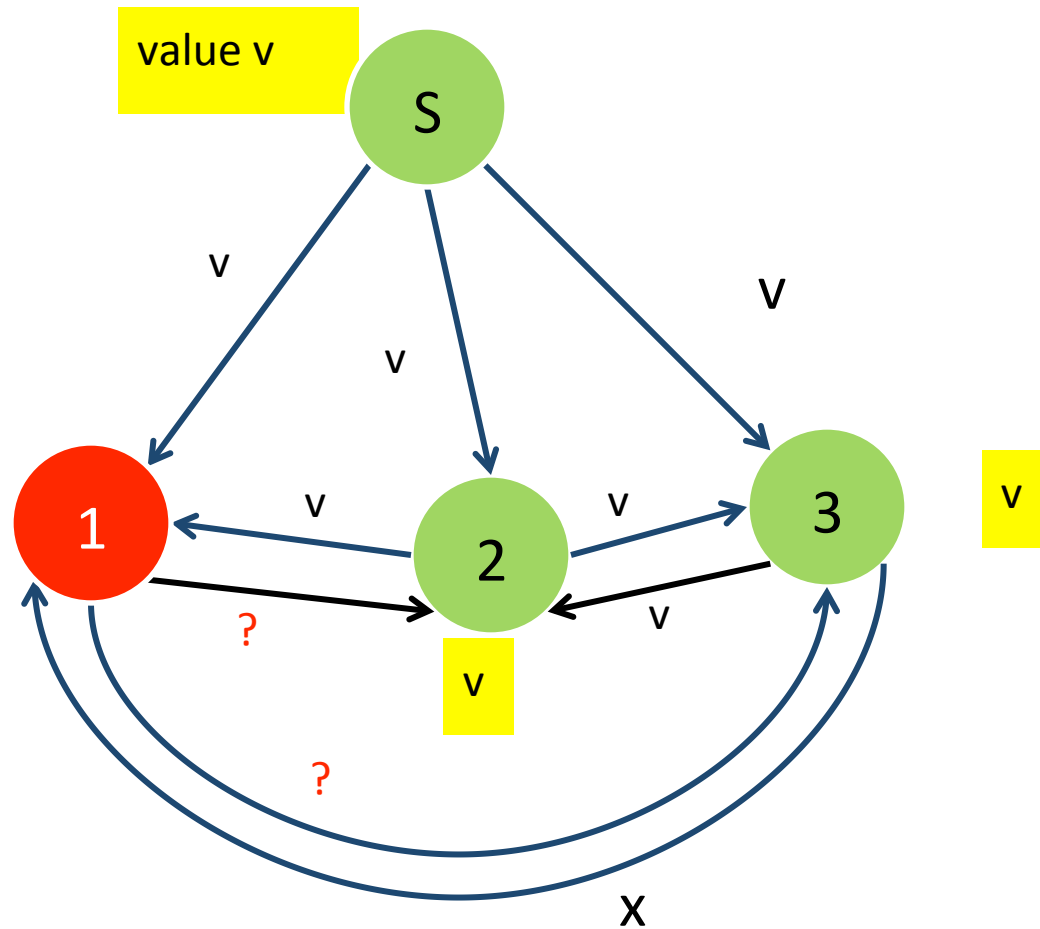# Byzantine Generals Algorithm

# Byzantine Generals Algorithm

# Byzantine Generals Algorithm



S

[v,w,x]

