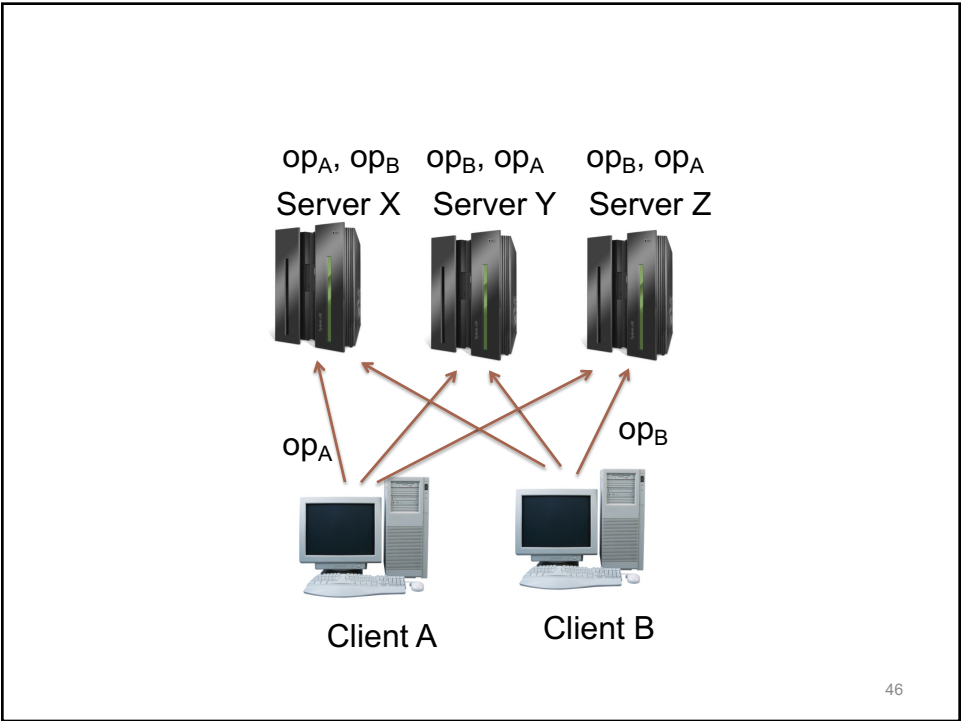# VIEWSTAMP REPLICATION: ORDERING UPDATES
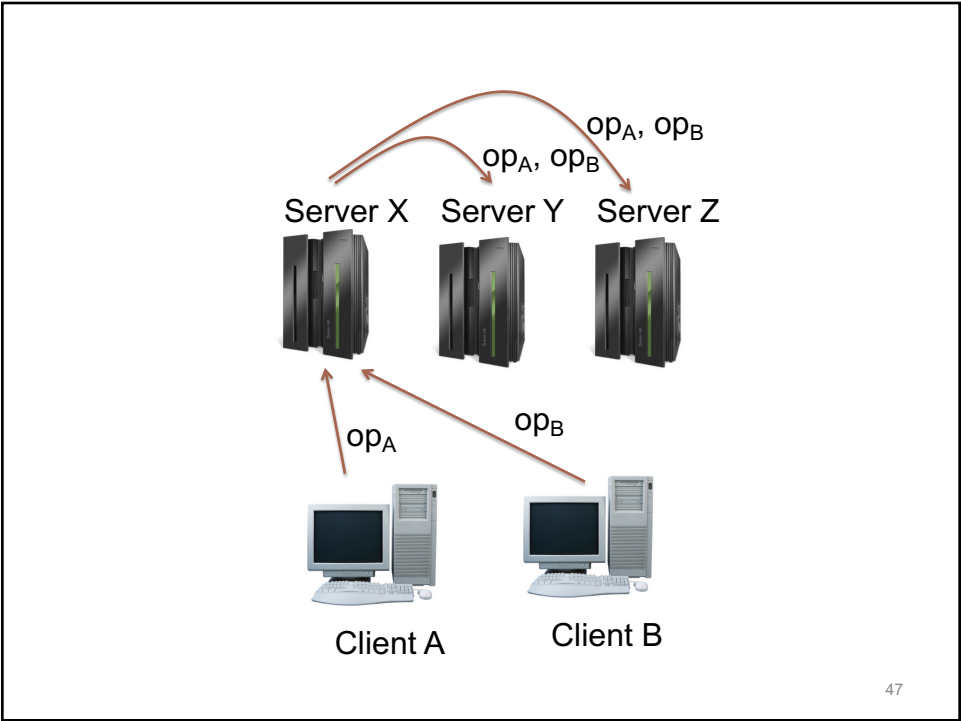
44

# Quorum Consensus

- Crash-stop failures

- Requires 2F+1 replicas
  - Operations must intersect for at least one replica
  - Want availability for both reads and writes
  - Read and write quorums of F+1 nodes

