

Addendum to “Near-Optimal Latency Versus Cost Tradeoffs in Geo-Distributed Storage”

1 Introduction

Geo-distributed storage systems that provide strong consistency are fundamentally limited in their ability to trade off latency for cost. This technical report supplements the paper [2] we have written on this topic in three ways. First, we describe a set of constraints that limit which read or write latencies can be achieved given a storage budget, and use them to compute a lower bound on the optimal three-way trade-off between read latency, write latency, and storage costs (§2). Second, we provide a proof of correctness for the PANDO consensus protocol (§3). Finally, we provide a TLA+ specification (§4) that we have used to model check PANDO’s consistency guarantees under a variety of settings. In our paper [2], we show that PANDO achieves a near-optimal trade-off.

2 Lower bound on the latency–cost tradeoff

We begin by describing the setting and assumptions under which we optimize latency and cost:

- **Strong Consistency.** All reads and writes must be linearizable; i.e., all writes are totally ordered and every read returns the last successful write.
- **Data centers are the unit of failure.** We ignore failures within a data center and ignore the cost of providing durability within a data center.
- **Reads and writes are executed at a known set of data centers.** We assume that clients issue requests through a small number of “front-end” data centers running application code. We refer to the set of front-end data centers that access an object as the object’s *access set* and the set of data centers storing the data for that object as its *data sites*.
- **No common-case optimizations.** We do not consider common-case optimizations that improve only throughput or median latency (e.g. read leases [4, 3, 5, 8], follow the sun [10]). Our focus is on what the lowest achievable tail latencies are.

- **Wide-area network latency dominates and has low variance.** We assume that latency between data centers is high enough that the difference between median and tail latency in the intra-data center storage (< 5 ms in existing SSD-based systems [1]) is insignificant and that the network itself has low variance.
- **Data is incompressible.** Data compression is an orthogonal tool that can be used to reduce the size of data stored locally at a data site or can be used at a front-end before it issues a write. Therefore, we assume that data cannot be or has already been compressed.
- **Writes are not commutative.** We do not assume that write operations on an object commute and therefore must resolve any conflicts during write execution.

Although the latency–cost trade-off achieved by an optimal approach may not be bound by limitations of existing protocols, there are still constraints that limit what can be achieved. While an optimal approach may not use existing data distribution techniques such as replication or erasure coding, it still must fetch (and distribute) the equivalent of one copy’s worth of data. Furthermore, even though there may exist consensus protocols that have yet to be developed, any approach that guarantees linearizability must be able to *detect updates* in order to react to them when reading or writing data. Following this line of reasoning, we have assembled a list of properties that have to be maintained by any approach (optimal or otherwise):

1. In the event that arbitrary f data sites fail, the remaining data sites *in every write quorum* must contain the equivalent of one copy worth of data.

Justification: If not met, f failures can lead to unavailability or data loss.

2. Every read quorum must contain at least $f + 1$ data sites.

Justification: If a read quorum has $\leq f$ data sites, any write that is completed while all data sites in this quorum are unavailable will not be read when a front-end subsequently reads using this quorum. Hence, it is not possible to guarantee linearizability.

3. Any read quorum must contain at least one copy worth of data in total.

Justification: For a front-end to be able to reconstruct an object’s data when it reads the object, it must receive at least one full copy of the object’s data.

4. All read–write and write–write quorum pairs must have at least one data site in their intersection.

Justification: Without this, it is possible for a write to succeed after contacting just Data Sites A and B followed by a read that contacts only C and D. The read operation may return stale data which would violate linearizability. Similarly, write quorums must overlap to resolve conflicts as they occur, since it is not possible (in general) to resolve them later on.

5. There must be at least $2f + 1$ data sites overall.

$A.ppn$	Promised proposal number stored at acceptor A
$A.apn$	Accepted proposal number stored at acceptor A
$A.vid$	Accepted value id stored at acceptor A
$A.vlen$	Accepted value length stored at acceptor A
$A.vsplit$	Accepted value split stored at acceptor A
vid_v	Value id for v , typically a hash or random number
$vlen_v$	Length of v , used to remove padding from coding
$\text{Split}(v, A)$	(Computed on proposers) The erasure-coded split associated with acceptor A

Figure 1: Summary of notation.

Justification: Based on Properties (1) and (2), each quorum must contain at least $f + 1$ data sites. If f sites fail, we have to have a read and write quorum available to continue servicing requests. Therefore, after f failures, we need an additional $f + 1$ sites remaining, leading to a minimum total of $2f + 1$.