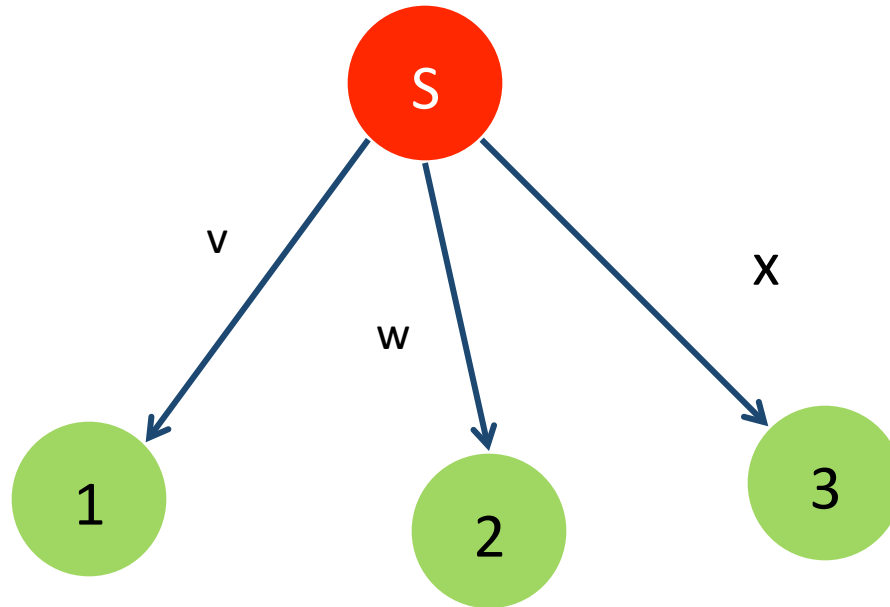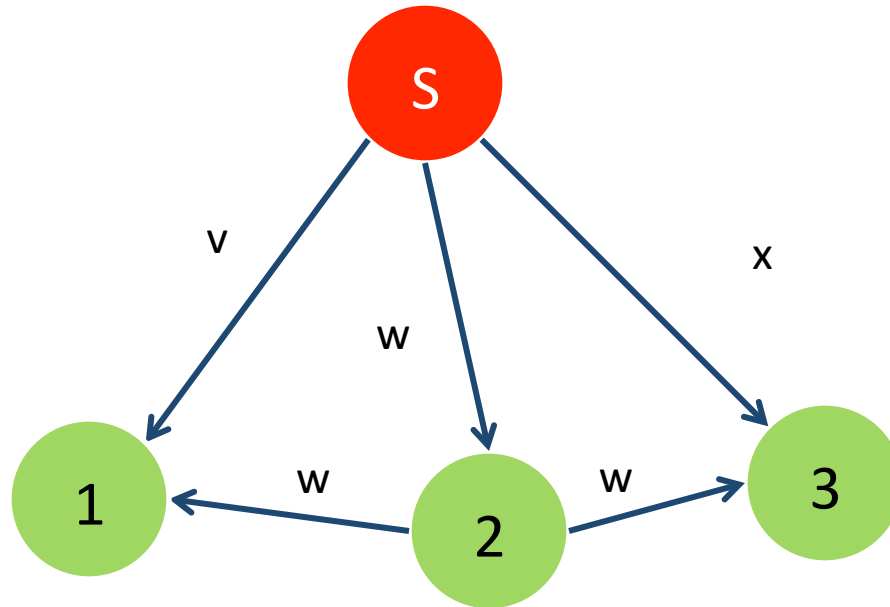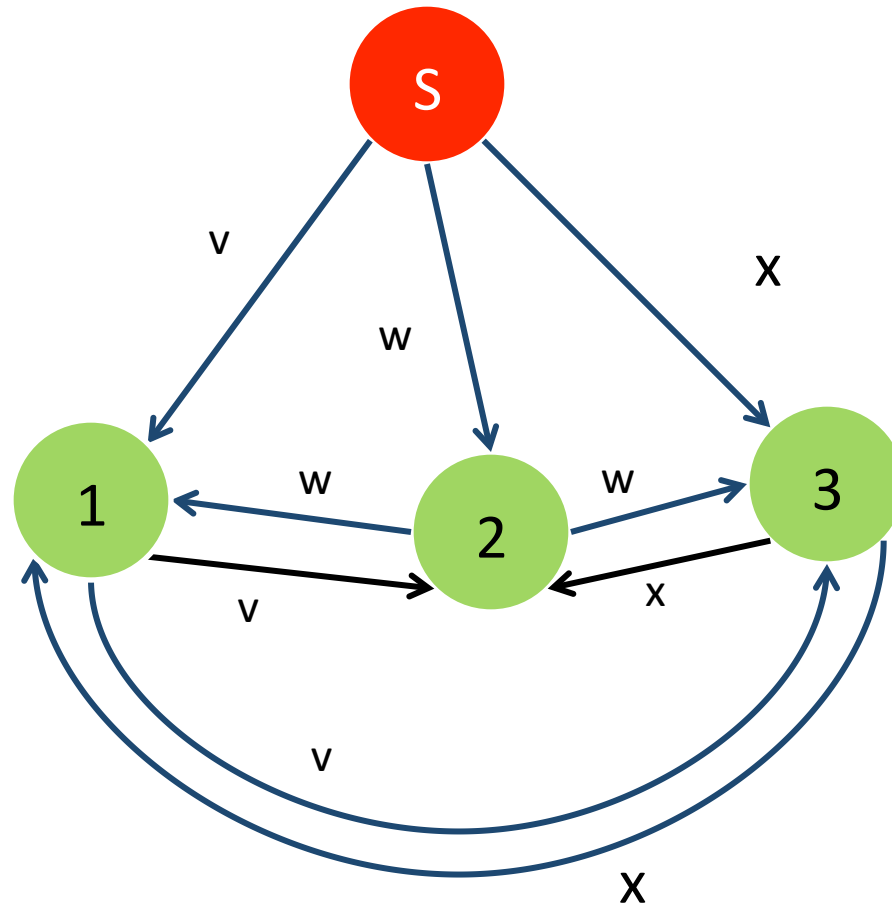1    w    2    w    3

