

An Efficient and Fault-Tolerant Update Commitment Protocol for Weakly Connected Replicas ^{*}

João Barreto^{**} and Paulo Ferreira^{***}

INESC-ID / IST

Rua Alves Redol, 9, 1000 Lisboa, Portugal
{joao.barreto, paulo.ferreira}@inesc-id.pt

Abstract. Mobile and other loosely-coupled environments call for decentralized optimistic replication protocols that provide highly available access to shared objects, while ensuring eventual convergence towards a strongly consistent state. In this paper we propose a novel epidemic weighted voting protocol for achieving such goal. Epidemic weighted voting approaches eliminate the single point of failure limitation of primary commit approaches. Our protocol introduces a significant improvement over other epidemic weighted voting solutions by allowing multiple, *happened-before* related updates to be committed at a single distributed election round. We demonstrate that our protocol is especially advantageous with the weak connectivity levels that characterize mobile and other loosely-coupled networks. We support such claims by presenting comparison results obtained from side-by-side execution of reference protocols in a simulated environment.

1 Introduction

Data replication is a fundamental mechanism for most distributed systems for performance, scalability and fault tolerance reasons. In particular, optimistic replication protocols [1] are of extreme importance in mobile and other loosely-coupled network environments. The nature of these environments calls for decentralized replication protocols that are able to provide highly available full access to shared objects. Such requirement is accomplished by optimistic replication strategies, which, in contrast to their pessimistic counterparts, enable updates to be issued at any one replica regardless of the availability of other replicas.

As a trade-off, the issue of consistency in optimistic replication is problematic. Since replicas are allowed to be updated at any time and circumstance, updates may conflict if issued concurrently at distinct replicas. Some optimistic replication protocols ensure that, from such a possibly inconsistent *tentative* state, replicas evolve towards an eventual consistent *stable* state. For this end,

^{*} This is a preprint of an article accepted for publication in *Concurrency and Computation: Practice and Experience* Copyright (c) 2006.

^{**} Funded by FCT Grant SFRH/BD/13859

^{***} Funded by FCT Project UbiRep (POSI/CHS/47832/2002)

a distributed consensus algorithm is executed so as to reach an agreement on a common order in which tentative updates should be committed.

There are many scenarios where users, in order to benefit from high availability, are willing to work with temporarily tentative data, provided that a commitment agreement regarding such data will eventually be reached. Consider, for instance, a laptop user that becomes disconnected from his corporate file server after leaving his office. If necessary, he may expect to be able to modify a report that is currently replicated at his laptop, even if tentatively.

Furthermore, such worker may meet other mobile team colleagues carrying their replicas and, in an ad-hoc fashion, establish a short term work group to collaboratively work on the report. A set of *happened-before* related [2] tentative updates will result from such activity. Hopefully, if no update is concurrently issued from outside the group, such tentative work will be eventually committed by the underlying consistency protocol. Hence, the high availability provided by an optimistic replication strategy is especially interesting in the above scenarios. However, the usefulness of such approach strongly depends on the ability of the underlying replication protocol to efficiently achieve a commitment decision concerning the tentatively issued data. Users are typically not inclined towards working on tentative data unless they trust the protocol to rapidly provide the with strong consistency guarantees on such data.

Aiming at such central objective, this paper proposes a novel optimistic replication protocol, called *Version Vector Weighted Voting (VVWV)*, for efficient and highly available update commitment through the use of an epidemic weighted voting protocol based on version vectors [3]. The use of a voting approach eliminates the single point of failure of primary commit approaches [4]. Hence, the unavailability of any individual replica does not prohibit the progress of the update commitment process. Moreover, commitment agreement is accomplished without the need for a plurality quorum of replica hosts to be simultaneously accessible: voting information flows epidemically between replicas and update commitment is based solely on local information.

VVWV introduces a significant improvement over basic epidemic weighted voting solutions [5] by allowing multiple update candidates to participate in an election. By using version vectors, candidates consisting of one or more *happened-before* related updates may be voted and committed by running a single distributed election round. As a result, the overall number of anti-entropy sessions required to commit updates is decreased when compared to a basic weighted voting protocol [5]. Hence, update commitment delay is minimized; thus eventual strong consistency guarantees are more rapidly delivered to applications. Namely, such reduction is substantial in scenarios where frequent *happened-before* related updates are tentatively generated by applications. The examples presented above are representative of such update patterns. In worst case scenarios, VVWV behaves similarly to basic weighted voting protocols.