45

op$_A$, op$_B$    op$_B$, op$_A$    op$_B$, op$_A$
Server X    Server Y    Server Z

op$_A$

op$_B$

Client A        Client B

46

op$_A$, op$_B$

op$_A$, op$_B$

Server X    Server Y    Server Z

op$_A$

op$_B$

Client A        Client B

47

# Viewstamp Replication

- Primary-backup
- System moves through a sequence of views
  - Primary runs the protocol
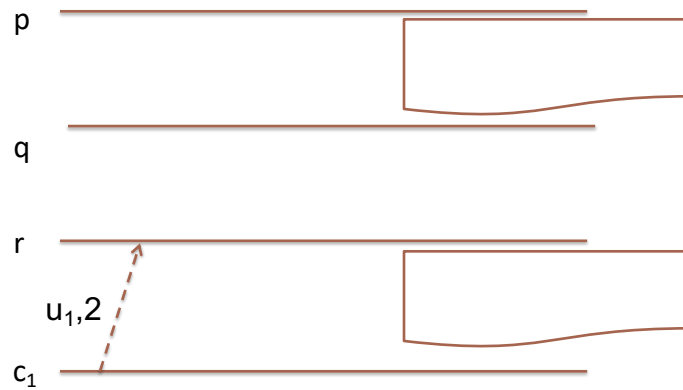  - Replicas do a view change if it fails

# Replica state

- A replica id i (between 0 and N-1)
  - Replica 0, replica 1, …
- A view number v#, initially 0
- Primary is the replica with id
  - i = v# mod N
- A log of <op, op#, status> entries
  - Status = prepared or committed

Client knows current view #
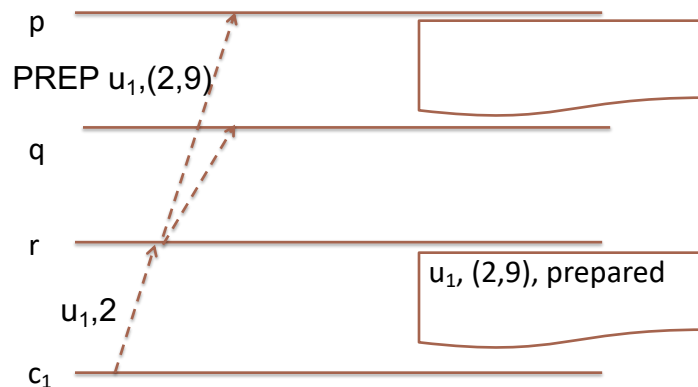
View # = 2 $\Rightarrow$ primary pid = 2 (i.e., r)

p

q

r

$u_1, 2$

$c_1$

50

---

Primary chooses logical timestamp for update, e.g. $LT(u_1) = 9$

Op # = (View #, time)

p

PREP $u_1, (2,9)$

q

r

$u_1, 2$

$c_1$

$u_1, (2,9)$, prepared

51

Primary waits for $\geq f$ replicas to respond ($\geq f+1$ total)

$u_1$, (2,9), prepared

p

$u_1$, (2,9), prepared

PREP $u_1$,(2,9)

q

r

$u_1$, (2,9), prepared

$u_1$,2

$c_1$

52

Client notified immediately after acks

$u_1$, (2,9), prepared

p

$u_1$, (2,9), prepared

PREP $u_1$,(2,9)

q

COMMIT (2,9)

r

$u_1$, (2,9), committed

$u_1$,2

$c_1$

53

**VIEWSTAMP REPLICATION:
VIEW CHANGE**

# View Changes

- Used to mask primary failures
- Replicas monitor the primary
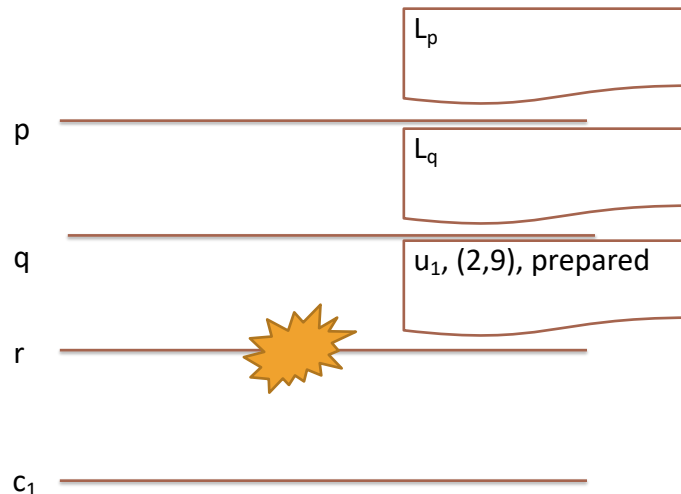- Replica requests next primary to do a view change

