

A Proof of Correctness for Egalitarian Paxos

Iulian Moraru¹, David G. Andersen¹, Michael Kaminsky²

¹ Carnegie Mellon University, ² Intel Labs

CMU-PDL-13-111

August 2013

Parallel Data Laboratory
Carnegie Mellon University
Pittsburgh, PA 15213-3890

Abstract

This paper presents a proof of correctness for Egalitarian Paxos (EPaxos), a new distributed consensus algorithm based on Paxos. EPaxos achieves three goals: (1) availability without interruption as long as a simple majority of replicas are reachable—its availability is not interrupted when replicas crash or fail to respond; (2) uniform load balancing across all replicas—no replicas experience higher load because they have special roles; and (3) optimal commit latency in the wide-area when tolerating one and two failures, under realistic conditions. Egalitarian Paxos is to our knowledge the first distributed consensus protocol to achieve all of these goals efficiently: requiring only a simple majority of replicas to be non-faulty, using a number of messages linear in the number of replicas to choose a command, and committing commands after just one communication round (one round trip) in the common case or after at most two rounds in any case.

Acknowledgements: This work was supported by Intel Science & Technology Center for Cloud Computing, Google, and the National Science Foundation under award CCF-0964474. We would like to thank the members and companies of the PDL Consortium (including Actifio, American Power Conversion, EMC Corporation, Emulex, Facebook, Fusion-io, Google, Hewlett-Packard Labs, Hitachi, Huawei Technologies Co., Intel Corporation, Microsoft Research, NEC Laboratories, NetApp, Inc., Oracle Corporation, Panasas, Riverbed, Samsung Information Systems America, Seagate Technology, STEC, Inc., Symantec Corporation, VMware, Inc., and Western Digital) for their interest, insights, feedback, and support.

Keywords: Distributed Consensus, Paxos

1 Introduction

Today’s clusters use fault-tolerant, highly available coordination engines like Chubby [1], Boxwood [7], or ZooKeeper[5] for activities such as operation sequencing, coordination, leader election, and resource discovery. An important limitation in these systems is that during efficient, normal operation, all the clients communicate with a single master (or leader) server at all times. This optimization, sometimes termed “Multi-Paxos,” is important to achieving high throughput in practical systems [2]. Changing the leader requires invoking additional consensus mechanisms that substantially reduce throughput.

This algorithmic limitation has several important consequences. First, it can impair scalability by placing a disproportionately high load on the master, which must process more messages than the other replicas [8]. Second, it can harm availability: if the master fails, the system cannot service requests until a new master is elected. Finally, as we show in this paper, traditional Paxos variants are sensitive to both long-term and transient load spikes and network delays that increase latency at the master. Previously proposed solutions such as partitioning or using proxy servers are undesirable because they restrict the type of operations the cluster can perform. For example, a partitioned cluster cannot perform atomic operations across partitions without using additional techniques.

Egalitarian Paxos (EPaxos) has no designated leader process. Instead, clients can choose, at every step, which replica to submit a command to, and in most cases the command will be committed without interfering with other concurrent commands. This allows the system to evenly distribute the load to all replicas, eliminating the first bottleneck identified above (having one server that must be on the critical path for all communication). The system can provide higher availability because there is no transient interruption because of leader election: there is no leader, and hence, no need for leader election, as long as more than half of the replicas are available. Finally, EPaxos’s flexible load distribution is better able to handle permanently or transiently slow nodes, substantially reducing both the median and tail commit latency.

2 Preliminaries

We begin by stating assumptions, definitions, and introducing our notation.

Messages exchanged by processes (clients and replicas) are asynchronous. Failures are non-Byzantine (a machine can fail by stopping to respond for an indefinite amount of time). The replicated state machine comprises N replicas. For every replica R there is an unbounded sequence of numbered instances $R.1$, $R.2$, $R.3$, ... that that replica is said to *own*. At most one command will be adopted in an instance. The ordering of the instances is *not pre-determined*—it is determined dynamically by the protocol, as commands are chosen.

It is important to understand that *committing* and *executing* commands are different actions, and that the commit and execution orders are not necessarily the same. A client of an EPaxos-based system will interact with the system through an interface of the following form:

To modify the replicated state, a client sends $Request(command)$ to a replica of its choice. A $RequestReply$ from that replica will notify the client that the command has been committed.

To read (a part of) the state, clients send $Read(objectID)$ messages and wait for $ReadReply$. A $Read$ is itself a special no-op command that interferes with updates to the object it is reading.

A client that receives a $RequestReply$ for a command knows only that the command has been committed, but has no information about whether the command has been executed or not. Only when the client reads the replicated state updated by its previously committed commands is it necessary for those commands to be executed.

3 Interfering Commands

Before we can describe Egalitarian Paxos in detail, we must define command *interference*.

Informally, two commands that interfere must be executed in the same order by all replicas.

Definition 1 (Interference). Two commands γ and δ *interfere* if there exists a sequence of commands Σ such that the serial execution Σ, γ, δ is not compatible (i.e., it produces different results than) the serial execution Σ, δ, γ .

As we explain in the following subsection, EPaxos guarantees that any two interfering commands will be executed in the same order with respect to each other on every replica. This is enough to guarantee that the executions on all replicas are compatible: the serial ordering of commands on a replica can be obtained from that on any other replica by commuting commutative commands.

Note that the interference relation is symmetric, but not necessarily transitive.

4 Protocol Guarantees

The formal guarantees that Egalitarian Paxos offers clients are similar to those provided by other Paxos variants:

Nontriviality Any command committed by any replica must have been proposed by a client.

Stability For any replica, the set of committed commands at any time is a subset of the committed commands at any later time. Furthermore, if at time t_1 a replica R has command γ committed at some instance $Q.i$, then R will have γ committed in $Q.i$ at any later time $t_2 > t_1$.

Consistency Two replicas can never have different commands committed for the same instance.

Execution consistency For any two commands γ and δ that interfere, if both γ and δ have been committed by any replicas, then γ and δ will be executed in the same order by every replica.

Execution linearizability If two interfering commands γ and δ are serialized by clients (i.e., δ is proposed only after γ is committed by any replica), then every replica will execute γ before δ .

Liveness A proposed command will eventually be committed by every non-faulty replica, as long as fewer than half the replicas are faulty and messages eventually reach their destination before their recipient times out. These are the same liveness guarantees provided by Paxos. By FLP [3], it is impossible to provide stronger guarantees for distributed consensus.

5 Simplified Egalitarian Paxos

In this section we describe the basic form of the Egalitarian Paxos protocol. In Section 6 we will show how to modify this protocol to reduce the quorum size.

5.1 The EPaxos Commit Protocol

As mentioned earlier, committing and executing commands are separate. Accordingly, EPaxos comprises two components: (1) the protocol for choosing (committing) commands and determining their ordering attributes in the process; and (2) the algorithm for executing commands based on these attributes.

We present a pseudocode description of the Egalitarian Paxos protocol for choosing commands below. The state of each replica is represented in the pseudocode by each replica's private *commands* array.

We split the description of the commit protocol into multiple phases. Not all phases are executed for every command: a command committed after the execution of Init, Phase 1, and Commit, is said to have

L, R, Q	Replicas (the command leader is usually denoted by L)
γ, δ	Commands
$R.i$, with $i = 1, 2, \dots$	Instances belonging to replica R
$epoch.i.R$, with $epoch, i = 0, 1, 2, \dots$	Ballot numbers generated by replica R ($epoch.0.R$ is the initial ballot for any instance $R.i$)
$deps_\gamma$	The list of dependencies for command γ
seq_γ	Approximate sequence number for command γ
$Q::commands[R][i]$	The state of replica Q at instance $R.i$

Figure 1: Summary of notation.

been executed on the *fast path*. The *slow path* involves the additional Paxos-Accept phase. The *Explicit Prepare* phase is only executed on failure recovery.

Phase 1 starts when a replica L receives a request (for a command γ) and becomes a command leader. L begins the process of choosing γ in the next available instance of its instance space. It also attaches what it believes are the correct attributes for that command:

deps is the list of all instances that contain commands (not necessarily committed) that interfere with γ ; we say that γ *depends* on those instances (and their corresponding commands);

seq is a sequence number used to break dependency cycles during the execution algorithm; *seq* is updated to be larger than the *seq* numbers of all commands in *deps*.

The command leader forwards the command and the initial attributes to at least a *fast quorum* of replicas as a *PreAccept* message. For now, we assume that a fast quorum contains $N - 1$ replicas, including the command leader. We will show in Section 5.7 that we can reduce the fast quorum size to $\lceil 3N/4 \rceil$ when $N > 3$.

Each replica, upon receiving the *PreAccept*, updates γ 's attributes according to the contents of its *commands* log, records γ and the new attributes in *commands*, and replies to the command leader.

If the command leader receives replies from enough replicas to constitute a fast quorum, and all the updated attributes are the same, it commits the command. If it doesn't receive enough replies, or the attributes in some replies have been updated differently than in others, then the command leader updates the attributes based on $\lfloor N/2 \rfloor + 1$ replies (taking the union of all *deps*, and the highest *seq*), and tells at least $\lfloor N/2 \rfloor + 1$ replicas to accept these attributes. This can be seen as running Paxos-Accept for choosing the triplet $(\gamma, deps_\gamma, seq_\gamma)$ in γ 's instance. At the end of this extra round, after replies from a majority (including itself), the command leader will reply to the client, and will send *Commit* messages asynchronously to all the other replicas.

Like classic Paxos, every message contains a ballot number (not presented explicitly in the pseudocode for phases other than Explicit Prepare). As in classic Paxos, the ballot number ensures message freshness: a replica will disregard any message with a smaller ballot than the largest it has seen for a certain instance. Ballot numbers used by a replica R have the form $epoch.i.R$, where *epoch* is the natural number id of the current configuration and i is a natural number (*epoch* and i take precedence when ordering ballots, while the replica id R ensures different replicas cannot generate the same ballot number). The initial *Prepare* phase is implicit for all instances $R.i$, for an initial ballot number $epoch.0.R$ and each replica is the default (i.e., initial) leader of its own instances. Whenever a command leader receives a NACK for one of its messages, indicating that some other replica has used a higher ballot in the same instance, that command leader will fallback on executing Explicit Prepare.

Init

Replica L on receiving $Request(\gamma)$ from a client:

1: start Phase 1 for γ at previously unused instance $L.i$

Phase 1

Replica L designated as leader for command γ (steps 2, 3 and 4 executed atomically):

- 2: $seq_\gamma \leftarrow \max(\{0\} \cup \{\text{seq. attribute of every command recorded in } L :: \text{commands that interferes w/ } \gamma\}) + 1$
- 3: $deps_\gamma \leftarrow \{(Q, j) \mid L :: \text{commands}[Q][j] \text{ interferes w/ } \gamma\}$
- 4: $L :: \text{commands}[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, PreAccepted)$
- 5: send $PreAccept(\gamma, seq_\gamma, deps_\gamma, L.i)$ to all other replicas in fast quorum that includes L

Replica R , on receiving $PreAccept(\gamma, seq_\gamma, deps_\gamma, L.i)$ from replica L (steps 6 through 10 executed atomically):

- 6: $max_seq \leftarrow \max(\{0\} \cup \{\text{seq. attribute of every command } \delta \text{ in } R :: \text{commands, s.t. } \gamma \text{ and } \delta \text{ interfere}\})$
- 7: update $seq_\gamma \leftarrow \max(\{seq_\gamma, max_seq + 1\})$
- 8: $deps_{local} \leftarrow \{(Q, j) \mid R :: \text{commands}[Q][j] \text{ interferes with } \gamma\}$
- 9: update $deps_\gamma \leftarrow deps_\gamma \cup deps_{local}$
- 10: $R :: \text{commands}[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, PreAccepted)$
- 11: reply $PreAcceptOK(\gamma, seq_\gamma, deps_\gamma, L.i)$ to L

Replica L (designated leader for command γ), on receiving at least $\lfloor N/2 \rfloor + 1$ $PreAcceptOK$ responses:

- 12: **if** received at least $N - 2$ $PreAcceptOK$'s with the same seq_γ and $deps_\gamma$ attributes **then**
- 13: reply $RequestReply(\gamma, L.i)$ to client
- 14: run Commit phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$
- 15: **else**
- 16: update $deps_\gamma \leftarrow \text{Union}(deps_\gamma \text{ from all replies})$
- 17: update $seq_\gamma \leftarrow \max(\{seq_\gamma \text{ of all replies}\})$
- 18: run Paxos-Accept phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$

Paxos-Accept

Designated leader replica L , for $(\gamma, seq_\gamma, deps_\gamma)$ at instance $L.i$

- 19: $L :: \text{commands}[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, Accepted)$
- 20: send $Accept(\gamma, seq_\gamma, deps_\gamma, L.i)$ to all other replicas in a slow quorum that includes L

Replica R , on receiving $Accept(\gamma, seq_\gamma, deps_\gamma, L.i)$:

- 21: $R :: \text{commands}[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, Accepted)$
- 22: reply $AcceptOK(\gamma, L.i)$ to L

Designated leader replica L , on receiving at least $\lfloor N/2 \rfloor$ $AcceptOK$ messages for instance $L.i$:

- 23: reply $RequestReply(\gamma, L.i)$ to client
- 24: run Commit phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$

Commit

Designated leader replica L , for $(\gamma, seq_\gamma, deps_\gamma)$ at instance $L.i$

- 25: $L :: \text{commands}[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, Committed)$
- 26: send $Commit(\gamma, seq_\gamma, deps_\gamma, L.i)$ to all other replicas

Replica R , on receiving $Commit(\gamma, seq_\gamma, deps_\gamma, L.i)$:

27: $R :: commands[L][i] \leftarrow (\gamma, seq_\gamma, deps_\gamma, Committed)$

Explicit Prepare Phase

Replica Q for instance $L.i$ of a potentially failed replica L

28: increment ballot number to $epoch.(b+1).Q$, (where $epoch.b.L$ was the default ballot number for instance $L.i$)

29: send $Prepare(epoch.(b+1).Q, L.i)$ to all replicas (including self)

30: wait for at least $\lfloor N/2 \rfloor + 1$ responses

31: let \mathcal{R} be set of replies w/ the highest ballot number

32: **if** \mathcal{R} contains a $(\gamma, seq_\gamma, deps_\gamma, Committed)$ **then**

33: run Commit phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$

34: **else if** \mathcal{R} contains a $(\gamma, seq_\gamma, deps_\gamma, Accepted)$ **then**

35: run Paxos-Accept phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$

36: **else if** \mathcal{R} contains at least $\lfloor N/2 \rfloor$ identical $(\gamma, seq_\gamma, deps_\gamma, PreAccepted)$ replies, and none is from L **then**

37: run Paxos-Accept phase for $(\gamma, seq_\gamma, deps_\gamma)$ at $L.i$

38: **else if** \mathcal{R} contains at least one $(\gamma, seq_\gamma, deps_\gamma, PreAccepted)$ **then**

39: start Phase 1 for γ at instance $L.i$, avoiding the fast path

40: **else**

41: start Phase 1 for a *no-op* at instance $L.i$, avoiding fast path

Replica R , on receiving $Prepare(epoch.b.Q, L.i)$ from Q

42: **if** $epoch.b.Q$ is larger than the most recent ballot number $epoch.x.Y$ for instance $L.i$ **then**

43: reply $PrepareOK(R :: commands[L][i], epoch.x.Y, L.i)$

44: **else**

45: reply NACK

5.2 The Execution Algorithm

To execute command γ committed in instance $R.i$, a replica will follow these steps:

1. Wait for $R.i$ to be committed (or run an explicit prepare phase to force it);
2. Build γ 's dependency graph by adding γ and all the commands in instances from γ 's dependency list as nodes, with directed edges from γ to these nodes, and then repeating this process recursively for all of γ 's dependencies (starting with step 1);
3. Find the strongly connected components, sort them topologically;
4. In decreasing topological order, for each strongly connected component, do:
 - 4.1 Sort all commands in the strongly connected component by their sequence number;
 - 4.2 Execute every command in increasing sequence number order (if it hasn't already been executed), and mark it as executed.

