

```

1  ┌────────────────────────── MODULE Cure ───────────────────────────┐
    │ See ICDCS2016: “Cure: Strong Semantics Meets High Availability and Low Latency”. │
5  └────────────────────────── EXTENDS Naturals, Sequences ───────────────────────────┘
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  └──────────────────────────┘
10  CONSTANTS
11  │ Key, the set of keys, ranged over by  $k \in Key$  │
12  │ Value, the set of values, ranged over by  $v \in Value$  │
13  │ Client, the set of clients, ranged over by  $c \in Client$  │
14  │ Partition, the set of partitions, ranged over by  $p \in Partition$  │
15  │ Datacenter, the set of datacenters, ranged over by  $d \in Datacenter$  │
16  │ KeySharding, the mapping from Key to Partition │
17  │ ClientAttachment the mapping from Client to Datacenter │
19  │  $NotVal \triangleq \text{CHOOSE } v : v \notin Value$  │
21  ASSUME
22  │  $\wedge KeySharding \in [Key \rightarrow Partition]$  │
23  │  $\wedge ClientAttachment \in [Client \rightarrow Datacenter]$  │
24  └──────────────────────────┘
25  VARIABLES
26  │ At the client side: │
27  │ cvc, cvc[c]: the vector clock of client  $c \in Client$  │
28  │ At the server side (each for partition  $p \in Partition$  in  $d \in Datacenter$ ): │
29  │ clock, clock[p][d]: the current clock │
30  │ pvc, pvc[p][d]: the vector clock │
31  │ css, css[p][d]: the stable snapshot │
32  │ store, store[p][d]: the kv store │
33  │ communication: │
34  │ msgs, the set of messages in transit │
35  │ incoming fifo[p][d]: incoming FIFO channel; for propagating updates and heartbeats │
37  │  $cVars \triangleq \langle cvc \rangle$  │
38  │  $sVars \triangleq \langle clock, pvc, css, store \rangle$  │
39  │  $mVars \triangleq \langle msgs, incoming \rangle$  │
40  │  $vars \triangleq \langle cvc, clock, pvc, css, store, msgs, incoming \rangle$  │
41  └──────────────────────────┘
42  │  $Clock \triangleq Nat$  │
43  │  $VC \triangleq [Datacenter \rightarrow Clock]$  vector clock with an entry per datacenter  $d \in Datacenter$  │
44  │  $VCInit \triangleq [d \in Datacenter \mapsto 0]$  │
45  │  $Merge(vc1, vc2) \triangleq [d \in Datacenter \mapsto Max(vc1[d], vc2[d])]$  │
46  │  $KVTuple \triangleq [key : Key, val : Value \cup \{NotVal\}, vc : VC]$  │
48  │  $Message \triangleq$  │

```

49  $[type : \{\text{"ReadRequest"}\}, key : Key, vc : VC, c : Client, p : Partition, d : Datacenter]$   
 50  $\cup [type : \{\text{"ReadReply"}\}, val : Value \cup \{\text{NotVal}\}, vc : VC, c : Client]$   
 51  $\cup [type : \{\text{"UpdateRequest"}\}, key : Key, val : Value, vc : VC, c : Client, p : Partition, d : Datacenter]$   
 52  $\cup [type : \{\text{"UpdateReply"}\}, ts : Clock, c : Client, d : Datacenter]$   
 53  $\cup [type : \{\text{"Replicate"}\}, d : Datacenter, kv : KVTuple]$   
 54  $\cup [type : \{\text{"Heartbeat"}\}, d : Datacenter, ts : Clock]$   
 56  $TypeOK \triangleq$   
 57  $\wedge cvc \in [Client \rightarrow VC]$   
 58  $\wedge clock \in [Partition \rightarrow [Datacenter \rightarrow Clock]]$   
 59  $\wedge pvc \in [Partition \rightarrow [Datacenter \rightarrow VC]]$   
 60  $\wedge css \in [Partition \rightarrow [Datacenter \rightarrow VC]]$   
 61  $\wedge store \in [Partition \rightarrow [Datacenter \rightarrow \text{SUBSET } KVTuple]]$   
 62  $\wedge msgs \subseteq Message$   
 63  $\wedge incoming \in [Partition \rightarrow [Datacenter \rightarrow Seq(Message)]]$   
 64  $\vdash$   
 65  $Init \triangleq$   
 66  $\wedge cvc = [c \in Client \mapsto VCInit]$   
