
1 | MODULE *TCS* |

The specification of the Transaction Certification Service (*TCS*) in *DISC'2018* “Multi-Shot Distributed Transaction Commit” by *Gregory Chockler* and *Alexey Gotsman*.

We have specified the multi-shot *2PC* protocol in Figure 1 of *DISC'2018*.

TODO: We plan

- to test *SER* using the Serializability Theorem
- to integrate *TCS* into a real distributed transaction protocol
- to implement certification functions for other isolation levels
- to specify the fault-tolerant commit protocol in Figure 5 of *DISC'2018*.

15 EXTENDS *Naturals, Integers, FiniteSets, Sequences, Functions, TLC,*
16 *FiniteSetsExt*

18 CONSTANTS

19 *Key*, the set of keys, ranged over by $k \in Key$

20 *Tid*, the set of transaction identifiers, ranged over by $t \in Tid$

21 *RSet*, *RSet*[t]: the read set of $t \in Tid$

22 *WSet*, *WSet*[t]: the write set of $t \in Tid$

23 *CVer*, *CVer*[t]: the commit version of $t \in Tid$

24 *Shard*, the set of shards, ranged over by $s \in Shard$

25 *Coord*, *Coord*[t]: the coordinator of $t \in Tid$

26 *KeySharding* *KeySharding*[k]: the shard that holds $k \in Key$

28 $NotTid \triangleq \text{CHOOSE } t : t \notin Tid$

30 $Ver \triangleq 0 \dots Cardinality(Tid)$ with a distinguished minimum version 0

31 $Slot \triangleq 0 \dots Cardinality(Tid) - 1$

33 $TShard(t) \triangleq \{KeySharding[k] : k \in (WSet[t] \cup \{kv[1] : kv \in RSet[t]\})\}$

35 ASSUME *TODO*: See Section 2 of *DISC'2018*

36 $\wedge RSet \in [Tid \rightarrow SUBSET (Key \times Ver)]$

37 $\wedge \forall t \in Tid: RSet[t] \setminus * \text{ } TODO: \text{ one version per object}$

38 $\wedge WSet \in [Tid \rightarrow SUBSET Key]$

39 $\wedge \setminus * \text{ } TODO: \text{ “no blind update” assumption}$

40 $\wedge CVer \in [Tid \rightarrow Ver]$

41 $\wedge \setminus * \text{ } TODO: \text{ higher than any of the versions read}$

42 $\wedge Coord \in [Tid \rightarrow Shard]$

43 $\wedge KeySharding \in [Key \rightarrow Shard]$

45 VARIABLES

46 *next*, *next*[s] $\in Z$ points to the last filled slot

47 *txn*, *txn*[s][i] is the transaction (identifier) to certify in the i -th slot

48 *vote*, *vote*[s][i] is the vote for *txn*[s][i]

49 *dec*, *dec*[s][i] is the decision for *txn*[s][i]

50 *phase*, *phase*[s][i] is the phase for *txn*[s][i]

