

# HotStuff: BFT Consensus with Linearity and Responsiveness

Yin, M., Malkhi, D., Reiter, M. K., Gueta, G. G., & Abraham, I.

PODC '19

# Contents

1. Contributions
2. Prior works
3. System model
4. Review: 2-phase paradigm (PBFT, Tendermint)
5. Basic HotStuff
6. Chained HotStuff

# Contributions

- **Linear view change**
  - Use a threshold signature scheme of which threshold is  $2f + 1$ .
- **Optimistic Responsiveness**
  - *Responsiveness* requires that a non-faulty leader can drive the protocol to consensus in time depending only on the actual message delays, independent of any known upper bound on message transmission delays
  - *Optimistic* means that *responsiveness* is required only in beneficial circumstances—here, after GST is reached.

# Prior works

Protocol	Authenticator complexity			Responsiveness
	<i>Correct leader</i>	<i>Leader failure (view-change)</i>	<i>f leader failures</i>	
DLS [25]	$O(n^4)$	$O(n^4)$	$O(n^4)$	
PBFT [20]	$O(n^2)$	$O(n^3)$	$O(fn^3)$	✓
SBFT [30]	$O(n)$	$O(n^2)$	$O(fn^2)$	✓
Tendermint [15] / Casper [17]	$O(n^2)$	$O(n^2)$	$O(fn^2)$	
Tendermint* / Casper*	$O(n)$	$O(n)$	$O(fn)$	
<b>HotStuff</b>	$O(n)$	$O(n)$	$O(fn)$	✓

Signatures can be combined using threshold signatures, though this optimization is not mentioned in their original works.

**Table 1: Performance of selected protocols after GST.**

# System Model (1)

- **Security model**
  - An *adversary* coordinates Byzantine replicas and learns all internal states held by these replicas
- **Network assumption**
  - Communication is p2p, authenticated and reliable.
  - “broadcast” involves the sender.
  - Partial synchrony model.

# System Model (2)

- **Threshold Signature**

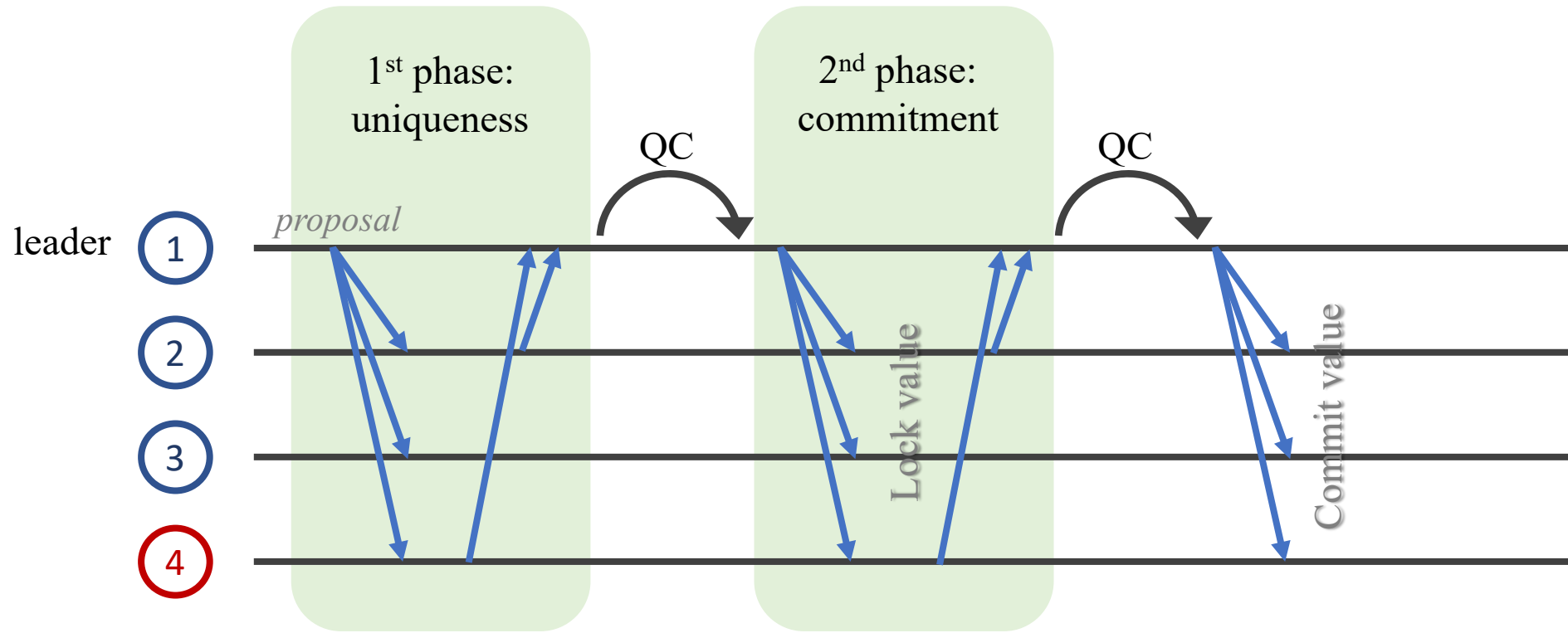
- A single public key held by all replicas, and each of the  $n$  replicas holds a distinct private key.
- The  $i$ -th replica can use its private key to contribute a partial signature  $\rho_i \leftarrow \mathbf{tsign}_i(m)$  on message  $m$ .
- $\sigma \leftarrow \mathbf{tcombine}(m, \{\rho_i\}_{i \in I})$  on message  $m$ , where  $|I| = \mathit{threshold} = 2f + 1$ .
- Any other replica can verify the signature using the public key and the function  $\mathbf{tverify}(\sigma, m)$

- **Authenticator complexity**

- an *authenticator* is either a partial signature or a signature.
- Sum of the number of *authenticators* received by replica  $i$  in the protocol to reach a consensus decision.