 67  $\wedge clock = [p \in Partition \mapsto [d \in Datacenter \mapsto 0]]$   
 68  $\wedge pvc = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$   
 69  $\wedge css = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$   
 70  $\wedge store = [p \in Partition \mapsto [d \in Datacenter \mapsto$   
 71  $\quad [key : \{k \in Key : KeySharding[k] = p\}, val : \{\text{NotVal}\}, vc : \{VCInit\}]]]$   
 72  $\wedge msgs = \{\}$   
 73  $\wedge incoming = [p \in Partition \mapsto [d \in Datacenter \mapsto \langle \rangle]]$   
 74  $\vdash$   
 75  $Send(m) \triangleq msgs' = msgs \cup \{m\}$   
 76  $SendAndDelete(sm, dm) \triangleq msgs' = (msgs \cup \{sm\}) \setminus \{dm\}$   
 78  $CanIssue(c) \triangleq \forall m \in msgs :$   
 79  $m.type \in \{\text{"ReadRequest"}, \text{"ReadReply"}, \text{"UpdateRequest"}, \text{"UpdateReply"}\} \Rightarrow m.c \neq c$   
 80  $\vdash$   
 81 **Client operations at client  $c \in Client$ .**  
 83  $Read(c, k) \triangleq$   $c \in Client$  reads from  $k \in Key$   
 84  $\wedge CanIssue(c)$   
 85  $\wedge Send([type \mapsto \text{"ReadRequest"}, key \mapsto k, vc \mapsto cvc[c],$   
 86  $\quad c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$   
 87  $\wedge \text{UNCHANGED } \langle cVars, sVars, incoming \rangle$   
 89  $ReadReply(c) \triangleq$   $c \in Client$  handles the reply to its read request  
 90  $\wedge \exists m \in msgs :$   
 91  $\wedge m.type = \text{"ReadReply"} \wedge m.c = c$  such  $m$  is unique due to well-formedness  
 92  $\wedge cvc' = [cvc \text{ EXCEPT } !c] = Merge(m.vc, @)$   
 93  $\wedge msgs' = msgs \setminus \{m\}$   
 94  $\wedge \text{UNCHANGED } \langle sVars, incoming \rangle$

```

96  $Update(c, k, v) \triangleq$   $c \in Client$  updates  $k \in Key$  with  $v \in Value$ 
97  $\wedge CanIssue(c)$ 
98  $\wedge Send([type \mapsto \text{"UpdateRequest"}, key \mapsto k, val \mapsto v,$ 
99  $vc \mapsto cvc[c], c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$ 
100  $\wedge UNCHANGED \langle cVars, sVars, incoming \rangle$ 

102  $UpdateReply(c) \triangleq$   $c \in Client$  handles the reply to its update request
103  $\wedge \exists m \in msgs :$ 
104  $\wedge m.type = \text{"UpdateReply"} \wedge m.c = c$  such  $m$  is unique due to well-formedness
105  $\wedge cvc' = [cvc \text{ EXCEPT } ![c][m.d] = m.ts]$ 
106  $\wedge msgs' = msgs \setminus \{m\}$ 
107  $\wedge UNCHANGED \langle sVars, incoming \rangle$ 

108 |-----|
109  $Server\ operations\ at\ partition\ p \in Partition\ in\ datacenter\ d \in Datacenter.$ 

111  $ReadRequest(p, d) \triangleq$  handle a "ReadRequest"
112  $\wedge \exists m \in msgs :$ 
113  $\wedge m.type = \text{"ReadRequest"} \wedge m.p = p \wedge m.d = d$ 
114  $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$ 
115  $\wedge LET\ kvs \triangleq \{kv \in store[p][d] :$ 
116  $\wedge kv.key = m.key$ 
117  $\wedge \forall dc \in Datacenter \setminus \{d\} : kv.vc[dc] \leq css'[p][d][dc]\}$ 
118  $lkv \triangleq CHOOSE\ kv \in kvs : \text{choose the latest one (Existence? Uniqueness?)}$ 
119  $\forall akv \in kvs, dc \in Datacenter : akv.vc[dc] \leq kv.vc[dc]$ 
120  $IN\ SendAndDelete([type \mapsto \text{"ReadReply"}, val \mapsto lkv.val, vc \mapsto lkv.vc, c \mapsto m.c], m)$ 
121  $\wedge UNCHANGED \langle cVars, clock, pvc, store, incoming \rangle$ 

123  $UpdateRequest(p, d) \triangleq$  handle a "UpdateRequest"
124  $\wedge \exists m \in msgs :$ 
125  $\wedge m.type = \text{"UpdateRequest"} \wedge m.p = p \wedge m.d = d$ 