# Correctness Requirement

- Operation order must be preserved by a view change

- For operations that are visible
  - executed by server
  - client received result

- An operation can be visible if it prepared at f+1 replicas
  - this is the commit point

---

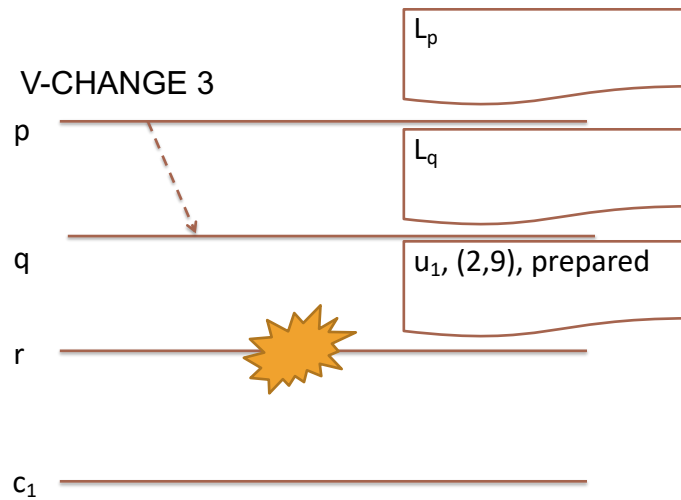Primary r has crashed after assigning $LT(u_1)=9$

$L_p$

p

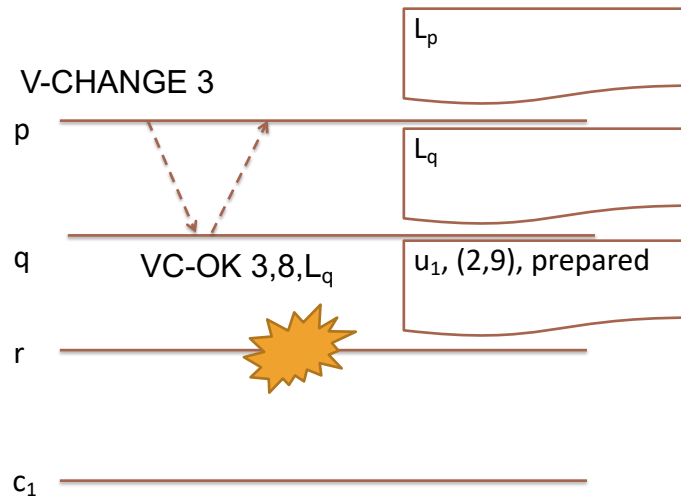$L_q$

q $\quad u_1, (2,9)$, prepared

r

$c_1$

## Slide 58

View # increments (to 3), new primary is p, start view change

V-CHANGE 3

$L_p$

p

$L_q$

q

$u_1, (2,9)$, prepared

r

$c_1$

58

58

## Slide 59

Replica q acks view change, returns last op# it knew of and its log ($L_q$ contains any operations that p never heard of)
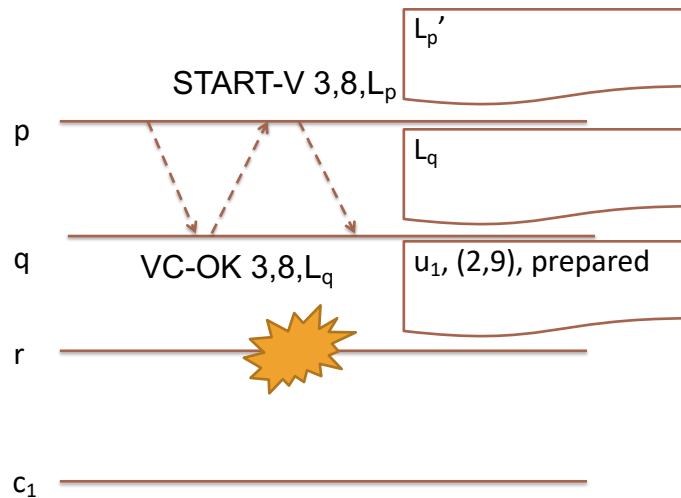
V-CHANGE 3

$L_p$

p

$L_q$

q    VC-OK $3,8,L_q$

$u_1, (2,9)$, prepared

r

$c_1$

59

59

8

Primary p lets replicas know (via its log $L_p'$) of ops that replicas were never informed of before old primary crashed