The paper is organized as follows. Section 2 describes related work, Section 3 introduces the protocol, which is evaluated in Section 4, and Section 5 concludes.

2 Related Work

The issue of optimistic data replication for mobile and loosely coupled environments has been addressed by a number of projects [1], with the common intent of offering high data availability. Most of the proposed solutions share the goal of our work by enforcing eventual convergence towards a strongly consistent stable form that is explicitly presented to applications.

Three main approaches can be distinguished. Firstly, Golding [6] proposes that each individual host commits an update when it is certain that it has been received by every replica. A main limitation is that the unavailability of any single replica stalls the entire commitment process. On the other hand, a primary commit strategy, such as the one adopted by Bayou [4], centralizes the commitment process in a single distinguished primary replica that establishes a total commit order over the updates it receives. Primary commit is able to rapidly commit updates, since it suffices for an update to be received by the primary replica to become committed, provided that no conflict is found. However, should the primary replica become unavailable, the commitment progress of updates generated by replicas other than the primary is inevitably halted.

Finally, a third approach uses voting so as to allow a plurality quorum to commit an update. In particular, Deno [5] relies on an epidemic voting protocol to support object replication in a transactional framework for loosely-connected environments. Deno requires one entire election round to be completed in order to commit each single update, if only non-commutable updates are considered. This is acceptable when applications are interested in knowing the commitment outcome of each tentatively issued update before issuing the next one. However, in the usage scenarios addressed by this paper, users and applications will often be interested in issuing multiple tentative updates before acknowledging their commitment. In such situations, the commitment delay imposed by Deno's voting protocol becomes unacceptably higher than that of primary commit.

3 Consistency Protocol

The following sections consider a model where a set of logical objects is replicated at N hosts. An object replica at a given host provides local applications with access to a version of the object contents, as stored by the replica. Such accesses may read or modify the object contents. In the case of the latter, an update is issued by the host and applied to the replica.

Updates issued at a given replica are propagated to other hosts in an epidemic fashion in order to eventually achieve object consistency. The local execution of an update is assumed to be recoverable, atomic and deterministic. The former means that a replica will not reach an inconsistent value if it fails before the update execution completes. It follows from the other two properties that the execution of the same ordered sequence of updates at two distinct replicas in the same initial consistent state will yield an identical final state. For simplicity and without loss of generality, we consider that each logical object is replicated at

every host in the system. For the sake of generality, the set of replicas may be dynamic, and thus change with the creation or removal of new hosts.

Hereafter, we assume an asynchronous system in which hosts can only fail silently. Network partitions may also occur, thus restricting connectivity between hosts that happen to be located in distinct partitions.

3.1 Overview

Due to the optimistic nature of VVWV, updates issued at a local replica are not immediately committed at every remaining replica. Instead, such updates are considered tentative since diverging sequences of updates may still be issued at other replicas. VVWV is responsible for eventually committing one of such diverging tentative sequences of updates at every replica, therefore ensuring eventual strong consistency.

Weighted Voting Commitment VVWV achieves this goal through a weighted voting approach [5], where concurrent tentative updates are regarded as rival candidates in an election. The hosts replicating a given logical object act as voters whose votes determine the outcome of each election between candidate updates to the object. A candidate update wins an election by collecting a plurality of votes, in which case it is committed and its rival candidates are discarded.

Elections consider a fixed per-object currency scheme, in which each voter is associated with a given amount of currency that determines its weight during voting rounds. The global currency of a logical object, distributed among its replica hosts, equals a fixed amount of 1. Currencies can be exchanged between hosts and the currency held by failed hosts can be recovered by running a *currency reevaluation* election, as discussed in [7].

Version Vector Candidates In some cases, applications will be interested in generating more than one tentative update prior to its commitment decision. These may include disconnected mobile applications and ad-hoc groups of mobile applications working cooperatively in the absence of a plurality quorum. Since the commitment decision may not be taken in the short-term, these applications may wish to issue a sequence of multiple, *happened-before* ordered tentative updates.

