

```

1  ─── MODULE CureKV ───
2  See ICDCS2016: “Cure: Strong Semantics Meets High Availability and Low Latency”.
3
4  EXTENDS Naturals, FiniteSets, TLC, SequenceUtils, RelationUtils, MathUtils
5
6  ───
7  CONSTANTS
8    Key,           the set of keys, ranged over by  $k \in Key$ 
9    Value,         the set of values, ranged over by  $v \in Value$ 
10   Client,        the set of clients, ranged over by  $c \in Client$ 
11   Partition,     the set of partitions, ranged over by  $p \in Partition$ 
12   Datacenter,    the set of datacenters, ranged over by  $d \in Datacenter$ 
13   KeySharding,   the mapping from Key to Partition
14   ClientAttachment the mapping from Client to Datacenter
15
16    $NotVal \triangleq \text{CHOOSE } v : v \notin Value$ 
17
18   ASSUME
19      $\wedge KeySharding \in [Key \rightarrow Partition]$ 
20      $\wedge ClientAttachment \in [Client \rightarrow Datacenter]$ 
21
22    $KeysOnPartition \triangleq [p \in Partition \mapsto \{k \in Key : KeySharding[k] = p\}]$ 
23
24   VARIABLES
25     At the client side:
26       cvc, cvc[c]: the vector clock of client  $c \in Client$ 
27     At the server side (each for partition  $p \in Partition$  in  $d \in Datacenter$ ):
28       clock, clock[p][d]: the current clock
29       pvc, pvc[p][d]: the vector clock
30       css, css[p][d]: the stable snapshot
31       store, store[p][d]: the kv store (represented by a map)
32       remote, remote[p][d]: the set of remote updates that have not been applied yet
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    $sVars \triangleq \langle clock, pvc, css, store, remote, L \rangle$ 
41    $mVars \triangleq \langle msgs, incoming \rangle$ 
42    $vars \triangleq \langle cvc, clock, pvc, css, store, remote, L, msgs, incoming \rangle$ 
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  $OpTuple \triangleq [type : \{“R”, “W”\}, kv : KVTuple, c : Client, cnt : Nat]$ 

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    $\wedge store \in [Partition \rightarrow [Datacenter \rightarrow [Key \rightarrow KVTuple]]]$ 
80    $\wedge remote \in [Partition \rightarrow [Datacenter \rightarrow \text{SUBSET } KVTuple]]$ 
81    $\wedge msgs \subseteq Message$ 
82    $\wedge incoming \in [Partition \rightarrow [Datacenter \rightarrow Seq(Message)]]$ 
83    $\wedge L \in [Client \rightarrow Seq(OpTuple)]$ 
84 |-----|
85  $Init \triangleq$ 
86    $\wedge cvc = [c \in Client \mapsto VCInit]$ 
87    $\wedge clock = [p \in Partition \mapsto [d \in Datacenter \mapsto 0]]$ 
88    $\wedge pvc = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$ 
89    $\wedge css = [p \in Partition \mapsto [d \in Datacenter \mapsto VCInit]]$ 
90    $\wedge store = [p \in Partition \mapsto [d \in Datacenter \mapsto$ 
91      $[k \in KeysOnPartition[p] \mapsto [key \mapsto k, val \mapsto NotVal, vc \mapsto VCInit]]]]$ 
92    $\wedge remote = [p \in Partition \mapsto [d \in Datacenter \mapsto \{\}]]$ 
93    $\wedge msgs = \{\}$ 
94    $\wedge incoming = [p \in Partition \mapsto [d \in Datacenter \mapsto \langle \rangle]]$ 

```

---

95  $\wedge L = [c \in Client \mapsto \langle \rangle]$   
 96 

---

  
 97 **Client operations at client  $c \in Client$ .**  
 98  
 99  $CanIssue(c) \triangleq \forall m \in msgs :$  **to ensure well-formedness of clients**  
 100  $m.type \in \{\text{"ReadRequest"}, \text{"ReadReply"}, \text{"UpdateRequest"}, \text{"UpdateReply"}\} \Rightarrow m.c \neq c$   
 101  
 102  $Read(c, k) \triangleq$   **$c \in Client$  reads from  $k \in Key$**   
 103  $\wedge CanIssue(c)$   
 104  $\wedge Send([type \mapsto \text{"ReadRequest"}, key \mapsto k, vc \mapsto cvc[c],$   
 105  $c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$   
 106  $\wedge UNCHANGED \langle cVars, sVars, incoming \rangle$   
 107  
 108  $ReadReply(c) \triangleq$   **$c \in Client$  handles the reply to its read request**  
 109  $\wedge \exists m \in msgs :$   
 110  $\wedge m.type = \text{"ReadReply"} \wedge m.c = c$  **such  $m$  is unique due to well-formedness**  
 111  $\wedge cvc' = [cvc \text{ EXCEPT } ![c] = Merge(m.vc, @)]$   
 112  $\wedge msgs' = msgs \setminus \{m\}$   
 113  $\wedge UNCHANGED \langle sVars, incoming \rangle$   
 114  
 115  $Update(c, k, v) \triangleq$   **$c \in Client$  updates  $k \in Key$  with  $v \in Value$**   
 116  $\wedge CanIssue(c)$   
 117  $\wedge Send([type \mapsto \text{"UpdateRequest"}, key \mapsto k, val \mapsto v,$   
 118  $vc \mapsto cvc[c], c \mapsto c, p \mapsto KeySharding[k], d \mapsto ClientAttachment[c]])$   
 119  $\wedge UNCHANGED \langle cVars, sVars, incoming \rangle$   
 120  
 121  $UpdateReply(c) \triangleq$   **$c \in Client$  handles the reply to its update request**  
 122  $\wedge \exists m \in msgs :$   
 123  $\wedge m.type = \text{"UpdateReply"} \wedge m.c = c$  **such  $m$  is unique due to well-formedness**  
 124  $\wedge cvc' = [cvc \text{ EXCEPT } ![c][m.d] = m.ts]$   
 125  $\wedge msgs' = msgs \setminus \{m\}$   
 126  $\wedge UNCHANGED \langle sVars, incoming \rangle$   
 127 