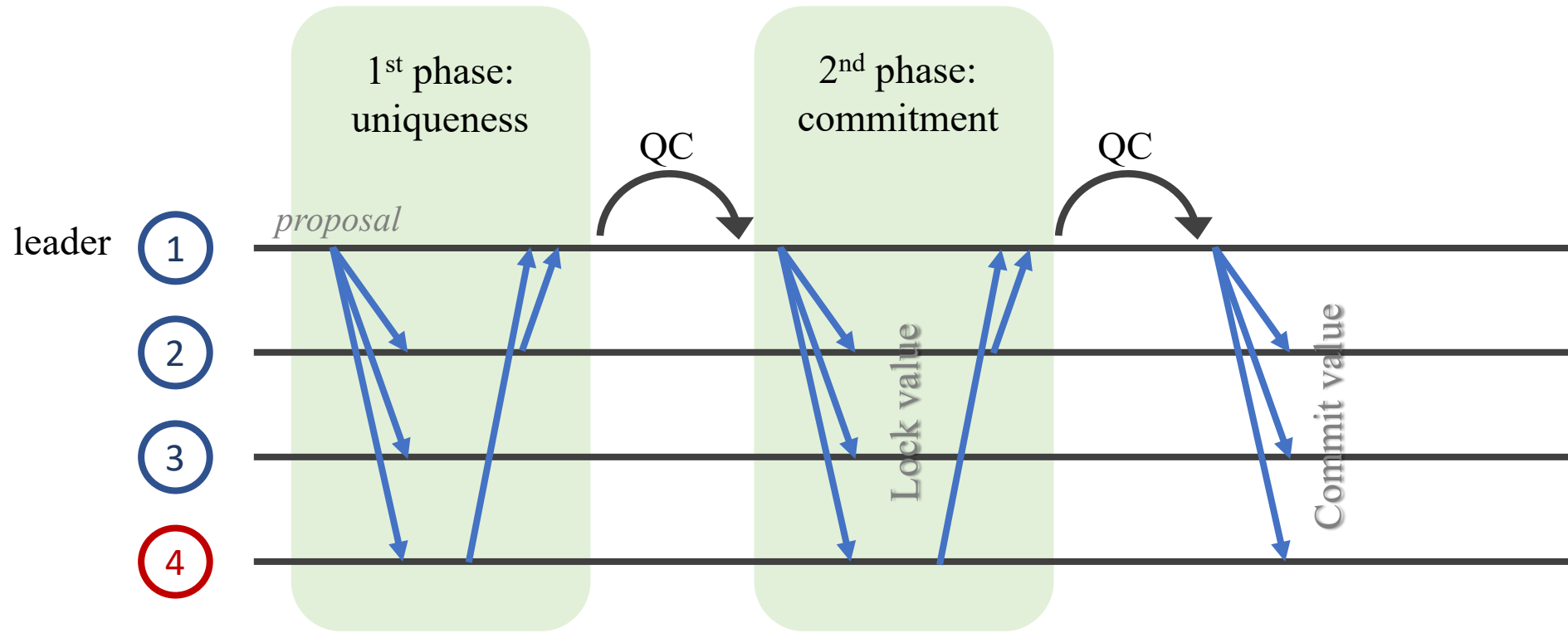
# Review: 2-phase paradigm



- **1<sup>st</sup> phase** : guarantees proposal uniqueness through the formation of a quorum certificate (QC) consisting of  $(n - f)$  votes. validity condition is important.
- **2<sup>nd</sup> phase** : guarantees that the next leader can convince replicas to vote for a safe proposal, using *commit-adopt*\*



# Review: 2-phase paradigm (cont.)

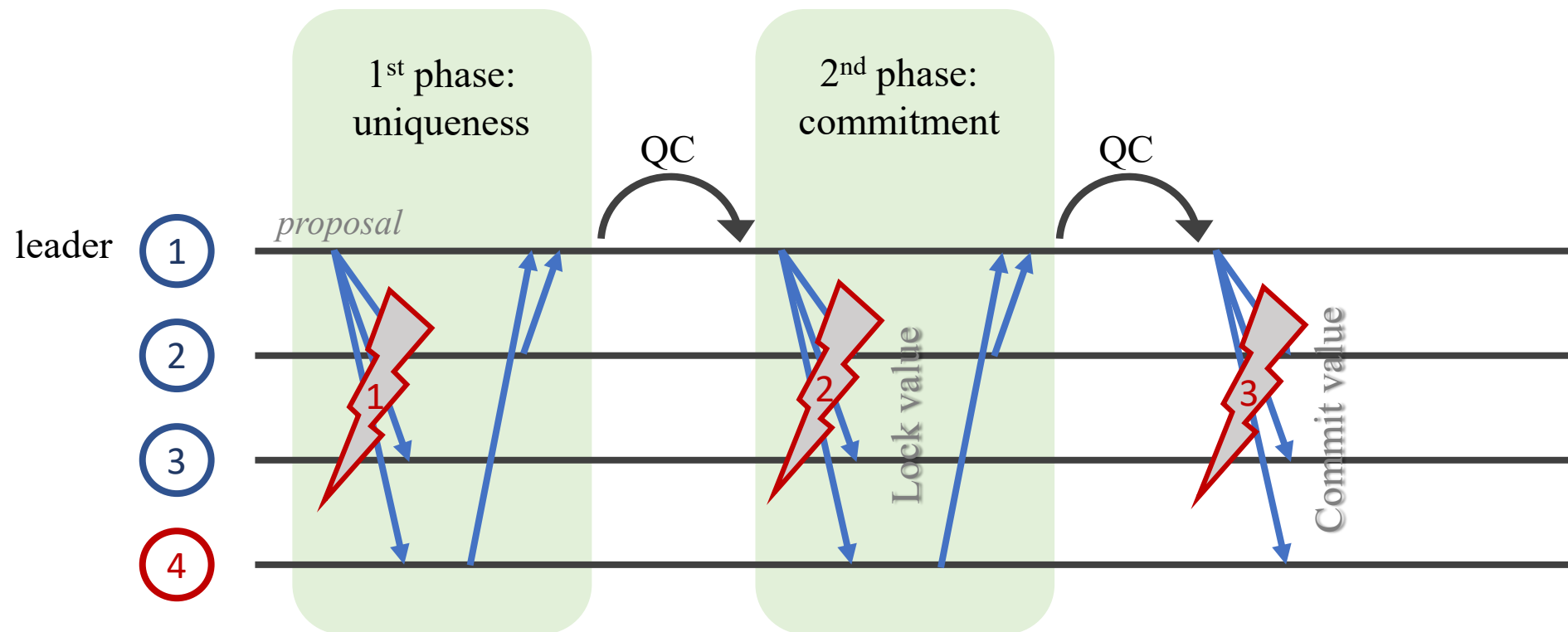


- ***Commit-Adopt***: if a correct replica commits a value  $v$ , at least  $2f + 1$  replicas adopt and protect  $v$  (lock the value  $v$ ). At least  $f + 1$  replicas are correct among those replicas, and the corresponding  $f + 1$  replicas guard safety of the committed value  $v$ . Thus, they would not vote on a conflicting value unless they are shown a proof that it is safe to do so.





