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

Synchronized consensus

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

6 CONSTANTS *N*, *FAILNUM*

7 ASSUME $N \leq 5 \wedge 0 \leq FAILNUM \wedge FAILNUM \leq 4$

8 *Nodes* $\triangleq 1 .. N$

11 **--algorithm** *syncCon3*

12 **{ variable** *FailNum* = *FAILNUM*,   Initialization block

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

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

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

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

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

19    **define {**

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

21     **}**

23    **macro** *MaybeFail*( ) **{**

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

25        **{ either**

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

27          **or skip ; } ;**   or not

28      **}**

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

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

32     **{**

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

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

35      $Q := Nodes$;

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

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

38        $mb[p] := mb[p] \cup \{d[self]\}$;      skip for attacking generals impossibility

39        $Q := Q \setminus \{p\}$;

40        *MaybeFail*() ;

41        **} ;**

42       **} ;**   end_while

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

44 *PR*:   **await** $(up[self] \wedge (\forall k \in Nodes : up[k] \Rightarrow pt[k] = pt[self]))$;   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      $mb[self] := \{\}$;

49    **} ;**   end_if

1

```
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          $\wedge mb = [n \in Nodes \mapsto \{\}]$
            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              $\wedge$ IF $up[self] \wedge \neg t[self]$
81                  THEN $\wedge$ IF $d[self] = -1$
82                              THEN $\wedge d' = [d$ EXCEPT $![self] = self]$
83                              ELSE $\wedge$ TRUE
84                                   $\wedge d' = d$
85                          $\wedge Q' = [Q$ EXCEPT $![self] = Nodes]$
86                          $\wedge pc' = [pc$ EXCEPT $![self] = \text{“PS”}]$
87                  ELSE $\wedge pc' = [pc$ EXCEPT $![self] = \text{“Done”}]$
88                       $\wedge$ UNCHANGED $\langle d, Q \rangle$
89              $\wedge$ UNCHANGED $\langle FailNum, up, pt, t, mb, pmb \rangle$

91   $PS(self) \triangleq \wedge pc[self] = \text{“PS”}$
92              $\wedge$ IF $up[self] \wedge Q[self] \neq \{\}$
93                  THEN $\wedge \exists p \in Q[self] :$
94                              $\wedge mb' = [mb$ EXCEPT $![p] = mb[p] \cup \{d[self]\}]$
95                              $\wedge Q' = [Q$ EXCEPT $![self] = Q[self] \setminus \{p\}]$
96                              $\wedge$ IF $FailNum > 0 \wedge up[self]$
97                                  THEN $\wedge \vee \wedge up' = [up$ EXCEPT $![self] = \text{FALSE}]$

2
```
```

$$98 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \wedge FailNum' = FailNum - 1$$
$$99 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \vee \wedge \text{TRUE}$$
$$100 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \wedge \text{UNCHANGED} \langle FailNum, up \rangle$$
$$101 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{ELSE} \quad \wedge \text{TRUE}$$
$$102 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \wedge \text{UNCHANGED} \langle FailNum, up \rangle$$
$$103 \qquad \qquad \qquad \qquad \qquad \qquad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{``PS''}]$$
$$104 \qquad \qquad \qquad \qquad \qquad \qquad \wedge pt' = pt$$
$$105 \qquad \qquad \qquad \qquad \text{ELSE} \quad \wedge \text{ IF } up[self]$$
$$106 \qquad \qquad \qquad \qquad \qquad \qquad \text{THEN} \quad \wedge pt' = [pt \text{ EXCEPT } ![self] = pt[self] + 1]$$
$$107 \qquad \qquad \qquad \qquad \qquad \qquad \text{ELSE} \quad \wedge \text{TRUE}$$
$$108 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \quad \wedge pt' = pt$$
$$109 \qquad \qquad \qquad \qquad \qquad \qquad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{``PR''}]$$
$$110 \qquad \qquad \qquad \qquad \qquad \qquad \wedge \text{UNCHANGED} \langle FailNum, up, mb, Q \rangle$$
$$111 \qquad \qquad \qquad \qquad \wedge \text{UNCHANGED} \langle t, d, pmb \rangle$$

$$113 \quad PR(self) \triangleq \wedge pc[self] = \text{``PR''}$$
$$114 \qquad \qquad \qquad \qquad \wedge (up[self] \wedge (\forall k \in Nodes : up[k] \Rightarrow pt[k] = pt[self]))$$
$$115 \qquad \qquad \qquad \qquad \wedge d' = [d \text{ EXCEPT } ![self] = SetMin(mb[self])]$$
$$116 \qquad \qquad \qquad \qquad \wedge \text{ IF } pmb[self] = mb[self]$$
$$117 \qquad \qquad \qquad \qquad \qquad \text{THEN} \quad \wedge t' = [t \text{ EXCEPT } ![self] = \text{TRUE}]$$
$$118 \qquad \qquad \qquad \qquad \qquad \text{ELSE} \quad \wedge \text{TRUE}$$
$$119 \qquad \qquad \qquad \qquad \qquad \qquad \qquad \wedge t' = t$$
$$120 \qquad \qquad \qquad \qquad \wedge pmb' = [pmb \text{ EXCEPT } ![self] = mb[self]]$$
$$121 \qquad \qquad \qquad \qquad \wedge mb' = [mb \text{ EXCEPT } ![self] = \{\}]$$
$$122 \qquad \qquad \qquad \qquad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{``P''}]$$
$$123 \qquad \qquad \qquad \qquad \wedge \text{UNCHANGED} \langle FailNum, up, pt, Q \rangle$$

$$125 \quad n(self) \triangleq P(self) \vee PS(self) \vee PR(self)$$

$$127 \quad Next \triangleq (\exists self \in Nodes : n(self))$$
$$128 \qquad \qquad \qquad \vee \boxed{\text{Disjunct to prevent deadlock on termination}}$$
$$129 \qquad \qquad \qquad \quad ((\forall self \in ProcSet : pc[self] = \text{``Done''}) \wedge \text{UNCHANGED } vars)$$

$$131 \quad Spec \triangleq \wedge Init \wedge \Box[Next]_{vars}$$
$$132 \qquad \qquad \qquad \wedge \forall self \in Nodes : \text{WF}_{vars}(n(self))$$

$$134 \quad Termination \triangleq \Diamond(\forall self \in ProcSet : pc[self] = \text{``Done''})$$

$$136 \quad \boxed{\text{END TRANSLATION}}$$

$$138 \quad Agreement \triangleq \forall i, j \in Nodes : t[i] \wedge t[j] \Rightarrow (d[i] = d[j] \wedge d[i] \neq -1)$$
$$139 \quad NoTerm \triangleq \neg \forall i \quad \in Nodes : up[i] \Rightarrow t[i]$$
$$140 \quad SyncTerm \triangleq \forall i, j \in Nodes : t[i] \wedge t[j] \Rightarrow pt[i] = pt[j]$$
$$141 \quad Term \triangleq \Diamond \forall i \in Nodes : up[i] \Rightarrow t[i]$$
$$142 \quad \boxed{Remember \triangleq \Box [ (\forall j \in Nodes: v'[p] \geq v[p]) ]_{vars}}$$
143

3

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.