

1 MODULE *CASpaxos*

This is an interesting extension of the Single-decree *Paxos* algorithm to a compare-and-swap type register. The algorithm is very similar to *Paxos*, but before starting the ACCEPT phase, proposers are free to mutate the value. The result is that the *Paxos* instance turns from a write-once register into a reusable register with atomic semantics, for example a compare-and-swap register.

10 EXTENDS *Integers, FiniteSets*

The data one has to define for the model. In this case, a set of possible *Values*, a set of acceptors, and a mutator which maps a ballot number and a value to a new value.

The *Mutator* is a good approximation of how compare-and-swap'ing proposers would choose the new values, abstracting away nondeterminism by using the ballot number to decide on the new value.

For model checking, infinite sets must be avoided. For convenience, we more or less explicitly assume that values can be compared.

24 CONSTANT *Values, Acceptors, Mutator*(_, _)

The set of quorums. We automatically construct this by taking all the subsets of *Acceptors* which are majorities, *i.e.* larger in number than the unchosen ones. This is for convenience; any definition for which *QuorumAssumption* below holds is valid.

32 $Quorums \triangleq \{S \in \text{SUBSET}(Acceptors) : \text{Cardinality}(S) > \text{Cardinality}(Acceptors \setminus S)\}$

The register in this algorithm can repeatedly change its value. For simplicity, we don't let it start from "no value" but explicitly specify a first value here. Choosing the smallest value is reasonable, but you can change this however you like.

40 $InitialValue \triangleq \text{CHOOSE } v \in Values : \forall vv \in Values : v \leq vv$

Sanity check that our defined quorums all have nontrivial pair-wise intersections. This is pretty clear with the above definition of *Quorums*, but note that you could specify any set of quorums (even minorities) and the algorithm should work the same way, as long as *QuorumAssumption* holds.

49 ASSUME $QuorumAssumption \triangleq \bigwedge \forall Q \in Quorums : Q \subseteq Acceptors$

50 $\bigwedge \forall Q1, Q2 \in Quorums : Q1 \cap Q2 \neq \{\}$

Ballot numbers are natural numbers, but it's good to have an alias so that you know when you're talking about ballots.

56 $Ballot \triangleq Nat$

 Now that we know what a ballot is, check that the *Mutator* maps *Ballots* and *Values* to *Values*.

61 ASSUME $\wedge \forall b \in \text{Ballot}, v \in \text{Values} : \text{Mutator}(b, v) \in \text{Values}$

 Define the set of all possible messages. In this specification, proposers are implicit. Messages originating from them are created “out of thin air” and not addressed to a specific acceptor. In practice they would be, though note that each acceptor would receive the same “message body”, and omitting the explicit recipient reduces the state space. Note also that messages are not explicitly rejected but simply not reacted to. In particular, the implicit proposer has no notion of which ballot to try next. The spec lets them try arbitrary ballots instead.
 A message is either a prepare request for a ballot, a prepare response, an accept request for a ballot with a new value, or an accept response.

78 $\text{Message} \triangleq$ $[type : \{ \text{"prepare-req"} \}, bal : \text{Ballot}]$
 79 \cup $[$
 80 $type : \{ \text{"prepare-rsp"} \}, acc : \text{Acceptors},$
 81 $promised : \text{Ballot},$ ballot for which promise is given
 82 $accepted : \text{Ballot},$ ballot at which *val* was accepted
 83 $val : \text{Values}$
 84 $]$
 85 $\cup [type : \{ \text{"accept-req"} \}, bal : \text{Ballot}, newVal : \text{Values}]$
 86 $\cup [type : \{ \text{"accept-rsp"} \}, acc : \text{Acceptors}, accepted : \text{Ballot}]$
 87 |-----|

 $\langle prepared[a], accepted[a], value[a] \rangle$ is the state of acceptor *a*: The ballot for which a promise has been made (*i.e.* no smaller ballot’s value will be accepted); the ballot at which the last value has been accepted; and the last accepted value.

95 VARIABLE *prepared*,
 96 *accepted*,
 97 *value*

 The set of all messages sent. In each state transition of the model, a message which could solicit a transaction may be reacted to. Note that this implicitly models that a message sent may be received multiple times, and that everything can arbitrarily reorder.

105 VARIABLE *msgs*

 An invariant which checks that the variables have values which make sense.

110 $\text{TypeOK} \triangleq \wedge prepared \in [\text{Acceptors} \rightarrow \text{Ballot}]$
 111 $\wedge accepted \in [\text{Acceptors} \rightarrow \text{Ballot}]$

112 $\wedge value \in [Acceptors \rightarrow Values]$
 113 $\wedge msgs \subseteq Message$

 The initial state of the model. Note that the state here has an initial committed value, ie the register doesn't start "empty". This is an inconsequential simplification.