[v,w,x]

v    x

[v,w,x]

v

x

Vote result identical at good peers

15

# Known Results

- Need 3f + 1 nodes to tolerate f failures

- Need $\Omega(n^2)$ messages in general

# $\Omega(n^2)$ Message Complexity

- Each message at least 1 bit

- $\Omega(n^2)$ bits "communication complexity" to agree on just 1 bit value

# Practical Byzantine Fault Tolerance

- Computer systems provide crucial services
- Computer systems fail
  - Crash-stop failure
  - Crash-recovery failure
  - Byzantine failure
- Example: natural disaster, malicious attack, hardware failure, software bug, etc.
- Need highly available service

Replicate to increase availability

# Challenges

Request A

Request B

Client

Client

# Requirements

- All replicas must handle same requests despite failure.
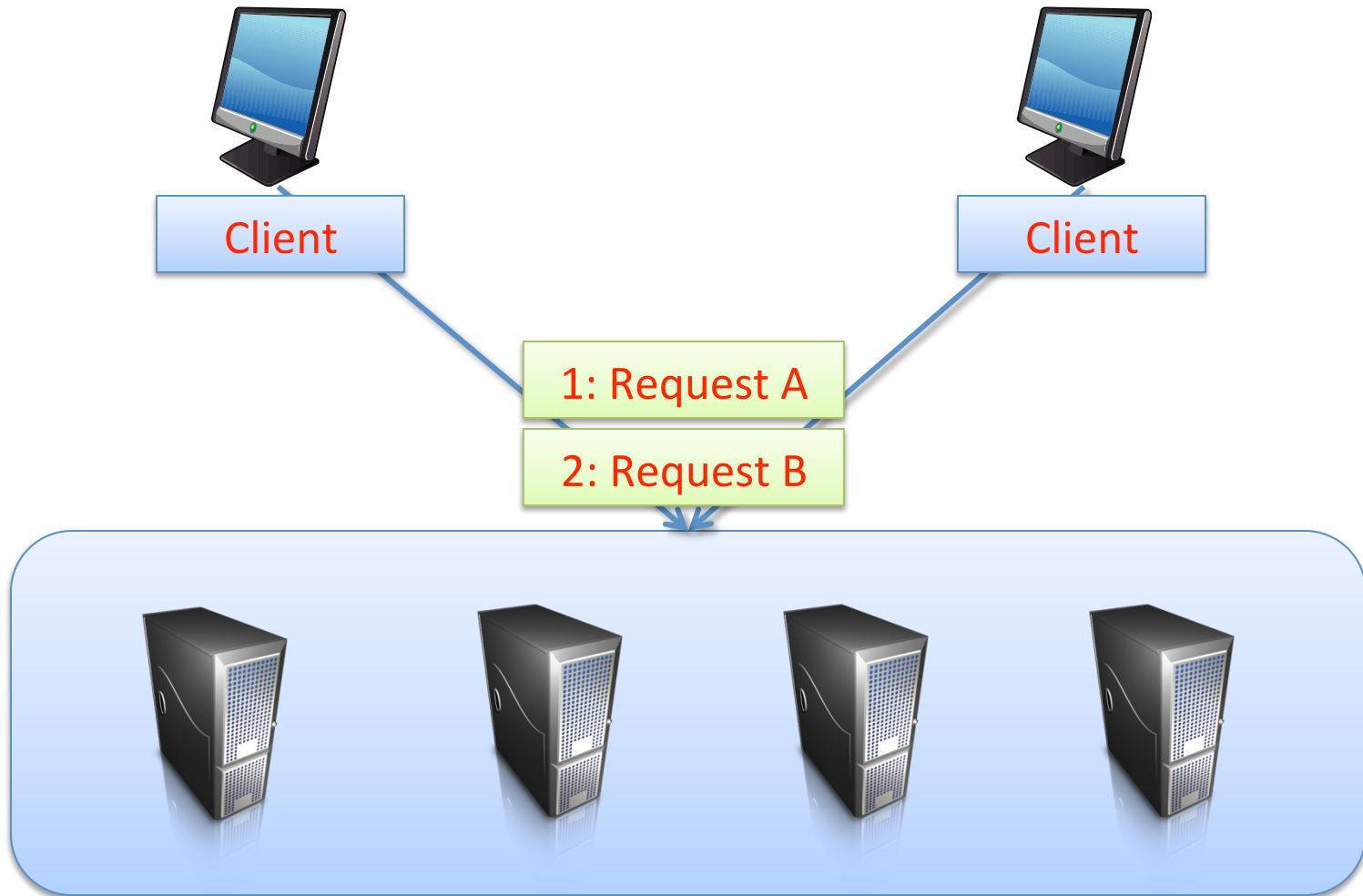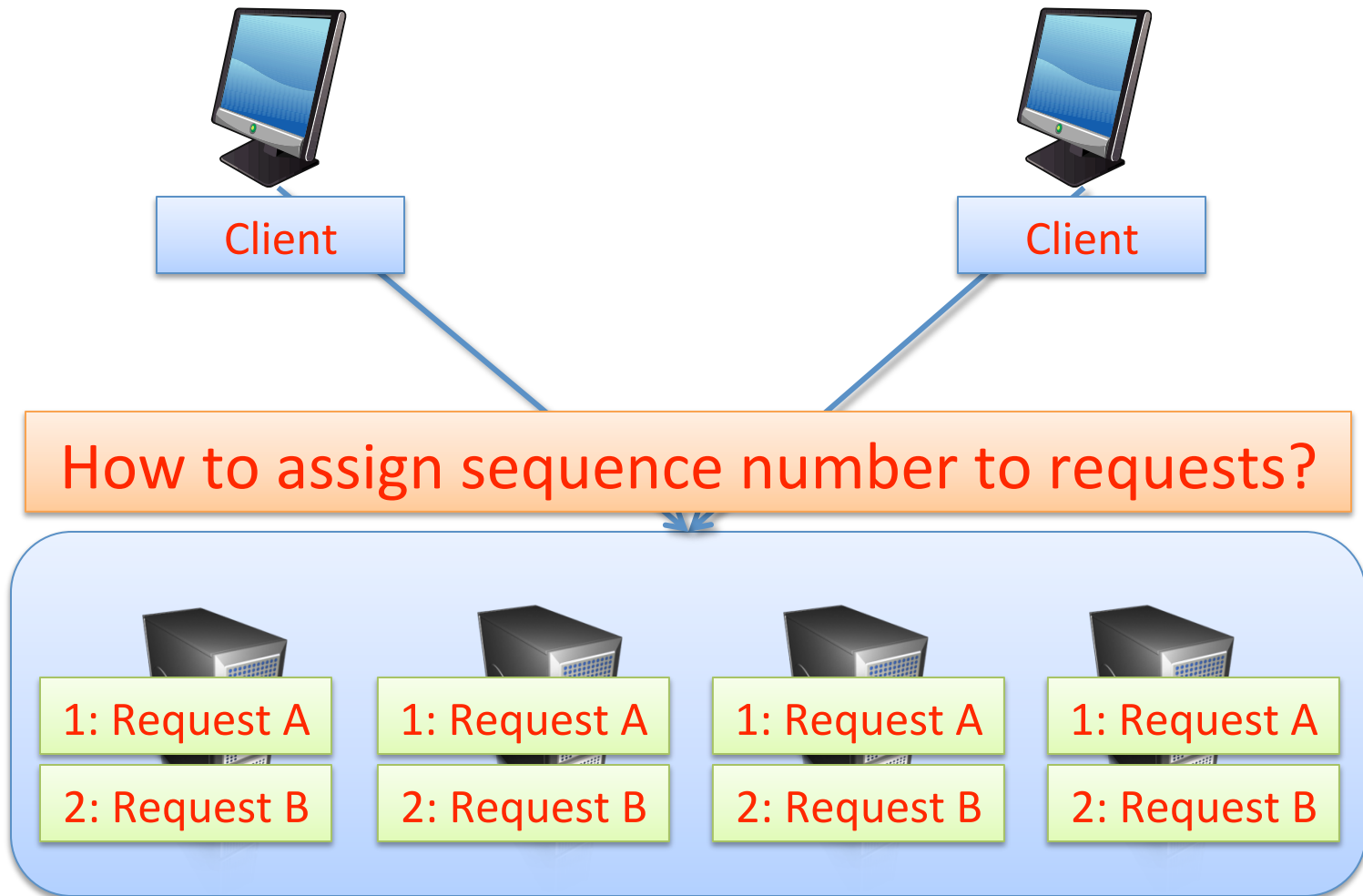
- Replicas must handle requests in identical order despite failure.