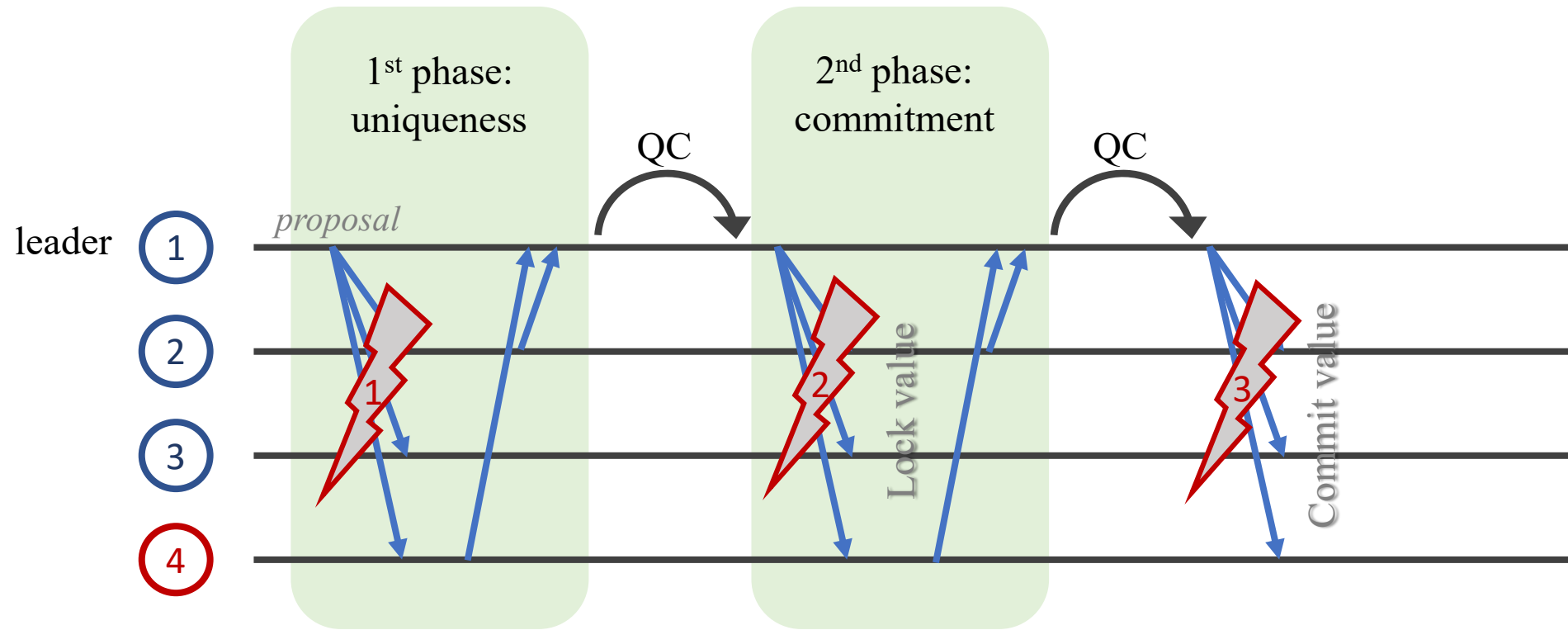
# Review: failure cases



1. No replica has a lock, no value is committed.
2. One or a few replicas are locked on a value, but no correct replica has committed.
3.  $2f + 1$  parties are locked, some correct replicas have committed but not all



# Review: failure cases

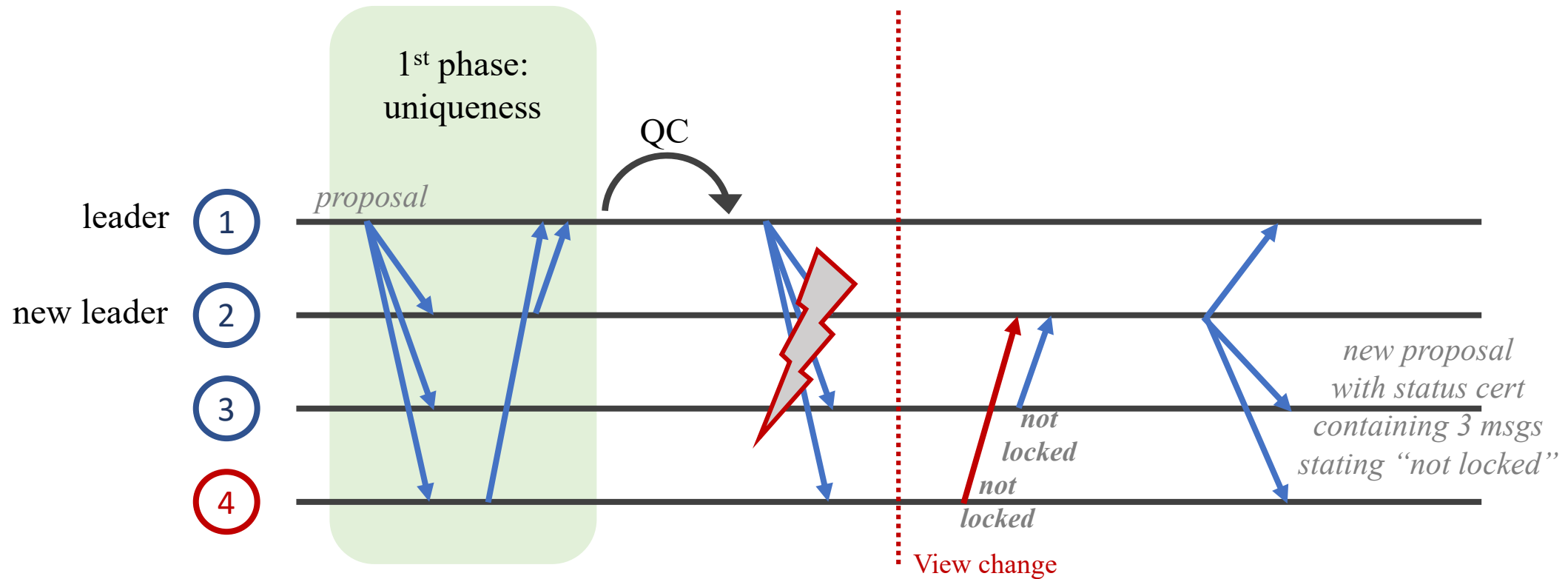


- ~~1. No replica has a lock, no value is committed.~~
2. One or a few replicas are locked on a value, but no correct replica has committed.
- ~~3.  $2f + 1$  parties are locked, some correct replicas have committed but not all~~

## Review: failure cases

- In case 2, some mechanism is required to ensure that these locked correct replicas vote for a safe proposal from a leader (perhaps conflicting with a locked value) after a view-change.
  1. What does the leader learn about the status of the system (the three scenarios) at the start of a view?
  2. How does the leader convince other replicas about the status of the system (and thus to vote for its proposal)?

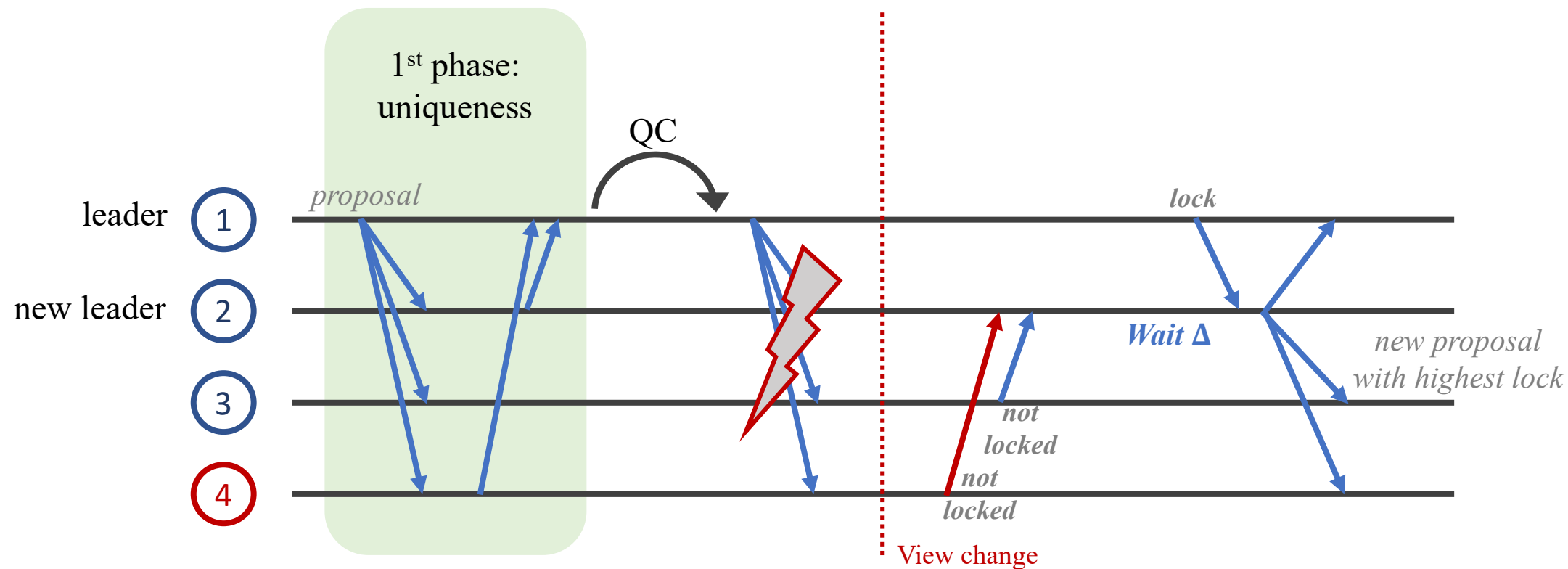
# Review: PBFT



- Leader sends the *status certificate* containing the  $2f + 1$  locks in its proposal.
- Quadratic complexity / Responsiveness