5.3 Keeping the Dependency List Small

Instead of including all interfering commands, we include only N dependencies in each list: the instance number $R.i$ with the highest i for which the current replica has seen an interfering command (not necessarily committed). If the interference relation is transitive (usually the case in practice) the most recent interfering command suffices, because its dependency graph will contain all commands committed in instances $R.j$, with $j < i$. Otherwise, every replica must assume that any unexecuted commands in previous instances $R.j$ ($j < i$) are possible dependencies and check them locally—a fast operation when commands are executed soon after being committed.

5.4 Recovering from Failures

A replica may need to find out the decision for an instance because it has commands to execute that depend on it. If the replica times out waiting for the commit for that instance, it will try to take ownership of it by running an *Explicit Prepare* phase, at the end of which it will either learn what command was being proposed in the problem instance (in which case it will finalize committing that command), or it will not learn any command (because no other replica has seen it), in which case it will commit a special no-op command to finalize the instance

Another failure-related situation is that where a client timed out waiting for a replica to reply and re-issues its command to a different replica. As a result, the same command can be proposed in two different instances, so every replica must be able to recognize duplicates, and only execute the command once. This situation is not specific to Egalitarian Paxos—it affects any replication protocol. An alternative solution is to make the application tolerant of re-executed commands.

5.5 Joining/Rejoining the Replica Set

We take an approach similar to a Vertical Paxos [6] system with majority read quorums: A new replica, or one that recovers without its memory, must receive a new ID and a new (higher) *epoch* number, e.g., from a configuration service or a human. It then sends *Join* messages to at least $F + 1$ live replicas (that are not themselves in the process of joining). Upon receiving a *Join*, a replica will first update its membership information and then the *epoch* part of each ballot number it uses or expects to receive for new instances—it will thus no longer acknowledge messages for instances initiated in older epochs, that it is not already aware of. It will then send the joining replica the list of committed or ongoing instances that the live replica is aware of. The joining replica becomes live only after receiving commits for all instances included in the replies to at least $F + 1$ *Join* messages. This strategy ensures that a joining replica and a replica that has been excluded from the new configuration cannot participate in voting commands at the same time—thus preserving the property that any two quorums overlap.

Once $F + 1$ live replicas have switched to the new configuration, a replica that has been excluded from the new configuration will not be able to act as command leader (its messages will not be acknowledged by a majority), nor can it participate in successful ballots. The excluded replica might still initiate instances, but they can be finalized only by live replicas that have committed to the new configuration. This eventually stops when all live replicas have committed to the new configuration, or when the excluded replica receives a *Join* message that makes it aware of its exclusion.

The strategy described in this section preserves the protocol guarantees. An instance is either finalized as if the joining replica has not joined yet (if the id of that instance was contained in one of the $F + 1$ replies to *Join* messages), or it is treated like a new command in the new configuration.

5.6 Proof of Properties

We prove that together, the commit protocol and execution algorithm guarantee the properties stated in Section 4.

Theorem 1 (Nontriviality). *Any command committed by any EPaxos replica must have been proposed by a client.*

Proof. For any command that reaches the Commit phase, a replica must have executed the Init phase. Init is only executed for commands proposed by clients. \square

Definition 2. If γ is a command with attributes seq_γ and $deps_\gamma$, we say that the tuple $(\gamma, seq_\gamma, deps_\gamma)$ is *safe* at instance $Q.i$ if $(\gamma, seq_\gamma, deps_\gamma)$ is the only tuple that is or will be committed at $Q.i$ by any replica.

Lemma 1. *EPaxos replicas commit only safe tuples.*

Proof.

1 The same ballot number cannot be used twice in the same instance.

PROOF:

1.1 No two different replicas can use the same ballot number.

PROOF: The ballot number chosen by a replica is based on its id, which is unique.

1.2 A replica never uses the same ballot number twice for the same instance

PROOF:

1.2.1 *Case:* If replicas store the command log in persistent memory, then a replica will never reinitiate the same instance twice with the same ballot number.

1.2.2 *Case:* If a crashed replica can forget the command log, it will be assigned a new id when it recovers.

1.2.3 *Q.E.D.*

Cases 1.2.1 and 1.2.2 are exhaustive.

1.3 *Q.E.D.*

Immediately from 1.1 and 1.2.

2 For any instance $Q.i$ there is at most one attempt (i.e., the *default* ballot $0.Q$) to choose a tuple without running Explicit Prepare first.

PROOF:

2.1 A replica Q starts an instance $Q.i$ at most once.

PROOF: A replica starts an instance only in the Init phase of the algorithm and it increments the instance number atomically every time it executes Init. The instance number never decreases. If a replica loses the content of its memory (e.g., after a crash), it will be assigned a previously unused replica id by a safe external configuration service—so the same instance can never be started twice.

2.2 No replica other than Q can start instance $Q.i$.

PROOF:

2.2.1 A replica with a different id $R \neq Q$ starts only instances $R.i \neq Q.i$

2.2.2 A new replica is never assigned the id of a previously started replica

2.2.3 *Q.E.D.*

Immediately from 2.2.1 and 2.2.2.

2.3 *Q.E.D.*

When not running Explicit Prepare, a replica tries to choose a command in an instance only if it starts that instance, and only for the default ballot. By 2.1 and 2.2, this can happen at most once per instance $Q.i$, in ballot $0.Q$.

3 Let $b_{smallest}$ be the smallest ballot number for which a tuple $(\gamma, seq_\gamma, deps_\gamma)$ has been committed at instance $Q.i$. Then any other commit at instance $Q.i$ commits the same tuple.

PROOF:

By induction on the ballot number b of all ballots committed for $Q.i$:

3.1 Base case: if $b = b_{smallest}$, then the same tuple is committed in both b and $b_{smallest}$.

PROOF:

By 1, b and $b_{smallest}$ must be the same ballots.

3.2 Induction step: if tuple $(\gamma, seq_\gamma, deps_\gamma)$ has been committed in ballot b_1 , then the next higher successful ballot $b > b_1$ will commit the same tuple.

PROOF:

Let b_2 be the next highest ballot number of a ballot attempted at instance $Q.i$. By 2, and since b_2 cannot be the default ballot for $Q.i$ (because there is a ballot b_1 smaller than it), b_2 is attempted after running Explicit Prepare. Furthermore, by the recovery procedure, any ballot attempted after Explicit Prepare must run the Paxos-Accept Phase.

3.2.1 *Case:* Ballot b_1 is committed directly after Phase 1.

Since b_1 is successful after Phase 1, then a fast quorum ($N - 1$ replicas) have recorded the same tuple $(\gamma, seq_\gamma, deps_\gamma)$ for instance $Q.i$. For b_2 to start, its leader must receive replies to *Prepare* messages from at least $\lfloor N/2 \rfloor + 1$ replicas. Therefore, at least $\lfloor N/2 \rfloor$ replicas will see a *Prepare* for b_2 after they have recorded $(\gamma, seq_\gamma, deps_\gamma)$ for ballot b_1 (if they had seen the larger ballot b_2 first, they would not have acknowledged any message for ballot b_1). b_2 's leader will therefore receive at least $\lfloor N/2 \rfloor$ *PrepareReply*'s with tuple $(\gamma, seq_\gamma, deps_\gamma)$ marked as pre-accepted.

If the leader of b_1 is among the replicas that reply to the *Prepare* of ballot b_2 , then it must have replied after the end of Phase 1 (otherwise it couldn't have completed the smaller ballot b_1), so it will have committed tuple $(\gamma, seq_\gamma, deps_\gamma)$ by then. The leader of b_2 will then know it is safe to commit the same tuple.

Below, we assume that the leader of b_1 is *not* among the replicas that reply to the *Prepare* of ballot b_2 .

3.2.1.1 *Subcase:* $N > 3$

The $\lfloor N/2 \rfloor$ replies with tuple $(\gamma, seq_\gamma, deps_\gamma)$ constitute a majority among the first $\lfloor N/2 \rfloor + 1$ *PrepareReply*'s. The leader of ballot b_2 , will therefore be able to identify tuple $(\gamma, seq_\gamma, deps_\gamma)$ as potentially committed, and use it in a Paxos-Accept Phase.

3.2.1.2 Subcase: $N = 3$

$\lfloor N/2 \rfloor = 1$ is not a majority among the first $\lfloor N/2 \rfloor + 1 = 2$ *PrepareReply*'s. However, for $N = 3$, a command leader commits a tuple after Phase 1 only if a *PreAcceptReply* matched the attributes in the initial *PreAccept*. The acceptor that has sent such a *PreAcceptReply* in ballot b_1 will convey this information in a *PrepareReply* for ballot b_2 . The leader of ballot b_2 will therefore use the correct tuple $(\gamma, seq_\gamma, deps_\gamma)$ in a Paxos-Accept Phase.

For ballots higher than b_2 to start, their leaders will follow the recovery procedure, and will receive either the same type of replies received by the leader of b_2 (as above), or it will receive at least one *PrepareReply* from a replica whose highest ballot is b_2 and has marked $(\gamma, seq_\gamma, deps_\gamma)$ as accepted. In either case, by the recovery procedure, the replica trying to take over instance $Q.i$ will have to use tuple $(\gamma, seq_\gamma, deps_\gamma)$ in a Paxos-Accept Phase. By simple induction, any ballot higher than b_1 will use tuple $(\gamma, seq_\gamma, deps_\gamma)$ in a Paxos-Accept Phase, including successful ballots.

3.2.2 Case: Ballot b_1 is committed after the Paxos-Accept Phase.

The tuple $(\gamma, seq_\gamma, deps_\gamma)$ is safe by the guarantees of classic Paxos.

3.2.3 Q.E.D.

Cases 2.2.1 and 2.2.2 are exhaustive.

3.3 Q.E.D.

The induction is complete.

4 Q.E.D.

Immediately from 3.

□

Theorem 2 (Consistency). *Two replicas can never have different commands committed for the same instance.*

Proof. We have already proved a stronger property: by Lemma 1, two replicas can never have different tuples (i.e., commands along with their commit attributes) committed for the same instance. □

Theorem 3 (Stability). *For any replica, the set of committed commands at any time is a subset of the committed commands at any later time. Furthermore, if at time t_1 a replica R has command γ committed at some instance $Q.i$, then R will have γ committed in $Q.i$ at any later time $t_2 > t_1$.*

Proof. By Theorem 2 and the extra assumption that committed commands are recorded in persistent memory. □

So far, we have shown that tuples are committed consistently across replicas. They are also stable, as long as they are recorded in persistent memory. We now show that having consistent attributes committed across all replicas is sufficient to guarantee that all interfering commands are executed in the same order on every replica:

Theorem 4 (Execution consistency). *If two interfering commands γ and δ are successfully committed (not necessarily by the same replica), they will be executed in the same order by every replica.*

Proof.

1 If γ and δ are successfully committed and $\gamma \sim \delta$, then either γ has δ in its dependency list when γ is committed (more precisely, γ has δ 's instance in its dependency list, but, for simplicity of notation, we use a command name to denote the pair comprising the command and the specific instance in which it has been committed), or δ has γ in its dependency list when δ is committed.

PROOF:

1.1 The attributes with which a command c is committed are the union of at least $\lfloor N/2 \rfloor + 1$ sets of attributes computed by as many replicas.

PROOF:

1.1.1 *Case:* c is committed immediately after Phase 1.

$N - 1$ replicas have input their attributes for c .

1.1.2 *Case:* c is committed after the Paxos-Accept phase.

1.1.2.1 *Subcase:* The Paxos-Accept phase starts after the execution of Phase 1.

Phase 1 ends after $\lfloor N/2 \rfloor$ replicas have replied to a *PreAccept* with the command leader's updated attributes (so the attributes are the union of $\lfloor N/2 \rfloor + 1$ sets of attributes, from as many replicas, including the command leader).

1.1.2.2 *Subcase:* The Paxos-Accept phase starts after $\lfloor N/2 \rfloor$ *PrepareReply*'s in the recovery phase, none of which is from the initial command leader.

Then $\lfloor N/2 \rfloor$ replicas, plus the initial command leader ($\lfloor N/2 \rfloor + 1$ replicas in total), have contributed to the set of attributes used for the subsequent Paxos-Accept phase.

1.1.2.3 *Subcase:* The Paxos-Accept phase starts after a *PrepareReply* from a replica R that had marked c as accepted.

Then some replica has to have previously initiated the Paxos-Accept phase that resulted in R receiving an *Accept*, so this subcase is reducible to one of the previous subcases.

1.1.2.4 *Q.E.D.*

The subcases enumerated above describe all possible circumstances in which a command is committed after the Paxos-Accept Phase.

1.1.3 *Case:* γ is committed after the current replica receives a *Commit* for γ from another replica.

The replica that initiates the *Commit* must be in one of the previous two cases.

1.1.4 *Q.E.D.*

The cases enumerated above are exhaustive.

1.2 *Q.E.D.*

By 1.1, at least one replica R contributes for both γ 's and δ 's final attributes. Because R records every command that it sees in its command log, and because $\gamma \sim \delta$, R will include the command it sees first in the dependency list of the command it sees second.

2 *Q.E.D.*

By 1, the final dependency *graphs* of γ and δ are in one of three cases:

2.1 *Case:* γ and δ are both in each other's dependency graph.

Then, by the execution algorithm, γ and δ have each other in their dependency graphs, and moreover, they are in the same strongly connected components of their respective graphs. By the execution algorithm, whenever one command is executed, the other is also executed. Since the execution algorithm is deterministic, and since, by Lemma 1, every replica builds the same dependency graphs for γ and δ , every replica will execute the commands in the same order.

2.2 *Case:* γ is in δ 's dependency graph, but δ is not in γ 's dependency graph.

The commands are in different strongly connected components in δ 's graph, and δ 's component is ordered after γ 's component in reversed topological order.

We show that γ is executed before δ by every replica:

2.2.1 *Subcase:* A replica tries to execute γ first.

The replica will execute γ without having executed δ .

2.2.2 *Subcase:* A replica tries to execute δ first.

By the execution algorithm, the replica will build δ 's dependency graph, which also contains γ in a strongly connected component that is ordered before δ 's component in reversed topological order. Then γ is executed before δ is executed.

2.3 *Case:* δ is in γ 's dependency graph, but γ is not in δ 's dependency graph.

Just like the previous case, with γ and δ interchanged.

2.4 *Q.E.D.*

The above three cases are exhaustive. In all cases, the commands are executed in the same order by every replica.

□

Theorem 5 (Execution linearizability). *If two interfering commands γ and δ are serialized by clients (i.e., δ is proposed only after γ is committed by any replica), then every replica will execute γ before δ .*

Proof.

1 γ will be in δ 's dependency graph.

PROOF:

By the time δ is proposed, γ will have been pre-accepted by at least $\lfloor N/2 \rfloor + 1$ replicas. For δ to be committed, it too has to be pre-accepted by at least $\lfloor N/2 \rfloor + 1$ replicas. Therefore, at least one replica R whose pre-accept is taken into account when establishing δ 's dependency list pre-accepts δ after it has pre-accepted γ . Since $\gamma \sim \delta$, R will put γ in δ 's dependency list.

2 The sequence number with which δ is committed will be higher than that with which γ is committed.

PROOF:

2.1 By the time any replica receives a request for δ from a client, at least $\lfloor N/2 \rfloor + 1$ replicas will have logged the final sequence number for γ .

