EXTENDS *FiniteSets*, *Naturals*, *Sequences*

Constants:
   *commands* : *set of records like*[*key* $\mapsto$ "key", *value* $\mapsto$ "value"].
   *replicas* : *set of replicas*.

CONSTANTS *commands*, *replicas*

ASSUME *IsFiniteSet*(*replicas*)

*Variables* :
  *leader* : *records the leader of each epoch*.
  *epoch* : *current epoch*(*the number of leader changes*).
  *proposedCmds* : *the set of proposed commands*.
  *proposeRequests* : *the set of propose could be received by each replica*.
  *proposeResponses* : *the set of responses in each epoch to each proposed*
          *command*.
  *specPools* : *the speculative pool of each replica*.
  *uncommittedCmds* : *the sequence of back* − *end protocol uncommitted commands*.
  *committedCmds* : *the sequence of back* − *end protocol committed commands*.
  *commitMsgs* : *the set of commit messages could be received by each replica*.
  *specExecPrevCmd* : *the index of the last same* − *key command in the committed*
          *sequence at the time the leader responds to the proposal*.

VARIABLES *leader*, *epoch*, *proposedCmds*, *proposeRequests*,
         *proposeResponses*, *specPools*, *uncommittedCmds*,
         *committedCmds*, *commitMsgs*, *specExecPrevCmd*

*The epoch space*.

$epoches \triangleq Nat$

*Special noLeader value*, *for future epoches*.

$noLeader \triangleq$ CHOOSE $r : r \notin replicas$

*In* $N = 2 * f + 1$ *replicas* :
  *quorum* : *a set of replicas that contains at least* $f + 1$ *replicas*.
  *superQuorum* : *a set of replicas that contains at least* $f + (f + 1)/2 + 1$
      *replicas*.
  *recoverQuorum* : *a set of replicas that contains at least*$(f + 1)/2 + 1$
      *replicas*.

$quorums \triangleq$
    LET $f \triangleq Cardinality(replicas) \div 2$
        $size \triangleq f + 1$
    IN   $\{q \in$ SUBSET $replicas : Cardinality(q) \geq size\}$

$superQuorums \triangleq$
    LET $f \triangleq Cardinality(replicas) \div 2$
        $size \triangleq f + (f + 1) \div 2 + 1$

IN $\{q \in$ SUBSET $replicas : Cardinality(q) \geq size\}$

$recoverQuorums \triangleq$
    LET $f \triangleq Cardinality(replicas) \div 2$
        $size \triangleq (f+1) \div 2 + 1$
    IN $\{q \in$ SUBSET $replicas : Cardinality(q) \geq size\}$

---

*Helper function for converting a set to a set containing all sequences*
*containing the elements of the set exactly once and no other elements.*

$SetToSeqs(set) \triangleq$
    LET $len \triangleq 1 .. Cardinality(set)$IN
        $\{f \in [len \to set] : \forall i, j \in len : i \neq j \Rightarrow f[i] \neq f[j]\}$

---

*Helper function for getting the index of the last element in a sequence*
*satisfying the predicate.*

$GetIdxInSeq(seq, Pred(\_)) \triangleq$
    LET $I \triangleq \{i \in 1 .. Len(seq) : Pred(seq[i])\}$IN
        IF $I \neq \{\}$ THEN CHOOSE $i \in I : \forall j \in I : j \leq i$ ELSE $0$

---

*Propose a command.*
*This is done by the client sending a proposeRequest to all replicas.*

$Propose(cmd) \triangleq$
    $\wedge \quad proposedCmds' = proposedCmds \cup \{cmd\}$
    $\wedge \quad proposeRequests' =$
        $[r \in replicas \mapsto proposeRequests[r] \cup \{cmd\}]$
    $\wedge \quad$ UNCHANGED $\langle leader, epoch, specPools, proposeResponses,$
                      $uncommittedCmds, committedCmds, commitMsgs,$
                      $specExecPrevCmd\rangle$