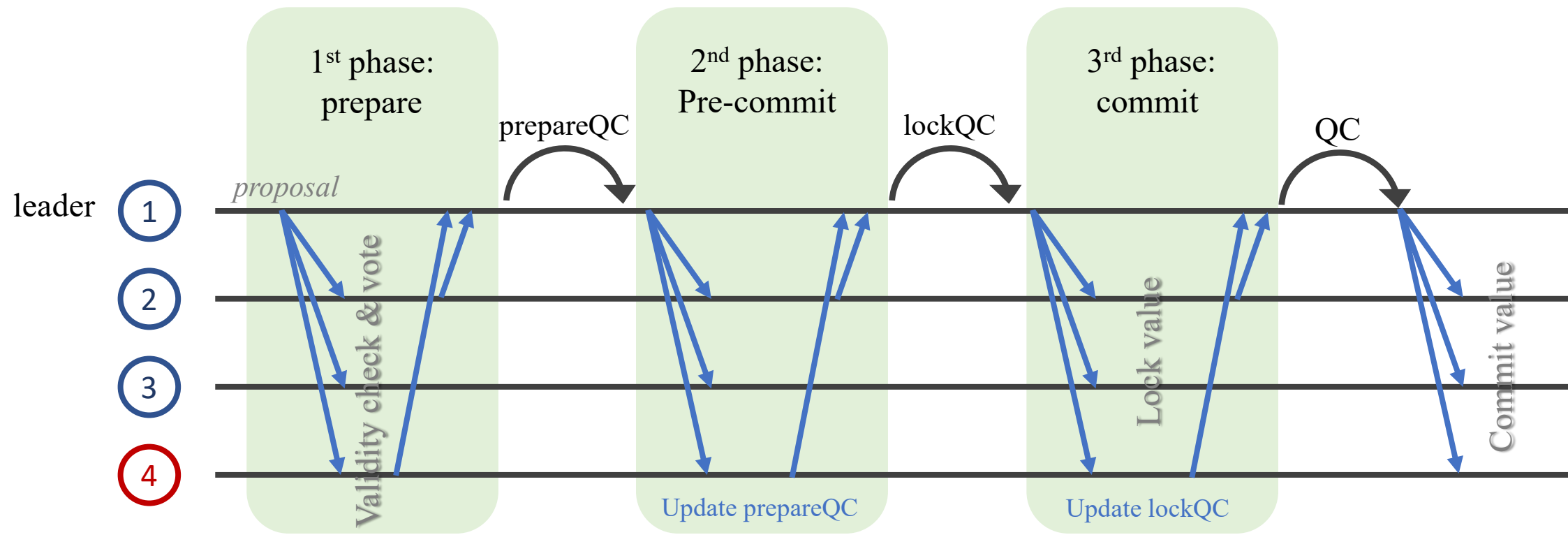


# Review: Tendermint



- Avoid sending a  $(2f + 1)$ -sized status certificate → linear view-change
- During  $\Delta$ , the leader collects all locks and make new proposal with highest lock.
- Linear view-change / Non-responsiveness

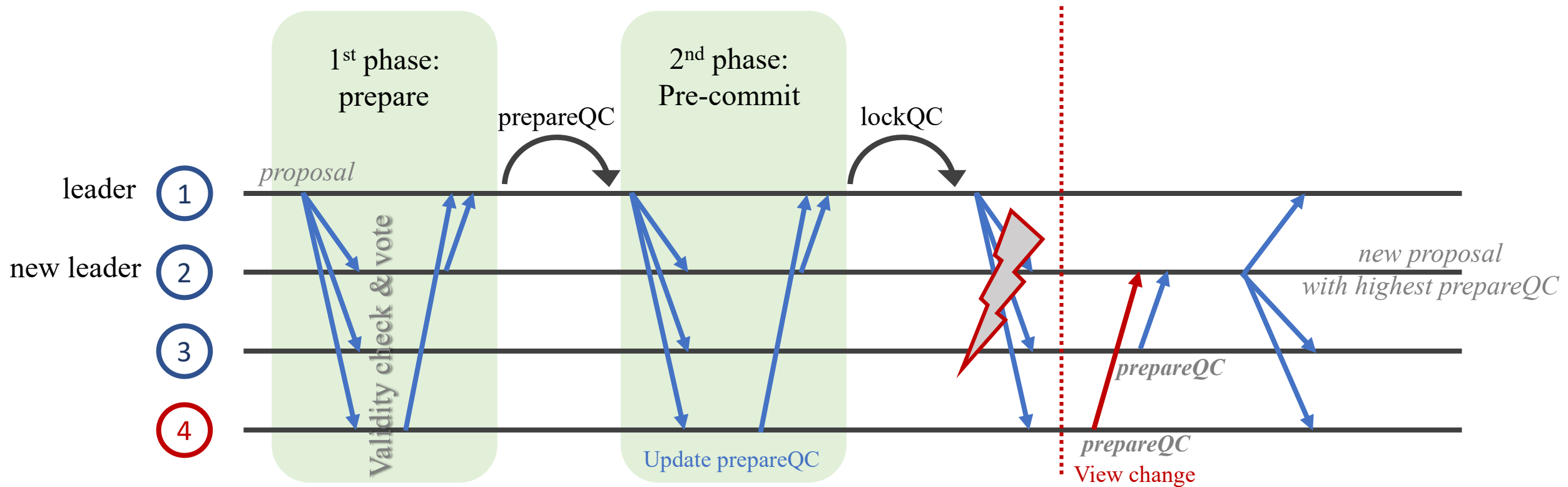
# Basic HotStuff



- prepareQC : qualified certificate for  $2f + 1$  validity votes.
- lockQC: qualified certificate for  $2f + 1$  locks on a proposal.

