

This is a specification of a variant of the classic *Paxos* consensus algorithm described in

AUTHOR = “*Leslie Lamport*”, *TITLE* = “The Part-Time Parliament”, *journal* = *ACM Transactions on Computing Systems*,
volume = 16,
Number = 2, *Month* = may, *Year* = 1998, *pages* = “133–169”

This algorithm was also described without proof in *Brian Oki’s Ph.D. thesis*.

It describes the actions that can be performed by leaders, but does not introduce explicit leader processes. More precisely, the specification is written as if there were a separate leader for each ballot.

This variant of the classic *Paxos* algorithm is an abstraction of an algorithm that is used in

AUTHOR = “*Leslie Lamport and Dahlia Malkhi and Lidong Zhou*”, *TITLE* = “Vertical *Paxos* and Primary-Backup Replication”, *Journal* = “*ACM SIGACT News (Distributed Computing Column)*”,
editor = {*Srikanta Tirthapura and Lorenzo Alvisi*},
booktitle = {*PODC*},
publisher = {*ACM*}, *YEAR* = 2009, *PAGES* = “312–313”

and in

Cheap paxos United States Patent 7249280 Inventors: *Lamport, Leslie B.*
Massa, Michael T.
Filing Date: 06/18/2004

In the classic *Paxos* algorithm, the leader sends a phase 2a message for a ballot *b* and value *v* that instructs acceptors to vote for *v* in ballot *b*. In terms of implementing the voting algorithm of module *VoteProof*, that 2a message serves two functions:

- It asserts that value *v* is safe at ballot *b*, so the acceptor can vote for it without violating invariant *VInv2*
- It tells the acceptors which single safe value they can vote for in ballot *b*, so they can vote for that value without violating *VInv3*.

The variant of the algorithm we specify here introduces phase 1c messages that perform the first function. The phase 2a message serves only the first function, being sent only if a 1c message had been sent for the value.

This variant of the algorithm is useful when reconfiguration is performed by using different sets of acceptors for different ballots. The leader propagates knowledge of what values are safe at ballot *b* so that the acceptors in the current configuration are no longer needed to determine that information. If the ballot *b* leader determines that all values are safe at *b*, then it sends a 1c message for every value and sends a phase 2a message only when it has a value to propose. The presence of the 1c messages removes dependency on the acceptors of ballots numbered *b* or lower for progress. (If the leader determines that only a single value is safe at *b*, then it sends the 1c and 2a messages together.)

In the algorithm described here, we do not include reconfiguration. Therefore, the sending of a 1c message serves only as a precondition for the sending of a 2a message with that value.

Classic *Paxos* and its variants maintain consensus in the presence of omission faults—faults in which a process fails to perform some enabled action or a message that is sent fails to be received. The safety specification, which is given by the *PlusCal* code, does not require that any action need ever be performed. A process need not execute an enabled action. Receipt of a message is modeled by a process performing the action enabled by that message having been sent, so message loss is also represented by a process not performing an enabled action. Thus, failures are never mentioned in the description of the algorithm.

87 EXTENDS *Integers, TLAPS*

88

The constant parameters and the set *Ballots* are the same as in the voting algorithm.

93 CONSTANT *Value, Acceptor, Quorum*

95 ASSUME $QA \triangleq \bigwedge \forall Q \in Quorum : Q \subseteq Acceptor$
 96 $\bigwedge \forall Q1, Q2 \in Quorum : Q1 \cap Q2 \neq \{\}$

98 $Ballot \triangleq Nat$

We are going to have a leader process for each ballot and an acceptor process for each acceptor. So we can use the ballot numbers and the acceptors themselves as the identifiers for these processes, we assume that the set of ballots and the set of acceptors are disjoint. For good measure, we also assume that -1 is not an acceptor, although that is probably not necessary.

108 ASSUME $BallotAssump \triangleq (Ballot \cup \{-1\}) \cap Acceptor = \{\}$

We define *None* to be an unspecified value that is not in the set *Value*.

113 $None \triangleq \text{CHOOSE } v : v \notin Value$

This is a message-passing algorithm, so we begin by defining the set *Message* of all possible messages. The messages are explained below with the actions that send them. A message *m* with *m.type* = “1a” is called a 1a message, and similarly for the other message types.