---

*How the leader process a proposeRequest.*

$ProcessProposeLeader(r, cmd) \triangleq$
    LET $specPoolHasConflict \triangleq$
        $\exists specCmd \in specPools[r] : specCmd.key = cmd.key$
        $uncommittedCmdsHasConflict \triangleq$
            $GetIdxInSeq(uncommittedCmds,$ LAMBDA $e : e.key = cmd.key) \neq 0$
    IN
        $\wedge proposeRequests' =$
        $[proposeRequests$ EXCEPT $![r] = @ \setminus \{cmd\}]$
        $\wedge specPools' =$
        $[specPools$ EXCEPT $![r] =$
            IF $\neg specPoolHasConflict$ THEN $@ \cup \{cmd\}$ ELSE $@]$
        $\wedge uncommittedCmds' = Append(uncommittedCmds, cmd)$
        $\wedge proposeResponses' =$
        $[proposeResponses$ EXCEPT $![cmd][epoch] =$
            IF $\neg specPoolHasConflict \wedge \neg uncommittedCmdsHasConflict$

THEN @ $\cup \{r\}$
                    ELSE @]
        $\wedge\ specExecPrevCmd' =$
            $[specExecPrevCmd$ EXCEPT $![cmd] =$
                IF $\neg specPoolHasConflict \wedge \neg uncommittedCmdsHasConflict$
                THEN $GetIdxInSeq(committedCmds,$ LAMBDA $e : e.key = cmd.key)$
                ELSE @]
        $\wedge$ UNCHANGED $\langle leader,\ epoch,\ proposedCmds,\ committedCmds,$
                        $commitMsgs\rangle$

$ProcessProposeNonLeader(r,\ cmd) \triangleq$
    LET $specPoolHasConflict \triangleq$
            $\exists\, specCmd \in specPools[r] : specCmd.key = cmd.key$
    IN
        $\wedge\ proposeRequests' =$
            $[proposeRequests$ EXCEPT $![r] = @ \setminus \{cmd\}]$
        $\wedge\ specPools' =$
            $[specPools$ EXCEPT $![r] =$
                IF $\neg specPoolHasConflict$ THEN @ $\cup \{cmd\}$ ELSE @]
        $\wedge\ proposeResponses' =$
            $[proposeResponses$ EXCEPT $![cmd][epoch] =$
                IF $\neg specPoolHasConflict$ THEN @ $\cup \{r\}$ ELSE @]
        $\wedge$ UNCHANGED $\langle leader,\ epoch,\ proposedCmds,\ uncommittedCmds,$
                        $committedCmds,\ commitMsgs,\ specExecPrevCmd\rangle$

$Commit \triangleq$
    $\wedge$   $committedCmds' = Append(committedCmds,\ Head(uncommittedCmds))$
    $\wedge$   $commitMsgs' =$
        $[r \in replicas \mapsto commitMsgs[r] \cup \{Head(uncommittedCmds)\}]$
    $\wedge$   $uncommittedCmds' = Tail(uncommittedCmds)$
    $\wedge$   UNCHANGED $\langle leader,\ epoch,\ specPools,\ proposedCmds,\ proposeRequests,$
                        $proposeResponses,\ specExecPrevCmd\rangle$

$ProcessCommitMsg(r,\ cmd) \triangleq$
    $\wedge\ commitMsgs' =$
        $[commitMsgs$ EXCEPT $![r] = @ \setminus \{cmd\}]$
    $\wedge\ specPools' = [specPools$ EXCEPT $![r] = @ \setminus \{cmd\}]$
    $\wedge$ UNCHANGED $\langle leader,\ epoch,\ proposedCmds,\ proposeRequests,$
                        $proposeResponses,\ uncommittedCmds,\ committedCmds,$

$$specExecPrevCmd\rangle$$

*Leader Change Action*

*The new leader should gather at least a quorum of replicas specPool to recover the commands.*

Commands existed in the *specPool* of a *RecoverQuorum* of replicas need to be recovered.

$LeaderChange(l) \triangleq$
    $\wedge\ leader' = [e \in epoches \mapsto \text{IF } e = epoch + 1 \text{ THEN } l \text{ ELSE } leader[e]]$
    $\wedge\ epoch' = epoch + 1$
    $\wedge\ \exists\, q \in quorums :$
       $\text{LET } specCmds \triangleq \text{UNION } \{specPools[r] : r \in q\}$
          $newSpecPool \triangleq$
             $\{cmd \in specCmds :$
                  $\{r \in q : cmd \in specPools[r]\} \in recoverQuorums\}$
       $\text{IN}$
          $\wedge\ specPools' = [specPools \text{ EXCEPT } ![l] = newSpecPool]$
          $\wedge\ uncommittedCmds' \in SetToSeqs(newSpecPool)$
    $\wedge\ \text{UNCHANGED } \langle proposedCmds, proposeRequests, proposeResponses,$
                   $committedCmds, commitMsgs, specExecPrevCmd\rangle$