120 $Init \triangleq \wedge prepared = [a \in Acceptors \mapsto 0]$
 121 $\wedge accepted = [a \in Acceptors \mapsto 0]$
 122 $\wedge value = [a \in Acceptors \mapsto InitialValue]$
 123 $\wedge msgs = \{\}$

 Sending a message just means adding it to the set of all messages.

128 $Send(m) \triangleq msgs' = msgs \cup \{m\}$

 A ballot is started by sending a prepare request (with the hope that responses will be received from a quorum). We could allow multiple prepare requests for a single ballot, but since they are all identical and we already model multiple-receipt for all messages, this adds only state space complexity. So a ballot will only be prepared once in this model.

138 $BallotActive(b) \triangleq \exists m \in msgs :$
 139 $\wedge m.type = \text{"prepare-req"}$
 140 $\wedge m.bal = b$
 141 $PrepareReq(b) \triangleq$
 142 $\wedge \neg BallotActive(b)$
 143 $\wedge Send([$
 144 $type \mapsto \text{"prepare-req"},$
 145 $bal \mapsto b$
 146 $])$
 147 $\wedge UNCHANGED (\langle prepared, accepted, value \rangle)$

 A prepare response can be sent if by an acceptor if a) a response was demanded via a prepare request and b) the acceptor has not already prepared that or any larger ballot. On success, the acceptor remembers that it has prepared the new ballot, and sends to response.

155 $PrepareRsp(a) \triangleq$
 156 $\wedge \exists m \in msgs :$
 157 $\wedge m.type = \text{"prepare-req"}$
 158 $\wedge m.bal > prepared[a]$
 159 $\wedge prepared' = [prepared \text{ EXCEPT } ![a] = m.bal]$
 160 $\wedge Send([$
 161 $acc \mapsto a,$
 162 $type \mapsto \text{"prepare-rsp"},$

```

163         promised  $\mapsto m.bal$ ,
164         accepted  $\mapsto accepted[a]$ ,
165         val  $\mapsto value[a]$ 
166     ])
167      $\wedge$  UNCHANGED  $\langle accepted, value \rangle$ 

```

An accept request can only be sent (*i.e.* be fabricated from thin air) if a) once; b) if prepare responses for the ballot have been received from a quorum; c) with a new value based on the most recently accepted value from the prepare responses.

```

175 AcceptReq(b, v)  $\triangleq$ 
176      $\wedge \neg \exists m \in msgs : m.type = \text{"accept-req"} \wedge m.bal = b$ 
177      $\wedge \exists Q \in Quorums :$ 
178         LET  $M \triangleq \{m \in msgs : \wedge m.type = \text{"prepare-rsp"} \wedge m.promised = b \wedge m.acc \in Q\}$ 
179          $\wedge m.promised = b$ 
180          $\wedge m.acc \in Q\}$ 
181     IN  $\wedge \forall a \in Q : \exists m \in M : m.acc = a$ 
182      $\wedge \exists m \in M :$ 
183          $\wedge m.val = v$ 
184          $\wedge \forall mm \in M : mm.accepted \leq m.accepted$ 
185      $\wedge$  LET  $newVal \triangleq Mutator(b, v)$  crucial difference from Paxos
186     IN Send([
187         type  $\mapsto \text{"accept-req"} ,$ 
188         bal  $\mapsto b,$ 
189         newVal  $\mapsto newVal$ 
190     ])
191      $\wedge$  UNCHANGED  $\langle \langle accepted, value, prepared \rangle \rangle$ 

```

An acceptor can reply to an accept request only if it hasn't yet prepared a higher ballot. Before replying, it makes sure it marks the ballot as prepared (as the particular acceptor may not have received the associated prepare request earlier), and updates its accepted ballot and the new value.

```

200 AcceptRsp(a)  $\triangleq$ 
201      $\wedge \exists m \in msgs :$ 
202          $\wedge m.type = \text{"accept-req"}$ 
203          $\wedge m.bal \geq prepared[a]$ 
204          $\wedge prepared' = [prepared \text{ EXCEPT } ![a] = m.bal]$ 
205          $\wedge accepted' = [accepted \text{ EXCEPT } ![a] = m.bal]$ 
206          $\wedge value' = [value \text{ EXCEPT } ![a] = m.newVal]$ 
207          $\wedge$  Send([
208             acc  $\mapsto a,$ 
209             type  $\mapsto \text{"accept-rsp"} ,$ 
210             accepted  $\mapsto m.bal$ 
211         ])

```

 Next is true if and only if the new state (*i.e.* that with primed variables) is valid. This is used to simulate the model by constructing new states until we run out of options. Concretely, the below means that either we prepare a ballot, or can react successfully to prepare request, or there is an acceptor which can find a message it can react to.