121 $Message \triangleq$
 122 $\cup [type : \{“1a”\}, bal : Ballot]$
 123 $\cup [type : \{“1b”\}, acc : Acceptor, bal : Ballot,$
 124 $mbal : Ballot \cup \{-1\}, mval : Value \cup \{None\}]$
 125 $\cup [type : \{“1c”\}, bal : Ballot, val : Value]$
 126 $\cup [type : \{“2a”\}, bal : Ballot, val : Value]$
 127 $\cup [type : \{“2b”\}, acc : Acceptor, bal : Ballot, val : Value]$

The algorithm is easiest to understand in terms of the set *msgs* of all messages that have ever been sent. A more accurate model would use one or more variables to represent the messages actually in transit, and it would include actions representing message loss and duplication as well as message receipt.

In the current spec, there is no need to model message loss explicitly. The safety part of the spec says only what messages may be received and does not assert that any message actually is received. Thus, there is no difference between a lost message and one that is never received. The liveness property of the spec will make it clear what messages must be received (and hence either not lost or successfully retransmitted if lost) to guarantee progress.

Another advantage of maintaining the set of all messages that have ever been sent is that it allows us to define the state function *votes* that implements the variable of the same name in the voting algorithm without having to introduce a history variable.

In addition to the variable *msgs*, the algorithm uses four variables whose values are arrays indexed by acceptor, where for any acceptor *a* :

maxBal[*a*] The largest ballot number in which *a* has participated

maxVBal[*a*] The largest ballot number in which *a* has voted, or -1 if it has never voted.

maxVVal[*a*] If *a* has voted, then this is the value it voted for in ballot *maxVBal*; otherwise it equals *None*.

As in the voting algorithm, an execution of the algorithm consists of an execution of zero or more ballots. Different ballots may be in progress concurrently, and ballots may not complete (and need not even start). A ballot *b* consists of the following actions (which need not all occur in the indicated order).

Phase1a : The leader sends a *1a* message for ballot *b*

Phase1b : If *maxBal*[*a*] < *b*, an acceptor *a* responds to the *1a* message by setting *maxBal*[*a*] to *b* and sending a *1b* message to the leader containing the values of *maxVBal*[*a*] and *maxVVal*[*a*].

Phase1c : When the leader has received ballot-*b* *1b* messages from a quorum, it determines some set of values that are safe at *b* and sends *1c* messages for them.

Phase2a : The leader sends a *2a* message for some value for which it has already sent a ballot-*b* *1c* message.

Phase2b : Upon receipt of the *2a* message, if *maxBal*[*a*] ≤ *b*, an acceptor *a* sets *maxBal*[*a*] and *maxVBal*[*a*] to *b*, sets *maxVVal*[*a*] to the value in the *2a* message, and votes for that value in ballot *b* by sending the appropriate *2b* message.

Here is the *PlusCal* code for the algorithm, which we call *PCon*.

```

191 --algorithm PCon{
192   variables maxBal = [a ∈ Acceptor ↦ -1],
193             maxVBal = [a ∈ Acceptor ↦ -1],
194             maxVVal = [a ∈ Acceptor ↦ None],
195             msgs = {}
196   define {
197     sentMsgs(t, b) ≜ {m ∈ msgs : (m.type = t) ∧ (m.bal = b)}

    We define ShowsSafeAt so that ShowsSafeAt(Q, b, v) is true for a quorum Q iff msgs contain
    ballot-b 1b messages from the acceptors in Q show that v is safe at b.

204     ShowsSafeAt(Q, b, v) ≜
205       LET Q1b ≜ {m ∈ sentMsgs("1b", b) : m.acc ∈ Q}
206       IN  ∧ ∀ a ∈ Q : ∃ m ∈ Q1b : m.acc = a
207           ∧ ∨ ∀ m ∈ Q1b : m.mbal = -1
208           ∨ ∃ m1c ∈ msgs :
209             ∧ m1c = [type ↦ "1c", bal ↦ m1c.bal, val ↦ v]
210             ∧ ∀ m ∈ Q1b : ∧ m1c.bal ≥ m.mbal
211             ∧ (m1c.bal = m.mbal) ⇒ (m.mval = v)
213   }
```

