

```

1  ┌────────────────────────── MODULE Cure ───────────────────────────┐
    │ See ICDCS2016: “Cure: Strong Semantics Meets High Availability and Low Latency”. │
5  └────────────────────────── EXTENDS Naturals, Sequences, TLC ───────────────────────────┘
6  ┌──────────────────────────┐
7  │  $Max(a, b) \triangleq \text{IF } a < b \text{ THEN } b \text{ ELSE } a$  │
8  │  $Min(S) \triangleq \text{CHOOSE } a \in S : \forall b \in S : a \leq b$  │
9  │  $Range(f) \triangleq \{f[x] : x \in \text{DOMAIN } f\}$  │
10 │  $Last(seq) \triangleq seq[Len(seq)]$  │
11 └──────────────────────────┘
12 CONSTANTS
13   Key,           the set of keys, ranged over by  $k \in Key$ 
14   Value,         the set of values, ranged over by  $v \in Value$ 
15   Client,        the set of clients, ranged over by  $c \in Client$ 
16   Partition,     the set of partitions, ranged over by  $p \in Partition$ 
17   Datacenter,    the set of datacenters, ranged over by  $d \in Datacenter$ 
18   KeySharding,   the mapping from Key to Partition
19   ClientAttachment the mapping from Client to Datacenter
21  $NotVal \triangleq \text{CHOOSE } v : v \notin Value$ 
23 ASSUME
24    $\wedge KeySharding \in [Key \rightarrow Partition]$ 
25    $\wedge ClientAttachment \in [Client \rightarrow Datacenter]$ 
26 ┌──────────────────────────┐
27 │ VARIABLES │
28 │ At the client side: │
29 │   cvc, cvc[c]: the vector clock of client  $c \in Client$  │
30 │ At the server side (each for partition  $p \in Partition$  in  $d \in Datacenter$ ): │
31 │   clock, clock[p][d]: the current clock │
32 │   pvc, pvc[p][d]: the vector clock │
33 │   css, css[p][d]: the stable snapshot │
34 │   store, store[p][d]: the kv store │
35 │   updates, updates[p][d]: the buffer of updates │
36 │ Clock management │
37 │   tick, tick[p][d]: toggle on clock ticks │
38 │ Client-server communication │
39 │   msgs the set of messages in transit │
41  $cVars \triangleq \langle cvc \rangle$ 
42  $sVars \triangleq \langle clock, pvc, css, store, updates, tick \rangle$ 
43  $mVars \triangleq \langle msgs \rangle$ 
44  $vars \triangleq \langle cvc, clock, pvc, css, store, updates, tick, msgs \rangle$ 
45 ┌──────────────────────────┐
46 │  $Clock \triangleq Nat$  │
47 │  $VC \triangleq [Datacenter \rightarrow Clock]$  vector clock with an entry per datacenter  $d \in Datacenter$  │
48 │  $VCInit \triangleq [d \in Datacenter \mapsto 0]$  │