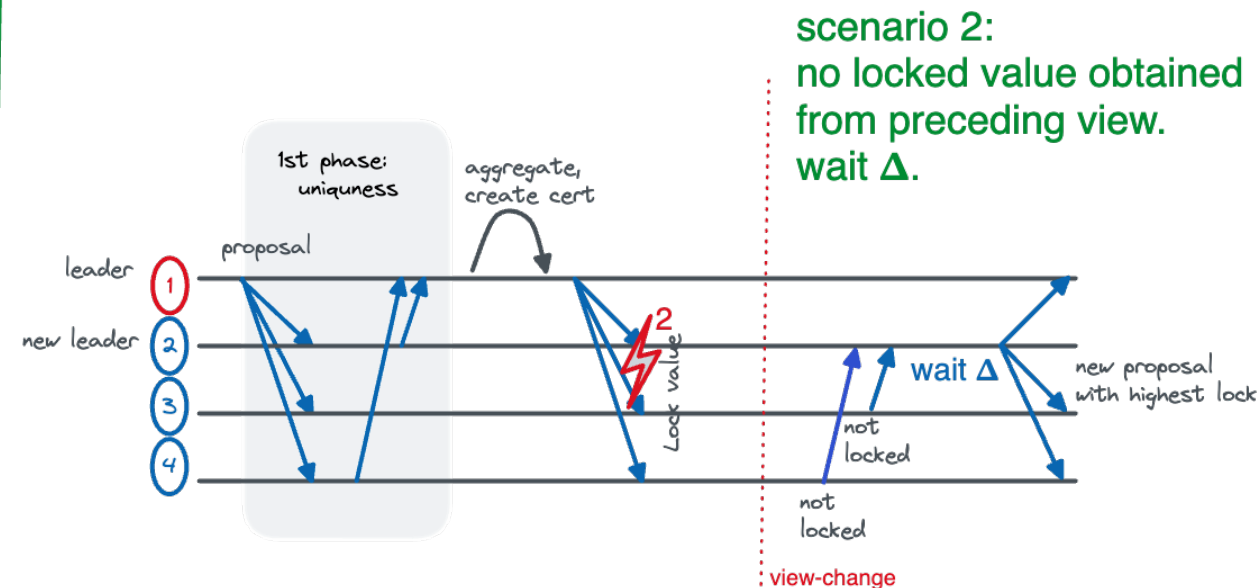
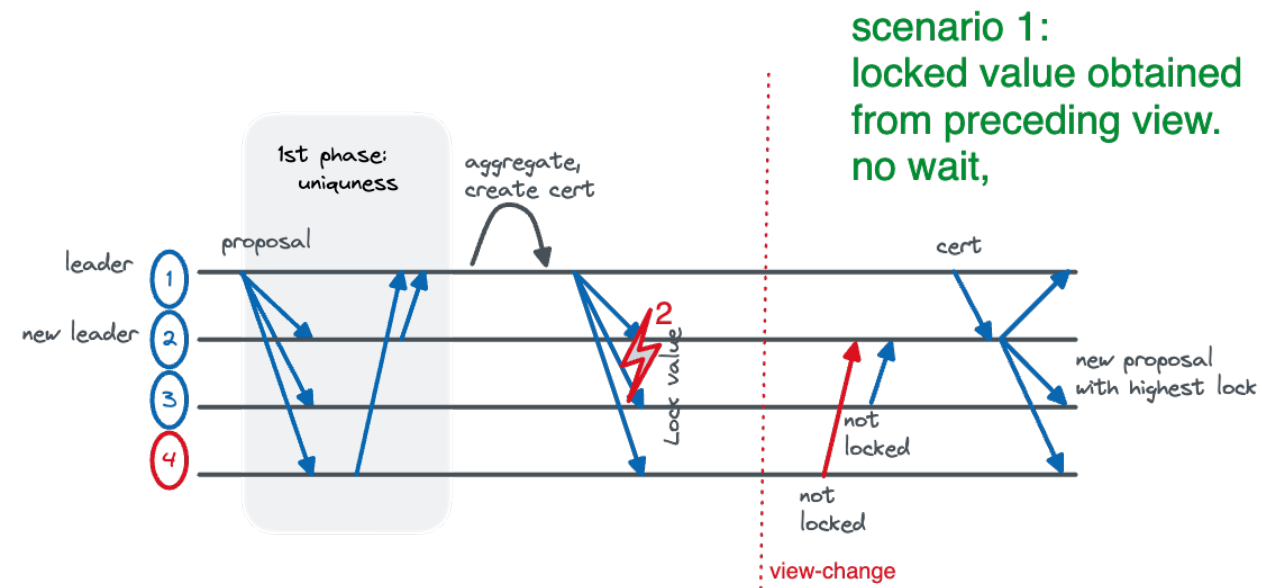


# Basic HotStuff \_ mitigation of case-2



- Each replica sends its own prepareQC (replica3 is not locked but has a prepareQC).
- Upon receiving  $2f+1$  prepareQCs, leader makes new proposal with *highQC* among those prepareQCs.
- Locked correct replicas vote to new proposal when  $\text{highQC} > \text{lockedQC}$ .

# Cf) HotStuff-2



- Highest  $qc$  must be found in 2-phase BFT.
- If the leader obtains a  $qc$  from the preceding view, the  $qc$  is maximal value that possibly exists in the system. In this case, it proceeds with a proposal in a responsive manner.
- Otherwise, must waits  $\Delta$ . This can happen when the leader from the previous view is malicious or the previous view was before GST.



# Basic HotStuff \_ pseudocode(1)

## Algorithm 1 Utilities (for replica $r$ ).

```
1: function Msg( $type, node, qc$ )
2:    $m.type \leftarrow type$ 
3:    $m.viewNumber \leftarrow curView$ 
4:    $m.node \leftarrow node$ 
5:    $m.justify \leftarrow qc$ 
6:   return  $m$ 
7: function VOTEMSG( $type, node, qc$ )
8:    $m \leftarrow \text{Msg}(type, node, qc)$ 
9:    $m.partialSig \leftarrow tsign_r((m.type, m.viewNumber, m.node))$ 
10:  return  $m$ 
11: procedure CREATELEAF( $parent, cmd$ )
12:    $b.parent \leftarrow parent$ 
13:    $b.cmd \leftarrow cmd$ 
14:   return  $b$ 
15: function QC( $V$ )
16:    $qc.type \leftarrow m.type : m \in V$ 
17:    $qc.viewNumber \leftarrow m.viewNumber : m \in V$ 
18:    $qc.node \leftarrow m.node : m \in V$ 
19:    $qc.sig \leftarrow tcombine(\langle qc.type, qc.viewNumber, qc.node \rangle,$ 
20:      $\{m.partialSig \mid m \in V\})$ 
21:   return  $qc$ 
22: function MATCHINGMSG( $m, t, v$ )
23:   return  $(m.type = t) \wedge (m.viewNumber = v)$ 
24: function MATCHINGQC( $qc, t, v$ )
25:   return  $(qc.type = t) \wedge (qc.viewNumber = v)$ 
26: function SAFENODE( $node, qc$ )
27:   return  $(node \text{ extends from } locked\ QC.node) \vee$  // safety rule
28:    $(qc.viewNumber > locked\ QC.viewNumber)$  // liveness rule
```

## Algorithm 2 Basic HotStuff protocol (for replica $r$ ).

```
1: for  $curView \leftarrow 1, 2, 3, \dots$  do
2:   ▷ PREPARE phase
3:   as a leader //  $r = \text{LEADER}(curView)$ 
4:     // we assume special NEW-VIEW messages from view 0
5:     wait for  $(n - f)$  NEW-VIEW messages:  $M \leftarrow \{m \mid \text{MATCHINGMSG}(m, \text{NEW-VIEW}, curView - 1)\}$ 
6:      $high\ QC \leftarrow \left( \arg \max_{m \in M} \{m.justify.viewNumber\} \right).justify$ 
7:      $curProposal \leftarrow \text{CREATELEAF}(high\ QC.node, \text{client's command})$ 
8:     broadcast MSG(PREPARE,  $curProposal$ ,  $high\ QC$ )
9:   as a replica
10:    wait for message  $m : \text{MATCHINGMSG}(m, \text{PREPARE}, curView)$  from  $\text{LEADER}(curView)$ 
11:    if  $m.node$  extends from  $m.justify.node \wedge$ 
12:      SAFENODE( $m.node, m.justify$ ) then
13:        send VOTEMSG(PREPARE,  $m.node, \perp$ ) to  $\text{LEADER}(curView)$ 
14:   ▷ PRE-COMMIT phase
15:   as a leader
16:     wait for  $(n - f)$  votes:  $V \leftarrow \{v \mid \text{MATCHINGMSG}(v, \text{PREPARE}, curView)\}$ 
17:      $prepare\ QC \leftarrow \text{QC}(V)$ 
18:     broadcast MSG(PRE-COMMIT,  $\perp$ ,  $prepare\ QC$ )
19:   as a replica
20:    wait for message  $m : \text{MATCHINGQC}(m.justify, \text{PREPARE}, curView)$  from  $\text{LEADER}(curView)$ 
21:     $prepare\ QC \leftarrow m.justify$ 
22:    send VOTEMSG(PRE-COMMIT,  $m.justify.node, \perp$ ) to  $\text{LEADER}(curView)$ 
```