---

  
 128 **Server operations at partition  $p \in Partition$  in datacenter  $d \in Datacenter$ .**  
 129  
 130  $Apply(p, d) \triangleq$  **apply remote updates which become stable due to up-to-date  $css[p][d]$**   
 131  $LET keys \triangleq \{kv.key : kv \in remote[p][d]\}$   
 132  $kvs \triangleq [k \in keys \mapsto \{kv \in remote[p][d] :$   
 133  $\wedge kv.key = k$   
 134  $\wedge \forall dc \in Datacenter \setminus \{d\} : kv.vc[dc] \leq css'[p][d][dc]\}]$   
 135  $IN \wedge store' = [store \text{ EXCEPT } ![p][d] = [k \in KeysOnPartition[p] \mapsto$   
 136  $IF k \in keys \wedge kvs[k] \neq \{\}$   
 137  $THEN CHOOSE kv \in kvs[k] : \forall akv \in kvs[k] : LTE(akv.vc, kv.vc)$   
 138  $ELSE @[k]]]$   
 139  $\wedge remote' = [remote \text{ EXCEPT } ![p][d] = @ \setminus (UNION (Range(kvs)))]$   
 140  
 141  $ReadRequest(p, d) \triangleq$  **handle a "ReadRequest"**

```

142  $\wedge \exists m \in msgs :$ 
143    $\wedge m.type = \text{"ReadRequest"} \wedge m.p = p \wedge m.d = d$ 
144    $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$ 
145    $\wedge LET \ kv \triangleq store[p][d][m.key]$ 
146   IN  $\wedge SendAndDelete([type \mapsto \text{"ReadReply"}, val \mapsto kv.val, vc \mapsto kv.vc, c \mapsto m.c], m)$ 
147    $\wedge Apply(p, d)$ 
148    $\wedge L' = [L \text{ EXCEPT } ![m.c] = Append(@, [type \mapsto \text{"R"}, kv \mapsto kv, c \mapsto m.c, cnt \mapsto Len(@) + 1])]$ 
149  $\wedge \text{UNCHANGED } \langle cVars, clock, pvc, incoming \rangle$ 

151  $UpdateRequest(p, d) \triangleq$  handle a "UpdateRequest"
152  $\wedge \exists m \in msgs :$ 
153    $\wedge m.type = \text{"UpdateRequest"} \wedge m.p = p \wedge m.d = d$ 
154    $\wedge m.vc[d] < clock[p][d]$  waiting condition; (" $\leq$ " strengthened to " $<$ ")
155    $\wedge css' = [css \text{ EXCEPT } ![p][d] = Merge(m.vc, @)]$ 
156    $\wedge LET \ kv \triangleq [key \mapsto m.key, val \mapsto m.val,$ 
157      $vc \mapsto [m.vc \text{ EXCEPT } ![d] = clock[p][d]]]$ 
158   IN  $\wedge store' = [store \text{ EXCEPT } ![p][d][m.key] = kv]$  TODO: Apply???
159    $\wedge SendAndDelete([type \mapsto \text{"UpdateReply"}, ts \mapsto clock[p][d], c \mapsto m.c, d \mapsto d], m)$ 
160    $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
161     IF  $dc = d$  THEN  $@[dc]$  ELSE  $Append(@[dc], [type \mapsto \text{"Replicate"}, d \mapsto d, kv \mapsto kv])]$ 
162    $\wedge L' = [L \text{ EXCEPT } ![m.c] = Append(@, [type \mapsto \text{"W"}, kv \mapsto kv, c \mapsto m.c, cnt \mapsto Len(@) + 1])]$ 
163  $\wedge \text{UNCHANGED } \langle cVars, clock, pvc, remote \rangle$ 

165  $Replicate(p, d) \triangleq$  handle a "Replicate"
166  $\wedge incoming[p][d] \neq \langle \rangle$ 
167  $\wedge LET \ m \triangleq Head(incoming[p][d])$ 
168 IN  $\wedge m.type = \text{"Replicate"}$ 
169  $\wedge remote' = [remote \text{ EXCEPT } ![p][d] = @ \cup \{m.kv\}]$ 
170  $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.kv.vc[m.d]]$ 
171  $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
172  $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, store, L, msgs \rangle$ 

174  $Heartbeat(p, d) \triangleq$  handle a "Heartbeat"
175  $\wedge incoming[p][d] \neq \langle \rangle$ 
176  $\wedge LET \ m \triangleq Head(incoming[p][d])$ 
177 IN  $\wedge m.type = \text{"Heartbeat"}$ 
178  $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][m.d] = m.ts]$ 
179  $\wedge incoming' = [incoming \text{ EXCEPT } ![p][d] = Tail(@)]$ 
180  $\wedge \text{UNCHANGED } \langle cVars, cvc, clock, css, store, remote, L, msgs \rangle$ 

181  $\vdash$  Clock management at partition  $p \in Partition$  in datacenter  $d \in Datacenter$ 
182  $Tick(p, d) \triangleq$   $clock[p][d]$  ticks
183  $\wedge clock' = [clock \text{ EXCEPT } ![p][d] = @ + 1]$ 
184  $\wedge pvc' = [pvc \text{ EXCEPT } ![p][d][d] = clock'[p][d]]$ 
185  $\wedge incoming' = [incoming \text{ EXCEPT } ![p] = [dc \in Datacenter \mapsto$ 
186   IF  $dc = d$  THEN  $@[dc]$  ELSE  $Append(@[dc], [type \mapsto \text{"Heartbeat"}, d \mapsto d, ts \mapsto pvc'[p][d][d]])]$ 

