

```

1  |----- MODULE Cure -----|
2  | See ICDCS2016: “Cure: Strong Semantics Meets High Availability and Low Latency”. |
3  |-----|
4  |
5  | EXTENDS Naturals, Sequences, FiniteSets, TLC, RelationUtils |
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 |
18 |
19 |  $NotVal \triangleq \text{CHOOSE } v : v \notin Value$  |
20 |
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 |   history: |
34 |     L, L[c]: local history at client  $c \in Client$  |
35 |   communication: |
36 |     msgs, the set of messages in transit |
37 |     incoming incoming[p][d]: incoming FIFO channel for propagating updates and heartbeats |
38 |
39 |  $cVars \triangleq \langle cvc \rangle$  |
40 |  $sVars \triangleq \langle clock, pvc, css, store, L \rangle$  |
41 |  $mVars \triangleq \langle msgs, incoming \rangle$  |
42 |  $vars \triangleq \langle cvc, clock, pvc, css, store, L, msgs, incoming \rangle$  |
43 |-----|
44 |  $VC \triangleq [Datacenter \rightarrow Nat]$  vector clock with an entry per datacenter  $d \in Datacenter$  |
45 |  $VCInit \triangleq [d \in Datacenter \mapsto 0]$  |
46 |  $Merge(vc1, vc2) \triangleq [d \in Datacenter \mapsto Max(vc1[d], vc2[d])]$  |
47 |
48 |  $DC \triangleq Cardinality(Datacenter)$ 

```

```

49  $DCIndex \triangleq \text{CHOOSE } f \in [1 \dots DC \rightarrow \text{Datacenter}] : \text{Injective}(f)$ 
50  $LTE(vc1, vc2) \triangleq$  less-than-or-equal-to comparator for vector clocks
51   LET RECURSIVE  $LTEHelper(-, -, -)$ 
52      $LTEHelper(vc1h, vc2h, index) \triangleq$ 
53       IF  $index > DC$  THEN TRUE  $EQ$ 
54       ELSE LET  $d \triangleq DCIndex[index]$ 
55         IN CASE  $vc1h[d] < vc2h[d] \rightarrow$  TRUE  $LT$ 
56           □  $vc1h[d] > vc2h[d] \rightarrow$  FALSE  $GT$ 
57           □ OTHER  $\rightarrow LTEHelper(vc1h, vc2h, index + 1)$ 
58   IN  $LTEHelper(vc1, vc2, 1)$ 

60  $KVTuple \triangleq [key : Key, val : Value \cup \{NotVal\}, vc : VC]$ 
61  $OpTuple \triangleq [type : \{“R”, “W”\}, kv : KVTuple]$ 

63  $Message \triangleq$ 
64    $[type : \{“ReadRequest”\}, key : Key, vc : VC, c : Client, p : Partition, d : Datacenter]$ 
65    $\cup [type : \{“ReadReply”\}, val : Value \cup \{NotVal\}, vc : VC, c : Client]$ 
66    $\cup [type : \{“UpdateRequest”\}, key : Key, val : Value, vc : VC, c : Client, p : Partition, d : Datacenter]$ 
67    $\cup [type : \{“UpdateReply”\}, ts : Nat, c : Client, d : Datacenter]$ 
68    $\cup [type : \{“Replicate”\}, d : Datacenter, kv : KVTuple]$ 
69    $\cup [type : \{“Heartbeat”\}, d : Datacenter, ts : Nat]$ 

71  $Send(m) \triangleq msgs' = msgs \cup \{m\}$ 
72  $SendAndDelete(sm, dm) \triangleq msgs' = (msgs \cup \{sm\}) \setminus \{dm\}$ 

74  $TypeOK \triangleq$ 
75    $\wedge cvc \in [Client \rightarrow VC]$ 
76    $\wedge clock \in [Partition \rightarrow [Datacenter \rightarrow Nat]]$ 
77    $\wedge pvc \in [Partition \rightarrow [Datacenter \rightarrow VC]]$ 
78    $\wedge css \in [Partition \rightarrow [Datacenter \rightarrow VC]]$ 
79    $\wedge store \in [Partition \rightarrow [Datacenter \rightarrow \text{SUBSET } KVTuple]]$ 
80    $\wedge msgs \subseteq Message$ 
81    $\wedge incoming \in [Partition \rightarrow [Datacenter \rightarrow Seq(Message)]]$ 
82    $\wedge L \in [Client \rightarrow Seq(OpTuple)]$ 
83 |-----|

84  $Init \triangleq$ 
85    $\wedge cvc = [c \in Client \mapsto VCInit]$ 
86    $\wedge clock = [p \in Partition \mapsto [d \in Datacenter \mapsto 0]]$ 
87    $\wedge pvc = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$ 
88    $\wedge css = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$ 
89    $\wedge store = [p \in Partition \mapsto [d \in Datacenter \mapsto$ 
90      $[key : \{k \in Key : KeySharding[k] = p\}, val : \{NotVal\}, vc : \{VCInit\}]]]$ 
91    $\wedge msgs = \{\}$ 
92    $\wedge incoming = [p \in Partition \mapsto [d \in Datacenter \mapsto \langle \rangle]]$ 
93    $\wedge L = [c \in Client \mapsto \langle \rangle]$ 
94 |-----|