# Basic HotStuff \_ pseudocode(2)

## Algorithm 1 Utilities (for replica $r$ ).

```
1: function Msg( $type, node, qc$ )
2:    $m.type \leftarrow type$ 
3:    $m.viewNumber \leftarrow curView$ 
4:    $m.node \leftarrow node$ 
5:    $m.justify \leftarrow qc$ 
6:   return  $m$ 
7: function VOTEMSG( $type, node, qc$ )
8:    $m \leftarrow \text{Msg}(type, node, qc)$ 
9:    $m.partialSig \leftarrow \text{tsign}_r((m.type, m.viewNumber, m.node))$ 
10:  return  $m$ 
11: procedure CREATELEAF( $parent, cmd$ )
12:    $b.parent \leftarrow parent$ 
13:    $b.cmd \leftarrow cmd$ 
14:   return  $b$ 
15: function QC( $V$ )
16:    $qc.type \leftarrow m.type : m \in V$ 
17:    $qc.viewNumber \leftarrow m.viewNumber : m \in V$ 
18:    $qc.node \leftarrow m.node : m \in V$ 
19:    $qc.sig \leftarrow \text{tcombine}(\langle qc.type, qc.viewNumber, qc.node \rangle,$ 
     $\{m.partialSig \mid m \in V\})$ 
20:  return  $qc$ 
21: function MATCHINGMSG( $m, t, v$ )
22:   return  $(m.type = t) \wedge (m.viewNumber = v)$ 
23: function MATCHINGQC( $qc, t, v$ )
24:   return  $(qc.type = t) \wedge (qc.viewNumber = v)$ 
25: function SAFENODE( $node, qc$ )
26:   return  $(node \text{ extends from } lockedQC.node) \vee$  // safety rule
27:    $(qc.viewNumber > lockedQC.viewNumber)$  // liveness rule
```

## ▷ COMMIT phase

```
19:   as a leader
20:     wait for  $(n - f)$  votes:  $V \leftarrow \{v \mid \text{MATCHINGMSG}(v, \text{PRE-COMMIT}, curView)\}$ 
21:      $precommitQC \leftarrow QC(V)$ 
22:     broadcast MSG(COMMIT,  $\perp$ ,  $precommitQC$ )
23:   as a replica
24:     wait for message  $m : \text{MATCHINGQC}(m.justify, \text{PRE-COMMIT}, curView)$  from LEADER( $curView$ )
25:      $lockedQC \leftarrow m.justify$ 
26:     send VOTEMSG(COMMIT,  $m.justify.node, \perp$ ) to LEADER( $curView$ )
```

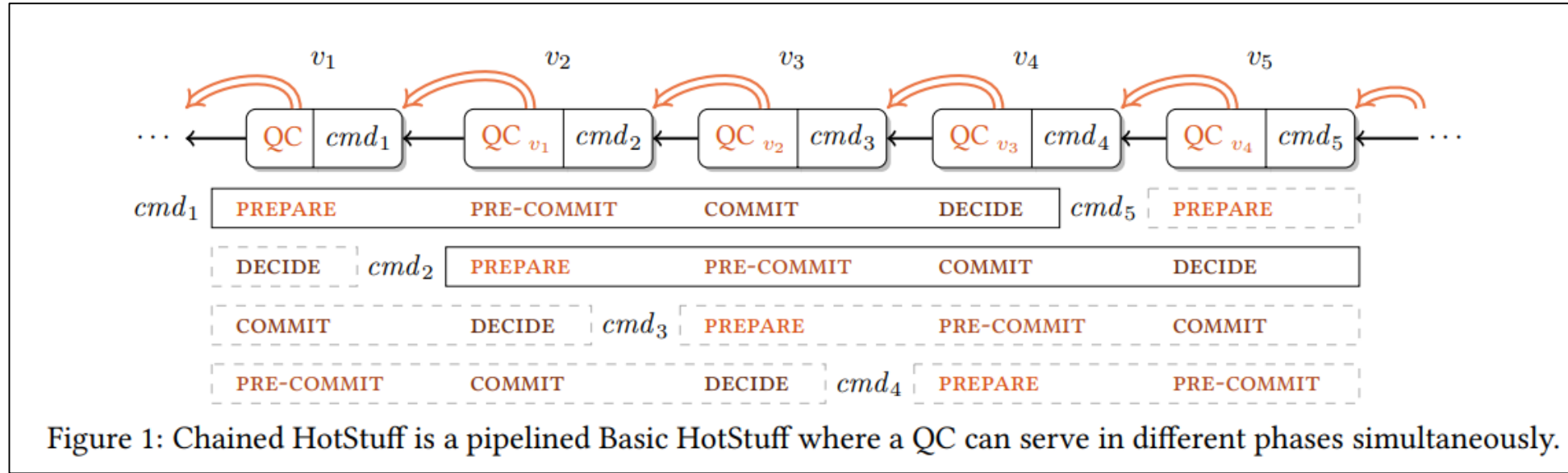
## ▷ DECIDE phase

```
27:   as a leader
28:     wait for  $(n - f)$  votes:  $V \leftarrow \{v \mid \text{MATCHINGMSG}(v, \text{COMMIT}, curView)\}$ 
29:      $commitQC \leftarrow QC(V)$ 
30:     broadcast MSG(DECIDE,  $\perp$ ,  $commitQC$ )
31:   as a replica
32:     wait for message  $m$  from LEADER( $curView$ )
33:     wait for message  $m : \text{MATCHINGQC}(m.justify, \text{COMMIT}, curView)$  from LEADER( $curView$ )
34:     execute new commands through  $m.justify.node$ , respond to clients
```

## ▷ Finally

```
35:   NEXTVIEW interrupt: goto this line if NEXTVIEW( $curView$ ) is called during “wait for” in any phase
36:   send MSG(NEW-VIEW,  $\perp$ ,  $prepareQC$ ) to LEADER( $curView + 1$ )
```