In order to efficiently accommodate for such update scenarios, VVWV employs version vectors to identify candidate updates in a weighted voting protocol. The flexibility brought by version vectors allows a sequence of one or more updates to run for the current election as a whole. In this case, the candidate is represented by the version vector corresponding to the tentative version obtained if the entire update sequence was applied to the replica. As the next sections explain, the voting protocol relies on the expressiveness of version vectors to decide if the update sequence or a prefix of it are to become committed. Consequently, candidates consisting of one or more *happened-before* related updates may be committed on a single distributed election round. In weakly connected

network environments, where such update patterns are expectably dominant, a substantial reduction of the update commitment delay is therefore achievable.

Hereafter, we use the notation $x \leq y$ to denote that $\forall i, x[i] \leq y[i]$, and $x < y$ to denote that $x \leq y$ and $\exists j : x[j] < y[j]$. Further, $x \parallel y$ means that x and y are conflicting, that is, neither $x \leq y$ nor $y \leq x$.

As the next sections describe in greater detail, voting information flows epidemically among hosts and the decision to commit an update is based only on local replica information. These are important properties for operation under mobile and loosely-coupled environments. Section 3.2 describes the state maintained at each replica and the two distinctly consistent views that the protocol offers of each replica. Section 3.3 addresses the storage of tentative updates and their corresponding commitment upon a replica value. Then, Section 3.4 describes the epidemic flow of consistency information and Section 3.5 finally defines how candidates are elected.

3.2 Replica State and Access to Stable and Tentative Views

Each replica r maintains the following state:

- $stable_r$, which consists of a version vector that identifies the most recent stable version that r is currently aware of;
- $votes_r[1..N]$, which stores, for each host $k = 1, 2, \dots, N$, the version vector corresponding to the candidate voted for by k , as known by r ; or \perp , if the vote of such host has not yet been known to r ;
- $cur_r[1..N]$, which stores, for each host $k = 1, 2, \dots, N$ whose vote replica r has knowledge of, the currency associated with such vote;
- $committed_r[1..c_r]$, which stores each committed update at r according to the agreed commitment order (where c_r denotes the number of committed updates at r);

Each host is able to offer two possibly distinct views over the value of a replica r to its applications and users: the stable and tentative views. The first view reflects a strongly consistent value of the replicated object that is obtained by the ordered application of the updates in $committed_r$. On the other hand, the tentative view exposes a weakly consistent value that corresponds to the candidate version that is currently voted by the local host, $votes_r[r]$. Update requests from applications are performed upon the tentative view. Hence, issuing a tentative update u on a replica r changes the state of r as follows.

1. If $votes_r[r] = \perp$, then $votes_r[r] \leftarrow adv_r(stable_r)$ ¹ and $cur_r[r] = currency_r$;
2. Otherwise, $votes_r[r] \leftarrow adv_r(votes_r[r])$;

In case r has not voted yet, then it casts a vote for the candidate representing u , issued upon the current stable version. Otherwise, it means that u is issued upon the value resulting from the ordered application of the tentative updates of the current vote of r . In this case, its vote is extended to a candidate that represents u and the tentative updates that precede it.

¹ adv_r advances the counter corresponding to r in the supplied version vector by one.

3.3 Update Commitment

The protocol proposed hereafter is orthogonal to the issues of actual transference and storage of tentative updates. In particular, VVWV does not impose the decision of whether to transfer and store, at each individual replica, (1) the tentative updates belonging to every candidate in the current election or, alternatively, (2) only those concerning the replica's own candidate.

This means that, at the time a host determines that a given candidate has won the election and, thus, its updates should be committed, such updates may not be immediately available. Instead, they will be eventually collected through succeeding anti-entropy sessions with other hosts. Consequently, there may occur a discrepancy between the most recent stable version identified by VVWV at a given replica r and the actual stable value that is locally accessible at r . In fact, the number of updates in *committed_r*, denoted by c_r , may be lower than the number of updates that have actually been determined by VVWV as belonging to the stable path. In such a case, the replica's stable value does not yet reflect the most recent stable version r is aware of.

As a consequence, VVWV is flexible enough to support hosts with differing memory limitations. On one hand, hosts with rich memory resources may store every update associated with each candidate, hence being able to immediately gain access to the most recent known stable value as each new stable version is determined by VVWV. On the other hand, memory-constrained devices may opt to restrict themselves to storing only the updates of their own candidate and, thus, allow for occasional delays in the availability of the most recent stable value when rival candidates win an election. In either case, however, the efficiency of VVWV in taking commitment decisions is not affected. Both strategies may transparently co-exist in a system of replicas of the same logical object.