```

95 Client operations at client $c \in Client$.

97 $CanIssue(c) \triangleq \forall m \in msgs : \text{to ensure well-formedness of clients}$
98 $m.type \in \{\text{"ReadRequest"}, \text{"ReadReply"}, \text{"UpdateRequest"}, \text{"UpdateReply"}\} \Rightarrow m.c \neq c$

100 $Read(c, k) \triangleq c \in Client \text{ reads from } k \in Key$
101 $\wedge CanIssue(c)$
102 $\wedge Send([type \mapsto \text{"ReadRequest"}, key \mapsto k, vc \mapsto cvc[c],$
103 $c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$
104 $\wedge UNCHANGED \langle cVars, sVars, incoming \rangle$

106 $ReadReply(c) \triangleq c \in Client \text{ handles the reply to its read request}$
107 $\wedge \exists m \in msgs :$
108 $\wedge m.type = \text{"ReadReply"} \wedge m.c = c \text{ such } m \text{ is unique due to well-formedness}$
109 $\wedge cvc' = [cvc \text{ EXCEPT } ![c] = Merge(m.vc, @)]$
110 $\wedge msgs' = msgs \setminus \{m\}$
111 $\wedge UNCHANGED \langle sVars, incoming \rangle$

113 $Update(c, k, v) \triangleq c \in Client \text{ updates } k \in Key \text{ with } v \in Value$
114 $\wedge CanIssue(c)$
115 $\wedge Send([type \mapsto \text{"UpdateRequest"}, key \mapsto k, val \mapsto v,$
116 $vc \mapsto cvc[c], c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$
117 $\wedge UNCHANGED \langle cVars, sVars, incoming \rangle$

119 $UpdateReply(c) \triangleq c \in Client \text{ handles the reply to its update request}$
120 $\wedge \exists m \in msgs :$
121 $\wedge m.type = \text{"UpdateReply"} \wedge m.c = c \text{ such } m \text{ is unique due to well-formedness}$
122 $\wedge cvc' = [cvc \text{ EXCEPT } ![c][m.d] = m.ts]$
123 $\wedge msgs' = msgs \setminus \{m\}$
124 $\wedge UNCHANGED \langle sVars, incoming \rangle$