START-V 3,8,$L_p$

$L_p'$

p

$L_q$

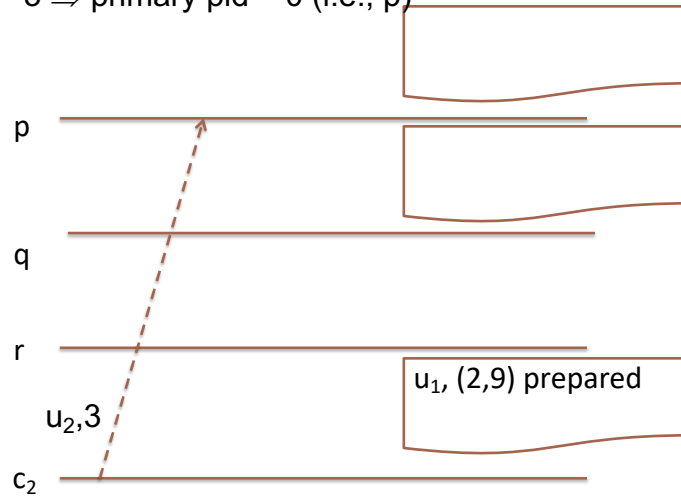q    VC-OK 3,8,$L_q$     $u_1$, (2,9), prepared

r

$c_1$

60

# View Change

- New primary may not know of all updates from old primary
- New primary asks replicas for their logs
- Any committed operation was acked by a quorum, so must be in log of a surviving replica
  - Primary takes the max of the logs returned
  - That log has most recent updates
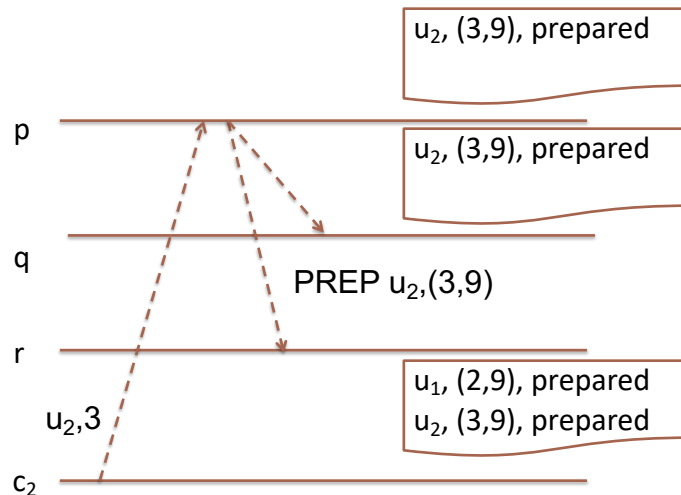
61

Second client $c_2$ knows current view #

View# = 3 $\Rightarrow$ primary pid = 0 (i.e., p)

$u_1$, (2,9) prepared
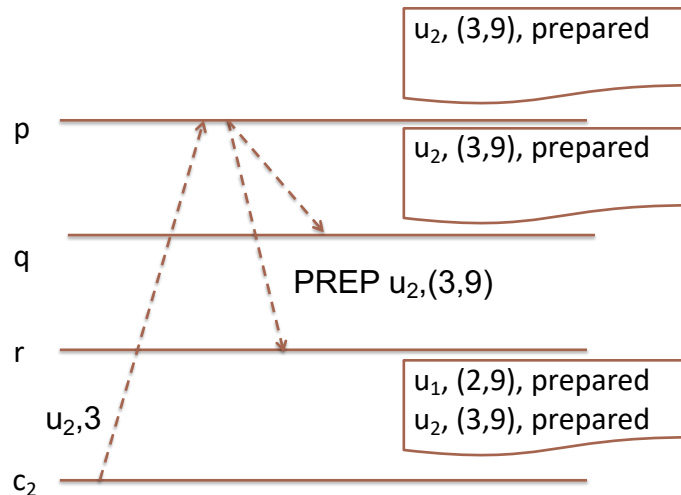
p

q

r

$u_2$,3

$c_2$

62

---

Primary p gets the other replicas to prepare to commit $u_2$

$u_2$, (3,9), prepared

p

$u_2$, (3,9), prepared

q

PREP $u_2$,(3,9)

r

$u_1$, (2,9), prepared
$u_2$, (3,9), prepared

$u_2$,3

$c_2$

63

Viewstamp (3,9) avoids confusion between $u_1$ and $u_2$ at r
(Replica r will not think p is committing $u_1$)

p

$u_2$, (3,9), prepared

$u_2$, (3,9), prepared

q

PREP $u_2$,(3,9)

r

$u_1$, (2,9), prepared
$u_2$, (3,9), prepared

$u_2$,3

$c_2$

64

64

# Persistent State

- Voting protocol: votes must survive failure
  - Save queue of pending updates on disk
- Viewstamp: primary can respond to client without recording commit on disk
  - View change: recover commitment from logs of surviving replicas
- Only need to persist state after view change
  - So if we crash and recover, we know view# when we crashed
  - Even that unnecessary with more expensive recovery protocol

65

# Summary

- Primary-Backup
  - Order updates at the primary
  - Order preserved by view change
- Quorum Consensus
- Split brain?
- What about FLP?

66

66