51 *msg*, the set of messages in transit

```

52      submitted      the set of  $t \in Tid$  that have been submitted to  $TCS$ 

54       $sVars \triangleq \langle next, txn, vote, dec, phase \rangle$ 
55       $vars \triangleq \langle next, txn, vote, dec, phase, msg, submitted \rangle$ 
56  |-----|
57      TODO:
58      "COMMIT/ABORT": using TRUE/FALSE (initially, FALSE???)
59      "PREPARE/PREPARE_ACK/DECISION": using CONSTANTS
60  |-----|
61       $Message \triangleq [type : \{ "PREPARE" \}, t : Tid, s : Shard]$ 
62       $\cup [type : \{ "PREPARE\_ACK" \}, s : Shard, n : Int, t : Tid, v : \{ "COMMIT", "ABORT" \}]$ 
63       $\cup [type : \{ "DECISION" \}, p : Int, d : \{ "COMMIT", "ABORT" \}, s : Shard]$ 
64  |-----|
65       $Send(m) \triangleq msg' = msg \cup m$ 
66       $Delete(m) \triangleq msg' = msg \setminus m$ 
67       $SendAndDelete(sm, dm) \triangleq msg' = (msg \cup sm) \setminus dm$ 
68  |-----|
69       $TypeOK \triangleq$ 
70       $\wedge next \in [Shard \rightarrow Int]$ 
71       $\wedge txn \in [Shard \rightarrow [Slot \rightarrow Tid \cup \{ NotTid \}]]$ 
72       $\wedge vote \in [Shard \rightarrow [Slot \rightarrow \{ "COMMIT", "ABORT", "NULL" \}]]$ 
73       $\wedge dec \in [Shard \rightarrow [Slot \rightarrow \{ "COMMIT", "ABORT", "NULL" \}]]$ 
74       $\wedge phase \in [Shard \rightarrow [Slot \rightarrow \{ "START", "PREPARED", "DECIDED" \}]]$ 
75       $\wedge msg \subseteq Message$ 
76       $\wedge submitted \subseteq Tid$ 
77  |-----|
78       $Init \triangleq$ 
79       $\wedge next = [s \in Shard \mapsto -1]$ 
80       $\wedge txn = [s \in Shard \mapsto [i \in Slot \mapsto NotTid]]$ 
81       $\wedge vote = [s \in Shard \mapsto [i \in Slot \mapsto "NULL"]]$ 
82       $\wedge dec = [s \in Shard \mapsto [i \in Slot \mapsto "NULL"]]$ 
83       $\wedge phase = [s \in Shard \mapsto [i \in Slot \mapsto "START"]]$ 
84       $\wedge msg = \{ \}$ 
85       $\wedge submitted = \{ \}$ 
86  |-----|
87       $KeyOnShard(s) \triangleq \{ k \in Key : KeySharding[k] = s \}$ 
88  |-----|
89       $ComputeVote(t, s, n) \triangleq$ 
90      LET  $cs \triangleq \{ k \in Slot :$  committed slots before position  $n$ 
91       $\wedge k < n$ 
92       $\wedge phase[s][k] = "DECIDED"$ 
93       $\wedge dec[s][k] = "COMMIT" \}$ 
94       $ct \triangleq \{ txn[s][k] : k \in cs \}$  committed transactions
95       $fv \triangleq$  IF  $\forall k \in KeyOnShard(s), v \in Ver :$ 
96       $\langle k, v \rangle \in RSet[t] \Rightarrow (\forall c \in ct : k \in WSet[c] \Rightarrow CVer[c] \leq v)$ 
97      THEN "COMMIT" ELSE "ABORT"
98       $ps \triangleq \{ k \in Slot :$  "prepared to commit" slots before position  $n$ 