From the viewpoint of VVWV, the procedure for committing a sequence of updates u_1, \dots, u_m is therefore comprised of the following steps:

1. For each update, $u_k (k = 1..m)$, $committed_r[c_r + k] \leftarrow u_k$;
2. $c_r \leftarrow c_r + m$;

3.4 Anti-entropy

Voting information is propagated through the system by anti-entropy sessions established between pairs of accessible hosts. An anti-entropy session is an unidirectional pull-based interaction in which a requesting host, holding replica A , updates its local election knowledge with information obtained from another host, holding replica B . In case B has more up-to-date election information, it transfers such information to A . Furthermore, if A has not yet voted for a candidate that is concurrent to the one voted for by B , A accepts the latter, thus contributing to its election.

Each anti-entropy session is carried out according to the following procedure, which should be executed atomically:

1. If $stable_A < stable_B$ then
 - (a) $stable_A \leftarrow stable_B$;

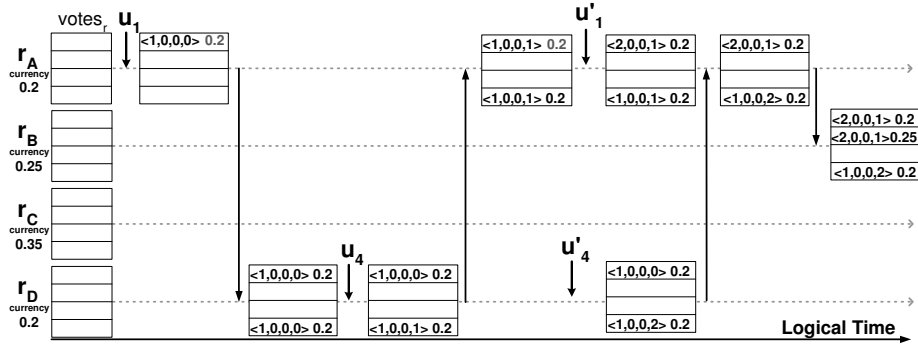


Fig. 1. Example of update generation and anti-entropy: four replicas with unevenly distributed currencies start from a common initial stable version $stable_r = \langle 0, 0, 0, 0 \rangle$.

- (b) Let $toCommit_A = u_1, \dots, u_m$ be the totally ordered sequence of uncommitted updates that, starting from the currently committed value at A , produces the new $stable_A$ version. If a prefix of $toCommit_A$, u_1, \dots, u_k ($k \leq m$) is locally available at A , then commit such updates in that order;
- (c) If $c_A < c_B$ then commit update sequence $committed_B[c_A+1], \dots, committed_B[c_B]$.
- (d) $\forall k$ s.t. $votes_A[k] \parallel stable_A$ or $votes_A[k] \leq stable_A$, then $votes_A[k] \leftarrow \perp$;
2. If $(votes_A[A] = \perp$ and $stable_A < votes_B[B])$ or $votes_A[A] < votes_B[B]$ then $votes_A[A] \leftarrow votes_B[B]$ and $cur_A[A] \leftarrow currency_A$;
3. $\forall k \neq A$ s.t. $(votes_A[k] = \perp$ and $stable_A < votes_B[k])$ or $votes_A[k] < votes_B[k]$ then $votes_A[k] \leftarrow votes_B[k]$ and $cur_A[k] \leftarrow cur_B[k]$.

The first step ensures that, in case B knows about a more recent stable version, A will adopt it. This means that A will regard the elections that originated such new stable version as completed (1a) and, therefore, commit the winner updates that are available at the moment.

In (1b), any winner updates that are already locally stored (tentatively) are committed. After that, the set of committed updates held by B that have not yet been committed at replica A are collected and committed by the latter (1c). Finally, further elections are prepared by keeping only the voting information that will still be meaningful for their outcome (1d). Namely, these are the votes on candidates that succeed (by *happened-before*) the stable version.

As a second step of anti-entropy, A is persuaded to vote for the same candidate as the one voted by B , provided that A has not yet voted for a concurrent candidate (2). Subsequently, A updates its current knowledge of the current election with relevant voting information that may be held by B (3). Namely, A stores each vote that it is not yet aware of or whose candidate is more complete than the one it currently has knowledge of.