126 Server operations at partition $p \in Partition$ in datacenter $d \in Datacenter$.

128 $ReadRequest(p, d) \triangleq \text{handle a "ReadRequest"}$
129 $\wedge \exists m \in msgs :$
130 $\wedge m.type = \text{"ReadRequest"} \wedge m.p = p \wedge m.d = d$
131 $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$
132 $\wedge LET \ kvs \triangleq \{kv \in store[p][d] :$
133 $\wedge kv.key = m.key$
134 $\wedge \forall dc \in Datacenter \setminus \{d\} : kv.vc[dc] \leq css'[p][d][dc]\}$
135 $lkv \triangleq CHOOSE \ kv \in kvs : \forall akv \in kvs : LTE(akv.vc, kv.vc)$
136 $IN \wedge SendAndDelete([type \mapsto \text{"ReadReply"}, val \mapsto lkv.val, vc \mapsto lkv.vc, c \mapsto m.c], m)$
137 $\wedge L' = [L \text{ EXCEPT } ![m.c] = Append(@, [type \mapsto \text{"R"}, kv \mapsto lkv])]$
138 $\wedge UNCHANGED \langle cVars, clock, pvc, store, incoming \rangle$

140 $UpdateRequest(p, d) \triangleq \text{handle a "UpdateRequest"}$
141 $\wedge \exists m \in msgs :$

```

142       $\wedge m.type = \text{"UpdateRequest"} \wedge m.p = p \wedge m.d = d$ 
143       $\wedge m.vc[d] < clock[p][d]$  waiting condition; ("≤" strengthened to "<")
144       $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$ 
145       $\wedge \text{LET } kv \triangleq [key \mapsto m.key, val \mapsto m.val,$ 
146         $vc \mapsto [m.vc \text{ EXCEPT } ![d] = clock[p][d]]]$ 
147      IN  $\wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup \{kv\}]$ 
148         $\wedge SendAndDelete([type \mapsto \text{"UpdateReply"}, ts \mapsto clock[p][d], c \mapsto m.c, d \mapsto d], m)$ 
149         $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
150          IF  $dc = d$  THEN  $@[dc]$  ELSE  $Append(@[dc], [type \mapsto \text{"Replicate"}, d \mapsto d, kv \mapsto kv])]$ 
151           $\wedge L' = [L \text{ EXCEPT } ![m.c] = Append(@, [type \mapsto \text{"R"}, kv \mapsto kv])]$ 
152       $\wedge \text{UNCHANGED } \langle cVars, clock, pvc \rangle$ 

154  $Replicate(p, d) \triangleq$  handle a "Replicate"
155    $\wedge incoming[p][d] \neq \langle \rangle$ 
156    $\wedge \text{LET } m \triangleq Head(incoming[p][d])$ 
157   IN  $\wedge m.type = \text{"Replicate"}$ 
158      $\wedge store' = [store \text{ EXCEPT } ![p][d] = @ \cup \{m.kv\}]$ 
159      $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.kv.vc[m.d]]$ 
160      $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
161    $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, L, msgs \rangle$ 

163  $Heartbeat(p, d) \triangleq$  handle a "Heartbeat"
164    $\wedge incoming[p][d] \neq \langle \rangle$ 
165    $\wedge \text{LET } m \triangleq Head(incoming[p][d])$ 
166   IN  $\wedge m.type = \text{"Heartbeat"}$ 
167      $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.ts]$ 
168      $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
169    $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, store, L, msgs \rangle$ 


---


171 Clock management at partition  $p \in Partition$  in datacenter  $d \in Datacenter$ 
172  $Tick(p, d) \triangleq$   $clock[p][d]$  ticks
173    $\wedge clock' = [clock \text{ EXCEPT } ![p][d] = @ + 1]$ 
174    $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = clock'[p][d]]$ 
175    $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
176     IF  $dc = d$  THEN  $@[dc]$  ELSE  $Append(@[dc], [type \mapsto \text{"Heartbeat"}, d \mapsto d, ts \mapsto pvc'[p][d][d]])]$ 
177    $\wedge \text{UNCHANGED } \langle cVars, cvc, css, store, L, msgs \rangle$ 

179  $UpdateCSS(p, d) \triangleq$  update  $css[p][d]$ 
180    $\wedge css' = [css \text{ EXCEPT } ![p][d] =$ 
181      $[dc \in Datacenter \mapsto Min(\{pvc[pp][d][dc] : pp \in Partition\})]$ 
182    $\wedge \text{UNCHANGED } \langle cVars, mVars, clock, pvc, store, L \rangle$ 


---


184  $Next \triangleq$ 
185    $\vee \exists c \in Client, k \in Key : Read(c, k)$ 
186    $\vee \exists c \in Client, k \in Key, v \in Value : Update(c, k, v)$ 
187    $\vee \exists c \in Client : ReadReply(c) \vee UpdateReply(c)$ 