# Challenges



Client

Client

1: Request A

2: Request B

# State Machine Replication



Client

Client

How to assign sequence number to requests?

1: Request A

1: Request A

1: Request A

1: Request A

2: Request B

2: Request B

2: Request B

2: Request B

# Primary Backup Mechanism



Client

Client

**What if the primary is faulty?**
Agreeing on sequence number
Agreeing on changing the primary (view change)

1: Request A

2: Request B

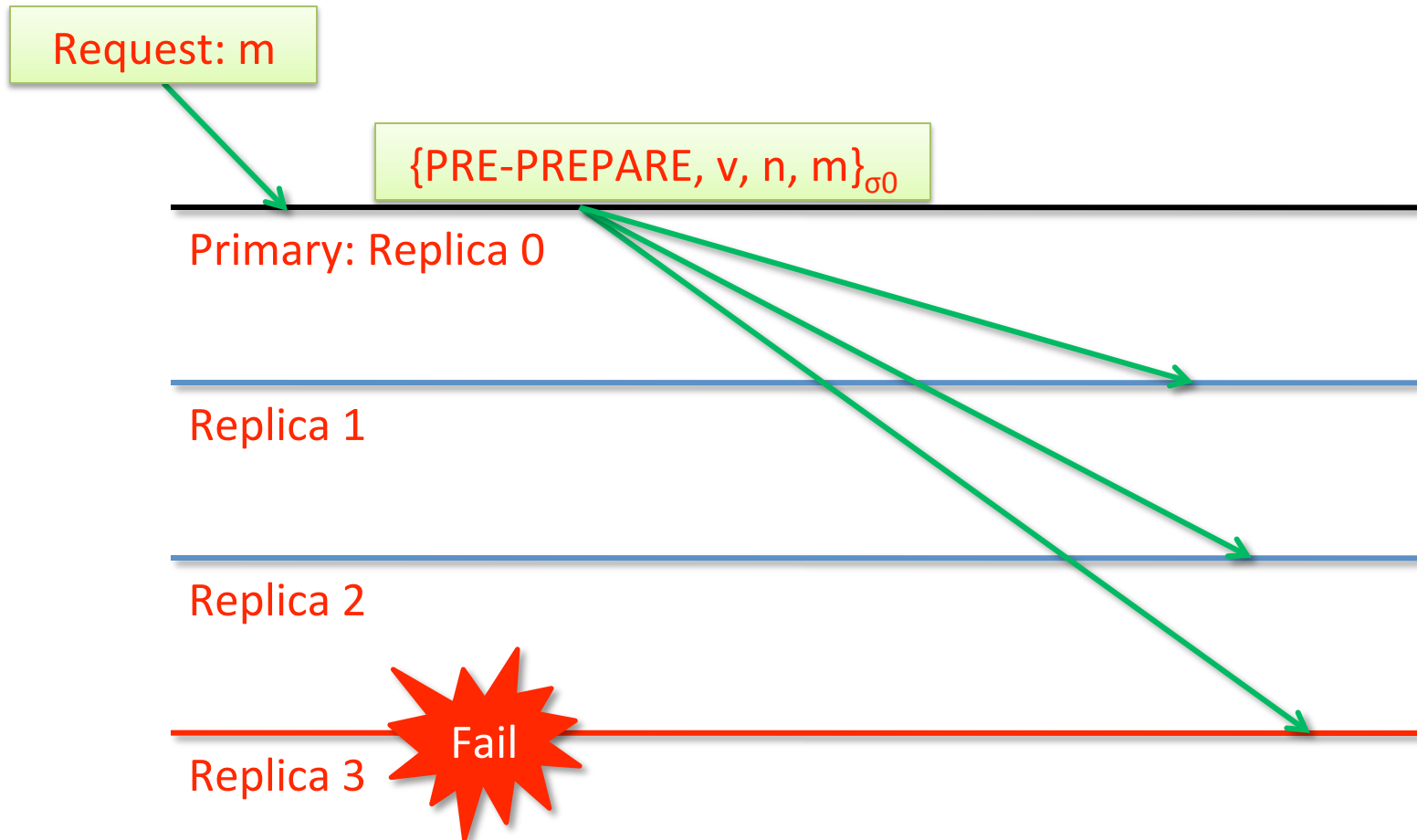View 0

# Normal Case Operation
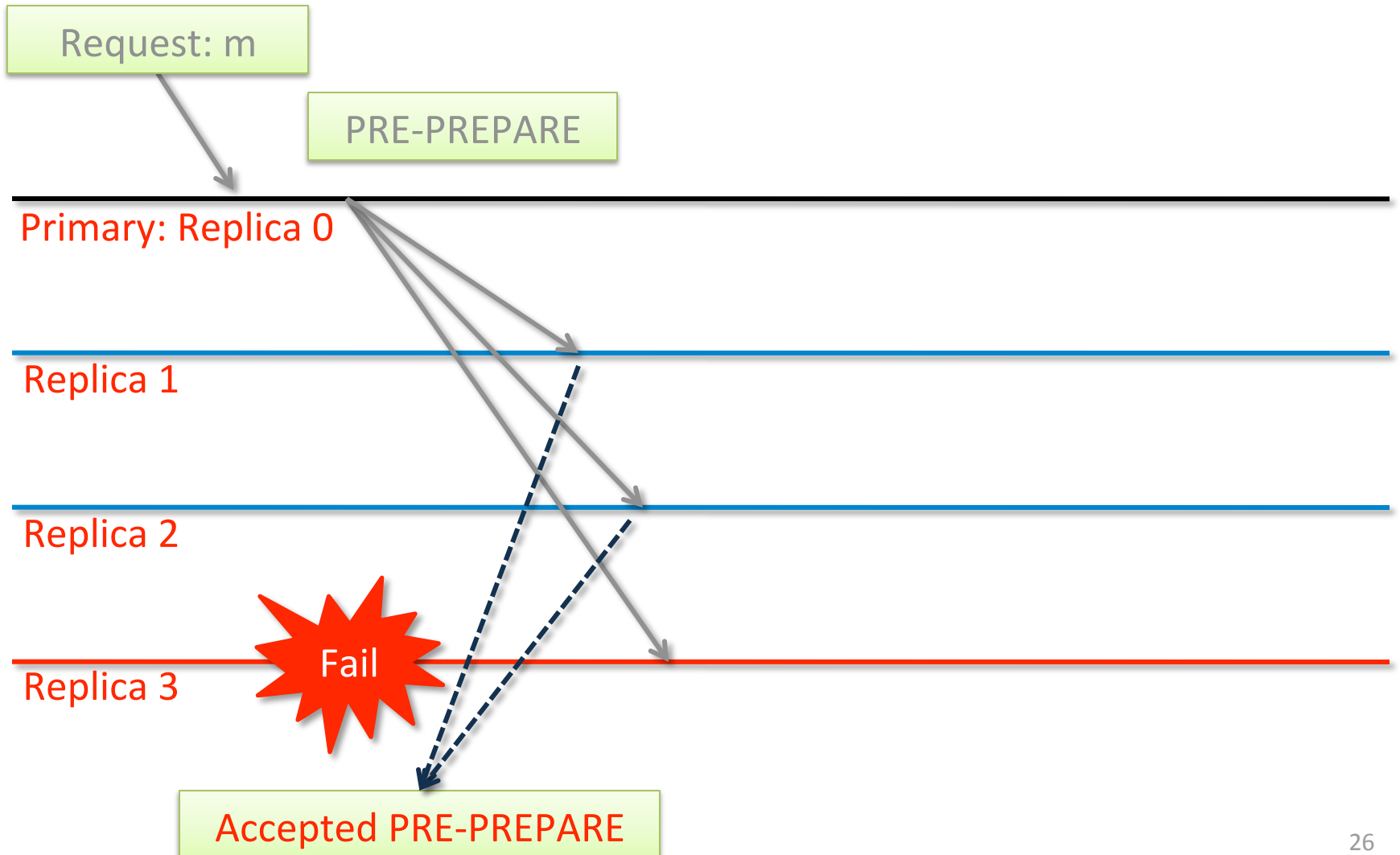
- Three phase algorithm:
  - PRE-PREPARE picks order of requests
  - PREPARE ensures order within views
  - COMMIT ensures order across views
- Replicas remember messages in log
- Messages are authenticated
  - $\{.\}_{\sigma k}$ denotes a message sent by k