These constraints are fundamental and we have used them to compute a lower bound on the optimal read-write-storage trade-off. As mentioned earlier, we assume that the set of front-end data centers is known which allows us to optimize data site selection. We do this for the lower bound by selecting one read quorum and one write quorum to use per front-end. We are not implying that these are the only quorums that can be used, but that, ignoring tail latency effects, a front-end will have *a* best read or write quorum to use. By enforcing the above constraints on just these quorums, we can compute a lower bound on the optimal trade-off.

3 The PANDO write protocol

PANDO bears a strong resemblance to Classic Paxos [7], but diverges in a few important ways. First, like Flexible Paxos [6], PANDO uses separate Phase 1 and Phase 2 quorums. Second, like RS-Paxos [9], PANDO erasure codes data which requires acceptors to store a few additional fields. Lastly, PANDO separates the notion of conflict detection and recovery by offering a *fast path* for Phase 1 when no conflicts occur. This last optimization enables PANDO to use smaller quorum intersections than a straightforward application of erasure coding.

As described in this paper, the PANDO write protocol allows a set of nodes to agree on the contents of a single value. Support for multiple values can be achieved by using PANDO to decide the value of each entry of a distributed log as is typically done with Paxos. In order to maintain consistency, we require that the quorums used with the PANDO write protocol have the following properties:

1. The intersection of any Phase 1a and Phase 2 quorums contains at least 1 split.
2. The intersection of any Phase 1b and Phase 2 quorums contains at least k splits.

In addition, if Phase 1a quorums are relied upon for reads, they must contain a minimum of k splits, or else the reader will have to contact additional nodes to recover data. Lastly, PANDO provides availability so long as both a Phase 1b and a Phase 2 quorum of nodes are available. Below is pseudocode for the PANDO write protocol.

Phase 1 (Prepare-Promise)

Proposer P initiates a write for value v :

1. Select a unique proposal number p (typically done using Lamport clocks).
2. Broadcast Prepare(p) messages to all acceptors.

Acceptor A , upon receiving Prepare(p) message from Proposer P :

3. If $p > A.ppn$ then set $A.ppn \leftarrow p$ and reply Promise($A.apn, A.vid, A.vlen, A.vsplit$).
4. Else reply NACK.

Proposer P , upon receiving Promise messages from a Phase 1a quorum:

5. If the values in all Promise responses are NULL, then skip to Phase 2 with $v' \leftarrow v$.

Proposer P , upon receiving Promise messages from a Phase 1b quorum:

6. Iterate over all Promise responses by their apn sorted in decreasing order.
 - (a) If there are at least k splits for value w associated with apn , recover the value w (using the associated $vlen$ and $vsplits$) and continue to Phase 2 with $v' \leftarrow w$.
7. If no value was recovered, continue to Phase 2 with $v' \leftarrow v$.

Phase 2 (Propose-Accept)

Proposer P , initiating Phase 2 to write value v' with proposal number p :

8. Broadcast Propose($p, vid_{v'}, vlen_{v'}, Split(v', A)$) messages to all acceptors.

Acceptor A , upon receiving Propose($p, vid, vlen, vsplit$) from a Proposer P :

9. If $p < A.ppn$ reply NACK
10. $A.ppn \leftarrow p$
11. $A.apn \leftarrow p$
12. $A.vid \leftarrow vid$
13. $A.vlen \leftarrow vlen$
14. $A.vsplit \leftarrow vsplit$
15. Reply Accept(p)

Proposer P , upon receiving Accept(p) messages from a Phase 2 quorum:

16. P now knows that v' was chosen, and can check whether the chosen value v' differs from the initial value v or not.

3.1 Proof of correctness

Definitions. We let \mathcal{A} refer to the set of all acceptors and use \mathcal{Q}_a , \mathcal{Q}_b , and \mathcal{Q}_2 refer to the sets of Phase 1a, Phase 1b, and Phase 2 quorums, respectively. Using this notation,

we restate our quorum assumptions:

$$Q \subseteq \mathcal{A} \quad \forall Q \in \mathcal{Q}_a \cup \mathcal{Q}_b \cup \mathcal{Q}_2 \quad (1)$$

$$|Q_a \cap Q_2| \geq 1 \quad \forall Q_a \in \mathcal{Q}_a, Q_2 \in \mathcal{Q}_2 \quad (2)$$

$$|Q_b \cap Q_2| \geq k \quad \forall Q_b \in \mathcal{Q}_b, Q_2 \in \mathcal{Q}_2 \quad (3)$$

Definition 1. A value is **chosen** if there exists a quorum of acceptors that all agree on the identity of the value and store splits that correspond to that value.

We now show that the PANDO write protocol provides the same guarantees as Paxos:

- **Nontriviality.** Any chosen value must have been proposed by a proposer.
- **Liveness.** A value will eventually be chosen as long as RPCs complete before they time out, and all acceptors in a Phase 1b and Phase 2 quorum are available.
- **Consistency.** At most one value can be chosen.
- **Stability.** Once a value is chosen, no other value may be chosen.