Figure 1 depicts an example with four replicas of a given object. After some update activity, a conflict occurs between u'_1 and u'_4 ; hence, two rival candidates end up running for the election.

3.5 Election decision

The candidates being voted in an election represent update paths that traverse through one of more versions beyond the initial point defined by the stable version, *stable*. These possibly divergent candidate update paths may share common prefix sub-paths. The following definition expresses such notion.

Definition 1: Given two version vectors, v_1 and v_2 , their *maximum common version* is given by a version vector, $mcv(v_1, v_2)$, s.t. $\forall k, mcv(v_1, v_2)[k] = \min(v_1[k], v_2[k])$. For simplicity, we represent $mcv(v_1, v_2, \dots, v_m)$ as the result of $mcv(mcv(mcv(v_1, v_2)), \dots, v_m)$.

Theorem 1: Let $v_1, \dots, v_m \in votes_r$, be one or more candidate versions known by replica r , each connoting a tentative update path starting from the stable version, $stable_r$. Their maximum common version, $mcv(v_1, \dots, v_m)$, constitutes the farthest version of an update sub-path that is mutually traversed by the update paths of v_1, \dots, v_m . Complementarily, the total currency voted for such common sub-path is obtained by $voted_r(mcv(v_1, \dots, v_m)) = cur_r[1] + \dots + cur_r[N]$. *Sketch of Proof:* Assume, by absurd, that $m = mcv(v_1, \dots, v_m)$, is not the farthest version of an update sub-path that is shared by the update paths of v_1, \dots, v_m ; instead, such version is given by $m' \neq m$. By definition of mcv , $m' \not\geq m$; otherwise, $m' \leq v_i, \forall i$, wouldn't be verified. So, $m' < m$, which implies that $\exists k$ s.t. $m'[k] < m[k] \leq v_1[k]$ and $m'[k] < m[k] \leq v_2[k]$. Since m' identifies the farthest common sub-path, the differences between $m'[k]$ and $v_1[k]$, as well as $m'[k]$ and $v_2[k]$, must have respectively resulted from concurrent tentative updates generated by replica k . However, replica k is not allowed to issue concurrent updates when holding the same $stable_k$: if $votes_k[k] = \perp$ (Section 3.1), it means that $\nexists i \neq k$ s.t. $votes_k[i]$ represents any tentative update from k . By anti-entropy, this is also verified at any other replica.

VVWV is responsible for progressively determining common sub-paths of candidate versions that manage to obtain a plurality of votes. Before describing this decision, we define the notion of collectable currency by a candidate c : $collectable_r(c) = \sum(cur[k] : votes_r[k] = \perp \vee votes_r[k] < c)$. It is obtained from the currencies of the votes that, despite not supporting c , may evolve to such condition; namely, votes on candidates happening-before c or \perp .

Definition 2: Let w be a version vector s.t. $w = mcv(w_1, \dots, w_m)$ where $w_1, \dots, w_m \in votes_r$ and $1 \leq m \leq N$. w wins an election when:

1. $voted_r(w) > 0.5$, or
2. $\forall l$ s.t. $l = mcv(l_1, \dots, l_k), l_1, \dots, l_k \in votes_r, 1 \leq k \leq N$, and $l \parallel w$, and also $\forall l$ s.t. $l = mcv(l_i : l_i \in votes_r \wedge l_i < w)$,
 - (a) $voted_r(w) > voted_r(l) + collectable_r(l)$, or
 - (b) $voted_r(w) = voted_r(l) + collectable_r(l)$ and $w \prec_r l$.

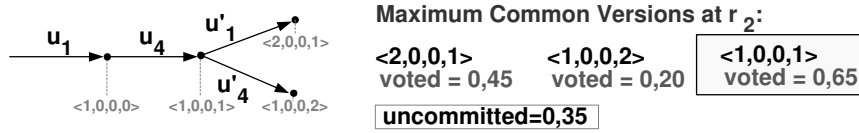


Fig. 2. Election decision for replica r_2 at the final state in Figure 1. Candidate $\langle 1, 0, 0, 1 \rangle$ has collected a plurality of votes and, thus, u_1 and u_4 will be committed in that order.