As before, we describe each action as a macro.

```
222 macro Phase1a() $\{ msgs := msgs \cup \{ [type \mapsto \text{"1a"}, bal \mapsto self] \}; \}$ 
```

```

230 macro Phase1b(b){
231   when (b > maxBal[self]) ∧ (sentMsgs("1a", b) ≠ {});
232   maxBal[self] := b;
233   msgs := msgs ∪ [{type ↦ "1b", acc ↦ self, bal ↦ b,
234                      mbal ↦ maxVVal[self], mval ↦ maxVVal[self]}];
235 }

```

```

249 macro Phase1c(S){
250   when  $\forall v \in S : \exists Q \in Quorum : ShowsSafeAt(Q, self, v)$ ;
251   msgs := msgs  $\cup \{[type \mapsto "1c", bal \mapsto self, val \mapsto v] : v \in S\}$ 
252 }

```

```

260 macro Phase2a(v) {
261   when  $\wedge sentMsgs("2a", self) = \{\}$ 
262      $\wedge [type \mapsto "1c", bal \mapsto self, val \mapsto v] \in msgs;$ 
263      $msgs := msgs \cup \{[type \mapsto "2a", bal \mapsto self, val \mapsto v]\}$ 
264 }

```

```

272 macro Phase2b(b){
273   when b ≥ maxBal[self];
274   with (m ∈ sentMsgs("2a", b)){
275     maxBal[self] := b;
276     maxVBal[self] := b;
277     maxVVal[self] := m.val;

```

```

278   msgs := msgs ∪ {[type ↦ "2b", acc ↦ self,
279                      bal ↦ b, val ↦ m.val]}
280   }
281 }

```

An acceptor performs the body of its *while* loop as a single atomic action by nondeterministically choosing a ballot in which its *Phase1b* or *Phase2b* action is enabled and executing that enabled action. If no such action is enabled, the acceptor does nothing.

```

289 process (acceptor ∈ Acceptor){
290   acc: while (TRUE){
291     with (b ∈ Ballot){either Phase1b(b)or Phase2b(b)
292   }
293 }
294 }

```

The leader of a ballot nondeterministically chooses one of its actions that is enabled (and the argument for which it is enabled) and performs it atomically. It does nothing if none of its actions is enabled.

```

302 process (leader ∈ Ballot){
303   ldr: while (TRUE){
304     either Phase1a()
305     or    with (S ∈ SUBSET Value){Phase1c(S)}
306     or    with (v ∈ Value){Phase2a(v)}
307   }
308 }
310 }

```

The translator produces the following TLA+ specification of the algorithm. Some blank lines have been deleted.

```

315 BEGIN TRANSLATION
316 VARIABLES maxBal, maxVVal, maxVVal, msgs

318 define statement
319 sentMsgs(t, b) ≜ {m ∈ msgs : (m.type = t) ∧ (m.bal = b)}

321 ShowsSafeAt(Q, b, v) ≜
322   LET Q1b ≜ {m ∈ sentMsgs("1b", b) : m.acc ∈ Q}
323   IN   ∧ ∀ a ∈ Q : ∃ m ∈ Q1b : m.acc = a
324       ∧ ∀ m ∈ Q1b : m.mbal = -1
325       ∧ ∃ m1c ∈ msgs :
326         ∧ m1c = [type ↦ "1c", bal ↦ m1c.bal, val ↦ v]
327         ∧ ∀ m ∈ Q1b : ∧ m1c.bal ≥ m.mbal
328         ∧ (m1c.bal = m.mbal) ⇒ (m.mval = v)

330 vars ≜ ⟨maxBal, maxVVal, maxVVal, msgs⟩