Theorem 1. (Nontriviality) PANDO will only choose values that have been proposed.

Proof. By definition, a value can only be chosen if it is present at a Phase 2 quorum of acceptors. Values are only stored at acceptors in response to Propose messages initiated by proposers. \square

Theorem 2. (Liveness) PANDO will eventually choose a value as long as RPCs complete before they time out, and all acceptors in a Phase 1b and Phase 2 quorum are available.

Proof. Let t refer to the (maximum) network and execution latency for an RPC. Since PANDO consists of two rounds of execution, a write can complete within $2t$ as long as a requested is uncontended. If all proposers retry RPCs using randomized exponential backoff, a time window of length $\geq 2t$ will eventually open where only Proposer P is executing. Since no other proposer is sending any RPCs during this time, both Phase 1 and Phase 2 will succeed for Proposer P . \square

Following the precedence of [6], we will show that PANDO provides both consistency and stability by proving that it provides a stronger guarantee.

Lemma 1. If a value v is chosen with proposal number p , then for any proposal with proposal number $p' > p$ and value v' , $v' = v$.

Proof. Recall that PANDO proposers use globally unique proposal numbers (Line 1); this makes it impossible for two different proposals to share a proposal number p . Therefore, if two proposals are both chosen, they must have different proposal numbers. Without loss of generality, we will prove Lemma 1 for the smallest p' such that $p' > p$. If $v' = v$ then we trivially have the desired property. Therefore, assume $v' \neq v$.

We will show that this case always results in a contradiction: either the prepare messages for p' will fail (and thus no propose messages will ever be sent) or the proposer will adopt and re-propose value v . Let $Q_{2,p}$ be the Phase 2 quorum used for proposal number p , and $Q_{a,p'}$ be the Phase 1a quorum used for p' . By Quorum Property 2, we know that $|Q_{2,p} \cap Q_{a,p'}|$ is non-empty. We will now look at the possible ordering of events at each acceptor A in the intersection of these two quorums ($Q_{2,p}$ and $Q_{a,p'}$):

- Case 1: A receives $\text{Prepare}(p')$ before $\text{Propose}(p, \dots)$.
The highest proposal number at A would be $p' > p$ by the time $\text{Propose}(p, \dots)$ was processed, and so A would reject $\text{Propose}(p, \dots)$. However, we know that this is not the case since $A \in Q_{2,p}$, so this is a contradiction.
- Case 2: A receives $\text{Propose}(p, \dots)$ before $\text{Prepare}(p')$.
The last promised proposal number at A is q such that $p \leq q < p'$ ($q > p'$ would be a contradiction since $\text{Prepare}(p')$ would fail even though $A \in Q_{a,p'}$). By our minimality assumption, we know that all proposals z such that $p \leq z < p'$ fail or re-propose v . Therefore, the acceptor A responds with $\text{Promise}(q, \text{val}_v, \dots)$.

At this point, the proposer has received at least one Promise message with a non-empty value. Therefore, it does not take the Phase 1 fast path and waits until it has a Phase 1b quorum (denoted $Q_{b,p'}$). Using the same logic as above, the proposer for p' will receive a minimum of k Promise messages each referencing value v since there are k acceptors in $Q_{b,p'} \cap Q_{2,p}$ (Quorum Property 3). Since the proposer has a minimum of k responses for v , it can reconstruct value v . Let q denote the highest proposal number among all k responses.

Besides those in $Q_{b,p'} \cap Q_{2,p}$, other acceptors in $Q_{b,p'}$ may return values that differ from v . We consider the proposal number q' for each of these accepted values:

- Case 1: $q' < q$. The proposer for p' will ignore the value for q' since it uses the highest proposal number for which it has k splits.
- Case 2: $p' < q'$. Not possible since $\text{Prepare}(p')$ would have failed.
- Case 3: $p < q' < p'$. This implies that a $\text{Propose}(q', v'')$ was issued where $v'' \neq v$. This violates our minimality assumption.

Therefore, the proposer will adopt value v since it can reconstruct it (the proposer has k splits from the acceptors in $Q_{b,p'} \cap Q_{2,p}$ alone) and the highest returned proposal number references it. This contradicts our assumption that $v' \neq v$. \square

Theorem 3. (*Consistency*) PANDO will choose at most one value.