PROOF:

2.1.1 *Case:* γ is committed directly after Phase 1.

Then $N - 1$ replicas have logged the same sequence number for γ , and this is the sequence number with which γ is committed.

2.1.2 *Case:* γ is committed after the Paxos-Accept Phase.

Then at least $\lfloor N/2 \rfloor + 1$ replicas have logged γ as accepted with its final attributes, including its sequence number.

2.1.3 *Q.E.D.*

Cases 2.1.1 and 2.1.2 are exhaustive.

2.2 Q.E.D.

By 2.1, at least one of the replicas that pre-accepts δ , whose *PreAcceptReply* is taken into account when establishing δ 's final attributes, will update δ 's sequence number to be higher than γ 's final sequence number.

3 Q.E.D.

At any replica R , there are two possible cases:

3.1 Case: R tries to execute γ before it tries to execute δ .

3.1.1 Subcase: δ is in γ 's dependency graph.

Then, by 1, δ and γ are in the same strongly connected component. By the execution algorithm and by 2, γ will be executed before δ .

3.1.2 Subcase: δ is not in γ 's dependency graph.

Then, by the execution algorithm, γ will be executed (at a moment when δ won't have been executed).

3.2 Case: R tries to execute δ before it tries to execute γ .

3.2.1 Subcase: δ is in γ 's dependency graph.

Then, by 1, δ and γ are in the same strongly connected component. By the execution algorithm and by 2, γ will be executed before δ .

3.2.2 Subcase: δ is not in γ 's dependency graph.

Then, by 1, γ is in a different strongly connected component than δ , and γ 's component is first in reversed topological order. By the execution algorithm, γ is executed before δ .

3.3 Q.E.D.

The above cases are exhaustive. In all cases γ is always executed before δ .

□

Finally, **liveness** is guaranteed with high probability as long as a majority of replicas are non-faulty: clients and replicas use time-outs to resend messages, and a client keeps retrying a command until a replica succeeds in committing that command.

5.7 Fast Egalitarian Paxos

We can use the Fast Paxos optimization in EPaxos to decrease the commit latency by one message delay by letting clients broadcast commands to all replicas. We do not explore this because of two main drawbacks: (1) the fast-path quorum size will be $\lceil 3N/4 \rceil$ (as in Generalized Paxos), which is by at least one replica larger than that in Optimized Egalitarian Paxos (which we describe in Section 6); and (2) when building *deps*, we can no longer identify commands by their instance numbers—we must use unique identifiers set by the clients instead.

6 Optimized Egalitarian Paxos

We now describe how Egalitarian Paxos can be enhanced by reducing the fast-path quorum size for increasing its performance: higher throughput and lower latency—including optimal commit latency in the wide area for setups with 3 and 5 replicas.

6.1 Preferred Fast-Path Quorums

Instead of sending *PreAccept* messages to every replica, a command leader sends *PreAccepts* to only those replicas in a fast-path quorum that includes itself. We call this mode of operation *thrifty*. The fast-path quorum can be static per command leader, or it can change for every new command—depending on inter-replica communication latency and dynamic load assessment.

Using this optimization has the immediate benefit of decreasing the overall number of messages processed by the system for each command, thus increasing the system throughput.

Another, less obvious consequence is that we can decrease the fast-path quorum size from $2F$ to $F + \lfloor \frac{F+1}{2} \rfloor$, where F is the maximum number of failures the system can tolerate (the total number of replicas is therefore $N = 2F + 1$). To achieve this, we make two additional modifications to simplified EPaxos:

1. We modify the fast path condition in Phase 2: the command leader commits a command on the fast path if both of the following conditions are fulfilled:

FP-quorum: The command leader receives $F + \lfloor \frac{F+1}{2} \rfloor - 1$ *PreAcceptReply*'s with identical *deps* and *seq* attributes, and

FP-deps-committed: For every command in *deps*, at least one of the replicas in the quorum (including the command leader itself) has recorded that command as Committed—acceptors pass this information to the command leader with at most one bit per each command included in *deps*.

The last condition is necessary for ensuring that the *seq* attribute for every command in *deps* is final (it will not change), and will aid in recovering from failures, as explained in the next subsection.

Alternatively, for up to seven total replicas, we can eliminate the second condition if we make the following modification to the protocol:

Accept-Deps: Every sender of an *Accept* or *AcceptReply* message will attach a dependency list updated right before the *Accept* / *AcceptReply* was sent; every receiver of an *Accept* or *AcceptReply* will store the message in its log permanently.

The updated dependency information in these recorded *Accept* and *AcceptReply* messages is only used during the recovery procedure, and has no role in the execution algorithm. Using this alternative protocol modification has the important benefit that it does not affect the chance of committing on the fast-path.

When using the *Accept-Deps* variant of the protocol, it is important for replicas that have not received a command to not set dependencies on its corresponding instance. This may happen in implementations that use the optimization described in Section 5.3 (which involves setting implicit dependencies on all instances with IDs smaller than a given instance ID) and may result in the approximate sequence number not being updated correctly. There are two ways an implementation can avoid this problem: (a) acceptors attach the ranges of instances that they have not seen to *PreAcceptReplies* or (b) acceptors do not reply to *PreAccepts* before receiving messages for all instances that they would set implicit dependencies on.

2. We modify the recovery procedure (i.e., the Explicit Prepare Phase), which we describe in the next subsection.

The rest of the algorithm remains the same as described in the previous section.

Note that for 3 and 5 replicas, the fast-path quorum sizes become 2 and 3, respectively, which is optimal (just like for classic Paxos).

Another consequence is that for 3 replicas, there is no chance of conflicts, even when all commands interfere. This is because the only reply that the command leader waits for the sole *PreAccept* it send does not have another *PreAcceptReply* to conflict with. As long as there are no failures and replicas reply timely, a 3-replica thrifty EPaxos state machine will commit every command after just one round of communication.

6.2 Failure Recovery in Optimized Egalitarian Paxos

We now describe the new recovery procedure (i.e., the new Explicit Prepare Phase) that allows us to use smaller fast-path quorums.

The recovery procedure guarantees that a command committed on the fast path will be committed even if its command leader and $F - 1$ other replicas have since failed.

Let R be a replica trying to decide instance $Q.i$ of a potentially failed replica Q :

1. R sends *Prepare* messages to all other replicas, with a higher ballot number than the initial ballot number for $Q.i$.

Each replica replies with the information recorded for $Q.i$, if any. R waits for at least $F + 1$ replies (including itself). If R does not receive $F + 1$ ACKS (because some replicas have received messages with higher ballots, and reply with NACKS), R increases the ballot number and retries.

2. **If** no replica has any information about $Q.i$, R exits recovery and starts the process of choosing a *no-op* at $Q.i$ by proposing it in the Paxos-Accept Phase.
3. **If** at least one replica has committed command γ in $Q.i$ (there is at most one such command), with attributes $deps_\gamma$ and seq_γ , R commits γ locally, sends *Commit*($Q.i, \gamma, deps_\gamma, seq_\gamma$) to every other replica, and exits recovery.
4. **If** at least one replica has accepted command $(\gamma, deps_\gamma, seq_\gamma)$ in $Q.i$, R exits recovery and starts a Paxos-Accept Phase for this tuple at $Q.i$ (it will choose the one accepted with the highest ballot number, if there are multiple different accepted tuples at $Q.i$).
5. **If** at least $\lfloor \frac{F+1}{2} \rfloor$ replicas have pre-accepted γ with the same attributes $(\gamma, deps_\gamma, seq_\gamma)$, **in** $Q.i$'s **default ballot** then **goto** 6.

Else R exits recovery and starts the process of choosing γ at $Q.i$, on the slow path (i.e., Phase 1, Phase 2, Paxos-Accept, Commit).

6. R sends *TentativePreAccept*($Q.i, \gamma, deps_\gamma, seq_\gamma$) to all the respondents that have not pre-accepted γ .

When receiving a *TentativePreAccept*($Q.i, \gamma, deps_\gamma, seq_\gamma$) a replica pre-accepts $(\gamma, deps_\gamma, seq_\gamma)$ at $Q.i$ if it has not already recorded an interfering command with conflicting attributes—i.e., any command δ such that:

- i. $\gamma \sim \delta$, and
- ii. $\gamma \notin deps_\delta$, and
- iii. (a) $\delta \notin deps_\gamma$, or

- (b) $\delta \in \text{deps}_\gamma$, but $\text{seq}_\delta \geq \text{seq}_\gamma$ (this subcase has one exception that does not constitute a conflict: δ and γ have the same initial command leader, and δ is recorded as pre-accepted).

Otherwise, if such a command δ with conflicting attributes exists, the receiver of the *TentativePreAccept* replies with NACK, δ 's instance, the identity of the command leader that has sent δ , and the status of δ (pre-accepted, accepted or committed).

7. (a) **If** the total number of replicas that have pre-accepted or tentatively pre-accepted $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$ is at least $F + 1$ (and we can count Q here too, even if it does not reply), R exits recovery and starts a Paxos-Accept Phase for this tuple at $Q.i$.
- (b) **Else if** a *TentativePreAccept* NACK returns a status of committed, R exits recovery and starts the process of choosing γ at $Q.i$, on the slow path.
- (c) **Else if** a *TentativePreAccept* NACK returns an instance not in γ 's dependency list, with a command leader that must have been part of γ 's fast quorum for γ to have been committed on the fast path, then R exits recovery and starts the process of choosing γ at $Q.i$, on the slow path.
- (d) **Else if** There exists command γ_0 such that R has deferred the recovery of γ_0 because of a conflict with γ , and γ_0 's initial command leader must have been part of γ 's fast quorum for γ to have been committed on the fast path, then R exits the recovery of γ and starts the process of choosing γ at $Q.i$, on the slow path.
- (e) **Else** R defers γ 's recovery, and tries to decide one of the uncommitted commands that conflicts with γ .

If $F \leq 3$ and we implement the **Accept-Depts** protocol modification, then step 7 becomes:

- 7'. (a) **If** the total number of replicas that have pre-accepted or tentatively pre-accepted $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$ is at least $F + 1$ (and we can count Q here too, even if it does not reply), R exits recovery and starts a Paxos-Accept Phase for this tuple at $Q.i$.
- (b) **Else if** a *TentativePreAccept* NACK returns a status of committed for a command δ such that $\delta \in \text{deps}_\gamma$ and $\text{seq}_\delta \geq \text{seq}_\gamma$, R checks the additional *Accept* and *AcceptReply* dependencies recorded at F other replicas. If there exists an *Accept* or *AcceptReply* for the tuple with which δ has been committed, such that γ is not part of its additional dependencies, and the sender of that message is part of γ 's fast quorum (i.e., must be part of the quorum for the fast-path hypothesis to hold), then R exits the recovery procedure and starts the process of choosing γ at $Q.i$ on the slow path. If, on the other hand, none of the F replicas has recorded such an *Accept* or *AcceptReply* message, then R resends the *TentativePreAccept* specifying that δ is no longer a conflict for $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$.
- (c) **Else if** a *TentativePreAccept* NACK returns a status of committed, R exits recovery and starts the process of choosing γ at $Q.i$, on the slow path.
- (d) **Else if** a *TentativePreAccept* NACK returns an instance not in γ 's dependency list, with a command leader that must have been part of γ 's fast quorum for γ to have been committed on the fast path, then R exits recovery and starts the process of choosing γ at $Q.i$, on the slow path.
- (e) **Else if** There exists command γ_0 such that R has deferred the recovery of γ_0 because of a conflict with γ , and γ_0 's initial command leader must have been part of γ 's fast quorum for γ to have been committed on the fast path, then R exits the recovery of γ and starts the process of choosing γ at $Q.i$, on the slow path.

- (f) **Else** R defers γ 's recovery, and tries to decide one of the uncommitted commands that conflicts with γ .

This decision process is depicted in Figure 2.

We are now ready to explain why the fast-path quorum must be $F + \lfloor \frac{F+1}{2} \rfloor$: so that the following lemma holds:

Lemma 2. *The recovery procedure for Thrifty Egalitarian Paxos can always make progress (as long as the system is live).*

Proof.

The recovery procedure blocks only if there exist commands c_1, c_2, \dots, c_n such that the recovery for c_i defers to c_{i+1} for any $i = 1..n - 1$, and c_n defers to c_1 . The recovery procedure must assume about every one of these commands that it might have been committed on the fast path.

I Case 1: The chain contains two consecutive commands γ and δ such that $\gamma \notin \text{deps}_\delta$ and $\delta \notin \text{deps}_\gamma$. Let R be a replica trying to recover γ . R must believe that γ may have been committed on the fast path. Eventually, R will defer γ and try to decide δ , and, by our initial assumption, it must believe that δ too may have been committed on the fast path. R must be aware of the following sets and their properties:

1. $RESP_\gamma$, the set of all the replicas in γ 's fast quorum ($QUOR_\gamma$) that have responded to R 's prepare messages, does not include L_γ , the initial command leader for γ (otherwise γ could be decided);
2. $RESP_\delta$, the set of all the replicas in δ 's quorum ($QUOR_\delta$) that have responded to R 's prepare messages, does not include L_δ , the initial command leader for δ ;
3. $|RESP_\gamma| \geq \lfloor \frac{F+1}{2} \rfloor$;
4. $|RESP_\delta| \geq \lfloor \frac{F+1}{2} \rfloor$;
5. $RESP_\gamma \cap RESP_\delta = \emptyset$ (because a replica cannot pre-accept both commands with conflicting attributes);
6. R infers that $L_\gamma \notin QUOR_\delta$ —otherwise δ could not have been committed on the fast path, since L_γ has set conflicting attributes for γ .
7. R infers that $L_\delta \notin QUOR_\gamma$ —otherwise γ could not have been committed on the fast path, since L_δ has set conflicting attributes for δ .
8. Since there are at most F replicas that do not reply to R , and L_γ (the possibly failed command leader for γ) must be one of them (otherwise R could decide γ), by 6, there are at most $F - 1$ replicas that may be part of $QUOR_\delta$ (we denote this superset by $\overline{QUOR_\delta}$) and that R does not receive replies from. Then, for R to believe δ may have been committed on the fast path, it must be the case that $|RESP_\delta| \geq \lfloor \frac{F+1}{2} \rfloor + 1$

By 2, 5 and 7, R must infer that the following sets are disjoint: $\overline{QUOR_\gamma}$ (i.e., the set of replicas that may be part of $QUOR_\gamma$), $RESP_\delta$, and $\{L_\delta\}$. By 8 and our fast-path quorum requirement, the cardinality of the union of these sets must be at least $F + \lfloor \frac{F+1}{2} \rfloor + \lfloor \frac{F+1}{2} \rfloor + 1 + 1 > 2F + 1$. But this is impossible, because this union must be a subset of the replica set, and its cardinality is $2F + 1$. Therefore, some of these sets overlap, so R cannot be simultaneously uncertain about γ and δ . Our assumption that the recovery procedure could deadlock is false.