```

```

332  $ProcSet \triangleq (Acceptor) \cup (Ballot)$ 
334  $Init \triangleq$  Global variables
335  $\wedge maxBal = [a \in Acceptor \mapsto -1]$ 
336  $\wedge maxVVal = [a \in Acceptor \mapsto -1]$ 
337  $\wedge maxVVal = [a \in Acceptor \mapsto None]$ 
338  $\wedge msgs = \{\}$ 
340  $acceptor(self) \triangleq \exists b \in Ballot :$ 
341  $\vee \wedge (b > maxBal[self]) \wedge (sentMsgs("1a", b) \neq \{\})$ 
342  $\wedge maxBal' = [maxBal \text{ EXCEPT } ![self] = b]$ 
343  $\wedge msgs' = (msgs \cup \{[type \mapsto "1b", acc \mapsto self, bal \mapsto b,$ 
344  $mbal \mapsto maxVVal[self], mval \mapsto maxVVal[self]]\})$ 
345  $\wedge \text{UNCHANGED } \langle maxVVal, maxVVal \rangle$ 
346  $\vee \wedge b \geq maxBal[self]$ 
347  $\wedge \exists m \in sentMsgs("2a", b) :$ 
348  $\wedge maxBal' = [maxBal \text{ EXCEPT } ![self] = b]$ 
349  $\wedge maxVVal' = [maxVVal \text{ EXCEPT } ![self] = b]$ 
350  $\wedge maxVVal' = [maxVVal \text{ EXCEPT } ![self] = m.val]$ 
351  $\wedge msgs' = (msgs \cup \{[type \mapsto "2b", acc \mapsto self,$ 
352  $bal \mapsto b, val \mapsto m.val]\})$ 
354  $leader(self) \triangleq \wedge \vee \wedge msgs' = (msgs \cup \{[type \mapsto "1a", bal \mapsto self]\})$ 
355  $\vee \wedge \exists S \in \text{SUBSET } Value :$ 
356  $\wedge \forall v \in S : \exists Q \in Quorum : ShowsSafeAt(Q, self, v)$ 
357  $\wedge msgs' = (msgs \cup \{[type \mapsto "1c", bal \mapsto self, val \mapsto v] : v \in S\})$ 
358  $\vee \wedge \exists v \in Value :$ 
359  $\wedge \wedge sentMsgs("2a", self) = \{\}$ 
360  $\wedge [type \mapsto "1c", bal \mapsto self, val \mapsto v] \in msgs$ 
361  $\wedge msgs' = (msgs \cup \{[type \mapsto "2a", bal \mapsto self, val \mapsto v]\})$ 
362  $\wedge \text{UNCHANGED } \langle maxBal, maxVVal, maxVVal \rangle$ 
364  $Next \triangleq (\exists self \in Acceptor : acceptor(self))$ 
365  $\vee (\exists self \in Ballot : leader(self))$ 
367  $Spec \triangleq Init \wedge \Box [Next]_{vars}$ 
369 END TRANSLATION
370 
We now rewrite the next-state relation in a way that makes it easier to use in a proof. We start by
defining the formulas representing the individual actions. We then use them to define the formula
 $TLANext$ , which is the next-state relation we would have written had we specified the algorithm
directly in TLA+ rather than in PlusCal.

378  $Phase1a(self) \triangleq$ 
379  $\wedge msgs' = (msgs \cup \{[type \mapsto "1a", bal \mapsto self]\})$ 
380  $\wedge \text{UNCHANGED } \langle maxBal, maxVVal, maxVVal \rangle$ 
382  $Phase1c(self, S) \triangleq$ 