221 $Next \triangleq \vee \exists b \in Ballot : \vee PrepareReq(b)$
 222 $\vee \exists v \in Values : AcceptReq(b, v)$
 223 $\vee \exists a \in Acceptors : PrepareRsp(a) \vee AcceptRsp(a)$

 Spec is the (default) entry point for the TLA+ model runner. The below formula is a temporal formula and means that the valid behaviors of the specification are those which initially satisfy *Init*, and from each step to the following the formula *Next* is satisfied, unless none of the variables changes (which is called a “stuttering step”). The model runner uses this to expand all possible behaviors.

233 $Spec \triangleq Init \wedge \Box [Next]_{(prepared, accepted, value, msgs)}$

 Equipped with a model, what invariants do we want to hold? Or, in other words, what exactly is it that we think the algorithm guarantees? Naively, each newly committed value should be in some relation to a previous value, so when you’re not thinking too hard about it, you could hope that when you take a committed value *A* and the previously committed value *B*, then *B* was created by mutating *A*. It’s not quite as simple, but going down that wrong track highlights how to use the model checker to find interesting histories.

If you have a minute, figure out why the above assumption is wrong. You don’t need more than three ballots and two concurrent proposals.

 For a given ballot, find the acceptors which (at some point in time) accepted that ballot.

252 $AcceptedBy(b) \triangleq \{a \in Acceptors :$
 253 $\exists m \in msgs : \wedge m.type = \text{“accept-rsp”}$
 254 $\wedge m.acc = a$
 255 $\wedge m.accepted = b\}$

 For a given ballot, find out whether the ballot was ever accepted by a quorum.

261 $AcceptedByQuorum(b) \triangleq \exists Q \in Quorums : AcceptedBy(b) \cap Q = Q$

 The set of committed ballot numbers. Note that 0 is trivially committed thanks to the initialization of the state.

267 $CommittedBallots \triangleq \{b \in Ballot : AcceptedByQuorum(b)\} \cup \{0\}$

 For a given ballot $b > 0$, find the next lowest ballot number which was committed (note that this doesn't have to be $b-1$).

273 $BallotCommittedBefore(b) \triangleq$ CHOOSE $c \in CommittedBallots :$
 274 $\quad \wedge \quad c < b$
 275 $\quad \wedge \quad \forall cc \in CommittedBallots :$
 276 $\quad \quad cc \geq b \vee cc \leq c$

 For a given ballot, collect all the values which an acceptor requested. This set will always either be empty or a singleton, but that's not something you can tell from this definition, though we assert it below.

283 $ValuesAt(b) \triangleq$ IF $b = 0$ THEN $\{InitialValue\}$
 284 \quad ELSE $\{v \in Values :$
 285 $\quad \quad \exists m \in msgs :$
 286 $\quad \quad \quad \wedge m.type = \text{"accept-rsp"}$
 287 $\quad \quad \quad \wedge m.accepted = b$
 288 $\quad \quad \quad \wedge \exists mm \in msgs :$
 289 $\quad \quad \quad \quad \wedge mm.type = \text{"accept-req"}$
 290 $\quad \quad \quad \quad \wedge mm.bal = b$
 291 $\quad \quad \quad \quad \wedge mm.newVal = v$
 292 $\quad \quad \quad \}$
 293 $OnlyOneValuePerBallot \triangleq \forall b \in Ballot : Cardinality(ValuesAt(b)) \leq 1$

MutationsLineUp is the main (ill-fated) assertion that each new value of the register is based on a previously committed value.

299 $UnwrapSingleton(s) \triangleq$ CHOOSE $v \in s : \text{TRUE} \quad \{x\} \mapsto x$
 300 $MutationsLineUp \triangleq$
 301 $\quad \forall b \in CommittedBallots \setminus \{0\} :$
 302 $\quad \quad \text{LET } newVal \triangleq UnwrapSingleton(ValuesAt(b))$
 303 $\quad \quad \quad prevCommitBallot \triangleq BallotCommittedBefore(b)$
 304 $\quad \quad \quad oldVal \triangleq UnwrapSingleton(ValuesAt(prevCommitBallot))$
 305 $\quad \quad \text{IN } newVal = Mutator(b, oldVal)$

DesiredProperties is a formula that we will tell the model checker to verify for each state in all valid behaviors. When it is something that actually holds, folks usually call it *Inv* (meaning an inductive invariant), and may try to prove it mechanically using the TLA proof checker, but it is a lot of work and often nontrivial.

In our case, there will be behaviors that violate *MutationsLineUp*.

316 $DesiredProperties \triangleq \wedge TypeOK$

```

317                                      $\wedge$  OnlyOneValuePerBallot
318                                      $\wedge$  MutationsLineUp
319 |_____|
    \ * Modification History
    \ * Last modified Fri Apr 07 02:26:16 EDT 2017 by tshottdorf
    \ * Created Thu Apr 06 02:12:06 EDT 2017 by tshottdorf

```