126  $\wedge m.vc[d] < clock[p][d]$  waiting condition; (" $\leq$ " strengthened to " $<$ ")
127  $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$ 
128  $\wedge LET\ kv \triangleq [key \mapsto m.key, val \mapsto m.val,$ 
129  $vc \mapsto [m.vc \text{ EXCEPT } ![d] = clock[p][d]]]$ 
130  $IN\ \wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup \{kv\}]$ 
131  $\wedge SendAndDelete([type \mapsto \text{"UpdateReply"}, ts \mapsto clock[p][d], c \mapsto m.c, d \mapsto d], m)$ 
132  $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
133  $IF\ dc = d\ THEN\ @[dc]\ ELSE\ Append(@[dc], [type \mapsto \text{"Replicate"}, d \mapsto d, kv \mapsto kv])]]$ 
134  $\wedge UNCHANGED \langle cVars, clock, pvc \rangle$ 

136  $Replicate(p, d) \triangleq$  handle a "Replicate"
137  $\wedge incoming[p][d] \neq \langle \rangle$ 
138  $\wedge LET\ m \triangleq Head(incoming[p][d])$ 
139  $IN\ \wedge m.type = \text{"Replicate"}$ 
140  $\wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup \{m.kv\}]$ 
141  $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.kv.vc[m.d]]$ 

```

```

142       $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
143       $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, msgs \rangle$ 

145   $Heartbeat(p, d) \triangleq$  handle a "Heartbeat"
146       $\wedge incoming[p][d] \neq \langle \rangle$ 
147       $\wedge \text{LET } m \triangleq Head(incoming[p][d])$ 
148       $\text{IN } \wedge m.type = \text{"Heartbeat"}$ 
149       $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.ts]$ 
150       $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
151       $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, store, msgs \rangle$ 
152  ────────────────────────────────────────────────────────────────────────────────────────────────────

153  Clock management at partition  $p \in Partition$  in datacenter  $d \in Datacenter$ 
154   $Tick(p, d) \triangleq$   $clock[p][d]$  ticks
155       $\wedge clock' = [clock \text{ EXCEPT } ![p][d] = @ + 1]$ 
156       $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = clock'[p][d]]$ 
157       $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
158           $\text{IF } dc = d \text{ THEN } @[dc] \text{ ELSE } Append(@[dc], [type \mapsto \text{"Heartbeat"}, d \mapsto d, ts \mapsto pvc'[p][d][d]])]]$ 
159       $\wedge \text{UNCHANGED } \langle cVars, cvc, css, store, msgs \rangle$ 

161   $UpdateCSS(p, d) \triangleq$  update  $css[p][d]$ 
162       $\wedge css' = [css \text{ EXCEPT } ![p][d] =$ 
163           $[dc \in Datacenter \mapsto Min(\{pvc[pp][d][dc] : pp \in Partition\})]]$ 
164       $\wedge \text{UNCHANGED } \langle cVars, mVars, clock, pvc, store \rangle$ 
165  ────────────────────────────────────────────────────────────────────────────────────────────────────

166   $Next \triangleq$ 
167       $\vee \exists c \in Client, k \in Key : Read(c, k)$ 
168       $\vee \exists c \in Client, k \in Key, v \in Value : Update(c, k, v)$ 
169       $\vee \exists c \in Client : ReadReply(c) \vee UpdateReply(c)$ 
170       $\vee \exists p \in Partition, d \in Datacenter :$ 
171           $\vee ReadRequest(p, d)$ 
172           $\vee UpdateRequest(p, d)$ 
173           $\vee Replicate(p, d)$ 
174           $\vee Heartbeat(p, d)$ 
175           $\vee Tick(p, d)$ 
176           $\vee UpdateCSS(p, d)$ 

178   $Spec \triangleq Init \wedge \Box [Next]_{vars}$ 
179  ────────────────────────────────────────────────────────────────────────────────────────────────────

```