```

```

383   $\wedge \forall v \in S : \exists Q \in \text{Quorum} : \text{ShowsSafeAt}(Q, \text{self}, v)$ 
384   $\wedge \text{msgs}' = (\text{msgs} \cup \{[type \mapsto "1c", bal \mapsto self, val \mapsto v] : v \in S\})$ 
385   $\wedge \text{UNCHANGED} \langle \text{maxBal}, \text{maxVVal}, \text{maxVVal} \rangle$ 

387   $\text{Phase2a}(\text{self}, v) \triangleq$ 
388   $\wedge \text{sentMsgs}("2a", \text{self}) = \{\}$ 
389   $\wedge [type \mapsto "1c", bal \mapsto self, val \mapsto v] \in \text{msgs}$ 
390   $\wedge \text{msgs}' = (\text{msgs} \cup \{[type \mapsto "2a", bal \mapsto self, val \mapsto v]\})$ 
391   $\wedge \text{UNCHANGED} \langle \text{maxBal}, \text{maxVVal}, \text{maxVVal} \rangle$ 

393   $\text{Phase1b}(\text{self}, b) \triangleq$ 
394   $\wedge b > \text{maxBal}[\text{self}]$ 
395   $\wedge \text{sentMsgs}("1a", b) \neq \{\}$ 
396   $\wedge \text{maxBal}' = [\text{maxBal} \text{ EXCEPT } ![self] = b]$ 
397   $\wedge \text{msgs}' = \text{msgs} \cup \{[type \mapsto "1b", acc \mapsto self, bal \mapsto b,$ 
398   $\text{mbal} \mapsto \text{maxVVal}[\text{self}], \text{mval} \mapsto \text{maxVVal}[\text{self}]]\}$ 
399   $\wedge \text{UNCHANGED} \langle \text{maxVVal}, \text{maxVVal} \rangle$ 

401   $\text{Phase2b}(\text{self}, b) \triangleq$ 
402   $\wedge b \geq \text{maxBal}[\text{self}]$ 
403   $\wedge \exists m \in \text{sentMsgs}("2a", b) :$ 
404   $\quad \wedge \text{maxBal}' = [\text{maxBal} \text{ EXCEPT } ![self] = b]$ 
405   $\quad \wedge \text{maxVVal}' = [\text{maxVVal} \text{ EXCEPT } ![self] = b]$ 
406   $\quad \wedge \text{maxVVal}' = [\text{maxVVal} \text{ EXCEPT } ![self] = m.val]$ 
407   $\quad \wedge \text{msgs}' = (\text{msgs} \cup \{[type \mapsto "2b", acc \mapsto self,$ 
408   $\quad \quad \quad \text{bal} \mapsto b, \text{val} \mapsto m.val]\})$ 

410   $\text{TLANext} \triangleq$ 
411   $\quad \vee \exists \text{self} \in \text{Acceptor} :$ 
412   $\quad \quad \exists b \in \text{Ballot} : \vee \text{Phase1b}(\text{self}, b)$ 
413   $\quad \quad \quad \vee \text{Phase2b}(\text{self}, b)$ 
414   $\quad \vee \exists \text{self} \in \text{Ballot} :$ 
415   $\quad \quad \vee \text{Phase1a}(\text{self})$ 
416   $\quad \quad \vee \exists S \in \text{SUBSET Value} : \text{Phase1c}(\text{self}, S)$ 
417   $\quad \quad \vee \exists v \in \text{Value} : \text{Phase2a}(\text{self}, v)$ 

```

The following theorem specifies the relation between the next-state relation *Next* obtained by translating the *PlusCal* code and the next-state relation *TLANext*.

```

424  THEOREM  $\text{NextDef} \triangleq (\text{Next} \equiv \text{TLANext})$ 
425   $\langle 1 \rangle 2.$  ASSUME NEW  $\text{self} \in \text{Acceptor}$ 
426  PROVE  $\text{acceptor}(\text{self}) \equiv \text{TLANext}!1!(\text{self})$ 
427  BY  $\langle 1 \rangle 2, \text{BallotAssump}$  DEF  $\text{acceptor}, \text{ProcSet}, \text{Phase1b}, \text{Phase2b}$ 
428   $\langle 1 \rangle 3.$  ASSUME NEW  $\text{self} \in \text{Ballot}$ 
429  PROVE  $\text{leader}(\text{self}) \equiv \text{TLANext}!2!(\text{self})$ 
430  BY  $\langle 1 \rangle 3, \text{BallotAssump}$  DEF  $\text{leader}, \text{ProcSet}, \text{Phase1a}, \text{Phase1c}, \text{Phase2a}$ 
431   $\langle 1 \rangle 4.$  QED

```

432 BY $\langle 1 \rangle 2, \langle 1 \rangle 3$ DEF *Next*, *TLANext*

433

The type invariant.

437 $TypeOK \triangleq \wedge maxBal \in [Acceptor \rightarrow Ballot \cup \{-1\}]$
 438 $\wedge maxVVal \in [Acceptor \rightarrow Ballot \cup \{-1\}]$
 439 $\wedge maxVVal \in [Acceptor \rightarrow Value \cup \{None\}]$
 440 $\wedge msgs \subseteq Message$

Here is the definition of the state-function *chosen* that implements the state-function of the same name in the voting algorithm.