# Chained HotStuff



# Chained HotStuff \_ pseudocode

```
6: for  $curView \leftarrow 1, 2, 3, \dots$  do
  ▷ GENERIC phase
7:   as a leader //  $r = \text{LEADER}(curView)$ 
    //  $M$  is the set of messages collected at the end of previous view by the leader of this view
8:      $highQC \leftarrow \left( \arg \max_{m \in M} \{m.justify.viewNumber\} \right).justify$ 
9:     if  $highQC.viewNumber > genericQC.viewNumber$  then  $genericQC \leftarrow highQC$ 
10:     $curProposal \leftarrow \text{CREATELEAF}(genericQC.node, \text{client's command}, genericQC)$ 
    // PREPARE phase
11:    broadcast MSG(GENERIC,  $curProposal, \perp$ )
12:  as a replica
13:    wait for message  $m : \text{MATCHINGMSG}(m, \text{GENERIC}, curView)$  from  $\text{LEADER}(curView)$ 
14:     $b^* \leftarrow m.node; b'' \leftarrow b^*.justify.node; b' \leftarrow b''.justify.node; b \leftarrow b'.justify.node$ 
15:    if  $\text{SAFE\_NODE}(b^*, b^*.justify)$  then
16:      send  $\text{VOTMSG}(\text{GENERIC}, b^*, \perp)$  to  $\text{LEADER}(curView + 1)$ 
    // start PRE-COMMIT phase on  $b^*$ 's parent
17:    if  $b^*.parent = b''$  then
18:       $genericQC \leftarrow b^*.justify$ 
    // start COMMIT phase on  $b^*$ 's grandparent
19:    if  $(b^*.parent = b'') \wedge (b''.parent = b')$  then
20:       $lockedQC \leftarrow b''.justify$ 
    // start DECIDE phase on  $b^*$ 's great-grandparent
21:    if  $(b^*.parent = b'') \wedge (b''.parent = b') \wedge (b'.parent = b)$  then
22:      execute new commands through  $b$ , respond to clients
23:  as the next leader
24:    wait for all messages:  $M \leftarrow \{m \mid \text{MATCHINGMSG}(m, \text{GENERIC}, curView)\}$ 
    until there are  $(n - f)$  votes:  $V \leftarrow \{v \mid v.partialSig \neq \perp \wedge v \in M\}$ 
25:     $genericQC \leftarrow QC(V)$ 
  ▷ Finally
26:  NEXTVIEW interrupt: goto this line if  $\text{NEXTVIEW}(curView)$  is called during “wait for” in any phase
27:  send  $\text{MSG}(\text{GENERIC}, \perp, genericQC)$  to  $\text{LEADER}(curView + 1)$ 
```

# Evaluation \_ setup

- EC2, c5.4xlarge (16 vCPUs supported by Intel Xeon Platinum 8000 processors @ 3.4GHz)
- The maximum TCP bandwidth measured by iperf was around 1.2 Gigabytes per second.
- The network latency between two machines was less than 1 ms. (LAN?)

# Evaluation \_ base performance

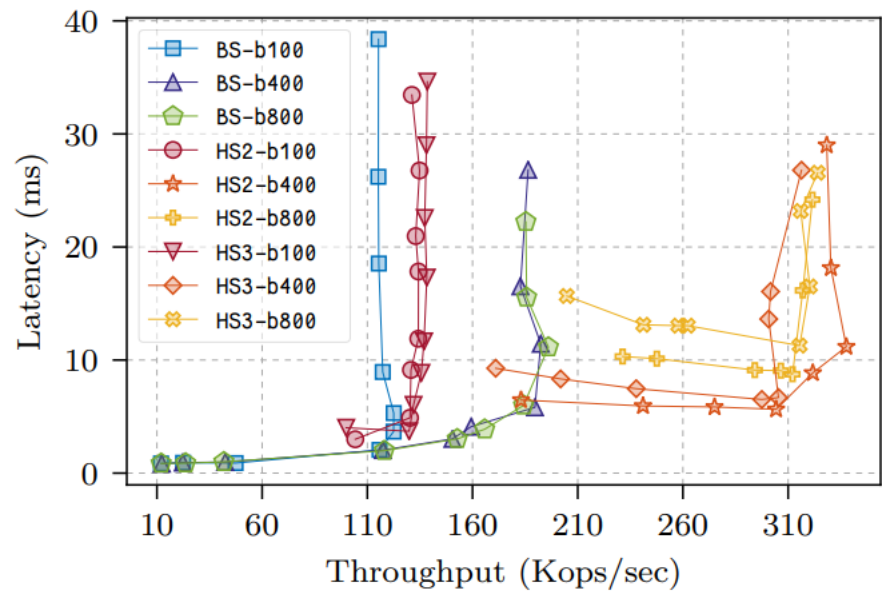


Figure 4: Throughput vs. latency with different choices of batch size, 4 replicas, 0/0 payload.

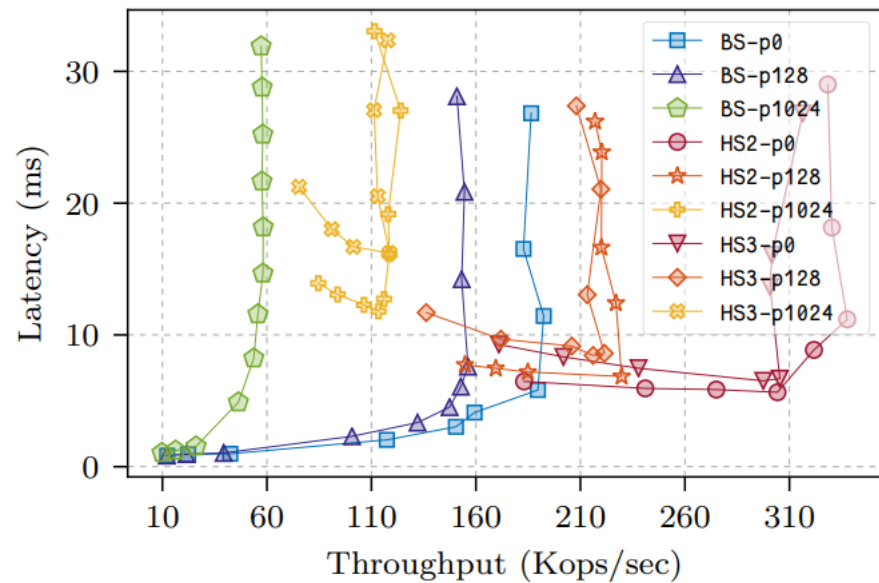
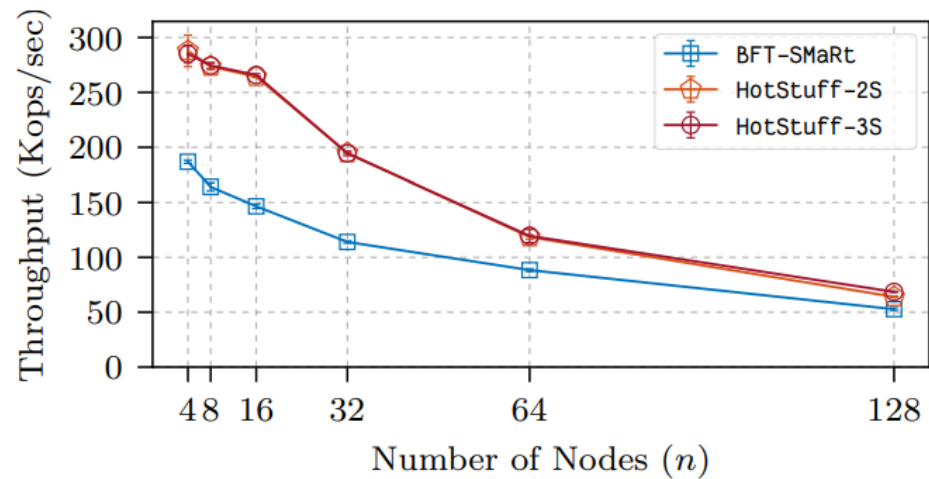


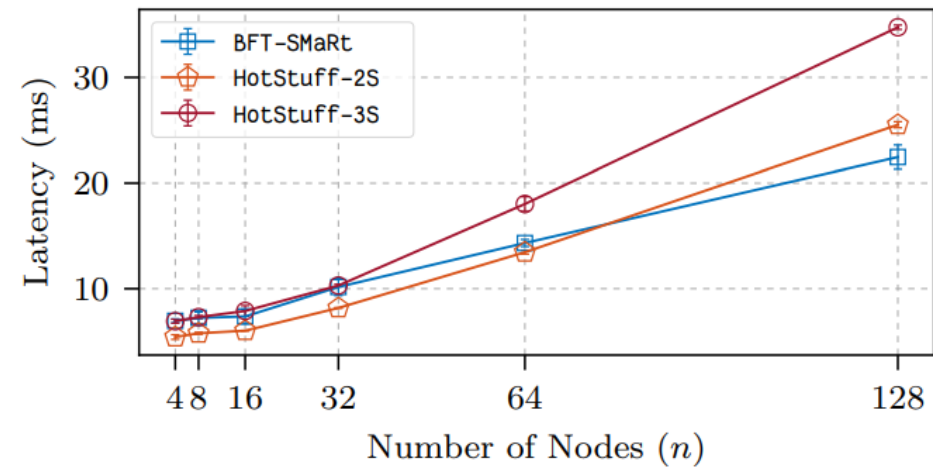
Figure 5: Throughput vs. latency with different choices of payload size, 4 replicas, batch size of 400.