# Pre-prepare Phase

Request: m

$\{\text{PRE-PREPARE}, v, n, m\}_{\sigma 0}$

Primary: Replica 0

Replica 1

Replica 2

Fail

Replica 3

# Prepare Phase

Request: m

PRE-PREPARE

Primary: Replica 0

Replica 1

Replica 2

Fail

Replica 3

Accepted PRE-PREPARE

# Prepare Phase



Request: m

PRE-PREPARE

Primary: Replica 0

Replica 1

{PREPARE, v, n, D(m), 1}$_{\sigma 1}$

Replica 2

Fail

Replica 3

Accepted PRE-PREPARE

# Prepare Phase



Request: m

Collect PRE-PREPARE + 2f matching PREPARE

PRE-PREPARE

Primary: Replica 0

$\{PREPARE, v, n, D(m), 1\}_{\sigma 1}$

Replica 1

Replica 2

Fail

Replica 3

Accepted PRE-PREPARE

# Commit Phase



Request: m

PRE-PREPARE

PREPARE

Primary: Replica 0

Replica 1

Replica 2

$\{COMMIT, v, n, D(m)\}_{\sigma 2}$

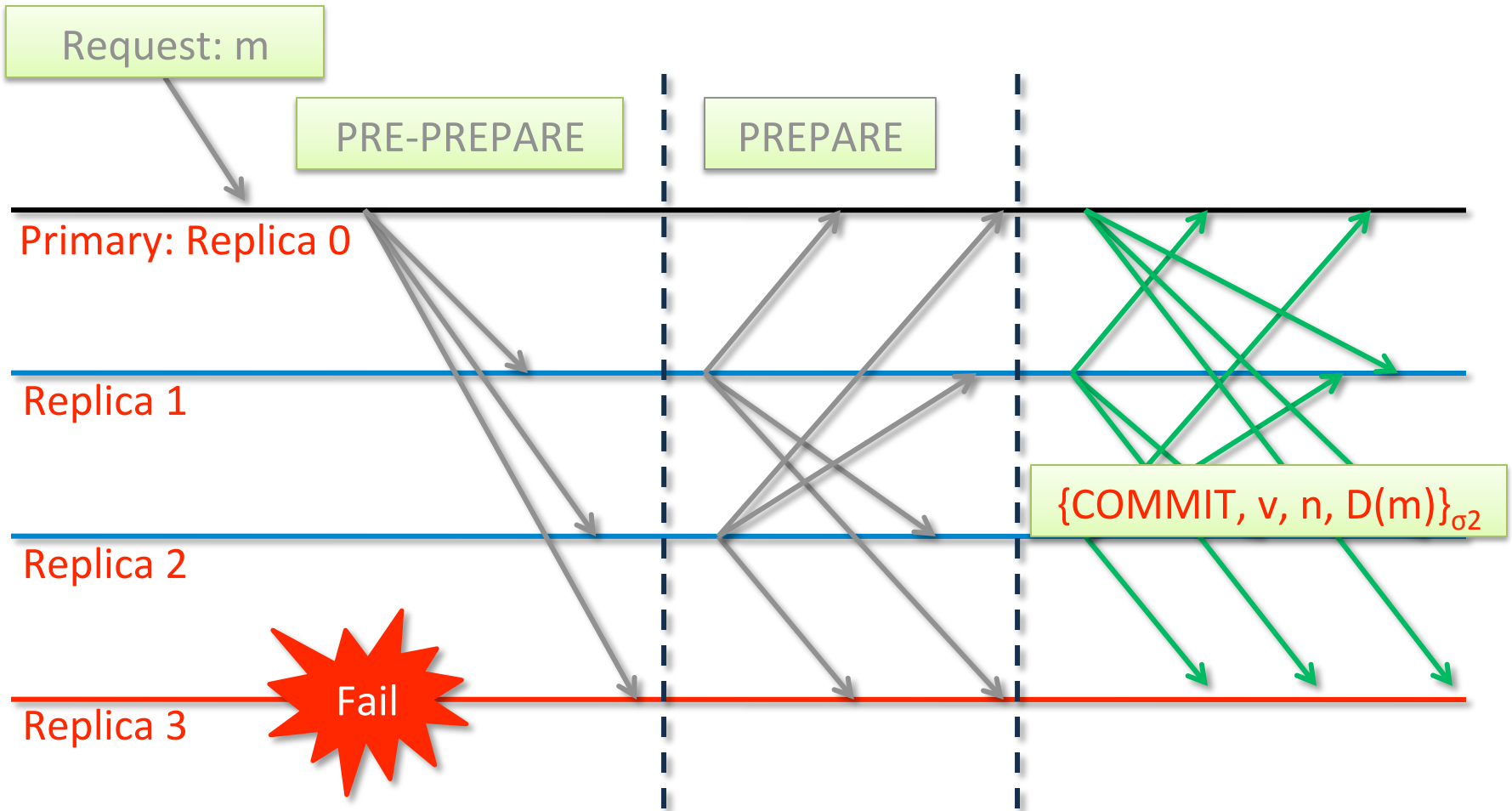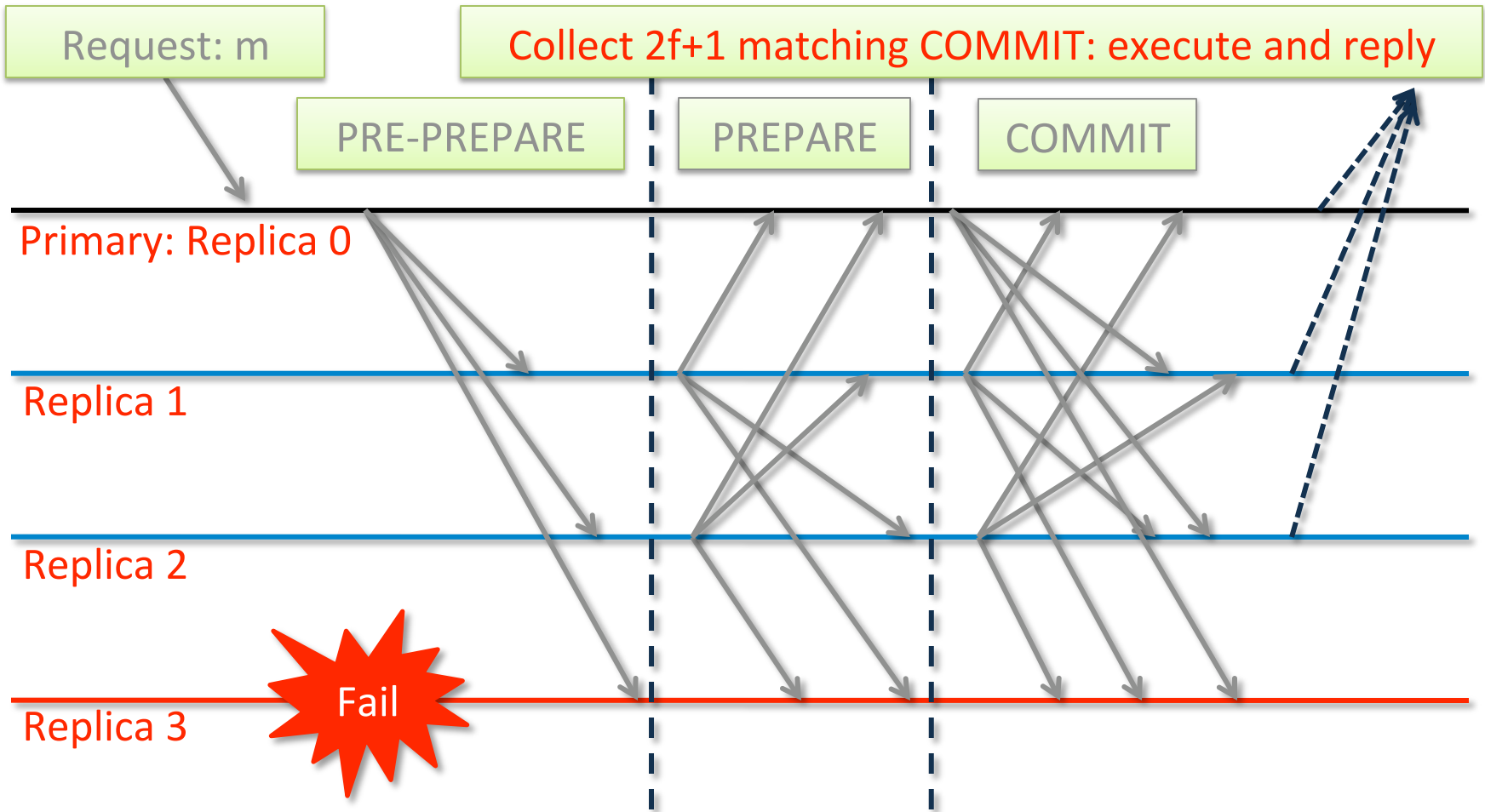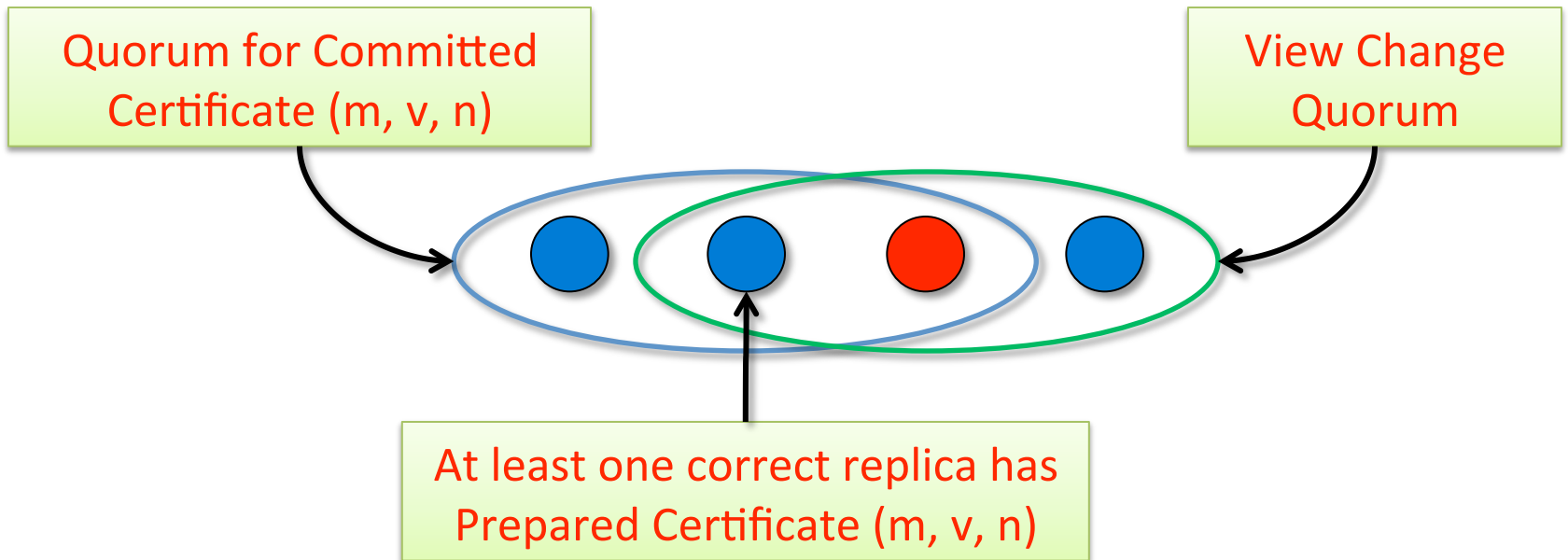Fail

Replica 3

# Commit Phase (2)

# View Change

- Provide liveness when primary fails
  - Timeouts trigger view changes
  - Select new primary (= view number mod 3f+1)
- Brief protocol
  - Replicas send VIEW-CHANGE message along with the requests they prepared so far
  - New primary collects *2f+1* VIEW-CHANGE messages
  - Constructs information about committed requests in previous views

# View Change Safety

- Goal: No two different committed request with same sequence number across views



Quorum for Committed Certificate (m, v, n)

View Change Quorum

At least one correct replica has Prepared Certificate (m, v, n)

# Related Works