Proof. Assume that two different proposals with proposal numbers p and q are chosen. Since proposers use globally unique proposal numbers, $p \neq q$. This implies that one of the proposal numbers is greater than the other, assume that $q > p$. By Lemma 1, the two proposals write the same value. \square

Theorem 4. (*Stability*) *Once a value is chosen by PANDO, no other value may be chosen.*

Proof. The proposal numbers used for any two chosen proposals will not be equal. Thus, with the additional assumption that acceptors store their state in durable storage, this follows immediately from Lemma 1. \square

4 TLA+ specification for PANDO reads and writes

MODULE *Pando*

EXTENDS *Integers*, *TLC*, *FiniteSets*

CONSTANTS *Acceptors*, *Ballots*, *Values*,
 Quorum1a, *Quorum1b*, *Quorum2*, *K*

ASSUME *QuorumAssumption* \triangleq
 \wedge *Quorum1a* \subseteq SUBSET *Acceptors*
 \wedge *Quorum1b* \subseteq SUBSET *Acceptors*
 \wedge *Quorum2* \subseteq SUBSET *Acceptors*
 Overlap of 1
 $\wedge \forall QA \in Quorum1a :$
 $\forall Q2 \in Quorum2 :$
 $Cardinality(QA \cap Q2) \geq 1$
 Overlap of *K*
 $\wedge \forall QB \in Quorum1b :$
 $\forall Q2 \in Quorum2 :$
 $Cardinality(QB \cap Q2) \geq K$

VARIABLES *msgs*, The set of messages that have been sent
 maxPBal, *maxPBal*[*a*] is the highest promised ballot (proposal number) at acceptor *a*
 maxABal, *maxABal*[*a*] is the highest accepted ballot (proposal number) at acceptor *a*
 maxVal, *maxVal*[*a*] is the value for *maxABal*[*a*] at acceptor *a*
 chosen, *chosen*[*a*] is the value that acceptor *a* heard was chosen (or else is *None*)
 readLog *readLog*[*b*] is the value that was read during ballot *b*

vars $\triangleq \langle msgs, maxPBal, maxABal, maxVal, chosen, readLog \rangle$
None \triangleq CHOOSE *v* : *v* \notin *Values*

Type invariants.

Messages \triangleq
 $[type : \{\text{"prepare"}\}, bal : Ballots]$
 \cup $[type : \{\text{"promise"}\}, bal : Ballots, maxABal : Ballots \cup \{-1\},$
 $maxVal : Values \cup \{None\}, acc : Acceptors,$
 $chosen : Values \cup \{None\}]$
 \cup $[type : \{\text{"propose"}\}, bal : Ballots, val : Values \cup \{None\},$
 $op : \{\text{"R"}, \text{"W"}\}]$
 \cup $[type : \{\text{"accept"}\}, bal : Ballots, val : Values, acc : Acceptors,$

$$\begin{aligned}
& \cup \quad \begin{aligned} & op : \{\text{"R"}, \text{"W"}\} \\ & [type : \{\text{"learn"}\}, bal : Ballots, val : Values] \end{aligned} \\
TypeOK & \triangleq \begin{aligned} & \wedge msgs \in \text{SUBSET Messages} \\ & \wedge maxABal \in [Acceptors \rightarrow Ballots \cup \{-1\}] \\ & \wedge maxPBal \in [Acceptors \rightarrow Ballots \cup \{-1\}] \\ & \wedge maxVal \in [Acceptors \rightarrow Values \cup \{None\}] \\ & \wedge chosen \in [Acceptors \rightarrow Values \cup \{None\}] \\ & \wedge readLog \in [Ballots \rightarrow Values \cup \{None\}] \\ & \wedge \forall a \in Acceptors : maxPBal[a] \geq maxABal[a] \end{aligned}
\end{aligned}$$

Initial state.

$$\begin{aligned}
Init & \triangleq \begin{aligned} & \wedge msgs = \{\} \\ & \wedge maxPBal = [a \in Acceptors \mapsto -1] \\ & \wedge maxABal = [a \in Acceptors \mapsto -1] \\ & \wedge maxVal = [a \in Acceptors \mapsto None] \\ & \wedge chosen = [a \in Acceptors \mapsto None] \\ & \wedge readLog = [b \in Ballots \mapsto None] \end{aligned}
\end{aligned}$$

Send message m .

$$Send(m) \triangleq msgs' = msgs \cup \{m\}$$

Prepare: The proposer chooses a ballot id and broadcasts prepare requests to all acceptors.
All writes start here.

$$\begin{aligned}
Prepare(b) & \triangleq \begin{aligned} & \wedge \neg \exists m \in msgs : (m.type = \text{"prepare"}) \wedge (m.bal = b) \\ & \wedge Send([type \mapsto \text{"prepare"}, bal \mapsto b]) \\ & \wedge \text{UNCHANGED} \langle maxPBal, maxABal, maxVal, chosen, readLog \rangle \end{aligned}
\end{aligned}$$