The above rules state the conditions that guarantee that a candidate has collected sufficient votes to win an election. The votes may constitute a majority, when the amount of currency voted on the winning candidate surpasses 0.5; or a plurality, when the voted currency is greater than the maximum potentially obtainable currency of any other rival candidate. This includes not only common sub-paths of *known* concurrent candidates; but also a hypothetical *unknown* concurrent candidate for which non-concurrent votes can still potentially evolve to. It should be noted that only the non-concurrent votes v_i such that $v_i < w$ may evolve to $v'_i > v_i$ where $v'_i \parallel w$; instead, this is not the case for votes v_j such that $v_j \geq w$, as w will happen-before any $v'_i > v_j$. Finally, the case of ties, decided by the \prec_r relation, is described in Section 3.5.

Determining if a candidate has won an election depends exclusively on information that is locally available at each replica. This means that, once having collected enough voting information, a given replica is able to decide, by its own, to commit a candidate version that locally fulfills the election winning conditions. Hence, update commitment is accomplished in a purely decentralized manner. Figure 2 illustrates the different maximum common versions that replica r_2 analyzes after the example in Figure 1 in order to determine the election outcome; r_2 is able to decide that updates u_1 and u_4 should be committed as soon as available at r_2 .

After finding a new winner version vector, w , a replica r atomically takes the following steps to accept the election decision and prepare for the next election:

1. $stable_r \leftarrow w$;
2. $\forall v_k \in votes_r$ s.t. $v_k \parallel w$ or $v_k \leq w$, $votes_r[k] \leftarrow \perp$;
3. Let $toCommit_r = u_1, \dots, u_m$ be the totally ordered sequence of uncommitted updates that, starting from the currently committed value at r , produces the new $stable_r$ version. If a prefix of $toCommit_r$, u_1, \dots, u_k ($k \leq m$) is locally available at r , then commit such updates in that order;

After accepting the election result by setting the winning version as the new stable version, the second step resets all the defeated candidates to \perp . Depending on the local availability of the updates that belong to the winning candidate, they may be committed into the replica's stable value; otherwise, further anti-entropy sessions will ensure that such updates are eventually collected and committed. A new election can then take place.

Theorem 2 (Correctness): After all elections have been completed at every replica and all updates belonging to the resulting stable path have been commit-

ted at every replica: $\forall r, t$, replica r has committed the same ordered sequence of updates as t .

Sketch of Proof: The proof is based on the one outlined for a basic weighted voting solution [5]. Assume that all replicas start with a common stable value, $stable_0$. It follows directly from the protocol that, if $votes_r[j] = k$ for any r, j, k , then for any l such that $stable_l = stable_r = stable_0$, $votes_l[j]$ will either be \perp or v , where v is comparable with k (i.e., $v \leq k$ or $v > k$). Let $S = v_1, \dots, v_m$, where $v_1, \dots, v_m \in votes_r$, and $w = mcv(v_1, \dots, v_m)$. Assume now that r decides that w wins the election, thus setting $stable_r \leftarrow w$; hence, the currency collected by w prevented any rival $c \parallel w$ to be declared winner. For each host m such that $stable_0 \leq stable_m < w$, it can be shown that, for each $s \in S$, $votes_m[s]$ must either be \perp or comparable with w ; this prevents any rival $c \parallel w$ to be decided winner at m . By anti-entropy, m will eventually receive enough voting information to determine that w has collected a plurality, or directly receive the outcome of the election from another replica p having $stable_p \geq w$. As the totally ordered set of updates that produce w is eventually propagated to every replica, they will be accordingly be committed in that order, the same at every replica.

Tie breaking relation A tie between candidates x and y is decided at replica r by choosing x if the partial relation $x \prec_r y$ holds. Intuitively, $x \prec_r y$ if a host v_x contributes with its vote to the candidate x and no other voter with a lower identifier than v_x contributes, or may ever contribute while x is still a candidate, with its vote to candidate y . Formally, \prec_r is defined as follows:

Definition 3: *Tie breaking relation.* Let x and y be two version vectors so that $x \parallel y$ and that, at replica r , $\exists i, j : stable_r < x \leq votes_r[i]$ and $stable_r < y \leq votes_r[j]$. We say that a tie between x and y at r is broken in favor of x , or $x \prec_r y$, if and only if $\exists v_x$ so that (1) $votes_r[v_x] > mcv(x, y)$ and $votes_r[v_x] \geq x$, and (2) $\forall v_y < v_x$, $votes_r[v_y] > mcv(x, y)$ and $votes_r[v_y] \parallel y$.