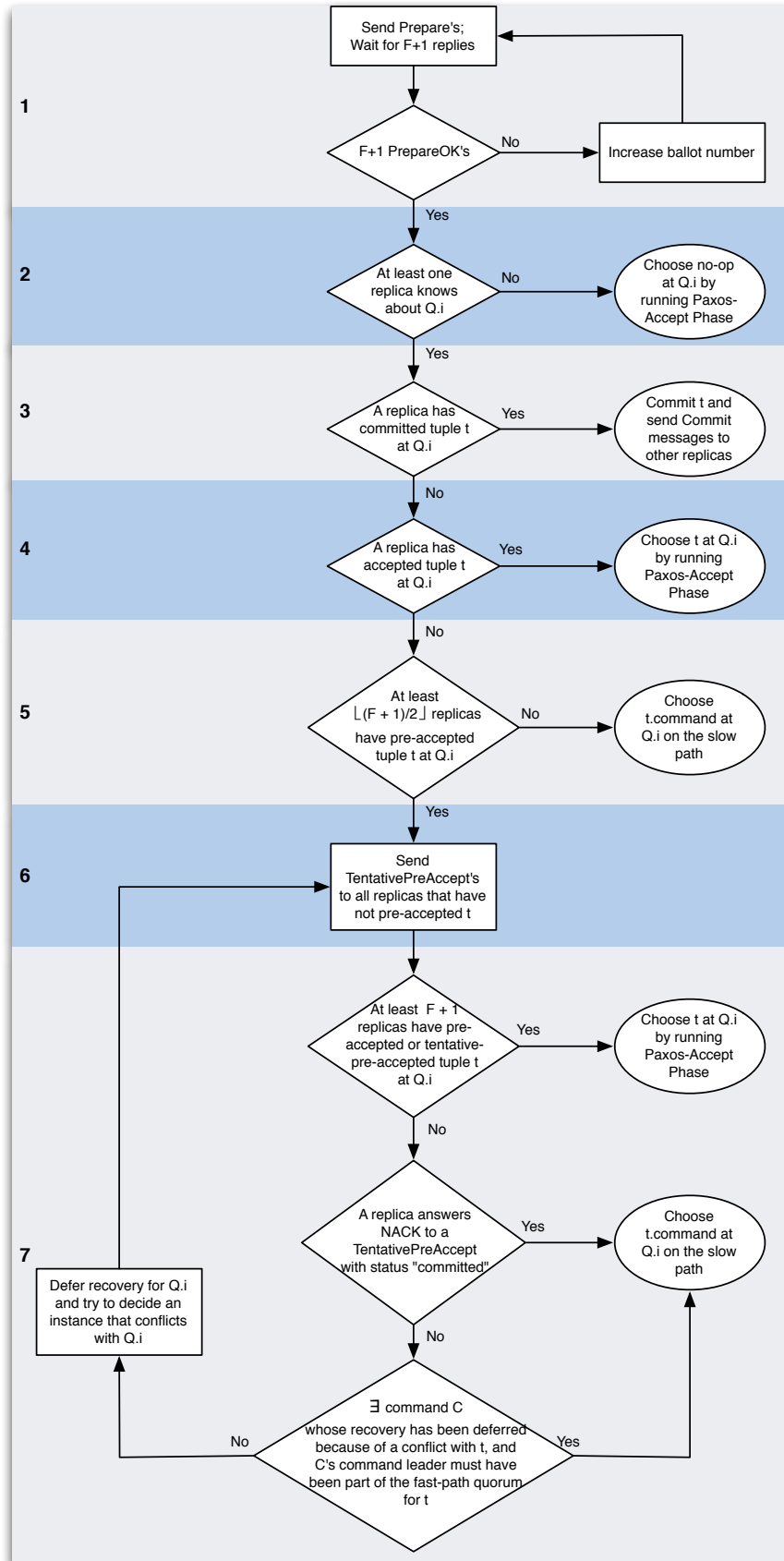


Figure 2: Decision process for recovery in optimized EPaxos.

II Case 2: $c_{i+1} \in \text{deps}_{c_i}$, for all $i = 1..n$ (with $c_{n+1} \equiv c_1$).

For recovery to defer, it must be the case that $c_i \sim c_{i+1}$, $c_i \notin \text{deps}_{c_{i+1}}$, $c_{i+1} \in \text{deps}_{c_i}$, and $\text{seq}_{c_i} \geq \text{seq}_{c_{i+1}}$ for any i (and $c_{n+1} \equiv c_1$). Then $\text{seq}_{c_1} = \text{seq}_{c_2} = \dots = \text{seq}_{c_n}$. Note also that this is only possible for $n \geq 3$.

1. **Subcase:** There exist c_i and c_{i+1} that have the same initial command leader.

Since c_{i+1} has not been decided, it must be the case that all the information available about c_{i+1} is that it has been pre-accepted by various replicas. By our definition of conflicts (step 6 of the recovery procedure, point iii.(b)), c_{i+1} will not conflict with c_i , which contradicts our assumption for this subcase.

2. **Subcase:** No two consecutive commands in the chain have the same command leader.

As noted earlier, it must be the case that $n \geq 3$, otherwise there is no conflict. Consider commands c_1, c_2 and c_3 . It is impossible that any replica that replies to *Prepare* messages has pre-accepted two consecutive commands in the chain (because then they would not have been pre-accepted with the same sequence numbers). Furthermore, at least $\lfloor \frac{F+1}{2} \rfloor$ responding replicas have pre-accepted each command in the chain (for c_i we denote this set of replicas as $RESP_{c_i}$). We show that the recovery procedure concludes that either c_1 's leader must have been part of c_2 's fast quorum or that c_2 's leader must have been part of c_3 's fast quorum:

2.1 **Sub-subcase:** F is odd.

Then $2F + 1 - \lfloor \frac{F+1}{2} \rfloor = F + \lfloor \frac{F+1}{2} \rfloor$, and since c_2 has not been pre-accepted by the $\lfloor \frac{F+1}{2} \rfloor$ replicas in $RESP_{c_1}$, all the other replicas, including c_1 's leader, must have been part of c_2 's fast quorum.

2.2 **Sub-subcase:** F is even.

Let *LIVE* be a set of $F + 1$ replicas that respond to *Prepare* messages (more replica may reply, we only consider $F + 1$ of them). No more than $\lfloor \frac{F+1}{2} \rfloor + 1$ of the replicas in *LIVE* can be part of any one command's fast quorum, because at least $\lfloor \frac{F+1}{2} \rfloor$ will be part of the fast quorum for the subsequent command in the chain.

If $|LIVE \cap RESP_{c_2}| = \lfloor \frac{F+1}{2} \rfloor$, then all replicas outside *LIVE* must have been part of c_2 's fast quorum, including c_1 's leader.

If $|LIVE \cap RESP_{c_2}| = 1 + \lfloor \frac{F+1}{2} \rfloor$, then $|LIVE \cap RESP_{c_3}| = \lfloor \frac{F+1}{2} \rfloor$, so all replicas outside *LIVE* must have been part of c_3 's fast quorum, including c_2 's leader.

The sub-subcases enumerated above are exhaustive. In all situations, the leader of a command c_i must have been part of the fast quorum for c_{i+1} . It is therefore impossible for c_{i+1} to have been committed on the fast-path, since c_i 's leader couldn't have pre-accepted c_{i+1} with the same sequence number as c_i without adding c_i to c_{i+1} 's dependency list. As per 7 (or 7') in the recovery procedure, the recovery procedure will eventually abandon the fast-path recovery for c_{i+1} .

The subcases enumerated above are exhaustive. □

Finally, we show that the recovery procedure is correct. We start by showing that it commits only safe tuples:

Theorem 6. *The Optimized Egalitarian Paxos recovery procedure commits only safe tuples.*

Proof.

Assume the recovery procedure is trying to recover instance $Q.i$. We show that the tuple that it commits at $Q.i$ is safe.

1 *Case*: No tuple is committed at instance $Q.i$ before the recovery procedure commits a tuple at $Q.i$.
 In all cases, the recovery procedure ends by choosing a tuple on the slow path, by running classic Paxos. The tuple is thus safe by the classic Paxos guarantees.

2 *Case*: A tuple $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$ has been committed at $Q.i$ before the recovery procedure terminates.

2.1 *Subcase*: $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$ has previously been committed on the slow path.

Then there must be at least $F + 1$ replicas that have accepted $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$. Since the recovery procedure terminates by running classic Paxos in all cases, it will use the same tuple in a Paxos-Accept Phase. By the guarantees of the classic Paxos algorithm, only this tuple can ever be committed at $Q.i$.

2.2 *Subcase*: $(\gamma, \text{deps}_\gamma, \text{seq}_\gamma)$ has previously been committed on the fast path.

Then there must be $F + \lfloor \frac{F+1}{2} \rfloor$ replicas that have pre-accepted this tuple at $Q.i$ before processing the *Prepares* of the recovery procedure (otherwise the initial command leader would have received NACKs for the initial *PreAccepts* and not taken the fast path). Since at most F replicas can be faulty, the recovery procedure will take into account the *PrepareReply*'s of at least $\lfloor \frac{F+1}{2} \rfloor$ of them, and by step 5 of the recovery procedure, it will try to obtain a quorum for this tuple. We show that it will succeed:

2.2.1 No interfering command $\delta \sim \gamma$, can be committed such that $\delta \notin \text{deps}_\gamma$ and $\gamma \notin \text{deps}_\delta$.

PROOF: δ must be pre-accepted by a majority of replicas, and that majority will intersect γ 's quorum (itself a majority) in at least one replica, which will ensure that at least one command will be in the other's *deps* set.

2.2.2 If the protocol does not implement **Accept-Depts**, but does implement **FP-deps-committed**, then no interfering command $\delta \sim \gamma$, $\delta \in \text{deps}_\gamma$, can be committed such that $\gamma \notin \text{deps}_\delta$ and $\text{seq}_\delta \geq \text{seq}_\gamma$.

PROOF:

We prove this by generalized induction. The relation that we run the induction on is $a \prec b \equiv$ "command a has been committed (in a particular instance) by the recovery procedure for the first time before command b has been committed (in a particular instance) by the recovery procedure for the first time".

2.2.2.1 *Base case*: Let γ_0 be the first command initially committed on the fast path and then committed again as a result of the recovery procedure (or one of the first, if multiple such commands are committed at the exact same time).

Assume there existed $\delta \sim \gamma_0$, $\delta \in \text{deps}_{\gamma_0}$, committed such that $\gamma_0 \notin \text{deps}_\delta$ and $\text{seq}_\delta \geq \text{seq}_{\gamma_0}$ at the time of γ_0 's recovery. Since γ_0 had been committed on the fast path, then by the *additional condition for the fast-path in optimized EPaxos (FP-deps-committed)*, all its dependencies, including δ must have been committed before seq_{γ_0} had been computed. Then, δ must have been committed again in the meantime with different attributes (thus breaking safety). But by Lemma 1, 1, 2.1, 2.2.1, and the recovery procedure, this could only have occurred if δ had been committed incorrectly by the recovery procedure (before γ_0), after initially having been committed on the fast-path—all other commit paths preserve safety. By our base case assumption, this is impossible, since $\gamma_0 \prec \delta$.

2.2.2.2 *Induction step*: The property holds for γ if it holds for every $\delta \prec \gamma$.

Assume there exists $\delta \sim \gamma$, $\delta \in \text{deps}_\gamma$, committed such that $\gamma \notin \text{deps}_\delta$ and $\text{seq}_\delta \geq \text{seq}_\gamma$. Since γ has been committed on the fast path, then, by the additional condition for the fast-path in optimized EPaxos, all its dependencies, including δ must have been committed

before seq_γ had been computed. Then, δ must have been committed again with different attributes (thus breaking safety). But by Lemma 1, 1, 2.1, 2.2.1 and the recovery procedure, this could only occur if δ has been committed incorrectly by the recovery procedure after initially having been committed on the fast-path—we have shown that all other commit paths preserve safety. Since γ has not been committed by the recovery procedure yet, $\delta \prec \gamma$. By the induction hypothesis and by 2.2.1, the recovery procedure would have exited δ 's recovery by correctly committing its initial fast-path attributes. Then seq_δ cannot be larger or equal to seq_γ , since seq_γ has been updated to be larger than seq_δ at δ 's initial commit time.

2.2.2.3 Q.E.D

The induction is complete.

2.2.3 If the protocol does not implement **FP-deps-committed**, but does implement **Accept-Deps**, the recovery procedure will not exit on the first Else case of step 7'.

PROOF: For the recovery procedure to exit on the first Else case of step 7', there must exist a committed tuple $(\delta, deps_\delta, seq_\delta)$, with $\delta \sim \gamma$, $\delta \in deps_\gamma$, $\gamma \notin deps_\delta$ and $seq_\delta \geq seq_\gamma$, and there must exist an *Accept* or *AcceptReply* message for this tuple sent by a replica R' in γ 's fast-path quorum such that the additional dependencies for this message do not include γ . Then R' must have accepted $(\delta, deps_\delta, seq_\delta)$ before pre-accepting γ . This is impossible, since then R' would not have pre-accepted γ with $seq_\gamma \leq seq_\delta$, as we know it must have for γ to be committed on the fast-path.

2.2.4 The recovery procedure will not exit on branches 7.c or 7d (7'.d or 7'.e, respectively): No replica in γ 's fast quorum can start instances for commands that interfere with γ and set conflicting attributes (as per the definition of conflicting attributes in step 6 of the recovery procedure), because all these replicas have pre-accepted γ with its attributes.

2.2.5 Q.E.D

By the recovery procedure, 2.2.1, 2.2.2, 2.2.3, 2.2.4 and Lemma 2 the recovery procedure will be successful in getting F replicas to pre-accept tuple $(\gamma, deps_\gamma, seq_\gamma)$ (not counting the implicit pre-accept of the initial command leader), and it will start the Paxos-Accept Phase for this tuple.

2.3 Q.E.D

Subcases 2.1 and 2.2 are exhaustive and safety is preserved in both.

3 Q.E.D.

Cases 1 and 2 are exhaustive and safety is preserved in both.

□

Next, we show that the recovery procedure preserves execution consistency:

Theorem 7. *The Optimized Egalitarian Paxos preserves execution consistency.*

Proof.

Let γ and δ be two commands that interfere and have been committed. We show that all replicas execute γ and δ in the same order.

1 *Case:* Both γ and δ have first been committed by their respective command leaders, without running the recovery procedure.

This is no different from simplified EPaxos: the different fast-path condition influences only the recovery path. By Theorem 6 and Theorem 4, γ and δ will be executed in the same order by every replica.

2 *Case*: γ is first committed as a result of the recovery procedure, while δ is first committed by its initial command leader without running the recovery procedure.

2.1 *Subcase*: γ is committed before step 7 of the recovery procedure, or after exiting one of the Else branches in step 7.

Then γ must have been pre-accepted by a majority of replicas and then committed after running the Paxos-Accept Phase. This too is reducible to the simple EPaxos case, so, by Theorem 4, γ and δ will be executed in a consistent order across all non-faulty replicas.

2.2 *Subcase*: γ is committed after exiting the recovery procedure on the If branch in step 7.

We show that either γ has δ as a dependency or δ has γ as a dependency:

2.2.1 *Sub-subcase*: γ had been pre-accepted with $\delta \in \text{deps}_\gamma$.

γ 's pre-accepted attributes as received in the recovery procedure at step 7 do not change, so γ will be committed with δ as a dependency.

2.2.2 *Sub-subcase*: γ had been pre-accepted with $\delta \notin \text{deps}_\gamma$.

Since the recovery procedure exits on the If branch of step 7, at least $F + 1$ replicas, including γ 's original command leader have pre-accepted γ as a result of a *PreAccept* or a *TentativePreAccept*. δ will also have been pre-accepted by a majority of replicas, so there is at least one replica that has pre-accepted both δ and γ , and whose replies are taken into account both when establishing δ 's commit attributes and in the recovery procedure for γ . Let this replica be R :

2.2.2.1 *Sub-sub-subcase*: R pre-accepts γ as a result of receiving a *PreAccept* from γ 's initial command leader.

Then R must have learned about γ before receiving a *PreAccept* for δ , so $\gamma \in \text{deps}_\delta$.

2.2.2.2 *Sub-sub-subcase*: R pre-accepts γ after receiving a *TentativePreAccept* during the recovery procedure.