Promise: If an acceptor receives a prepare request with ballot id greater than that of any prepare request which it has already responded to, then it responds to the request with a promise. The promise reply contains the proposal (if any) with the highest ballot id that it has accepted.

$$\begin{aligned}
Promise(a) & \triangleq \\ & \begin{aligned} & \exists m \in msgs : \\ & \quad \wedge m.type = \text{"prepare"} \\ & \quad \wedge m.bal > maxPBal[a] \\ & \quad \wedge Send([type \mapsto \text{"promise"}, acc \mapsto a, bal \mapsto m.bal, \\ & \quad \quad \quad maxABal \mapsto maxABal[a], maxVal \mapsto maxVal[a], \\ & \quad \quad \quad chosen \mapsto chosen[a]]) \\ & \quad \wedge maxPBal' = [maxPBal \text{ EXCEPT } ![a] = m.bal] \\ & \quad \wedge \text{UNCHANGED} \langle maxABal, maxVal, chosen, readLog \rangle \end{aligned}
\end{aligned}$$

Propose (fast path): The proposer waits until it collects promises from a Phase 1a quorum of acceptors. If no previous value is found, then the proposer can skip to Phase 2 with its own value.

$$\begin{aligned}
ProposeA(b) & \triangleq \\ & \begin{aligned} & \wedge \neg \exists m \in msgs : (m.type = \text{"propose"}) \wedge (m.bal = b) \\ & \wedge \exists v \in Values : \\ & \quad \wedge \exists Q \in Quorum1a : \end{aligned}
\end{aligned}$$

LET $Q1Msgs \triangleq \{m \in msgs : \wedge m.type = \text{"promise"} \\ \wedge m.bal = b \\ \wedge m.acc \in Q\}$
 IN
 Check for promises from all acc in Q
 $\wedge \forall a \in Q : \exists m \in Q1Msgs : m.acc = a$
 Make sure no previous vals have been returned in promises
 $\wedge \forall m \in Q1Msgs : m.maxABal = -1$
 $\wedge Send([type \mapsto \text{"propose"}, bal \mapsto b, val \mapsto v, op \mapsto \text{"W"}])$
 $\wedge \text{UNCHANGED} \langle maxPBal, maxABal, maxVal, chosen, readLog \rangle$

Propose (slow path): The proposer waits for promises from a Phase 1b quorum of acceptors. If no value is found accepted, then the proposer can pick its own value for the next phase. If any accepted coded split is found in one of the promises, the proposer detects whether there are at least K splits (for the particular value) in these promises. Next, the proposer picks up the recoverable value with the highest ballot, and uses it for next phase.

$ProposeB(b) \triangleq$
 $\wedge \neg \exists m \in msgs : (m.type = \text{"propose"}) \wedge (m.bal = b)$
 $\wedge \exists Q \in Quorum1b :$
 LET $Q1Msgs \triangleq \{m \in msgs : \wedge m.type = \text{"promise"} \\ \wedge m.bal = b \\ \wedge m.acc \in Q\}$
 $Q1Vals \triangleq [v \in Values \cup \{None\} \mapsto \\ \{m \in Q1Msgs : m.maxVal = v\}]$
 IN
 Check that all acceptors from Q responded
 $\wedge \forall a \in Q : \exists m \in Q1Msgs : m.acc = a$
 $\wedge \exists v \in Values :$
 \wedge No recoverable value, use anything
 $\vee \forall vv \in Values : Cardinality(Q1Vals[vv]) < K$
 Check if v is recoverable and of highest ballot
 \vee Use previous val if K splits exist
 $\wedge Cardinality(Q1Vals[v]) \geq K$
 $\wedge \exists m \in Q1Vals[v] :$
 Ensure no other recoverable value has a higher ballot
 $\wedge \forall mm \in Q1Msgs :$
 $\vee m.bal \geq mm.bal$
 $\vee Cardinality(Q1Vals[mm.maxVal]) < K$
 $\wedge Send([type \mapsto \text{"propose"}, bal \mapsto b, val \mapsto v, op \mapsto \text{"W"}])$
 $\wedge \text{UNCHANGED} \langle maxPBal, maxABal, maxVal, chosen, readLog \rangle$

Phase 2: If an acceptor receives an accept request with ballot i , it accepts the proposal unless it has already responded to a prepare request having a ballot greater than it does.

