1 ─────────────── MODULE *syncCon3* ───────────────

Synchronized consensus

5  EXTENDS *Integers*, *Sequences*, *FiniteSets*, *TLC*

6  CONSTANTS $N$, *FAILNUM*

7  ASSUME $N \leq 7 \wedge 0 \leq FAILNUM \wedge FAILNUM \leq 4$

8  $Nodes \triangleq 1 \mathinner{.\,.} N$

10  **--algorithm** *syncCon3*

11  **{ variable** $FailNum = FAILNUM$,   Initialization block

12       $up = [n \in Nodes \mapsto \text{TRUE}]$;   nodes are up

13       $pt = [n \in Nodes \mapsto 0]$;   nodes are at round 0

14       $t = [n \in Nodes \mapsto \text{FALSE}]$;   nodes are not terminated

15       $d = [n \in Nodes \mapsto -1]$;   nodes are not decided

16       $mb = [n \in Nodes \mapsto \{\}]$;   nodes have mailbox as emptyset

18    **define {**

19    $SetMin(S) \triangleq$ CHOOSE $i \in S : \forall j \in S : i \leq j$

20     **}**

22    **macro** *MaybeFail*( ) **{**

23       **if (** $FailNum > 0 \wedge up[self]$ **)**

24         **{ either**

25            **{** $up[self] := \text{FALSE}$; $FailNum := FailNum - 1$; **}**   Node may fail

26           **or skip; } ;**   or not

27       **}**

29    **fair process (** $n \in Nodes$ **)**

30    **variable** $pmb = \{\}$, $Q = \{\}$;

31     **{**

32  *P*: **while (** $up[self] \wedge \neg t[self]$ **) {**

33       **if (** $d[self] = -1$ **)** $d[self] := self$;   vote is set

34       $Q := Nodes$;   send message to up nodes

35  *PS*:  **while (** $up[self] \wedge Q \neq \{\}$ **) {**   send vote to $mb[p]$ one by one; this node can fail in between

36          **with (** $p \in Q$ **) {**

37            **if (** $pt[p] \geq pt[self] \vee \neg up[p]$ **) {**   send msgs for the same round

38             $mb[p] := mb[p] \cup \{d[self], self\}$;

39             $Q := Q \setminus \{p\}$; **} ;**   also down process with stale $pt$ should not stop progress

40            *MaybeFail*();

41          **} ;**

42        **} ;**   end_while

43       **if (** $up[self]$ **)** $pt[self] := pt[self] + 1$;   move to next round

44  *PR*:  **await** $(up[self] \wedge (\forall k \in Nodes : (up[k] \wedge \neg t[k]) \Rightarrow pt[k] \geq pt[self]))$;   wait for others to move

45       $d[self] := SetMin(mb[self])$;

46       **if (** $pmb = mb[self]$ **)** $t[self] := \text{TRUE}$;

47       $pmb := mb[self]$;

48       $mb[self] := \{\}$;

1

49      **}** ;    *end_if*
50    **}**    process
51  **}**

\ * *PR* label critical for nonblocking;
\ * Remove up in *PR* label, to show the *FLP* result with asynchronous rounds!

55    BEGIN TRANSLATION
56    VARIABLES *FailNum*, *up*, *pt*, *t*, *d*, *mb*, *pc*

58    define statement
59    $SetMin(S) \triangleq$ CHOOSE $i \in S : \forall j \in S : i \leq j$

61    VARIABLES *pmb*, *Q*

63    $vars \triangleq \langle FailNum, up, pt, t, d, mb, pc, pmb, Q \rangle$

65    $ProcSet \triangleq (Nodes)$

67    $Init \triangleq$    Global variables
68          $\wedge FailNum = FAILNUM$
69          $\wedge up = [n \in Nodes \mapsto \text{TRUE}]$
70          $\wedge pt = [n \in Nodes \mapsto 0]$
71          $\wedge t = [n \in Nodes \mapsto \text{FALSE}]$
72          $\wedge d = [n \in Nodes \mapsto -1]$
73          $\wedge mb = [n \in Nodes \mapsto \{\}]$
74          Process *n*
75          $\wedge pmb = [self \in Nodes \mapsto \{\}]$
76          $\wedge Q = [self \in Nodes \mapsto \{\}]$
77          $\wedge pc = [self \in ProcSet \mapsto \text{"P"}]$