```

49 $KVTuple \triangleq [key : Key, val : Value \cup \{NotVal\}, vc : VC]$
 51 $Message \triangleq$
 52 $[type : \{ "ReadRequest" \}, key : Key, vc : VC, c : Client, p : Partition, d : Datacenter]$
 53 $\cup [type : \{ "ReadReply" \}, val : Value \cup \{NotVal\}, vc : VC, c : Client]$
 54 $\cup [type : \{ "UpdateRequest" \}, key : Key, val : Value, vc : VC, c : Client, p : Partition, d : Datacenter]$
 55 $\cup [type : \{ "UpdateReply" \}, ts : Clock, c : Client, d : Datacenter]$
 56 $d \in Datacenter$: the source datacenter; $dd \in Datacenter$: the destination datacenter
 57 $\cup [type : \{ "Replicate" \}, d : Datacenter, kvs : Seq(KVTuple), p : Partition, dd : Datacenter]$
 58 $\cup [type : \{ "Heartbeat" \}, d : Datacenter, ts : Clock, p : Partition, dd : Datacenter]$
 60 $TypeOK \triangleq$
 61 $\wedge cvc \in [Client \rightarrow VC]$
 62 $\wedge clock \in [Partition \rightarrow [Datacenter \rightarrow Clock]]$
 63 $\wedge pvc \in [Partition \rightarrow [Datacenter \rightarrow VC]]$
 64 $\wedge css \in [Partition \rightarrow [Datacenter \rightarrow VC]]$
 65 $\wedge store \in [Partition \rightarrow [Datacenter \rightarrow SUBSET KVTuple]]$
 66 $\wedge updates \in [Partition \rightarrow [Datacenter \rightarrow Seq(KVTuple)]]$
 67 $\wedge tick \in [Partition \rightarrow [Datacenter \rightarrow BOOLEAN]]$
 68 $\wedge msgs \subseteq Message$
 69 \mid
 70 $Init \triangleq$
 71 $\wedge cvc = [c \in Client \mapsto VCInit]$
 72 $\wedge clock = [p \in Partition \mapsto [d \in Datacenter \mapsto 0]]$
 73 $\wedge pvc = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$
 74 $\wedge css = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$
 75 $\wedge store = [p \in Partition \mapsto [d \in Datacenter \mapsto$
 76 $[key : \{k \in Key : KeySharding[k] = p\}, val : \{NotVal\}, vc : \{VCInit\}]]]$
 77 $\wedge updates = [p \in Partition \mapsto [d \in Datacenter \mapsto \langle \rangle]]$
 78 $\wedge tick = [p \in Partition \mapsto [d \in Datacenter \mapsto FALSE]]$
 79 $\wedge msgs = \{\}$
 80 \mid
 81 $Send(m) \triangleq msgs' = msgs \cup \{m\}$
 82 $SendSet(ms) \triangleq msgs' = msgs \cup ms$
 83 $SendAndDelete(sm, dm) \triangleq msgs' = (msgs \cup \{sm\}) \setminus \{dm\}$
 85 $CanIssue(c) \triangleq \forall m \in msgs :$
 86 $m.type \in \{ "ReadRequest", "ReadReply", "UpdateRequest", "UpdateReply" \} \Rightarrow m.c \neq c$
 87 \mid
 88 Client operations at client $c \in Client$.
 90 $Read(c, k) \triangleq$ $c \in Client$ reads from $k \in Key$
 91 $\wedge CanIssue(c)$
 92 $\wedge Send([type \mapsto "ReadRequest", key \mapsto k, vc \mapsto cvc[c],$
 93 $c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$
 94 $\wedge UNCHANGED \langle cVars, sVars \rangle$

```

96  $ReadReply(c) \triangleq$   $c \in Client$  handles the reply to its read request
97  $\wedge \exists m \in msgs :$ 
98  $\wedge m.type = \text{"ReadReply"} \wedge m.c = c$  such  $m$  is unique
99  $\wedge cvc' = [cvc \text{ EXCEPT } ![c] = [d \in Datacenter \mapsto Max(m.vc[d], @[d])]]$ 
100  $\wedge msgs' = msgs \setminus \{m\}$ 
101  $\wedge \text{UNCHANGED } \langle sVars \rangle$ 

103  $Update(c, k, v) \triangleq$   $c \in Client$  updates  $k \in Key$  with  $v \in Value$ 
104  $\wedge CanIssue(c)$ 
105  $\wedge Send([type \mapsto \text{"UpdateRequest"}, key \mapsto k, val \mapsto v,$ 
106  $vc \mapsto cvc[c], c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$ 
107  $\wedge \text{UNCHANGED } \langle cVars, sVars \rangle$ 

109  $UpdateReply(c) \triangleq$   $c \in Client$  handles the reply to its update request
110  $\wedge \exists m \in msgs :$ 
111  $\wedge m.type = \text{"UpdateReply"} \wedge m.c = c$  such  $m$  is unique
112  $\wedge cvc' = [cvc \text{ EXCEPT } ![c][m.d] = m.ts]$ 
113  $\wedge msgs' = msgs \setminus \{m\}$ 
114  $\wedge \text{UNCHANGED } \langle sVars \rangle$ 

115 |-----|
116  $\text{Server operations at partition } p \in Partition \text{ in datacenter } d \in Datacenter.$ 

118  $ReadRequest(p, d) \triangleq$  handle a "ReadRequest"
119  $\wedge \exists m \in msgs :$ 
120  $\wedge m.type = \text{"ReadRequest"} \wedge m.p = p \wedge m.d = d$  such  $m$  may be not unique
121  $\wedge css' = [css \text{ EXCEPT } ![p][d] =$ 
122  $[dc \in Datacenter \mapsto \text{IF } dc = d \text{ THEN } @[dc] \text{ ELSE } Max(m.vc[dc], @[dc])]]$ 
123  $\wedge \text{LET } kvs \triangleq \{kv \in store[p][d] :$ 
124  $\wedge kv.key = m.key$ 
125  $\wedge \forall dc \in Datacenter \setminus \{d\} : kv.vc[dc] \leq css'[p][d][dc]\}$ 
126  $lkv \triangleq \text{CHOOSE } kv \in kvs : \text{choose the latest one (Existence? Uniqueness?)}$ 
127  $\forall akv \in kvs, dc \in Datacenter : akv.vc[dc] \leq kv.vc[dc]$ 
128  $\text{IN } SendAndDelete([type \mapsto \text{"ReadReply"}, val \mapsto lkv.val, vc \mapsto lkv.vc, c \mapsto m.c], m)$ 
129  $\wedge \text{UNCHANGED } \langle cVars, clock, pvc, updates, tick \rangle$ 

131  $UpdateRequest(p, d) \triangleq$  handle a "UpdateRequest"
132  $\wedge \exists m \in msgs :$ 
133  $\wedge m.type = \text{"UpdateRequest"} \wedge m.p = p \wedge m.d = d$  such  $m$  may be not unique
134  $\wedge m.vc[d] < clock[p][d]$  waiting condition; (" $\leq$ " strengthened to " $<$ ")
135  $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = clock[p][d]]$ 
136  $\wedge css' = [css \text{ EXCEPT } ![p][d] =$ 
137  $[dc \in Datacenter \mapsto \text{IF } dc = d \text{ THEN } @[dc] \text{ ELSE } Max(m.vc[dc], @[dc])]]$ 
138  $\wedge \text{LET } kv \triangleq [key \mapsto m.key, val \mapsto m.val,$ 
139  $vc \mapsto [m.vc \text{ EXCEPT } ![d] = clock[p][d]]]$ 
140  $\text{IN } \wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup \{kv\}]$ 
141  $\wedge updates' = [updates \text{ EXCEPT } ![p][d] = @ \circ \langle kv \rangle]$ 