$Init \triangleq$
    $\exists\, r \in replicas :$
      $\text{LET } initEpoch \triangleq 1\ initLeader \triangleq r\ \text{IN}$
          $\wedge\ leader = [e \in epoches \mapsto$
            $\text{IF } e = initEpoch \text{ THEN } initLeader \text{ ELSE } noLeader]$
          $\wedge\ epoch = initEpoch$
          $\wedge\ proposedCmds = \{\}$
          $\wedge\ proposeRequests = [replica \in replicas \mapsto \{\}]$
          $\wedge\ proposeResponses =$
            $[cmd \in commands \mapsto [e \in epoches \mapsto \{\}]]$
          $\wedge\ specPools = [replica \in replicas \mapsto \{\}]$
          $\wedge\ uncommittedCmds = \langle\rangle$
          $\wedge\ committedCmds = \langle\rangle$
          $\wedge\ commitMsgs = [replica \in replicas \mapsto \{\}]$
          $\wedge\ specExecPrevCmd = [cmd \in commands \mapsto 0]$

$Next \triangleq$
    $\vee\ \exists\, cmd \in (commands \setminus proposedCmds) : Propose(cmd)$
    $\vee\ \exists\, r \in replicas : \exists\, cmd \in proposeRequests[r] :$
      $\text{IF } leader[epoch] = r$
       $\text{THEN } ProcessProposeLeader(r, cmd)$
       $\text{ELSE } ProcessProposeNonLeader(r, cmd)$
    $\vee\ uncommittedCmds \neq \langle\rangle \wedge Commit$
    $\vee\ \exists\, r \in replicas : \exists\, cmd \in commitMsgs[r] : ProcessCommitMsg(r, cmd)$
    $\vee\ \exists\, l \in replicas : LeaderChange(l)$

4

$$Spec \;\triangleq\; Init \wedge \Box[Next]\langle leader,\ epoch,\ specPools,\ proposedCmds,$$
$$proposeRequests,\ proposeResponses,$$
$$uncommittedCmds,\ committedCmds,\ commitMsgs,$$
$$specExecPrevCmd\rangle$$

$TypeOK \;\triangleq$

$\quad \wedge \quad leader \in [epoches \to (replicas \cup \{noLeader\})]$

$\quad \wedge \quad epoch \in epoches$

$\quad \wedge \quad proposedCmds \subseteq commands$

$\quad \wedge \quad proposeRequests \in [replicas \to \text{SUBSET}\ commands]$

$\quad \wedge \quad proposeResponses \in [commands \to [epoches \to \text{SUBSET}\ replicas]]$

$\quad \wedge \quad specPools \in [replicas \to \text{SUBSET}\ commands]$

$\quad \wedge \quad uncommittedCmds \in \text{UNION}\ \{SetToSeqs(s) : s \in \text{SUBSET}\ commands\}$

$\quad \wedge \quad committedCmds \in \text{UNION}\ \{SetToSeqs(s) : s \in \text{SUBSET}\ commands\}$

$\quad \wedge \quad commitMsgs \in [replicas \to \text{SUBSET}\ commands]$

$\quad \wedge \quad specExecPrevCmd \in [commands \to 0\ ..\ Cardinality(commands)]$

$Stability \;\triangleq$

$\quad \forall\, cmd \in commands : \forall\, e \in epoches :$

$\quad\quad (\wedge\ leader[e] \in proposeResponses[cmd][e]$

$\quad\quad \wedge\ proposeResponses[cmd][e] \in superQuorums) \Rightarrow$

$\quad\quad\quad \text{LET}\ idx \;\triangleq\; GetIdxInSeq(committedCmds,\ \text{LAMBDA}\ t : t = cmd)$

$\quad\quad\quad\quad prevExecCmds \;\triangleq\; SubSeq(committedCmds,\ 1,\ idx)$

$\quad\quad\quad \text{IN}$

$\quad\quad\quad\quad \wedge\ idx \neq 0$

$\quad\quad\quad\quad \wedge\ GetIdxInSeq(prevExecCmds,\ \text{LAMBDA}\ t : t.key = cmd.key) =$

$\quad\quad\quad\quad\quad specExecPrevCmd[cmd]$

$\text{THEOREM}\ Spec \Rightarrow \Box TypeOK \wedge \Diamond Stability$