79    $P(self) \triangleq \wedge pc[self] = \text{"P"}$
80              $\wedge \text{IF } up[self] \wedge \neg t[self]$
81                  $\text{THEN } \wedge \text{IF } d[self] = -1$
82                          $\text{THEN } \wedge d' = [d \text{ EXCEPT } ![self] = self]$
83                          $\text{ELSE } \wedge \text{TRUE}$
84                                $\wedge d' = d$
85                      $\wedge Q' = [Q \text{ EXCEPT } ![self] = Nodes]$
86                      $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"PS"}]$
87                  $\text{ELSE } \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"Done"}]$
88                      $\wedge \text{UNCHANGED } \langle d, Q \rangle$
89              $\wedge \text{UNCHANGED } \langle FailNum, up, pt, t, mb, pmb \rangle$

91    $PS(self) \triangleq \wedge pc[self] = \text{"PS"}$
92              $\wedge \text{IF } up[self] \wedge Q[self] \neq \{\}$
93                  $\text{THEN } \wedge \exists p \in Q[self]:$
94                          $\wedge \text{IF } pt[p] \geq pt[self] \vee \neg up[p]$
95                              $\text{THEN } \wedge mb' = [mb \text{ EXCEPT } ![p] = mb[p] \cup \{d[self], self\}]$
96                                  $\wedge Q' = [Q \text{ EXCEPT } ![self] = Q[self] \setminus \{p\}]$

2

$$
\begin{array}{ll}
97 & \qquad\qquad\qquad\qquad\qquad \text{ELSE}\quad \wedge \text{TRUE} \\
98 & \qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge \text{UNCHANGED } \langle mb,\; Q \rangle \\
99 & \qquad\qquad\qquad\quad \wedge \text{IF } FailNum > 0 \wedge up[self] \\
100 & \qquad\qquad\qquad\qquad\quad \text{THEN}\quad \wedge \vee \wedge up' = [up \text{ EXCEPT } ![self] = \text{FALSE}] \\
101 & \qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge FailNum' = FailNum - 1 \\
102 & \qquad\qquad\qquad\qquad\qquad\qquad \vee \wedge \text{TRUE} \\
103 & \qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge \text{UNCHANGED } \langle FailNum,\; up \rangle \\
104 & \qquad\qquad\qquad\qquad\quad \text{ELSE}\quad \wedge \text{TRUE} \\
105 & \qquad\qquad\qquad\qquad\qquad\qquad \wedge \text{UNCHANGED } \langle FailNum,\; up \rangle \\
106 & \qquad\qquad\qquad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{``PS''}] \\
107 & \qquad\qquad\qquad \wedge pt' = pt \\
108 & \qquad\qquad \text{ELSE}\quad \wedge \text{IF } up[self] \\
109 & \qquad\qquad\qquad\qquad \text{THEN}\quad \wedge pt' = [pt \text{ EXCEPT } ![self] = pt[self] + 1] \\
110 & \qquad\qquad\qquad\qquad \text{ELSE}\quad \wedge \text{TRUE} \\
111 & \qquad\qquad\qquad\qquad\qquad\quad \wedge pt' = pt \\
112 & \qquad\qquad\qquad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{``PR''}] \\
113 & \qquad\qquad\qquad \wedge \text{UNCHANGED } \langle FailNum,\; up,\; mb,\; Q \rangle \\
114 & \qquad\quad \wedge \text{UNCHANGED } \langle t,\; d,\; pmb \rangle \\
\\
116 & PR(self) \;\triangleq\; \wedge pc[self] = \text{``PR''} \\
117 & \qquad\qquad\quad \wedge (up[self] \wedge (\forall\, k \in Nodes : (up[k] \wedge \neg t[k]) \Rightarrow pt[k] \geq pt[self])) \\
118 & \qquad\qquad\quad \wedge d' = [d \text{ EXCEPT } ![self] = SetMin(mb[self])] \\
119 & \qquad\qquad\quad \wedge \text{IF } pmb[self] = mb[self] \\
120 & \qquad\qquad\qquad\qquad \text{THEN}\quad \wedge t' = [t \text{ EXCEPT } ![self] = \text{TRUE}] \\
121 & \qquad\qquad\qquad\qquad \text{ELSE}\quad \wedge \text{TRUE} \\
122 & \qquad\qquad\qquad\qquad\qquad\quad \wedge t' = t \\
123 & \qquad\qquad\quad \wedge pmb' = [pmb \text{ EXCEPT } ![self] = mb[self]] \\
124 & \qquad\qquad\quad \wedge mb' = [mb \text{ EXCEPT } ![self] = \{\}] \\
125 & \qquad\qquad\quad \wedge pc' \; = [pc \text{ EXCEPT } ![self] \; = \text{``P''}] \\
126 & \qquad\qquad\quad \wedge \text{UNCHANGED } \langle FailNum,\; up,\; pt,\; Q \rangle \\
\\
128 & n(self) \;\triangleq\; P(self) \vee PS(self) \vee PR(self) \\
\\
130 & Next \;\triangleq\; (\exists\, self \in Nodes : n(self)) \\
131 & \qquad\quad \vee \;\; \boxed{\text{Disjunct to prevent deadlock on termination}} \\
132 & \qquad\qquad ((\forall\, self \in ProcSet : pc[self] = \text{``Done''}) \wedge \text{UNCHANGED } vars) \\
\\
134 & Spec \;\triangleq\; \wedge Init \wedge \Box[Next]_{vars} \\
135 & \qquad\qquad \wedge \forall\, self \in Nodes : \text{WF}_{vars}(n(self)) \\
\\
137 & Termination \;\triangleq\; \Diamond(\forall\, self \in ProcSet : pc[self] = \text{``Done''}) \\
\\
139 & \boxed{\text{END TRANSLATION}} \\
\\
141 & Agreement \;\triangleq\; \forall\, i,\, j \in Nodes : t[i] \wedge t[j] \Rightarrow (d[i] = d[j] \wedge d[i] \neq -1) \\
142 & NoTerm \;\triangleq\; \neg \forall\, i \quad \in Nodes : up[i] \Rightarrow t[i] \\
143 & SyncTerm \;\triangleq\; \forall\, i,\, j \in Nodes : t[i] \wedge t[j] \Rightarrow pt[i] = pt[j]
\end{array}
$$