Since \prec_r is a partial relation, it may happen that, at a given moment, two candidates x and y are not related by \prec_r . In this case, the tie breaking decision is deferred until sufficient voting information is received to decide either $x \prec_r y$ or $y \prec_r x$ (or until the tie ceases to be verified). The following theorem ensures the correctness of VVWV when election decisions are made by breaking ties using \prec_r . Intuitively, it states that, if a tie between two rival candidates x and y at a given election is broken at a replica r by deciding the victory of x , then no different tie breaking decision will ever be obtained at any replica.

Theorem 3: If, at some moment, $x \prec_r y$ is verified at a replica r , then $y \prec_i x$ will never be verified at any replica $i = 1..N$.

Sketch of Proof: Assume that, at a given moment, $x \prec_r y$ and that v_x is defined as in Definition 3. Then, it follows directly from the definition of \prec_r that $y \not\prec_r x$. Further, assume, by absurd, that $y \prec_{r'} x$, where r' denotes r in some point in the future. This means that either (a) $votes_{r'}[v_x] \not\geq x$, or (b) $\exists v_y < v_x :$

$votes_{r'}[v_y] > mcv(x, y)$ and $votes_{r'}[v_y]$ is comparable to y , i.e. $votes_{r'}[v_y] < y$ or $votes_{r'}[v_y] \geq y$ (where previously, by definition of \prec_r , $votes_r[v_y] \parallel y$ since $x \prec_r y$). According to the protocol, $\forall i, votes_r[i] \neq \perp$, $votes_r[i]$ may only be changed, at each single step, to $v > votes_r[i]$, or to \perp ; it follows from this that (a) and (b) imply that, at some step, $votes_r[v_x] \leftarrow \perp$, or $votes_r[v_y] \leftarrow \perp$, respectively. Such assignment is caused by a new $stable'_r$ version at r such that $stable'_r \not\leq mcv(x, y)$ (otherwise, neither $votes_r[v_x] > mcv(x, y)$ nor $votes_r[v_y] > mcv(x, y)$ would have been changed). So, one of two cases must have occurred: (1) $stable'_r \parallel mcv(x, y)$, and hence $stable'_r \parallel x$ and $stable'_r \parallel y$; or (2) $stable'_r > mcv(x, y)$, and thus (by a corollary of Theorem 1, and as $x \parallel y$) either $stable'_r \parallel x$ or $stable'_r \parallel y$ (or both). It can easily be proven that VVWV ensures that $\forall i, votes_r[i] > stable'_r$; so, neither both x and y will be candidates with $stable'_r$, as at least $x \parallel stable'_r$ or $y \parallel stable'_r$, and therefore $y \not\prec_r x$; this contradicts the initial hypothesis. A similar reasoning extends the proof to the case of other replicas.

4 Evaluation

C# implementations of the primary commit (*Primary*), basic weighted voting (*Basic WV*) and version vector weighted voting (*VVWV*) protocols were run side-by-side in a simulated environment. The simulator includes a collection of replicas of a common logical object, randomly distributed by a set of network partitions. Time is divided into logical time slices; at each time slice, each replica: (1) with a given *mobility probability*, migrates to a different, randomly chosen, network partition; (2) pulls anti-entropy information from a partner, randomly selected from the set replicas present in its current partition; and, (3) generates, with a given *update probability*, one tentative update. Each replica may be active or inactive; in the case of the latter, its update probability is null. An inactive replica exchanges, with a given *activation probability*, its status with an active replica after pulling anti-entropy information from it. The differentiation between active and inactive replicas allows for non-uniform update models to be simulated, namely the hot-spot model [8], which assumes, based on empirical evidence, that updates typically occur in a small set of replicas.

The protocols were evaluated against an increasing number of partitions. Since update contention is prone to arise in a partitioned system, the update commitment delay is not a sufficiently meaningful measure for our purposes, as it does not take into account the discarded updates. Instead, a better evaluation is provided by the update commitment ratio of each protocol, i.e. the percentage of issued updates that is committed at all replicas.