Then, according to the conditions in step 6 of the recovery procedure, either R had already pre-accepted δ such that $\gamma \in \text{deps}_\delta$, or δ reaches R after the *TentativePreAccept* for γ . In either case, $\gamma \in \text{deps}_\delta$ when δ commits.

In conclusion $\gamma \in \text{deps}_\delta$

2.2.3 *Q.E.D.*

Sub-subcases 2.2.1 and 2.2.3 are exhaustive.

By step 2 of the proof for Theorem 4, since at least one command is committed with the other in its dependency list, every replica will execute the commands in the same order.

3 *Case*: δ is first committed as a result of the recovery procedure, while γ is first committed by its initial command leader without running the recovery procedure.

Just like case 2, with γ and δ interchanged.

4 *Case*: Both γ and δ are first committed after the recovery procedure.

If at least one of the commands is committed before step 7 in the recovery procedure, or after exiting step 7 on one of the Else branches, the situation is reducible to one of the previous cases.

The only remaining subcase is that when both commands are committed after exiting step 7 on the If branch. Assume no command has the other in its dependency list when exiting step 7 of the recovery procedure. But each command has been pre-accepted by a majority of replicas (either as a result of *PreAccepts* or

TentativePreAccepts). Then there must be at least one replica R that pre-accepts both commands, and whose replies are taken into account when establishing each command’s commit attributes. If R pre-accepts γ before δ , then, by the conditions in step 6 of the recovery procedure, R will not acknowledge δ without a dependency for γ (and vice-versa). This contradicts our assumption.

Then at least one command is in the other’s dependency list, and by step 2 in the proof for Theorem 4, the commands will be executed in the same order on every replica.

5 Q.E.D

Cases 1, 2, 3 and 4 are exhaustive. □

Finally, we show that the recovery procedure preserves execution linearizability:

Theorem 8. *The Optimized Egalitarian Paxos preserves execution linearizability.*

Proof. Let γ and δ be two interfering commands serialized by clients: δ is proposed only after a replica has committed γ . We show that γ will always be executed before δ

By the time δ is proposed, a majority of replicas have either pre-accepted or accepted γ with its final (commit) attributes. At least one of these replicas will pre-accept δ as a result of receiving a *PreAccept* or a *TentativePreAccept*, and its reply will be considered in deciding δ ’s final attributes. Let this replica be R :

1 Case: R receives a *PreAccept* for δ .

Then R will put γ in $deps_\delta$ and it will increment seq_δ to be larger than seq_γ . Since R ’s reply is considered when deciding δ ’s final attributes, δ ’s dependency list will include γ and its sequence number will be larger than γ ’s at commit time. By the execution algorithm, γ will always be executed before δ .

2 Case: R receives a *TentativePreAccept* for δ at some point other than step 7’ in the recovery procedure.

Since R will ACK (otherwise δ would not be committed), and $\delta \notin deps_\gamma$ (because δ was proposed after γ was committed), by the conditions in step 6 of the recovery procedure, it must hold that $\gamma \in deps_\delta$ and $seq_\gamma < seq_\delta$. By the execution algorithm, γ will always be executed before δ .

3 Case: R receives a *TentativePreAccept* for δ at step 7’ in the recovery procedure.

In this case, there must exist a command γ' such that $\gamma' \sim \delta$, $\gamma' \in deps_\delta$, $\delta \notin deps_{\gamma'}$ and $seq_\delta \leq seq_{\gamma'}$, and R is instructed to ignore the conflict between the attributes of δ and those of γ' . We show that $\gamma \neq \gamma'$.

Assume $\gamma' = \gamma$. Then γ must have first been committed on the slow-path, because otherwise δ couldn’t have acquired a stale dependency on γ (i.e., a dependency where seq_δ hasn’t been updated to be larger than the committed seq_γ)—this remains true for implicit dependencies, as per the discussion in Section 6.1. Then there must exist a replica R' that has participated in the Paxos-Accept phase for γ , and that was also supposed to be part of δ ’s fast-path quorum. We first note that $L_\delta \neq R'$, where L_δ is the initial command leader for δ , because it wouldn’t have set the conflicting sequence number for δ after having accepted the commit-time attributes for γ .

3.1 Subcase $F = 2$.

The recovery procedure for δ must not have received a response from L_δ (otherwise it would have exited before step 7’). Then there is at most one more replica that fails to respond during the recovery procedure, and this replica must be R' (otherwise R' would have replied with non-conflicting attributes for δ , and the recovery procedure would have ended).

- 3.1.1 *Sub-subcase* R' was the leader of the Paxos-Accept phase that first committed γ . Then it would have sent *AcceptReply* messages to two other replicas, and these messages would have included additional dependencies (as per **Accept-deps**) that would not have included δ (because δ had not been proposed at this point). Since at most $F = 2$ replica are faulty, these acceptors must be part of the replicas that reply during the recovery for δ , so the recovery procedure will discover the recorded *AcceptReplies* and will not send the *TentativePreAccept* in step 7'.
- 3.1.2 *Sub-subcase* R' was an acceptor in the Paxos-Accept phase that first committed γ . Then the leader for this Paxos-Accept phase must have responded during the recovery for δ with the *Accept* message from R' . The additional dependencies for this message would not contain δ , so the recovery procedure will not send the *TentativePreAccept* in step 7'.

3.2 *Subcase* $F = 3$.

In this case, the recovery procedure for δ reaches the Else branches of step 7' only if exactly two replicas part of δ 's fast-path quorum reply to *Prepare* messages (otherwise the recovery procedure either exits before step 7', or on the If branch of step 7'). Then of the three un-responsive replicas, one is L_δ , and the other two must be part of both δ 's fast-path quorum and the Paxos-Accept quorum that has first committed γ . Then at least one of the two other replicas part of γ 's Paxos-Accept quorum must have received and recorded an *Accept* or *AcceptReply* message with additional dependencies that did not include δ . Since these replicas are both responsive during the recovery for δ , the recovery procedure could not have sent the *TentativePreAccept* in step 7'.

- 3.3 *Q.E.D.* Step 7' can be executed only if $F = 2$ or $F = 3$. Both subcases have ended in contradiction, so our assumption that $\gamma = \gamma'$ is false. Then R cannot receive a *TentativePreAccept* that forces it to pre-accept δ such that $seq_\delta \leq seq_\gamma$. By 3.1, 3.2 and the execution algorithm, γ will always be executed before δ .

3 *Q.E.D*

Cases 1 and 2 and 3 are exhaustive.

□

7 Linearizability and Serializability in Egalitarian Paxos

Egalitarian Paxos guarantees that interfering commands are executed in a linearizable order. We now show what the overall consistency properties of a system using Egalitarian Paxos are.

If the interference relation is transitive (which is equivalent with a setting where each command refers to the state of a single object) then execution linearizability implies the whole system is linearizable, as defined by Herlihy and Wing: linearizability is a *local* property, meaning that “a system is linearizable if each individual object is linearizable” [4].

The corresponding consistency property in a system where each command can update and read multiple objects is strict serializability. This is a setting where the interference relation is not necessarily transitive. EPaxos does not guarantee strict serializability without a simple modification, described later in this section.

For example, imagine a system with two distinct objects A and B, and two concurrent clients *client1* and *client2*. If each client issues the following commands sequentially (i.e., they wait for the a command to be committed before issuing the next)

client 1: *update A; update B*
 client 2: *read B; read A*

it is impossible for *client2* to see an updated *B* and an unmodified *A*. In this case, all operations are linearizable, because they each refer to a single object. However, if *client2* were to issue the composite command *read A and B* instead, it is possible that it will see an updated *B* and an unmodified *A*:

client 1: *update A; update B*
 client 2: *read A and B*

This is because the read may be concurrent with both updates, and since the updates do not interfere, the system can choose the following ordering: *update B; read A and B; update A*.

We can modify EPaxos to guarantee strict serializability: clients are sent the commit notification for a command only after all instances in the dependencies graph of the current command have been committed. We call this version *EPaxos-strict*. We show that it guarantees strict serializability.

Theorem 9 (Strict serializability). *Every execution in EPaxos-strict is equivalent to a serial execution where non-concurrent commands are executed in their temporal order.*

Proof.

- 1 We define $a <_i b$ to be true if commands a and b interfere, and a has been executed by a replica before b (by execution consistency, if $a <_i b$ then every replica will execute a before b). It is easy to see that $<_i$ is a partial order relation.
- 2 Because commands are executed serially, every execution in EPaxos-strict (as well as in EPaxos) defines a total order that is an extension of $<_i$ (by execution consistency).
- 3 We define $a <_{time} b$ to be true if any client has been notified that a has been committed (along with its entire dependency graph) before b is proposed.
- 4 We show that $<_i \cup <_{time}$ is a partial order relation. We define $<$ to be $(<_i \cup <_{time})$. Assume $<$ is not a partial order relation. Then there exists a sequence of commands c_1, \dots, c_n such that $c_1 < c_2 < \dots < c_n < c_1$. Consider the shortest such cycle.
 - 4.1 Because $<_i$ is a partial order relation, there must be at least one pair of consecutive commands c_{j-1} and c_j such that $c_{j-1} <_i c_j$ is not true. Then it must be true that $c_{j-1} <_{time} c_j$. By re-indexing, let these two commands be c_1 and c_2 .
 - 4.2 The cycle must have more than two commands: it is impossible for $c_2 <_{time} c_1$ since $c_1 <_{time} c_2$; $c_2 <_i c_1$ implies that c_2 is in c_1 's dependency graph which also contradicts $c_1 <_{time} c_2$, since c_1 could not have ended before its entire dependency graph was committed (by EPaxos-strict).
 - 4.3 There is no $j > 2$ such that $c_{j-1} <_i c_j$ does not hold: Assume there is such a j . Then $c_{j-1} <_{time} c_j$. If c_1 ends (i.e., its corresponding client is notified) before c_j starts, then $c_1 <_{time} c_j$, and we have the shorter cycle c_1, c_j, \dots, c_n . If c_1 ends after c_j starts, then c_2 must start after c_{j-1} ends, so $c_{j-1} <_{time} c_2$, and we have the shorter cycle c_2, \dots, c_{j-1} .
 - 4.4 By 4.3 $c_2 <_i \dots <_i c_n <_i c_1$. But $c_{j-1} <_i c_j$ implies that c_{j-1} is in c_j 's dependency graph. Since this relation is transitive, it must hold that c_2 is in c_1 's dependency graph. But, by EPaxos-strict, this implies that c_2 was committed before the commit notification for c_1 was sent to its corresponding client. This contradicts $c_1 <_{time} c_2$.
- 5 By the definition of *interference*, all serial executions that extend $<_i$ are compatible—i.e., they produce the same state and read results. In particular, any execution that extends the partial order $<_i \cup <_{time}$ will be compatible with any serial execution that extends $<_i$. Then by 2, all executions in EPaxos-strict are compatible with a serial execution that extends $<_i \cup <_{time}$. *Q.E.D*

Client notification	Interference must be transitive	Consistency guarantee
After commit	Yes	Linearizability
After commit	No	Per-object linearizability
After command and entire dependency graph have been committed	No	Strict serializability

Table 1: Consistency guarantees in EPaxos.

□

Table 1 summarizes the consistency guarantees of EPaxos.

Because the dependency graph of a command is complete (all the dependencies are committed) by the time the client receives a notification for it, EPaxos-strict does not require approximate sequence numbers to guarantee linearizability, and this simplifies the recovery procedure. The drawback of EPaxos-strict is that clients may experience higher latency for conflicting commands.

8 Conclusion

We have presented a proof of correctness for Egalitarian Paxos, a new state machine replication protocol based on Paxos. EPaxos’s decentralized and uncoordinated design, as well as its small fast-path quorum size, have important benefits for the availability, performance and performance stability of both local and wide area replication.

References

- [1] Michael Burrows. The Chubby lock service for loosely-coupled distributed systems. In *Proc. 7th USENIX OSDI*, Seattle, WA, November 2006.
- [2] Tushar D. Chandra, Robert Griesemer, and Joshua Redstone. Paxos made live: an engineering perspective. In *Proc. 26th ACM SOSP, PODC ’07*, pages 398–407, New York, NY, USA, 2007. ACM.
- [3] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *J. ACM*, 32(2):374–382, April 1985.
- [4] Maurice P. Herlihy and Jeannette M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.*, 12(3), July 1990.
- [5] Patrick Hunt, Mahadev Konar, Flavio P. Junqueira, and Benjamin Reed. ZooKeeper: wait-free coordination for internet-scale systems. In *Proc. USENIX ATC, USENIXATC’10*, Berkeley, CA, USA, 2010. USENIX Association.
- [6] Leslie Lamport, Dahlia Malkhi, and Lidong Zhou. Vertical Paxos and primary-backup replication. Technical report, Microsoft Research, 2009.
- [7] John MacCormick, Nick Murphy, Marc Najork, Chandramohan A. Thekkath, and Lidong Zhou. Boxwood: abstractions as the foundation for storage infrastructure. In *Proc. 6th USENIX OSDI*, San Francisco, CA, December 2004.
- [8] Yanhua Mao, Flavio P. Junqueira, and Keith Marzullo. Mencius: building efficient replicated state machines for WANs. In *Proc. 8th USENIX OSDI*, pages 369–384, San Diego, CA, December 2008.

APPENDIX

TLA+ specification for the Egalitarian Paxos commit protocol:

```

┌────────────────────────── MODULE EgalitarianPaxos ───────────────────────────┐
EXTENDS Naturals, FiniteSets
└────────────────────────────────────────────────────────────────────────────────┘

Max(S)  $\triangleq$  IF S = {} THEN 0 ELSE CHOOSE i ∈ S : ∀ j ∈ S : j ≤ i

Constant parameters:
  Commands: the set of all possible commands
  Replicas: the set of all EPaxos replicas
  FastQuorums(r): the set of all fast quorums where r is a command leader
  SlowQuorums(r): the set of all slow quorums where r is a command leader

CONSTANTS Commands, Replicas, FastQuorums(-), SlowQuorums(-), MaxBallot
ASSUME IsFiniteSet(Replicas)

Quorum conditions: (simplified)

ASSUME ∀ r ∈ Replicas :
  ∧ SlowQuorums(r) ⊆ SUBSET Replicas
  ∧ ∀ SQ ∈ SlowQuorums(r) :
    ∧ r ∈ SQ
    ∧ Cardinality(SQ) = (Cardinality(Replicas) ÷ 2) + 1

ASSUME ∀ r ∈ Replicas :
  ∧ FastQuorums(r) ⊆ SUBSET Replicas
  ∧ ∀ FQ ∈ FastQuorums(r) :
    ∧ r ∈ FQ
    ∧ Cardinality(FQ) = (Cardinality(Replicas) ÷ 2) +
      ((Cardinality(Replicas) ÷ 2) + 1) ÷ 2

Special none command

none  $\triangleq$  CHOOSE c : c ∉ Commands

The instance space

Instances  $\triangleq$  Replicas × (1 .. Cardinality(Commands))

The possible status of a command in the log of a replica.