$Accept(a) \triangleq$
 $\wedge \exists m \in msgs :$
 $\wedge m.type = \text{"propose"}$
 $\wedge m.bal \geq maxPBal[a]$

$$\begin{aligned}
&\wedge \text{maxABal}' = [\text{maxABal} \text{ EXCEPT } ![a] = m.\text{bal}] \\
&\wedge \text{maxPBal}' = [\text{maxPBal} \text{ EXCEPT } ![a] = m.\text{bal}] \\
&\wedge \text{maxVal}' = [\text{maxVal} \text{ EXCEPT } ![a] = m.\text{val}] \\
&\wedge \text{Send}([type \mapsto \text{"accept"}, \text{bal} \mapsto m.\text{bal}, \text{acc} \mapsto a, \text{val} \mapsto m.\text{val}, \\
&\quad \text{op} \mapsto m.\text{op}]) \\
&\wedge \text{UNCHANGED} \langle \text{chosen}, \text{readLog} \rangle
\end{aligned}$$

ProposerEnd: If the proposer receives acknowledgements from a Phase 2 quorum, then it knows that the value was chosen and broadcasts this.

$$\begin{aligned}
&\text{ProposerEnd}(b) \triangleq \\
&\quad \wedge \exists v \in \text{Values} : \\
&\quad \wedge \exists Q \in \text{Quorum2} : \\
&\quad \quad \text{LET } Q2\text{msgs} \triangleq \{m \in \text{msgs} : \wedge m.\text{type} = \text{"accept"} \\
&\quad \quad \quad \wedge m.\text{bal} = b \\
&\quad \quad \quad \wedge m.\text{val} = v \\
&\quad \quad \quad \wedge m.\text{acc} \in Q\} \\
&\quad \text{IN} \\
&\quad \quad \text{Check for accept messages from all members of } Q \\
&\quad \quad \wedge \forall a \in Q : \exists m \in Q2\text{msgs} : m.\text{acc} = a \\
&\quad \quad \text{If this was in response to a read, log the result} \\
&\quad \quad \wedge \text{Read: log the result} \\
&\quad \quad \quad \vee \wedge \exists m \in Q2\text{msgs} : m.\text{op} = \text{"R"} \\
&\quad \quad \quad \quad \wedge \text{readLog}' = [\text{readLog} \text{ EXCEPT } ![b] = v] \\
&\quad \quad \quad \text{Write: don't log the result} \\
&\quad \quad \quad \vee (\forall m \in Q2\text{msgs} : m.\text{op} = \text{"W"} \wedge \text{UNCHANGED} \langle \text{readLog} \rangle) \\
&\quad \quad \wedge \text{Send}([type \mapsto \text{"learn"}, \text{bal} \mapsto b, \text{val} \mapsto v]) \\
&\quad \wedge \text{UNCHANGED} \langle \text{maxABal}, \text{maxPBal}, \text{maxVal}, \text{chosen} \rangle
\end{aligned}$$

Learn: A proposer has announced that value v is chosen.

$$\begin{aligned}
&\text{Learn}(a) \triangleq \\
&\quad \wedge \exists m \in \text{msgs} : \\
&\quad \quad \wedge m.\text{type} = \text{"learn"} \\
&\quad \quad \text{Process accept before learn, needed for ReadInv, not the protocol} \\
&\quad \quad \wedge \text{maxABal}[a] \geq m.\text{bal} \\
&\quad \quad \wedge \text{chosen}' = [\text{chosen} \text{ EXCEPT } ![a] = m.\text{val}] \\
&\quad \wedge \text{UNCHANGED} \langle \text{msgs}, \text{maxPBal}, \text{maxABal}, \text{maxVal}, \text{readLog} \rangle
\end{aligned}$$

Count how many splits of v we have received.

$$\text{CountSplitsOf}(\text{resps}, v) \triangleq \text{Cardinality}(\{m \in \text{resps} : m.\text{maxVal} = v\})$$

FastRead: Check if any value returned from a Phase 1a quorum was chosen. If we have enough splits to reconstruct that value, then return immediately. If not, wait for Phase 1b quorum. If we have a value that was marked chosen, return. Otherwise, perform a write-back.

$$\begin{aligned}
&\text{FastRead}(b) \triangleq \\
&\quad \wedge \neg \exists m \in \text{msgs} : (m.\text{type} = \text{"propose"}) \wedge (m.\text{bal} = b) \\
&\quad \wedge \\
&\quad \quad \text{Fastest path: Phase 1a quorum has } k \text{ splits and the value is chosen}
\end{aligned}$$