```

$\wedge \text{SendAndDelete}([type \mapsto \text{"UpdateReply"}, ts \mapsto clock[p][d], c \mapsto m.c, d \mapsto d], m)$
 $\wedge \text{UNCHANGED} \langle cVars, clock, pvc, tick \rangle$
 $\text{PropagateUpdates}(p, d) \triangleq \text{propagate buffered updates to other datacenters}$
 $\wedge \text{IF } updates[p][d] \neq \langle \rangle$
 $\text{THEN } \wedge \text{SendSet}([type : \{\text{"Replicate"}\}, d : \{d\}, kvs : \{updates[p][d]\}, p : \{p\}, dd : Datacenter \setminus \{d\})$
 $\wedge updates' = [updates \text{ EXCEPT } ![p][d] = \langle \rangle]$
 $\wedge \text{UNCHANGED} \langle tick \rangle$
 $\text{ELSE } \wedge tick[p][d]$
 $\wedge \text{SendSet}([type : \{\text{"Heartbeat"}\}, d : \{d\}, ts : \{pvc[p][d][d]\}, p : \{p\}, dd : Datacenter \setminus \{d\})$
 $\wedge tick' = [tick \text{ EXCEPT } ![p][d] = \text{FALSE}]$
 $\wedge \text{UNCHANGED} \langle updates \rangle$
 $\wedge \text{UNCHANGED} \langle cVars, cvc, clock, pvc, css, store \rangle$
 $\text{Replicate}(p, d) \triangleq \text{handle a "Replicate"}$
 $\wedge \exists m \in msgs :$
 $\wedge m.type = \text{"Replicate"} \wedge m.p = p \wedge m.dd = d$
 $\wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup Range(m.kvs)]$
 $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = Last(m.kvs).vc[m.d]]$
 $\wedge msgs' = msgs \setminus \{m\}$
 $\wedge \text{UNCHANGED} \langle cVars, cvc, clock, css, updates, tick \rangle$
 $\text{Heartbeat}(p, d) \triangleq \text{handle a "Heartbeat"}$
 $\wedge \exists m \in msgs :$
 $\wedge m.type = \text{"Heartbeat"} \wedge m.p = p \wedge m.dd = d$
 $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.ts]$
 $\wedge msgs' = msgs \setminus \{m\}$
 $\wedge \text{UNCHANGED} \langle cVars, cvc, clock, css, store, updates, tick \rangle$

 $\text{Clock management at partition } p \in Partition \text{ in datacenter } d \in Datacenter$
 $\text{Tick}(p, d) \triangleq \text{clock}[p][d] \text{ ticks}$
 $\wedge clock' = [clock \text{ EXCEPT } ![p][d] = @ + 1]$
 $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = clock'[p][d]]$
 $\wedge tick' = [tick \text{ EXCEPT } ![p][d] = \text{TRUE}]$
 $\wedge \text{UNCHANGED} \langle cVars, mVars, cvc, css, store, updates \rangle$

 $\text{UpdateCSS}(p, d) \triangleq \text{update } css[p][d]$
 $\wedge css' = [css \text{ EXCEPT } ![p][d] =$
 $[dc \in Datacenter \mapsto \text{IF } dc = d \text{ THEN } @[dc]$
 $\text{ELSE } Min(\{pvc[pp][d][dc] : pp \in Partition\})]$
 $\wedge \text{UNCHANGED} \langle cVars, mVars, clock, pvc, store, updates, tick \rangle$

 $\text{Next} \triangleq$
 $\vee \exists c \in Client, k \in Key : Read(c, k)$
 $\vee \exists c \in Client, k \in Key, v \in Value : Update(c, k, v)$
 $\vee \exists c \in Client : ReadReply(c) \vee UpdateReply(c)$

```

188       $\vee \exists p \in Partition, d \in Datacenter :$ 
189           $\vee ReadRequest(p, d)$ 
190           $\vee UpdateRequest(p, d)$ 
191           $\vee PropagateUpdates(p, d)$ 
192           $\vee Replicate(p, d)$ 
193           $\vee Heartbeat(p, d)$ 
194           $\vee Tick(p, d)$ 
195           $\vee UpdateCSS(p, d)$ 

```

```

197   $Spec \triangleq Init \wedge \Box [Next]_{vars}$ 

```

```

198  ┌──────────────────────────────────────────────────────────────────────────────────┐
199  └──────────────────────────────────────────────────────────────────────────────────┘

```