```

$Status \triangleq \{ \text{"not-seen"}, \text{"pre-accepted"}, \text{"accepted"}, \text{"committed"} \}$

All possible protocol messages:

$Message \triangleq$

- $[type : \{ \text{"pre-accept"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas,$
 $cmd : Commands \cup \{ none \}, deps : \text{SUBSET } Instances, seq : Nat]$
- $\cup [type : \{ \text{"accept"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas,$
 $cmd : Commands \cup \{ none \}, deps : \text{SUBSET } Instances, seq : Nat]$
- $\cup [type : \{ \text{"commit"} \},$
 $inst : Instances, ballot : Nat \times Replicas,$
 $cmd : Commands \cup \{ none \}, deps : \text{SUBSET } Instances, seq : Nat]$
- $\cup [type : \{ \text{"prepare"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas]$
- $\cup [type : \{ \text{"pre-accept-reply"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas,$
 $deps : \text{SUBSET } Instances, seq : Nat, committed : \text{SUBSET } Instances]$
- $\cup [type : \{ \text{"accept-reply"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas]$
- $\cup [type : \{ \text{"prepare-reply"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas, prev_ballot : Nat \times Replicas,$
 $status : Status,$
 $cmd : Commands \cup \{ none \}, deps : \text{SUBSET } Instances, seq : Nat]$
- $\cup [type : \{ \text{"try-pre-accept"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas,$
 $cmd : Commands \cup \{ none \}, deps : \text{SUBSET } Instances, seq : Nat]$
- $\cup [type : \{ \text{"try-pre-accept-reply"} \}, src : Replicas, dst : Replicas,$
 $inst : Instances, ballot : Nat \times Replicas, status : Status \cup \{ \text{"OK"} \}]$

Variables:

$cmdLog$ = the commands log at each replica
 $proposed$ = command that have been proposed
 $executed$ = the log of executed commands at each replica
 $sentMsg$ = sent (but not yet received) messages
 $crtInst$ = the next instance available for a command leader
 $leaderOfInst$ = the set of instances each replica has started but not yet finalized
 $committed$ = maps commands to set of commit attributes tuples
 $ballots$ = largest ballot number used by any replica
 $preparing$ = set of instances that each replica is currently preparing (*i.e.* recovering)

VARIABLES $cmdLog, proposed, executed, sentMsg, crtInst, leaderOfInst,$

committed, ballots, preparing

$$\begin{aligned}
\text{TypeOK} &\triangleq \\
&\wedge \text{cmdLog} \in [\text{Replicas} \rightarrow \text{SUBSET} [\text{inst} : \text{Instances}, \\
&\hspace{15em} \text{status} : \text{Status}, \\
&\hspace{15em} \text{ballot} : \text{Nat} \times \text{Replicas}, \\
&\hspace{15em} \text{cmd} : \text{Commands} \cup \{\text{none}\}, \\
&\hspace{15em} \text{deps} : \text{SUBSET} \text{Instances}, \\
&\hspace{15em} \text{seq} : \text{Nat}]] \\
&\wedge \text{proposed} \in \text{SUBSET} \text{Commands} \\
&\wedge \text{executed} \in [\text{Replicas} \rightarrow \text{SUBSET} (\text{Nat} \times \text{Commands})] \\
&\wedge \text{sentMsg} \in \text{SUBSET} \text{Message} \\
&\wedge \text{crtInst} \in [\text{Replicas} \rightarrow \text{Nat}] \\
&\wedge \text{leaderOfInst} \in [\text{Replicas} \rightarrow \text{SUBSET} \text{Instances}] \\
&\wedge \text{committed} \in [\text{Instances} \rightarrow \text{SUBSET} ((\text{Commands} \cup \{\text{none}\}) \times \\
&\hspace{15em} (\text{SUBSET} \text{Instances}) \times \\
&\hspace{15em} \text{Nat})] \\
&\wedge \text{ballots} \in \text{Nat} \\
&\wedge \text{preparing} \in [\text{Replicas} \rightarrow \text{SUBSET} \text{Instances}] \\
\text{vars} &\triangleq \langle \text{cmdLog}, \text{proposed}, \text{executed}, \text{sentMsg}, \text{crtInst}, \text{leaderOfInst}, \\
&\hspace{10em} \text{committed}, \text{ballots}, \text{preparing} \rangle
\end{aligned}$$

Initial state predicate

$$\begin{aligned}
\text{Init} &\triangleq \\
&\wedge \text{sentMsg} = \{\} \\
&\wedge \text{cmdLog} = [r \in \text{Replicas} \mapsto \{\}] \\
&\wedge \text{proposed} = \{\} \\
&\wedge \text{executed} = [r \in \text{Replicas} \mapsto \{\}] \\
&\wedge \text{crtInst} = [r \in \text{Replicas} \mapsto 1] \\
&\wedge \text{leaderOfInst} = [r \in \text{Replicas} \mapsto \{\}] \\
&\wedge \text{committed} = [i \in \text{Instances} \mapsto \{\}] \\
&\wedge \text{ballots} = 1 \\
&\wedge \text{preparing} = [r \in \text{Replicas} \mapsto \{\}]
\end{aligned}$$

Actions

$$\begin{aligned}
\text{StartPhase1}(C, \text{cleader}, Q, \text{inst}, \text{ballot}, \text{oldMsg}) &\triangleq \\
&\text{LET } \text{newDeps} \triangleq \{\text{rec.inst} : \text{rec} \in \text{cmdLog}[\text{cleader}]\} \\
&\hspace{2em} \text{newSeq} \triangleq 1 + \text{Max}(\{t.\text{seq} : t \in \text{cmdLog}[\text{cleader}]\}) \\
&\hspace{2em} \text{oldRecs} \triangleq \{\text{rec} \in \text{cmdLog}[\text{cleader}] : \text{rec.inst} = \text{inst}\} \text{IN} \\
&\wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{cleader}] = (@ \setminus \text{oldRecs}) \cup \\
&\hspace{10em} \{\{\text{inst} \mapsto \text{inst}, \\
&\hspace{10em} \text{status} \mapsto \text{"pre-accepted"}\},
\end{aligned}$$

$$\begin{aligned}
& ballot \mapsto ballot, \\
& cmd \mapsto C, \\
& deps \mapsto newDeps, \\
& seq \mapsto newSeq]] \\
\wedge leaderOfInst' = [leaderOfInst \text{ EXCEPT } ![cleader] = @ \cup \{inst\}] \\
\wedge sentMsg' = (sentMsg \setminus oldMsg) \cup \\
& [type : \{\text{"pre-accept"}\}, \\
& src : \{cleader\}, \\
& dst : Q \setminus \{cleader\}, \\
& inst : \{inst\}, \\
& ballot : \{ballot\}, \\
& cmd : \{C\}, \\
& deps : \{newDeps\}, \\
& seq : \{newSeq\}]
\end{aligned}$$

$$\begin{aligned}
Propose(C, cleader) &\triangleq \\
\text{LET } newInst &\triangleq \langle cleader, crtInst[cleader] \rangle \\
newBallot &\triangleq \langle 0, cleader \rangle \\
\text{IN } \wedge proposed' &= proposed \cup \{C\} \\
&\wedge (\exists Q \in FastQuorums(cleader) : \\
&\quad StartPhase1(C, cleader, Q, newInst, newBallot, \{\}) \\
&\wedge crtInst' = [crtInst \text{ EXCEPT } ![cleader] = @ + 1] \\
&\wedge \text{UNCHANGED } \langle executed, committed, ballots, preparing \rangle
\end{aligned}$$

$$\begin{aligned}
Phase1Reply(replica) &\triangleq \\
\exists msg \in sentMsg : \\
&\wedge msg.type = \text{"pre-accept"} \\
&\wedge msg.dst = replica \\
&\wedge \text{LET } oldRec \triangleq \{rec \in cmdLog[replica] : rec.inst = msg.inst\} \text{IN} \\
&\quad \wedge (\forall rec \in oldRec : \\
&\quad\quad (rec.ballot = msg.ballot \vee rec.ballot[1] < msg.ballot[1])) \\
&\quad \wedge \text{LET } newDeps \triangleq msg.deps \cup \\
&\quad\quad (\{t.inst : t \in cmdLog[replica]\} \setminus \{msg.inst\}) \\
&\quad\quad newSeq \triangleq Max(\{msg.seq, \\
&\quad\quad\quad 1 + Max(\{t.seq : t \in cmdLog[replica]\})\}) \\
&\quad instCom \triangleq \{t.inst : t \in \{tt \in cmdLog[replica] : \\
&\quad\quad tt.status \in \{\text{"committed"}, \text{"executed"}\}\}\} \text{IN} \\
&\wedge cmdLog' = [cmdLog \text{ EXCEPT } ![replica] = (@ \setminus oldRec) \cup \\
&\quad\quad \{[inst \mapsto msg.inst, \\
&\quad\quad\quad status \mapsto \text{"pre-accepted"}, \\
&\quad\quad\quad ballot \mapsto msg.ballot, \\
&\quad\quad\quad cmd \mapsto msg.cmd, \\
&\quad\quad\quad deps \mapsto newDeps, \\
&\quad\quad\quad seq \mapsto newSeq]\}] \\
&\wedge sentMsg' = (sentMsg \setminus \{msg\}) \cup
\end{aligned}$$

$$\begin{aligned}
& \{[type \mapsto \text{"pre-accept-reply"}, \\
& \quad src \mapsto replica, \\
& \quad dst \mapsto msg.src, \\
& \quad inst \mapsto msg.inst, \\
& \quad ballot \mapsto msg.ballot, \\
& \quad deps \mapsto newDeps, \\
& \quad seq \mapsto newSeq, \\
& \quad committed \mapsto instCom]\} \\
\wedge \text{UNCHANGED } & \langle proposed, crtInst, executed, leaderOfInst, \\
& \quad committed, ballots, preparing \rangle
\end{aligned}$$

$$\begin{aligned}
Phase1Fast(cleader, i, Q) & \triangleq \\
& \wedge i \in leaderOfInst[cleader] \\
& \wedge Q \in FastQuorums(cleader) \\
& \wedge \exists record \in cmdLog[cleader] : \\
& \quad \wedge record.inst = i \\
& \quad \wedge record.status = \text{"pre-accepted"} \\
& \quad \wedge record.ballot[1] = 0 \\
& \quad \wedge \text{LET } replies \triangleq \{msg \in sentMsg : \\
& \quad \quad \wedge msg.inst = i \\
& \quad \quad \wedge msg.type = \text{"pre-accept-reply"} \\
& \quad \quad \wedge msg.dst = cleader \\
& \quad \quad \wedge msg.src \in Q \\
& \quad \quad \wedge msg.ballot = record.ballot\}IN \\
& \quad \wedge (\forall replica \in (Q \setminus \{cleader\}) : \\
& \quad \quad \exists msg \in replies : msg.src = replica) \\
& \quad \wedge (\forall r1, r2 \in replies : \\
& \quad \quad \wedge r1.deps = r2.deps \\
& \quad \quad \wedge r1.seq = r2.seq) \\
& \quad \wedge \text{LET } r \triangleq \text{CHOOSE } r \in replies : \text{TRUEIN} \\
& \quad \quad \wedge \text{LET } localCom \triangleq \{t.inst : \\
& \quad \quad \quad t \in \{tt \in cmdLog[cleader] : \\
& \quad \quad \quad \quad tt.status \in \{\text{"committed"}, \text{"executed"}\}\} \\
& \quad \quad \quad extCom \triangleq \text{UNION } \{msg.committed : msg \in replies\}IN \\
& \quad \quad \quad (r.deps \subseteq (localCom \cup extCom)) \\
& \quad \quad \wedge cmdLog' = [cmdLog \text{ EXCEPT } ![cleader] = (@ \setminus \{record\}) \cup \\
& \quad \quad \quad \{[inst \mapsto i, \\
& \quad \quad \quad \quad status \mapsto \text{"committed"}, \\
& \quad \quad \quad \quad ballot \mapsto record.ballot, \\
& \quad \quad \quad \quad cmd \mapsto record.cmd, \\
& \quad \quad \quad \quad deps \mapsto r.deps, \\
& \quad \quad \quad \quad seq \mapsto r.seq]\}] \\
& \quad \wedge sentMsg' = (sentMsg \setminus replies) \cup \\
& \quad \quad \{[type \mapsto \text{"commit"}, \\
& \quad \quad \quad inst \mapsto i,
\end{aligned}$$

$$\begin{aligned}
& ballot \mapsto record.ballot, \\
& cmd \mapsto record.cmd, \\
& deps \mapsto r.deps, \\
& seq \mapsto r.seq\} \\
\wedge leaderOfInst' = [leaderOfInst \text{ EXCEPT } ![cleader] = @ \setminus \{i\}] \\
\wedge committed' = [committed \text{ EXCEPT } ![i] = \\
\quad @ \cup \{\langle record.cmd, r.deps, r.seq \rangle\}] \\
\wedge \text{UNCHANGED } \langle proposed, executed, crtInst, ballots, preparing \rangle
\end{aligned}$$

$$\begin{aligned}
Phase1Slow(cleader, i, Q) \triangleq & \\
& \wedge i \in leaderOfInst[cleader] \\
& \wedge Q \in SlowQuorums(cleader) \\
& \wedge \exists record \in cmdLog[cleader] : \\
& \quad \wedge record.inst = i \\
& \quad \wedge record.status = \text{"pre-accepted"} \\
& \quad \wedge \text{LET } replies \triangleq \{msg \in sentMsg : \\
& \quad \quad \wedge msg.inst = i \\
& \quad \quad \wedge msg.type = \text{"pre-accept-reply"} \\
& \quad \quad \wedge msg.dst = cleader \\
& \quad \quad \wedge msg.src \in Q \\
& \quad \quad \wedge msg.ballot = record.ballot\} \text{IN} \\
& \quad \wedge (\forall replica \in (Q \setminus \{cleader\}) : \exists msg \in replies : msg.src = replica) \\
& \quad \wedge \text{LET } finalDeps \triangleq \text{UNION } \{msg.deps : msg \in replies\} \\
& \quad \quad finalSeq \triangleq \text{Max}(\{msg.seq : msg \in replies\}) \text{IN} \\
& \quad \wedge cmdLog' = [cmdLog \text{ EXCEPT } ![cleader] = (@ \setminus \{record\}) \cup \\
& \quad \quad \quad \{[inst \mapsto i, \\
& \quad \quad \quad \quad status \mapsto \text{"accepted"}, \\
& \quad \quad \quad \quad ballot \mapsto record.ballot, \\
& \quad \quad \quad \quad cmd \mapsto record.cmd, \\
& \quad \quad \quad \quad deps \mapsto finalDeps, \\
& \quad \quad \quad \quad seq \mapsto finalSeq]\}] \\
& \quad \wedge \exists SQ \in SlowQuorums(cleader) : \\
& \quad \quad (sentMsg' = (sentMsg \setminus replies) \cup \\
& \quad \quad \quad [type : \{\text{"accept"}\}, \\
& \quad \quad \quad \quad src : \{cleader\}, \\
& \quad \quad \quad \quad dst : SQ \setminus \{cleader\}, \\
& \quad \quad \quad \quad inst : \{i\}, \\
& \quad \quad \quad \quad ballot : \{record.ballot\}, \\
& \quad \quad \quad \quad cmd : \{record.cmd\}, \\
& \quad \quad \quad \quad deps : \{finalDeps\}, \\
& \quad \quad \quad \quad seq : \{finalSeq\}]) \\
& \quad \wedge \text{UNCHANGED } \langle proposed, executed, crtInst, leaderOfInst, \\
& \quad \quad \quad committed, ballots, preparing \rangle
\end{aligned}$$