```

```

188       $\wedge$  UNCHANGED  $\langle cVars, cvc, css, store, remote, L, msgs \rangle$ 
190  UpdateCSS( $p, d$ )  $\triangleq$  update  $css[p][d]$ 
191       $\wedge css' = [css \text{ EXCEPT } ![p][d] =$ 
192           $[dc \in Datacenter \mapsto SetMin(\{pvc[pp][d][dc] : pp \in Partition\})]$ 
193       $\wedge Apply(p, d)$ 
194       $\wedge$  UNCHANGED  $\langle cVars, mVars, clock, pvc, L \rangle$ 
195  |-----|
196  Next  $\triangleq$ 
197       $\vee \exists c \in Client, k \in Key : Read(c, k)$ 
198       $\vee \exists c \in Client, k \in Key, v \in Value : Update(c, k, v)$ 
199       $\vee \exists c \in Client : ReadReply(c) \vee UpdateReply(c)$ 
200       $\vee \exists p \in Partition, d \in Datacenter :$ 
201           $\vee ReadRequest(p, d)$ 
202           $\vee UpdateRequest(p, d)$ 
203           $\vee Replicate(p, d)$ 
204           $\vee Heartbeat(p, d)$ 
205           $\vee Tick(p, d)$ 
206           $\vee UpdateCSS(p, d)$ 
208  Spec  $\triangleq Init \wedge \Box [Next]_{vars}$ 
209  |-----|
210  Valid( $s$ )  $\triangleq$  Is  $s$  a valid serialization?
211      LET RECURSIVE ValidHelper( $-, -$ )
212          ValidHelper( $seq, kvs$ )  $\triangleq$ 
213              IF  $seq = \langle \rangle$  THEN TRUE
214              ELSE LET  $op \triangleq Head(seq)$ 
215                  IN IF  $op.type = "W"$  overwritten
216                      THEN ValidHelper( $Tail(seq), op.kv.key :> op.kv.vc @@ kvs$ )
217                      ELSE  $\wedge op.kv.vc = kvs[op.kv.key]$ 
218                           $\wedge ValidHelper(Tail(seq), kvs)$ 
219      IN ValidHelper( $s, [k \in Key \mapsto VInit]$ ) with initial values
221  CM  $\triangleq$  causal memory consistency model; see Ahamad@DC'1995
222      LET  $ops \triangleq$  UNION  $\{Range(L[c]) : c \in Client\}$ 
223       $rops \triangleq \{op \in ops : op.type = "R"\}$ 
224       $wops \triangleq \{op \in ops : op.type = "W"\}$ 
225       $so \triangleq$  UNION  $\{SeqToRel(L[c]) : c \in Client\}$  session order
226       $rf \triangleq \{\langle w, r \rangle \in wops \times rops : w.kv.key = r.kv.key \wedge w.kv.vc = r.kv.vc\}$ 
227       $co \triangleq TC(so \cup rf)$  causality order
228      IN  $\forall c \in Client :$ 
229           $\exists sc \in PermutationsOf(L[c] \circ SetToSeq(wops)) :$ 
230               $\wedge Valid(sc)$ 
231               $\wedge Respect(sc, co)$ 
233  THEOREM Spec  $\Rightarrow \Box CM$ 

```

234 |  
| \\* Modification History  
| \\* Last modified *Thu Jan 30 23:55:36 CST 2020* by *hengxin*  
| \\* Created *Tue Jan 28 21:28:25 CST 2020* by *hengxin*