$\vee \wedge \exists Q \in \text{Quorum1a} :$
 $\text{LET } R\text{Msgs} \triangleq \{m \in \text{msgs} : \wedge m.type = \text{"promise"} \\ \wedge m.bal = b \\ \wedge m.acc \in Q\}$
 IN $\text{Check that all acceptors from } Q \text{ responded}$
 $\wedge \forall a \in Q : \exists m \in R\text{Msgs} : m.acc = a$
 $\text{Check that we have } k \text{ splits of a chosen value}$
 $\wedge \exists m \in R\text{Msgs} :$
 $\wedge m.chosen \neq \text{None}$
 $\wedge \text{CountSplitsOf}(R\text{Msgs}, m.chosen) \geq K$
 $\wedge \text{readLog}' = [\text{readLog EXCEPT } ![b] = m.chosen]$
 $\wedge \text{UNCHANGED} \langle \text{msgs}, \text{maxPBal}, \text{maxABal}, \text{maxVal}, \text{chosen} \rangle$
 Fast path: Phase 1b quorum has k splits and the value is chosen
 $\vee \wedge \exists Q \in \text{Quorum1b} :$
 $\text{LET } R\text{Msgs} \triangleq \{m \in \text{msgs} : \wedge m.type = \text{"promise"} \\ \wedge m.bal = b \\ \wedge m.acc \in Q\}$
 IN $\text{Check that all acceptors from } Q \text{ responded}$
 $\wedge \forall a \in Q : \exists m \in R\text{Msgs} : m.acc = a$
 $\text{Check that we have } k \text{ splits of a chosen value}$
 $\wedge \exists m \in R\text{Msgs} :$
 $\wedge m.chosen \neq \text{None}$
 $\wedge \text{CountSplitsOf}(R\text{Msgs}, m.chosen) \geq K$
 $\wedge \text{readLog}' = [\text{readLog EXCEPT } ![b] = m.chosen]$
 $\wedge \text{UNCHANGED} \langle \text{msgs}, \text{maxPBal}, \text{maxABal}, \text{maxVal}, \text{chosen} \rangle$
 Slow path: Phase 1b recovery and write back
 $\vee \wedge \exists Q \in \text{Quorum1b} :$
 $\text{LET } Q1\text{Msgs} \triangleq \{m \in \text{msgs} : \wedge m.type = \text{"promise"} \\ \wedge m.bal = b \\ \wedge m.acc \in Q\}$
 $Q1\text{Vals} \triangleq [v \in \text{Values} \cup \{\text{None}\} \mapsto \\ \{m \in Q1\text{Msgs} : m.maxVal = v\}]$
 IN $\text{Check that all acceptors from } Q \text{ responded}$
 $\wedge \forall a \in Q : \exists m \in Q1\text{Msgs} : m.acc = a$
 $\wedge \exists v \in \text{Values} :$
 $\text{Check if } v \text{ is recoverable and of highest ballot}$
 $\text{Use previous val if } K \text{ splits exist}$
 $\wedge \text{Cardinality}(Q1\text{Vals}[v]) \geq K$
 $\wedge \exists m \in Q1\text{Vals}[v] :$
 $\text{Ensure no other recoverable value has a higher ballot}$
 $\wedge \forall mm \in Q1\text{Msgs} :$
 $\vee m.bal \geq mm.bal$
 $\vee \text{Cardinality}(Q1\text{Vals}[mm.maxVal]) < K$
 $\text{readLog will be updated in ProposerEnd}$
 $\wedge \text{Send}([type \mapsto \text{"propose"}, bal \mapsto b, val \mapsto v,$

$$\begin{aligned}
& op \mapsto \text{"R"}) \\
& \wedge \text{UNCHANGED} \langle \text{maxPBal}, \text{maxABal}, \text{maxVal}, \text{chosen}, \text{readLog} \rangle \\
& \text{No value recovered: Return None} \\
& \vee \wedge \text{readLog}' = [\text{readLog EXCEPT } ![b] = \text{None}] \\
& \wedge \text{UNCHANGED} \langle \text{msgs}, \text{maxPBal}, \text{maxABal}, \text{maxVal}, \text{chosen} \rangle
\end{aligned}$$

Next state.

$$\begin{aligned}
\text{Next} \triangleq & \vee \exists b \in \text{Ballots} : \vee \text{Prepare}(b) \\
& \vee \text{ProposeA}(b) \\
& \vee \text{ProposeB}(b) \\
& \vee \text{ProposerEnd}(b) \\
& \vee \text{FastRead}(b) \\
& \vee \exists a \in \text{Acceptors} : \text{Promise}(a) \vee \text{Accept}(a) \vee \text{Learn}(a)
\end{aligned}$$

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

Invariant helpers.