$$Phase2Reply(replica) \triangleq$$

$$\begin{aligned}
& \exists \text{msg} \in \text{sentMsg} : \\
& \quad \wedge \text{msg.type} = \text{"accept"} \\
& \quad \wedge \text{msg.dst} = \text{replica} \\
& \quad \wedge \text{LET } \text{oldRec} \triangleq \{ \text{rec} \in \text{cmdLog}[\text{replica}] : \text{rec.inst} = \text{msg.inst} \} \text{IN} \\
& \quad \quad \wedge (\forall \text{rec} \in \text{oldRec} : (\text{rec.ballot} = \text{msg.ballot} \vee \\
& \quad \quad \quad \text{rec.ballot}[1] < \text{msg.ballot}[1])) \\
& \quad \wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{replica}] = (@ \setminus \text{oldRec}) \cup \\
& \quad \quad \quad \{ [\text{inst} \mapsto \text{msg.inst}, \\
& \quad \quad \quad \text{status} \mapsto \text{"accepted"}, \\
& \quad \quad \quad \text{ballot} \mapsto \text{msg.ballot}, \\
& \quad \quad \quad \text{cmd} \mapsto \text{msg.cmd}, \\
& \quad \quad \quad \text{deps} \mapsto \text{msg.deps}, \\
& \quad \quad \quad \text{seq} \mapsto \text{msg.seq}] \}] \\
& \quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \{ \text{msg} \}) \cup \\
& \quad \quad \quad \{ [\text{type} \mapsto \text{"accept-reply"}, \\
& \quad \quad \quad \text{src} \mapsto \text{replica}, \\
& \quad \quad \quad \text{dst} \mapsto \text{msg.src}, \\
& \quad \quad \quad \text{inst} \mapsto \text{msg.inst}, \\
& \quad \quad \quad \text{ballot} \mapsto \text{msg.ballot}] \} \\
& \quad \wedge \text{UNCHANGED } \langle \text{proposed}, \text{crtInst}, \text{executed}, \text{leaderOfInst}, \\
& \quad \quad \quad \text{committed}, \text{ballots}, \text{preparing} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{Phase2Finalize}(\text{cleader}, i, Q) \triangleq \\
& \quad \wedge i \in \text{leaderOfInst}[\text{cleader}] \\
& \quad \wedge Q \in \text{SlowQuorums}(\text{cleader}) \\
& \quad \wedge \exists \text{record} \in \text{cmdLog}[\text{cleader}] : \\
& \quad \quad \wedge \text{record.inst} = i \\
& \quad \quad \wedge \text{record.status} = \text{"accepted"} \\
& \quad \wedge \text{LET } \text{replies} \triangleq \{ \text{msg} \in \text{sentMsg} : \\
& \quad \quad \quad \wedge \text{msg.inst} = i \\
& \quad \quad \quad \wedge \text{msg.type} = \text{"accept-reply"} \\
& \quad \quad \quad \wedge \text{msg.dst} = \text{cleader} \\
& \quad \quad \quad \wedge \text{msg.src} \in Q \\
& \quad \quad \quad \wedge \text{msg.ballot} = \text{record.ballot} \} \text{IN} \\
& \quad \wedge (\forall \text{replica} \in (Q \setminus \{ \text{cleader} \}) : \exists \text{msg} \in \text{replies} : \\
& \quad \quad \quad \text{msg.src} = \text{replica}) \\
& \quad \wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{cleader}] = (@ \setminus \{ \text{record} \}) \cup \\
& \quad \quad \quad \{ [\text{inst} \mapsto i, \\
& \quad \quad \quad \text{status} \mapsto \text{"committed"}, \\
& \quad \quad \quad \text{ballot} \mapsto \text{record.ballot}, \\
& \quad \quad \quad \text{cmd} \mapsto \text{record.cmd}, \\
& \quad \quad \quad \text{deps} \mapsto \text{record.deps}, \\
& \quad \quad \quad \text{seq} \mapsto \text{record.seq}] \}] \\
& \quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \text{replies}) \cup
\end{aligned}$$

$$\begin{aligned}
& \{[type \mapsto \text{"commit"}, \\
& \quad inst \mapsto i, \\
& \quad ballot \mapsto record.ballot, \\
& \quad cmd \mapsto record.cmd, \\
& \quad deps \mapsto record.deps, \\
& \quad seq \mapsto record.seq]\} \\
\wedge \text{committed}' &= [\text{committed EXCEPT } ![i] = @ \cup \\
& \quad \{ \langle record.cmd, record.deps, record.seq \rangle \}] \\
\wedge \text{leaderOfInst}' &= [\text{leaderOfInst EXCEPT } ![cleader] = @ \setminus \{i\}] \\
\wedge \text{UNCHANGED} &\langle proposed, executed, crtInst, ballots, preparing \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{Commit}(replica, cmsg) \triangleq \\
& \text{LET } oldRec \triangleq \{rec \in cmdLog[replica] : rec.inst = cmsg.inst\} \text{ IN} \\
& \wedge \forall rec \in oldRec : (rec.status \notin \{\text{"committed"}, \text{"executed"}\} \wedge \\
& \quad rec.ballot[1] \leq cmsg.ballot[1]) \\
& \wedge cmdLog' = [cmdLog EXCEPT ![replica] = (@ \setminus oldRec) \cup \\
& \quad \{[inst \mapsto cmsg.inst, \\
& \quad \quad status \mapsto \text{"committed"}, \\
& \quad \quad ballot \mapsto cmsg.ballot, \\
& \quad \quad cmd \mapsto cmsg.cmd, \\
& \quad \quad deps \mapsto cmsg.deps, \\
& \quad \quad seq \mapsto cmsg.seq]\}] \\
& \wedge \text{committed}' = [\text{committed EXCEPT } ![cmsg.inst] = @ \cup \\
& \quad \{ \langle cmsg.cmd, cmsg.deps, cmsg.seq \rangle \}] \\
& \wedge \text{UNCHANGED} \langle proposed, executed, crtInst, leaderOfInst, \\
& \quad sentMsg, ballots, preparing \rangle
\end{aligned}$$

Recovery actions

$$\begin{aligned}
& \text{SendPrepare}(replica, i, Q) \triangleq \\
& \wedge i \notin leaderOfInst[replica] \\
& \wedge i \notin preparing[replica] \\
& \wedge ballots \leq MaxBallot \\
& \wedge \neg(\exists rec \in cmdLog[replica] : \\
& \quad \wedge rec.inst = i \\
& \quad \wedge rec.status \in \{\text{"committed"}, \text{"executed"}\}) \\
& \wedge sentMsg' = sentMsg \cup \\
& \quad [type : \{\text{"prepare"}\}, \\
& \quad \quad src : \{replica\}, \\
& \quad \quad dst : Q, \\
& \quad \quad inst : \{i\}, \\
& \quad \quad ballot : \{\langle ballots, replica \rangle\}] \\
& \wedge ballots' = ballots + 1 \\
& \wedge preparing' = [preparing EXCEPT ![replica] = @ \cup \{i\}] \\
& \wedge \text{UNCHANGED} \langle cmdLog, proposed, executed, crtInst,
\end{aligned}$$

$\langle \text{leaderOfInst}, \text{committed} \rangle$

$$\begin{aligned} \text{ReplyPrepare}(\text{replica}) &\triangleq \\ &\exists \text{msg} \in \text{sentMsg} : \\ &\quad \wedge \text{msg.type} = \text{“prepare”} \\ &\quad \wedge \text{msg.dst} = \text{replica} \\ &\quad \wedge \forall \exists \text{rec} \in \text{cmdLog}[\text{replica}] : \\ &\quad\quad \wedge \text{rec.inst} = \text{msg.inst} \\ &\quad\quad \wedge \text{msg.ballot}[1] > \text{rec.ballot}[1] \\ &\quad\quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \{\text{msg}\}) \cup \\ &\quad\quad\quad \{ [\text{type} \mapsto \text{“prepare-reply”}, \\ &\quad\quad\quad\quad \text{src} \mapsto \text{replica}, \\ &\quad\quad\quad\quad \text{dst} \mapsto \text{msg.src}, \\ &\quad\quad\quad\quad \text{inst} \mapsto \text{rec.inst}, \\ &\quad\quad\quad\quad \text{ballot} \mapsto \text{msg.ballot}, \\ &\quad\quad\quad\quad \text{prev_ballot} \mapsto \text{rec.ballot}, \\ &\quad\quad\quad\quad \text{status} \mapsto \text{rec.status}, \\ &\quad\quad\quad\quad \text{cmd} \mapsto \text{rec.cmd}, \\ &\quad\quad\quad\quad \text{deps} \mapsto \text{rec.deps}, \\ &\quad\quad\quad\quad \text{seq} \mapsto \text{rec.seq}] \} \\ &\quad \wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{replica}] = (@ \setminus \{\text{rec}\}) \cup \\ &\quad\quad\quad \{ [\text{inst} \mapsto \text{rec.inst}, \\ &\quad\quad\quad\quad \text{status} \mapsto \text{rec.status}, \\ &\quad\quad\quad\quad \text{ballot} \mapsto \text{msg.ballot}, \\ &\quad\quad\quad\quad \text{cmd} \mapsto \text{rec.cmd}, \\ &\quad\quad\quad\quad \text{deps} \mapsto \text{rec.deps}, \\ &\quad\quad\quad\quad \text{seq} \mapsto \text{rec.seq}] \} \\ &\quad \wedge \text{IF } \text{rec.inst} \in \text{leaderOfInst}[\text{replica}] \text{ THEN} \\ &\quad\quad \wedge \text{leaderOfInst}' = [\text{leaderOfInst} \text{ EXCEPT } ![\text{replica}] = \\ &\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad @ \setminus \{\text{rec.inst}\}] \\ &\quad\quad \wedge \text{UNCHANGED } \langle \text{proposed}, \text{executed}, \text{committed}, \\ &\quad\quad\quad\quad \text{crtInst}, \text{ballots}, \text{preparing} \rangle \\ &\quad \text{ELSE UNCHANGED } \langle \text{proposed}, \text{executed}, \text{committed}, \text{crtInst}, \\ &\quad\quad\quad\quad \text{ballots}, \text{preparing}, \text{leaderOfInst} \rangle \\ \\ \vee \wedge \neg(\exists \text{rec} \in \text{cmdLog}[\text{replica}] : \text{rec.inst} = \text{msg.inst}) \\ &\quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \{\text{msg}\}) \cup \\ &\quad\quad\quad \{ [\text{type} \mapsto \text{“prepare-reply”}, \\ &\quad\quad\quad\quad \text{src} \mapsto \text{replica}, \\ &\quad\quad\quad\quad \text{dst} \mapsto \text{msg.src}, \\ &\quad\quad\quad\quad \text{inst} \mapsto \text{msg.inst}, \\ &\quad\quad\quad\quad \text{ballot} \mapsto \text{msg.ballot}, \\ &\quad\quad\quad\quad \text{prev_ballot} \mapsto \langle 0, \text{replica} \rangle, \\ &\quad\quad\quad\quad \text{status} \mapsto \text{“not-seen”}, \\ &\quad\quad\quad\quad \text{cmd} \mapsto \text{none}, \end{aligned}$$

$$\begin{aligned}
& \text{deps} \mapsto \{\}, \\
& \text{seq} \mapsto 0\} \\
\wedge \text{cmdLog}' = & [\text{cmdLog} \text{ EXCEPT } ![\text{replica}] = @ \cup \\
& \{[\text{inst} \mapsto \text{msg.inst}, \\
& \text{status} \mapsto \text{"not-seen"}, \\
& \text{ballot} \mapsto \text{msg.ballot}, \\
& \text{cmd} \mapsto \text{none}, \\
& \text{deps} \mapsto \{\}, \\
& \text{seq} \mapsto 0]\}] \\
\wedge \text{UNCHANGED} & \langle \text{proposed}, \text{executed}, \text{committed}, \text{crtInst}, \text{ballots}, \\
& \text{leaderOfInst}, \text{preparing} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{PrepareFinalize}(\text{replica}, i, Q) & \triangleq \\
& \wedge i \in \text{preparing}[\text{replica}] \\
& \wedge \exists \text{rec} \in \text{cmdLog}[\text{replica}] : \\
& \quad \wedge \text{rec.inst} = i \\
& \quad \wedge \text{rec.status} \notin \{\text{"committed"}, \text{"executed"}\} \\
& \quad \wedge \text{LET } \text{replies} \triangleq \{ \text{msg} \in \text{sentMsg} : \\
& \quad \quad \wedge \text{msg.inst} = i \\
& \quad \quad \wedge \text{msg.type} = \text{"prepare-reply"} \\
& \quad \quad \wedge \text{msg.dst} = \text{replica} \\
& \quad \quad \wedge \text{msg.ballot} = \text{rec.ballot} \} \text{IN} \\
& \quad \wedge (\forall \text{rep} \in Q : \exists \text{msg} \in \text{replies} : \text{msg.src} = \text{rep}) \\
& \quad \wedge \forall \exists \text{com} \in \text{replies} : \\
& \quad \quad \wedge (\text{com.status} \in \{\text{"committed"}, \text{"executed"}\}) \\
& \quad \quad \wedge \text{preparing}' = [\text{preparing} \text{ EXCEPT } ![\text{replica}] = @ \setminus \{i\}] \\
& \quad \quad \wedge \text{sentMsg}' = \text{sentMsg} \setminus \text{replies} \\
& \quad \quad \wedge \text{UNCHANGED} \langle \text{cmdLog}, \text{proposed}, \text{executed}, \text{crtInst}, \text{leaderOfInst}, \\
& \quad \quad \quad \text{committed}, \text{ballots} \rangle \\
& \quad \vee \wedge \neg(\exists \text{msg} \in \text{replies} : \text{msg.status} \in \{\text{"committed"}, \text{"executed"}\}) \\
& \quad \wedge \exists \text{acc} \in \text{replies} : \\
& \quad \quad \wedge \text{acc.status} = \text{"accepted"} \\
& \quad \quad \wedge (\forall \text{msg} \in (\text{replies} \setminus \{\text{acc}\}) : \\
& \quad \quad \quad (\text{msg.prev_ballot}[1] \leq \text{acc.prev_ballot}[1] \vee \\
& \quad \quad \quad \text{msg.status} \neq \text{"accepted"})) \\
& \quad \quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \text{replies}) \cup \\
& \quad \quad \quad [\text{type} : \{\text{"accept"}\}, \\
& \quad \quad \quad \text{src} : \{\text{replica}\}, \\
& \quad \quad \quad \text{dst} : Q \setminus \{\text{replica}\}, \\
& \quad \quad \quad \text{inst} : \{i\}, \\
& \quad \quad \quad \text{ballot} : \{\text{rec.ballot}\}, \\
& \quad \quad \quad \text{cmd} : \{\text{acc.cmd}\}, \\
& \quad \quad \quad \text{deps} : \{\text{acc.deps}\}, \\
& \quad \quad \quad \text{seq} : \{\text{acc.seq}\}] \\
& \quad \wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{replica}] = (@ \setminus \{\text{rec}\}) \cup
\end{aligned}$$