The measurements were obtained with the fixed settings of 10 replicas with mobility and activation probabilities of 20% and 40%, respectively, running for 2000 time slices on each experiment; we observed that the variation of such values does not have a relevant impact on obtained results. Three update models were tested with different numbers of active replicas: ten, five and just a single one; a global update probability of 5%, evenly divided by the active replicas, was considered.

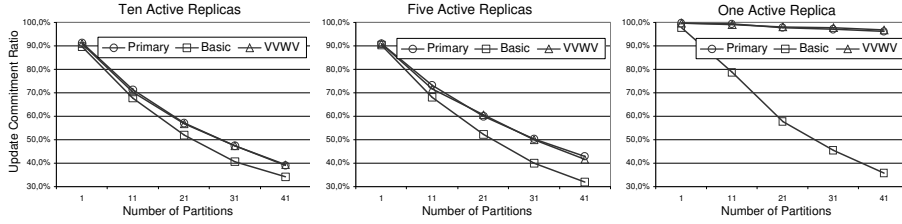


Fig. 3. Update commitment ratios vs. number of partitions, for different numbers of active replicas.

The commitment ratio is directly affected by the efficiency of each evaluated update commitment protocol, since if updates remain in their tentative state for longer periods, the probability of conflicts is higher; hence, lower commitment ratios reflect longer delays imposed by the update commitment process. Fig. 3 shows that, as expected, update commitment ratios decrease as the connectivity among replicas is weakened by an increasing number of partitions.

However, Primary and VVWV are able to ensure higher ratios than Basic WV as partitioning grows. Situations of multiple *happened-before* related tentative updates occur more frequently as updates remain tentative for longer periods. Hence, such results are explained by the efficiency of the former protocols in the commitment of multiple *happened-before* related updates, in contrast to Basic WV. Such situations are also increased as the global update probability is distributed by a smaller number of active replicas. Accordingly, the advantage of Primary and VVWV over Basic WV is accentuated as the number of active replicas decreases. It should be noted that higher update probabilities yielded equivalent, yet magnified, conclusions. On the other hand, Primary and VVWV have similar ratios; however, VVWV has the crucial advantage of not depending on a single point of failure.

Finally, similar experiments compared the two update storage alternatives of VVWV (see Section 3.3). A maximum improvement of 0.8% was attained by storing the updates of all candidates, which suggests that the more resource-efficient alternative of storing only the updates of a replica’s own candidate is acceptable.

5 Conclusions

We propose a novel epidemic weighted voting protocol, VVWV, for achieving the goal of optimistic update commitment that allows multiple *happened-before* related update candidates to be committed at a single election round. Simulation results show that, under weak connectivity conditions, VVWV is advantageous relatively to a basic weighted voting protocol, while attaining similar update commitment ratios to the less fault-tolerant primary commit protocol.

Acknowledgements

We acknowledge Pierre Sutra for the useful discussions that lead to corrections to the initial version of the protocol.

References

1. Saito, Y., Shapiro, M.: Optimistic replication. *ACM Comput. Surv.* **37** (2005) 42–81
2. Lamport, L.: Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM* **21** (1978) 558–565
3. Parker, D.S., et al: Detection of mutual inconsistency in distributed systems. *Distributed systems, Vol. II: distributed data base systems* (1986) 306–312
4. Petersen, K., et al: Flexible update propagation for weakly consistent replication. In: *Proceedings of the 16th ACM Symp. on Operating Systems Principles*. (1997)
5. Keleher, P.: Decentralized replicated-object protocols. In: *Proc. of the 18th Annual ACM Symp. on Principles of Distributed Computing (PODC'99)*. (1999)
6. Golding, R., Long, D.: Modeling replica divergence in a weak-consistency protocol for global scale distributed data bases. *Technical Report UCSC-CRL-93-09* (1993)
7. Cetintemel, U., Keleher, P.: Light-weight currency management mechanisms in mobile and weakly-connected environments. *Dist. Par. Databases* **11** (2002) 53–71
8. Ratner, D., Reiher, P., Popek, G.: Roam: A scalable replication system for mobile computing. In: *DEXA '99: Proceedings of the 10th International Workshop on Database & Expert Systems Applications, Washington, DC, USA, IEEE Computer Society* (1999) 96