$$\text{AllChosenWereAcceptedByPhase2} \triangleq$$

$$\begin{aligned}
& \forall a \in \text{Acceptors} : \\
& \vee \text{chosen}[a] = \text{None} \\
& \vee \exists Q \in \text{Quorum2} : \\
& \quad \forall a2 \in Q : \\
& \quad \quad \exists m \in \text{msgs} : \wedge m.\text{type} = \text{"accept"} \\
& \quad \quad \wedge m.\text{acc} = a2 \\
& \quad \quad \wedge m.\text{val} = \text{chosen}[a]
\end{aligned}$$

$$\text{OnlyOneChosen} \triangleq$$

$$\begin{aligned}
& \forall a, aa \in \text{Acceptors} : \\
& (\text{chosen}[a] \neq \text{None} \wedge \text{chosen}[aa] \neq \text{None}) \Rightarrow (\text{chosen}[a] = \text{chosen}[aa])
\end{aligned}$$

$$\begin{aligned}
\text{VotedForIn}(a, v, b) \triangleq & \exists m \in \text{msgs} : \wedge m.\text{type} = \text{"accept"} \\
& \wedge m.\text{val} = v \\
& \wedge m.\text{bal} = b \\
& \wedge m.\text{acc} = a
\end{aligned}$$

$$\begin{aligned}
\text{ProposedValue}(v, b) \triangleq & \exists m \in \text{msgs} : \wedge m.\text{type} = \text{"propose"} \\
& \wedge m.\text{val} = v \\
& \wedge m.\text{bal} = b \\
& \wedge m.\text{op} = \text{"W"}
\end{aligned}$$

$$\text{NoOtherFutureProposal}(v, b) \triangleq$$

$$\begin{aligned}
& \forall vv \in \text{Values} : \\
& \quad \forall bb \in \text{Ballots} : \\
& \quad (bb > b \wedge \text{ProposedValue}(vv, bb)) \Rightarrow v = vv
\end{aligned}$$

$$\text{ChosenIn}(v, b) \triangleq \exists Q \in \text{Quorum2} : \forall a \in Q : \text{VotedForIn}(a, v, b)$$

$$\text{ChosenBy}(v, b) \triangleq \exists b2 \in \text{Ballots} : (b2 \leq b \wedge \text{ChosenIn}(v, b2))$$

$$Chosen(v) \triangleq \exists b \in Ballots : ChosenIn(v, b)$$

Invariants.

$$LearnInv \triangleq AllChosenWereAcceptedByPhase2 \wedge OnlyOneChosen$$

$$ReadInv \triangleq \forall b \in Ballots : readLog[b] = None \vee ChosenBy(readLog[b], b)$$

$$ConsistencyInv \triangleq \forall v1, v2 \in Values : Chosen(v1) \wedge Chosen(v2) \Rightarrow (v1 = v2)$$

$$StabilityInv \triangleq$$

$$\forall v \in Values : \forall b \in Ballots : ChosenIn(v, b) \Rightarrow NoOtherFutureProposal(v, b)$$

$$AcceptorInv \triangleq$$

$$\forall a \in Acceptors :$$

$$\wedge (maxVal[a] = None) \equiv (maxABal[a] = -1)$$

$$\wedge maxABal[a] \leq maxPBal[a]$$

$$\wedge (maxABal[a] \geq 0) \Rightarrow VotedForIn(a, maxVal[a], maxABal[a])$$

$$\wedge \forall c \in Ballots :$$

$$c > maxABal[a] \Rightarrow \neg \exists v \in Values : VotedForIn(a, v, c)$$

References

- [1] Azure Cosmos DB - globally distributed database service — Microsoft Azure. <https://azure.microsoft.com/en-us/services/cosmos-db/>.
- [2] Anonymous. Near-optimal latency versus cost tradeoffs in geo-distributed storage.
- [3] J. Baker, C. Bond, J. C. Corbett, J. Furman, A. Khorlin, J. Larson, J.-M. Leon, Y. Li, A. Lloyd, and V. Yushprakh. Megastore: Providing scalable, highly available storage for interactive services. In *CIDR*, 2011.
- [4] M. Burrows. The chubby lock service for loosely-coupled distributed systems. In *OSDI*, 2006.
- [5] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rolig, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, and D. Woodford. Spanner: Google’s globally-distributed database. In *OSDI*, 2012.
- [6] H. Howard, D. Malkhi, and A. Spiegelman. Flexible Paxos: Quorum Intersection Revisited. In *OPODIS*, 2016.
- [7] L. Lamport. The part-time parliament. *ACM Transactions on Computer Systems (TOCS)*, 16(2):133–169, 1998.
- [8] I. Moraru, D. G. Andersen, and M. Kaminsky. Paxos quorum leases: Fast reads without sacrificing writes. In *SoCC*, 2014.
- [9] S. Mu, K. Chen, Y. Wu, and W. Zheng. When Paxos meets erasure code: Reduce network and storage cost in state machine replication. In *HPDC*, 2014.
- [10] Z. Shen, Q. Jia, G.-E. Sela, B. Rainero, W. Song, R. van Renesse, and H. Weatherspoon. Follow the sun through the clouds: Application migration for geographically shifting workloads. In *Proceedings of the Seventh ACM Symposium on Cloud Computing*, pages 141–154. ACM, 2016.