$$\begin{aligned}
& \{[inst \mapsto i, \\
& \quad status \mapsto \text{"accepted"}, \\
& \quad ballot \mapsto rec.ballot, \\
& \quad cmd \mapsto acc.cmd, \\
& \quad deps \mapsto acc.deps, \\
& \quad seq \mapsto acc.seq]\} \\
& \wedge preparing' = [preparing \text{ EXCEPT } ![replica] = @ \setminus \{i\}] \\
& \wedge leaderOfInst' = [leaderOfInst \text{ EXCEPT } ![replica] = @ \cup \{i\}] \\
& \wedge \text{UNCHANGED } \langle proposed, executed, crtInst, committed, ballots \rangle \\
\vee \wedge \neg(\exists msg \in replies : \\
& \quad msg.status \in \{\text{"accepted"}, \text{"committed"}, \text{"executed"}\}) \\
& \wedge \text{LET } preaccepts \triangleq \{msg \in replies : msg.status = \text{"pre-accepted"}\} \text{ IN} \\
& (\vee \wedge \forall p1, p2 \in preaccepts : \\
& \quad p1.cmd = p2.cmd \wedge p1.deps = p2.deps \wedge p1.seq = p2.seq \\
& \wedge \neg(\exists pl \in preaccepts : pl.src = i[1]) \\
& \wedge \text{Cardinality}(preaccepts) \geq \text{Cardinality}(Q) - 1 \\
& \wedge \text{LET } pac \triangleq \text{CHOOSE } pac \in preaccepts : \text{TRUEIN} \\
& \quad \wedge sentMsg' = (sentMsg \setminus replies) \cup \\
& \quad \quad [type : \{\text{"accept"}\}, \\
& \quad \quad src : \{replica\}, \\
& \quad \quad dst : Q \setminus \{replica\}, \\
& \quad \quad inst : \{i\}, \\
& \quad \quad ballot : \{rec.ballot\}, \\
& \quad \quad cmd : \{pac.cmd\}, \\
& \quad \quad deps : \{pac.deps\}, \\
& \quad \quad seq : \{pac.seq\}] \\
& \wedge cmdLog' = [cmdLog \text{ EXCEPT } ![replica] = (@ \setminus \{rec\}) \cup \\
& \quad \{[inst \mapsto i, \\
& \quad \quad status \mapsto \text{"accepted"}, \\
& \quad \quad ballot \mapsto rec.ballot, \\
& \quad \quad cmd \mapsto pac.cmd, \\
& \quad \quad deps \mapsto pac.deps, \\
& \quad \quad seq \mapsto pac.seq]\} \\
& \wedge preparing' = [preparing \text{ EXCEPT } ![replica] = @ \setminus \{i\}] \\
& \wedge leaderOfInst' = [leaderOfInst \text{ EXCEPT } ![replica] = @ \cup \{i\}] \\
& \wedge \text{UNCHANGED } \langle proposed, executed, crtInst, committed, ballots \rangle \\
\vee \wedge \forall p1, p2 \in preaccepts : p1.cmd = p2.cmd \wedge \\
& \quad p1.deps = p2.deps \wedge \\
& \quad p1.seq = p2.seq \\
& \wedge \neg(\exists pl \in preaccepts : pl.src = i[1]) \\
& \wedge \text{Cardinality}(preaccepts) < \text{Cardinality}(Q) - 1 \\
& \wedge \text{Cardinality}(preaccepts) \geq \text{Cardinality}(Q) \div 2 \\
& \wedge \text{LET } pac \triangleq \text{CHOOSE } pac \in preaccepts : \text{TRUEIN} \\
& \quad \wedge sentMsg' = (sentMsg \setminus replies) \cup \\
& \quad \quad [type : \{\text{"try-pre-accept"}\},
\end{aligned}$$

$$\begin{aligned}
& \text{src} : \{\text{replica}\}, \\
& \text{dst} : Q, \\
& \text{inst} : \{i\}, \\
& \text{ballot} : \{\text{rec.ballot}\}, \\
& \text{cmd} : \{\text{pac.cmd}\}, \\
& \text{deps} : \{\text{pac.deps}\}, \\
& \text{seq} : \{\text{pac.seq}\} \\
& \wedge \text{preparing}' = [\text{preparing} \text{ EXCEPT } ![\text{replica}] = @ \setminus \{i\}] \\
& \wedge \text{leaderOfInst}' = [\text{leaderOfInst} \text{ EXCEPT } ![\text{replica}] = @ \cup \{i\}] \\
& \wedge \text{UNCHANGED} \langle \text{cmdLog}, \text{proposed}, \text{executed}, \\
& \quad \text{crtInst}, \text{committed}, \text{ballots} \rangle \\
\vee & \wedge \vee \exists p1, p2 \in \text{preaccepts} : p1.\text{cmd} \neq p2.\text{cmd} \vee \\
& \quad p1.\text{deps} \neq p2.\text{deps} \vee \\
& \quad p1.\text{seq} \neq p2.\text{seq} \\
& \vee \exists pl \in \text{preaccepts} : pl.\text{src} = i[1] \\
& \vee \text{Cardinality}(\text{preaccepts}) < \text{Cardinality}(Q) \div 2 \\
& \wedge \text{preaccepts} \neq \{\} \\
& \wedge \text{LET } \text{pac} \triangleq \text{CHOOSE } \text{pac} \in \text{preaccepts} : \text{pac.cmd} \neq \text{noneIN} \\
& \quad \wedge \text{StartPhase1}(\text{pac.cmd}, \text{replica}, Q, i, \text{rec.ballot}, \text{replies}) \\
& \quad \wedge \text{preparing}' = [\text{preparing} \text{ EXCEPT } ![\text{replica}] = @ \setminus \{i\}] \\
& \quad \wedge \text{UNCHANGED} \langle \text{proposed}, \text{executed}, \text{crtInst}, \text{committed}, \text{ballots} \rangle \\
\vee & \wedge \forall \text{msg} \in \text{replies} : \text{msg.status} = \text{"not-seen"} \\
& \quad \wedge \text{StartPhase1}(\text{none}, \text{replica}, Q, i, \text{rec.ballot}, \text{replies}) \\
& \quad \wedge \text{preparing}' = [\text{preparing} \text{ EXCEPT } ![\text{replica}] = @ \setminus \{i\}] \\
& \quad \wedge \text{UNCHANGED} \langle \text{proposed}, \text{executed}, \text{crtInst}, \text{committed}, \text{ballots} \rangle
\end{aligned}$$

$\text{ReplyTryPreaccept}(\text{replica}) \triangleq$

$$\begin{aligned}
& \exists \text{tpa} \in \text{sentMsg} : \\
& \quad \wedge \text{tpa.type} = \text{"try-pre-accept"} \\
& \quad \wedge \text{tpa.dst} = \text{replica} \\
& \quad \wedge \text{LET } \text{oldRec} \triangleq \{\text{rec} \in \text{cmdLog}[\text{replica}] : \text{rec.inst} = \text{tpa.inst}\} \text{IN} \\
& \quad \quad \wedge \forall \text{rec} \in \text{oldRec} : \text{rec.ballot}[1] \leq \text{tpa.ballot}[1] \wedge \\
& \quad \quad \quad \text{rec.status} \notin \{\text{"accepted"}, \text{"committed"}, \text{"executed"}\} \\
& \quad \wedge \vee (\exists \text{rec} \in \text{cmdLog}[\text{replica}] \setminus \text{oldRec} : \\
& \quad \quad \wedge \text{tpa.inst} \notin \text{rec.deps} \\
& \quad \quad \wedge \vee \text{rec.inst} \notin \text{tpa.deps} \\
& \quad \quad \quad \vee \text{rec.seq} \geq \text{tpa.seq} \\
& \quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \{\text{tpa}\}) \cup \\
& \quad \quad \{[\text{type} \mapsto \text{"try-pre-accept-reply"}, \\
& \quad \quad \text{src} \mapsto \text{replica}, \\
& \quad \quad \text{dst} \mapsto \text{tpa.src}, \\
& \quad \quad \text{inst} \mapsto \text{tpa.inst}, \\
& \quad \quad \text{ballot} \mapsto \text{tpa.ballot}, \\
& \quad \quad \text{status} \mapsto \text{rec.status}]\} \\
& \quad \wedge \text{UNCHANGED} \langle \text{cmdLog}, \text{proposed}, \text{executed}, \text{committed}, \text{crtInst},
\end{aligned}$$

$$\begin{aligned}
& \text{ballots, leaderOfInst, preparing} \rangle \\
\vee \wedge (\forall \text{rec} \in \text{cmdLog}[\text{replica}] \setminus \text{oldRec} : \\
& \quad \text{tpa.inst} \in \text{rec.deps} \vee (\text{rec.inst} \in \text{tpa.deps} \wedge \\
& \quad \quad \text{rec.seq} < \text{tpa.seq})) \\
\wedge \text{sentMsg}' = (\text{sentMsg} \setminus \{\text{tpa}\}) \cup \\
& \quad \{[\text{type} \mapsto \text{"try-pre-accept-reply"}, \\
& \quad \text{src} \mapsto \text{replica}, \\
& \quad \text{dst} \mapsto \text{tpa.src}, \\
& \quad \text{inst} \mapsto \text{tpa.inst}, \\
& \quad \text{ballot} \mapsto \text{tpa.ballot}, \\
& \quad \text{status} \mapsto \text{"OK"}]\} \\
\wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{replica}]] = (@ \setminus \text{oldRec}) \cup \\
& \quad \{[\text{inst} \mapsto \text{tpa.inst}, \\
& \quad \text{status} \mapsto \text{"pre-accepted"}, \\
& \quad \text{ballot} \mapsto \text{tpa.ballot}, \\
& \quad \text{cmd} \mapsto \text{tpa.cmd}, \\
& \quad \text{deps} \mapsto \text{tpa.deps}, \\
& \quad \text{seq} \mapsto \text{tpa.seq}]\} \\
\wedge \text{UNCHANGED} \langle \text{proposed, executed, committed, crtInst, ballots,} \\
& \quad \text{leaderOfInst, preparing} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{FinalizeTryPreAccept}(\text{cleader}, i, Q) \triangleq \\
& \quad \exists \text{rec} \in \text{cmdLog}[\text{cleader}] : \\
& \quad \quad \wedge \text{rec.inst} = i \\
& \quad \quad \wedge \text{LET } \text{tprs} \triangleq \{ \text{msg} \in \text{sentMsg} : \text{msg.type} = \text{"try-pre-accept-reply"} \wedge \\
& \quad \quad \quad \text{msg.dst} = \text{cleader} \wedge \text{msg.inst} = i \wedge \\
& \quad \quad \quad \text{msg.ballot} = \text{rec.ballot} \} \text{IN} \\
& \quad \quad \wedge \forall r \in Q : \exists \text{tpr} \in \text{tprs} : \text{tpr.src} = r \\
& \quad \quad \wedge \forall \text{tpr} \in \text{tprs} : \text{tpr.status} = \text{"OK"} \\
& \quad \quad \wedge \text{sentMsg}' = (\text{sentMsg} \setminus \text{tprs}) \cup \\
& \quad \quad \quad [\text{type} : \{\text{"accept"}\}, \\
& \quad \quad \quad \text{src} : \{\text{cleader}\}, \\
& \quad \quad \quad \text{dst} : Q \setminus \{\text{cleader}\}, \\
& \quad \quad \quad \text{inst} : \{i\}, \\
& \quad \quad \quad \text{ballot} : \{\text{rec.ballot}\}, \\
& \quad \quad \quad \text{cmd} : \{\text{rec.cmd}\}, \\
& \quad \quad \quad \text{deps} : \{\text{rec.deps}\}, \\
& \quad \quad \quad \text{seq} : \{\text{rec.seq}\}] \\
& \quad \quad \wedge \text{cmdLog}' = [\text{cmdLog} \text{ EXCEPT } ![\text{cleader}]] = (@ \setminus \{\text{rec}\}) \cup \\
& \quad \quad \quad \{[\text{inst} \mapsto i, \\
& \quad \quad \quad \text{status} \mapsto \text{"accepted"}, \\
& \quad \quad \quad \text{ballot} \mapsto \text{rec.ballot}, \\
& \quad \quad \quad \text{cmd} \mapsto \text{rec.cmd}, \\
& \quad \quad \quad \text{deps} \mapsto \text{rec.deps},
\end{aligned}$$

$$\begin{aligned}
& seq \mapsto rec.seq\}}] \\
& \wedge \text{UNCHANGED } \langle proposed, executed, committed, crtInst, ballots, \\
& \quad leaderOfInst, preparing \rangle \\
\vee & \wedge \exists tpr \in tprs : tpr.status \in \{ \text{"accepted"}, \text{"committed"}, \text{"executed"} \} \\
& \wedge \text{StartPhase1}(rec.cmd, cleader, Q, i, rec.ballot, tprs) \\
& \wedge \text{UNCHANGED } \langle proposed, executed, committed, crtInst, ballots, \\
& \quad leaderOfInst, preparing \rangle \\
\vee & \wedge \exists tpr \in tprs : tpr.status = \text{"pre-accepted"} \\
& \wedge \forall tpr \in tprs : tpr.status \in \{ \text{"OK"}, \text{"pre-accepted"} \} \\
& \wedge sentMsg' = sentMsg \setminus tprs \\
& \wedge leaderOfInst' = [leaderOfInst \text{ EXCEPT } ![cleader] = @ \setminus \{i\}] \\
& \wedge \text{UNCHANGED } \langle cmdLog, proposed, executed, committed, crtInst, \\
& \quad ballots, preparing \rangle
\end{aligned}$$

Action groups

CommandLeaderAction \triangleq

$$\begin{aligned}
& \vee (\exists C \in (\text{Commands} \setminus \text{proposed}) : \\
& \quad \exists cleader \in \text{Replicas} : \text{Propose}(C, cleader)) \\
& \vee (\exists cleader \in \text{Replicas} : \exists inst \in leaderOfInst[cleader] : \\
& \quad \vee (\exists Q \in \text{FastQuorums}(cleader) : \text{Phase1Fast}(cleader, inst, Q)) \\
& \quad \vee (\exists Q \in \text{SlowQuorums}(cleader) : \text{Phase1Slow}(cleader, inst, Q)) \\
& \quad \vee (\exists Q \in \text{SlowQuorums}(cleader) : \text{Phase2Finalize}(cleader, inst, Q)) \\
& \quad \vee (\exists Q \in \text{SlowQuorums}(cleader) : \text{FinalizeTryPreAccept}(cleader, inst, Q)))
\end{aligned}$$

ReplicaAction \triangleq

$$\begin{aligned}
& \exists replica \in \text{Replicas} : \\
& \quad (\vee \text{Phase1Reply}(replica) \\
& \quad \vee \text{Phase2Reply}(replica) \\
& \quad \vee \exists cmsg \in sentMsg : (cmsg.type = \text{"commit"} \wedge \text{Commit}(replica, cmsg)) \\
& \quad \vee \exists i \in \text{Instances} : \\
& \quad \quad \wedge crtInst[i[1]] > i[2] \quad \text{This condition states that the instance has} \\
& \quad \quad \quad \text{been started by its original owner} \\
& \quad \quad \wedge \exists Q \in \text{SlowQuorums}(replica) : \text{SendPrepare}(replica, i, Q) \\
& \quad \vee \text{ReplyPrepare}(replica) \\
& \quad \vee \exists i \in preparing[replica] : \\
& \quad \quad \exists Q \in \text{SlowQuorums}(replica) : \text{PrepareFinalize}(replica, i, Q) \\
& \quad \vee \text{ReplyTryPreaccept}(replica))
\end{aligned}$$

Next action

Next \triangleq

$$\begin{aligned}
& \vee \text{CommandLeaderAction} \\
& \vee \text{ReplicaAction}
\end{aligned}$$

The complete definition of the algorithm

$$Spec \triangleq Init \wedge \Box[Next]_{vars}$$

Theorems

Nontriviality \triangleq

$$\forall i \in Instances : \\ \Box(\forall C \in committed[i] : C \in proposed \vee C = none)$$

Stability \triangleq

$$\forall replica \in Replicas : \\ \forall i \in Instances : \\ \forall C \in Commands : \\ \Box((\exists rec1 \in cmdLog[replica] : \\ \wedge rec1.inst = i \\ \wedge rec1.cmd = C \\ \wedge rec1.status \in \{“committed”, “executed”\}) \Rightarrow \\ \Box(\exists rec2 \in cmdLog[replica] : \\ \wedge rec2.inst = i \\ \wedge rec2.cmd = C \\ \wedge rec2.status \in \{“committed”, “executed”\})))$$

Consistency \triangleq

$$\forall i \in Instances : \\ \Box(Cardinality(committed[i]) \leq 1)$$

THEOREM $Spec \Rightarrow (\Box TypeOK) \wedge Nontriviality \wedge Stability \wedge Consistency$