```

```

99       $\wedge k < n$ 
100       $\wedge phase[s][k] = \text{"PREPARED"}$ 
101       $\wedge vote[s][k] = \text{"COMMIT"}$ 
102       $pt \triangleq \{txn[s][k] : k \in ps\}$  "prepared to commit" transactions
103       $gv \triangleq \text{IF } \forall k \in KeyOnShard(s), v \in Ver :$ 
104           $\wedge \langle k, v \rangle \in RSet[t] \Rightarrow (\forall p \in pt : k \notin WSet[p])$ 
105           $\wedge k \in WSet[t] \Rightarrow (\forall p \in pt : \langle k, v \rangle \notin RSet[p])$ 
106          THEN "COMMIT" ELSE "ABORT"
107      IN   IF  $fv = \text{"COMMIT"} \wedge gv = \text{"COMMIT"}$  THEN "COMMIT" ELSE "ABORT"

109   $ComputeDecision(vs) \triangleq$ 
110      IF  $\forall v \in vs : v = \text{"COMMIT"}$  THEN "COMMIT" ELSE "ABORT"
111  |-----|
112   $Certify(t) \triangleq$  Certify  $t \in Tid$ 
113       $\wedge t \in Tid \setminus submitted$ 
114       $\wedge Send([type : \{\text{"PREPARE"}\}, t : \{t\}, s : TShard(t))]$ 
115       $\wedge submitted' = submitted \cup \{t\}$ 
116       $\wedge \text{UNCHANGED } sVars$ 

118   $Prepare(t, s) \triangleq$  Prepare  $t \in Tid$  on  $s \in Shard$  when receive " $PREPARE(t)$ " message
119       $\wedge \exists m \in msg :$ 
120           $\wedge m = [type \mapsto \text{"PREPARE"}, t \mapsto t, s \mapsto s]$ 
121           $\wedge next' = [next \text{ EXCEPT } !s] @ + 1]$ 
122           $\wedge txn' = [txn \text{ EXCEPT } !s][next'[s]] = t]$ 
123           $\wedge vote' = [vote \text{ EXCEPT } !s][next'[s]] = ComputeVote(t, s, next'[s])]$ 
124           $\wedge phase' = [phase \text{ EXCEPT } !s][next'[s]] = \text{"PREPARED"}$ 
125           $\wedge SendAndDelete(\{[type \mapsto \text{"PREPARE\_ACK"},$ 
126               $s \mapsto s,$ 
127               $n \mapsto next'[s],$ 
128               $t \mapsto t,$ 
129               $v \mapsto vote'[s][next'[s]]\},$ 
130               $\{m\})$ 
131           $\wedge \text{UNCHANGED } \langle dec, submitted \rangle$ 

133   $PrepareAck(t, s) \triangleq$  PrepareAck for  $t \in Tid$  on shard  $s \in Shard$  when receive all " $PREPARE\_ACK$ " messages for  $t$ 
134       $\wedge s = Coord[t]$ 
135       $\wedge \text{LET } ms \triangleq \{m \in msg : m.type = \text{"PREPARE\_ACK"} \wedge m.t = t\}$ 
136           $vs \triangleq \{m.v : m \in ms\}$ 
137           $ss \triangleq \{m.s : m \in ms\}$ 
138      IN    $\wedge ss = TShard(t)$ 
139           $\wedge SendAndDelete(\{[type \mapsto \text{"DECISION"},$ 
140               $p \mapsto ChooseUnique(ms, \text{LAMBDA } m : m.s = shard).n,$ 
141               $d \mapsto ComputeDecision(vs),$ 
142               $s \mapsto shard] : shard \in ss\},$ 
143               $ms)$ 
144       $\wedge \text{UNCHANGED } \langle sVars, submitted \rangle$ 

```

```

146  $Decision(s) \triangleq$  Decide on shard  $s \in Shard$  when receive a "DECISION" message
147    $\wedge \exists m \in msg :$ 
148      $\wedge m.type = \text{"DECISION"}$ 
149      $\wedge m.s = s$ 
150      $\wedge dec' = [dec \text{ EXCEPT } ![s][m.p] = m.d]$ 
151      $\wedge phase' = [phase \text{ EXCEPT } ![s][m.p] = \text{"DECIDED"}]$ 
152      $\wedge Delete(\{m\})$ 
153    $\wedge \text{UNCHANGED } \langle next, txn, vote, submitted \rangle$ 
154 |-----|
155 | TODO: adding the two non-deterministic actions |
156 |-----|
158 |
159  $Next \triangleq$ 
160    $\vee \exists t \in Tid : Certify(t)$ 
161    $\vee \exists t \in Tid, s \in Shard :$ 
162      $\vee Prepare(t, s)$ 
163      $\vee PrepareAck(t, s)$ 
164    $\vee \exists s \in Shard :$ 
165      $\vee Decision(s)$ 
167  $Spec \triangleq Init \wedge \Box [Next]_{vars}$ 
168 |-----|
  
```

* Modification History
 * Last modified Sun Jun 13 19:27:55 CST 2021 by *hengxin*
 * Created Sat Jun 12 21:01:57 CST 2021 by *hengxin*