144    $Term \stackrel{\Delta}{=} \Diamond \forall\, i \in Nodes : up[i] \Rightarrow t[i]$

145      $Remember \stackrel{\Delta}{=} \Box\, [\ (\forall\, j \in Nodes:\ v'[p] \ge v[p])\ ]\_vars$

146

Agreement. Two correct processes can not commit to different decision variables. ($i,j$:$ti$ tj :$di = dj$) Validity (Nontriviality). If all initial values are equal, correct processes must decide on that value. ($k$:: ($i$:: $vi = k$)) ($i$ :$ti$ :$di = vi$) Termination. The system eventually terminates. true ($i$:: $ti$)

Synchronous consensus Every process broadcasts (to all other processes, including itself) its initial value $vi$. In a synchronous network, this can be done in a single "round" of messages. After this round, each process decides on the minimum value it received. If no faults occur, this algorithm is correct. In the presence of a crash fault, however, a problem can arise. In particular, a problem may occur if a process crashes during a round. When this happens, some processes may have received its (low) initial value, but others may not have.

To address this concern, consider this simplifying assumption: say that at most 1 process can crash. How can we modify the algorithm to handle such a failure? Answer: by using 2 rounds. In $1st$ round, processes broadcast their own initial value. In $2nd$ round, processes broadcast the minimum value they heard. Each process then decides on the min value among all the sets of values it received in $2nd$ round. If the one crash occurs during the first round, the second round ensures that all processes have the same set of values from which to decide. Else, if the one crash occurs during the second round, the first round must have completed without a crash and hence all processes have the same set of values from which to decide.

The key observation is that if no crash occurs during a round, all processes have the same set of values from which to decide and they correctly decide on the same minimum value. Thus, to tolerate multiple crashes, say $f$, the protocol is modified to have $f + 1$ rounds of synchronous communication. Of course, this requires knowing $f$, an upper bound on the number of possible crash faults.