446 $chosen \triangleq \{v \in Value : \exists Q \in Quorum, b \in Ballot :$
 447 $\quad \forall a \in Q : \exists m \in msgs : \wedge m.type = "2b"$
 448 $\quad \wedge m.acc = a$
 449 $\quad \wedge m.bal = b$
 450 $\quad \wedge m.val = v\}$

451

We now define the refinement mapping under which this algorithm implements the specification in module *Voting*.

As we observed, votes are registered by sending phase *2b* messages. So the array *votes* describing the votes cast by the acceptors is defined as follows.

462 $votes \triangleq [a \in Acceptor \mapsto$
 463 $\quad \{\langle m.bal, m.val \rangle : m \in \{mm \in msgs : \wedge mm.type = "2b"$
 464 $\quad \wedge mm.acc = a\}\}]$

We now instantiate module *Voting*, substituting:

- The constants *Value*, *Acceptor*, and *Quorum* declared in this module for the corresponding constants of that module *Voting*.
- The variable *maxBal* and the defined state function *votes* for the correspondingly-named variables of module *Voting*.

475 $V \triangleq \text{INSTANCE } VoteProof$

477

We now define *PInv* to be what I believe to be an inductive invariant and assert the theorems for proving that this algorithm implements the voting algorithm under the refinement mapping specified by the *INSTANCE* statement. Whether *PInv* really is an inductive invariant will be determined only by a rigorous proof.

485 $PAccInv \triangleq \forall a \in Acceptor :$
 486 $\quad \wedge maxBal[a] \geq maxVVal[a]$
 487 $\quad \wedge \forall b \in (maxVVal[a] + 1) .. (maxBal[a] - 1) : V!DidNotVoteIn(a, b)$
 488 $\quad \wedge (maxVVal[a] \neq -1) \Rightarrow V!VotedFor(a, maxVVal[a], maxVVal[a])$
 490 $P1bInv \triangleq \forall m \in msgs :$
 491 $\quad (m.type = "1b") \Rightarrow$
 492 $\quad \wedge (maxBal[m.acc] \geq m.bal) \wedge (m.bal > m.mbal)$
 493 $\quad \wedge \forall b \in (m.mbal + 1) .. (m.bal - 1) : V!DidNotVoteIn(m.acc, b)$

495 $P1cInv \triangleq \forall m \in msgs : (m.type = "1c") \Rightarrow V!SafeAt(m.bal, m.val)$

497 $P2aInv \triangleq \forall m \in msgs :$
 498 $(m.type = "2a") \Rightarrow \exists m1c \in msgs : \wedge m1c.type = "1c"$
 499 $\wedge m1c.bal = m.bal$
 500 $\wedge m1c.val = m.val$

The following theorem is interesting in its own right. It essentially asserts the correctness of the definition of *ShowsSafeAt*.

505 THEOREM $PT1 \triangleq TypeOK \wedge P1bInv \wedge P1cInv \Rightarrow$
 506 $\forall Q \in Quorum, b \in Ballot, v \in Value :$
 507 $ShowsSafeAt(Q, b, v) \Rightarrow V!SafeAt(b, v)$

509 $PInv \triangleq TypeOK \wedge PAccInv \wedge P1bInv \wedge P1cInv \wedge P2aInv$

511 THEOREM $Invariance \triangleq Spec \Rightarrow \Box PInv$

513 THEOREM $Implementation \triangleq Spec \Rightarrow V!Spec$

The following result shows that our definition of *chosen* is the correct one, because it implements the state-function *chosen* of the voting algorithm.

520 THEOREM $Spec \Rightarrow \Box(chosen = V!chosen)$

The four theorems above have been checked by *TLC* for a model with 3 acceptors, 2 values, and 3 ballot numbers. Theorem *PT1* was checked as an invariant, therefore checking only that it is true for all reachable states. This model is large enough that it would most likely have revealed any "coding" errors in the algorithm. We believe that the algorithm is well-enough understood that it is unlikely to contain any fundamental errors.

531

\ * Modification History
 \ * Last modified Sat Nov 16 22:19:39 CST 2019 by hengxin
 \ * Last modified Tue Feb 08 12:09:41 PST 2011 by lamport