# Evaluation \_ scalability



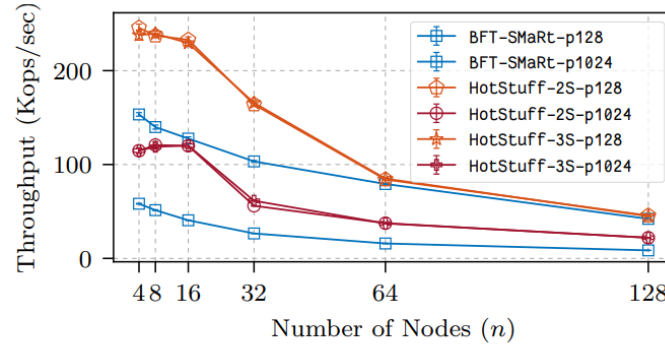
(a) Throughput



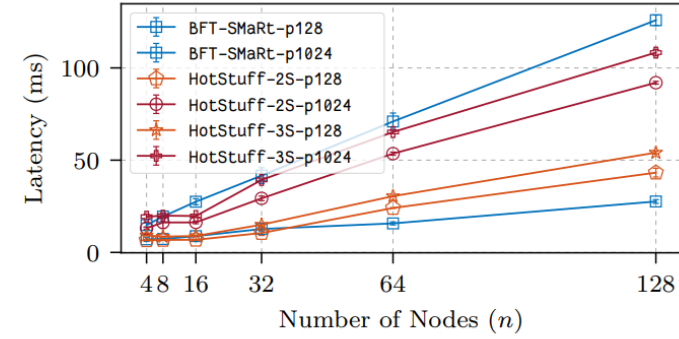
(b) Latency

Figure 6: Scalability with 0/0 payload, batch size of 400.

# Evaluation \_ scalability

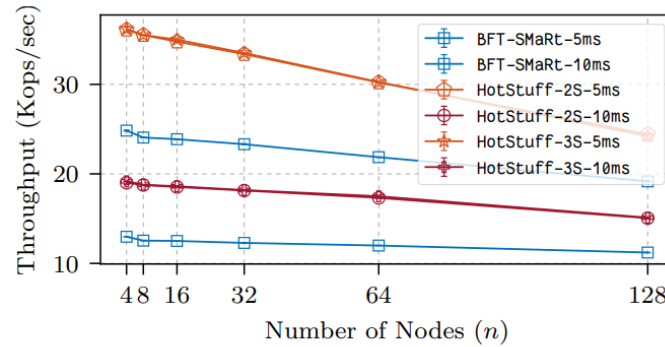


(a) Throughput

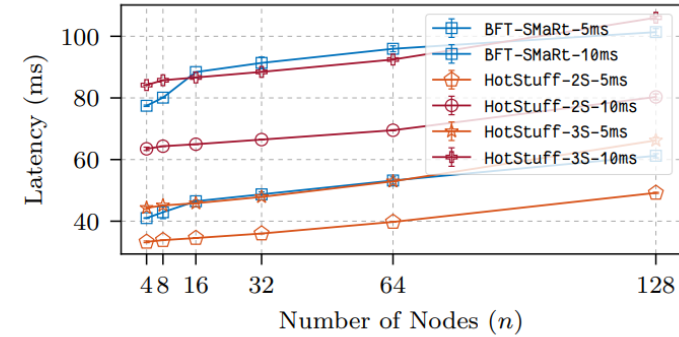


(b) Latency

Figure 7: Scalability for 128/128 payload or 1024/1024 payload, with batch size of 400.



(a) Throughput



(b) Latency

Figure 8: Scalability for inter-replica latency 5ms  $\pm$  0.5ms or 10ms  $\pm$  1.0ms, with 0/0 payload, batch size of 400.



# Evaluation \_ view change

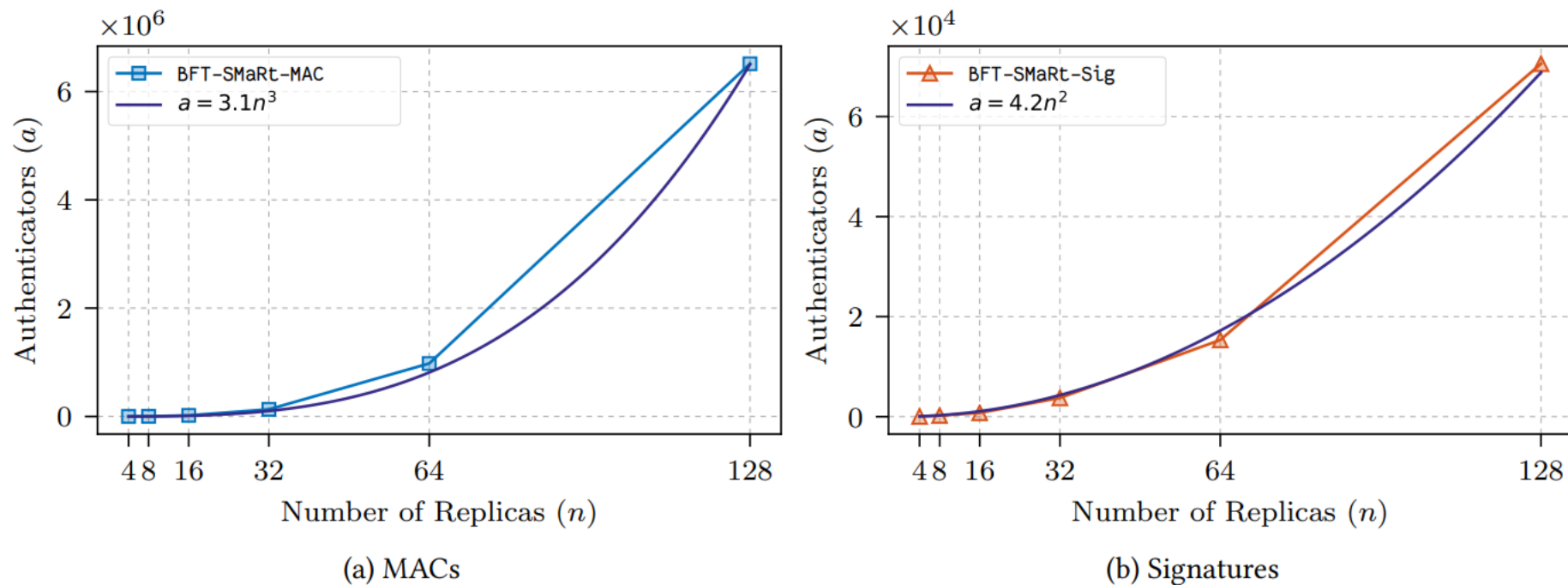


Figure 10: Number of extra authenticators used for each BFT-SMaRt view change.