```

```

188      $\vee \exists p \in \text{Partition}, d \in \text{Datacenter} :$ 
189          $\vee \text{ReadRequest}(p, d)$ 
190          $\vee \text{UpdateRequest}(p, d)$ 
191          $\vee \text{Replicate}(p, d)$ 
192          $\vee \text{Heartbeat}(p, d)$ 
193          $\vee \text{Tick}(p, d)$ 
194          $\vee \text{UpdateCSS}(p, d)$ 
196      $\text{Spec} \triangleq \text{Init} \wedge \Box[\text{Next}]_{\text{vars}}$ 
197 |-----|
198      $\text{so} \triangleq \text{UNION } \{ \text{SeqToRel}(L[c]) : c \in \text{Client} \}$  session order
200      $\text{rf} \triangleq$  read-from (or called writes-into) relation
201     LET  $\text{ops} \triangleq \text{UNION } \{ \text{Range}(L[c]) : c \in \text{Client} \}$ 
202      $\text{rops} \triangleq \{ op \in \text{ops} : op.type = \text{"R"} \}$ 
203      $\text{keys} \triangleq \{ op.kv.key : op \in \text{rops} \}$ 
204      $\text{wops} \triangleq \{ op \in \text{ops} : op.type = \text{"W"} \} \cup$  initial writes
205          $[type : \{ \text{"W"} \}, kv : [key : \text{keys}, val : \{ \text{NotVal} \}, vc : \{ \text{VCInit} \}]]$ 
206     IN  $\{ \langle w, r \rangle \in \text{wops} \times \text{rops} : w.kv.key = r.kv.key \wedge w.kv.vc = r.kv.vc \}$ 
208      $\text{co} \triangleq \text{TC}(\text{so} \cup \text{rf})$  causality order
210      $\text{Valid}(s) \triangleq$  Is  $s$  a valid serialization?
211     LET RECURSIVE  $\text{ValidHelper}(-, -)$ 
212          $\text{ValidHelper}(seq, kvs) \triangleq$ 
213             IF  $seq = \langle \rangle$  THEN TRUE
214             ELSE LET  $op \triangleq \text{Head}(seq)$ 
215                 IN IF  $op.type = \text{"W"}$ 
216                     THEN  $\text{ValidHelper}(\text{Tail}(seq), kvs @ @ op.kv.key :> op.kv.vc)$ 
217                     ELSE  $\wedge op.kv.vc = kvs[op.kv.key]$ 
218                          $\wedge \text{ValidHelper}(\text{Tail}(seq), kvs)$ 
219     IN  $\text{ValidHelper}(s, \{ \})$ 
221     TODO: Handling initial writes
222      $\text{CM} \triangleq$  causal memory consistency model; see  $DC'1995$ 
223     LET  $\text{ops} \triangleq \text{UNION } \{ \text{Range}(L[c]) : c \in \text{Client} \}$ 
224      $\text{wops} \triangleq \{ op \in \text{ops} : op.type = \text{"W"} \}$ 
225     IN  $\forall c \in \text{Client} :$ 
226          $\exists sc \in \text{Seq}(\text{Range}(L[c]) \cup \text{wops}) :$  TODO: performance?
227          $\wedge \text{Valid}(sc)$ 
228          $\wedge \text{Respect}(sc, co)$ 
229 |-----|

```