

```

1  ┌────────────────── MODULE cosmos_client ───────────────────┐
    │ Microsoft Azure Cosmos DB TLA+ specification for the five consistency levels the service offers.
    │ The spec focuses on the consistency guarantees Cosmos DB provides to the clients, without the
    │ details of the protocol implementation.
    └────────────────────────────────────────────────────────────────┘

9  EXTENDS Naturals, Integers, Reals, Sequences, FiniteSets, TLC

    Number of regions
14  CONSTANT NumRegions
15  CONSTANT NumWriteRegions

17  ASSUME NumRegions ∈ Nat
18  ASSUME NumWriteRegions ≥ 1 ∧ NumWriteRegions ≤ NumRegions

    Number of clients per region for modeling
23  CONSTANT NumClientsPerRegion

25  ASSUME NumClientsPerRegion ∈ Nat

    MaxNumOp max number of operations from client
30  CONSTANT MaxNumOp

    Consistency level
    (1) strong (Linearizability)
    (2) bounded (Bounded Staleness)
    (3) session
    (4) prefix (Consistent Prefix)
    (5) eventual
40  CONSTANT Consistency

42  ASSUME Consistency ∈ { "strong", "bounded_staleness", "session", "consistent_prefix", "eventual" }

44  The bounded version differences in Bounded Staleness consistency
45  CONSTANT K

47  ASSUME K ∈ Nat

    All regions in topology
50  Regions ≜ 1 .. NumRegions
51  All writable regions in topology
52  WriteRegions ≜ 1 .. NumWriteRegions
53  All clients with local region
54  Clients ≜ {⟨r, j⟩ : r ∈ Regions, j ∈ 1 .. NumClientsPerRegion}

    All possible operations in history
59  Operations ≜ [type : { "write" }, data : Nat, region : WriteRegions, client : Clients]
60  ∪ [type : { "read" }, data : Nat, region : Regions, client : Clients]

63 --algorithm cosmos_client

```

```

64 {
65
66   variables   Max staleness. Strong is a special case of bounded with  $K = 1$ 
67   Bound = CASE Consistency = "strong"  $\rightarrow 1$ 
68              $\square$  Consistency = "bounded_staleness"  $\rightarrow K$ 
69              $\square$  Consistency = "session"  $\rightarrow MaxNumOp$ 
70              $\square$  Consistency = "consistent_prefix"  $\rightarrow MaxNumOp$ 
71              $\square$  Consistency = "eventual"  $\rightarrow MaxNumOp$ ;
72
73   Client operation history
74   History =  $\langle \rangle$ ;
75
76   Latest data value in each region
77   Data =  $[r \in Regions \mapsto 0]$ ;
78
79   Tentative log in each region
80   Database =  $[r \in Regions \mapsto \langle \rangle]$ ;
81
82   Value used by clients
83   value = 0;
84
85   define
86   {
87     Removing duplicates from a sequence:
88     RECURSIVE RemDupRec(–, –)
89     RemDupRec(es, seen)  $\triangleq$  IF es =  $\langle \rangle$  THEN  $\langle \rangle$ 
90                          ELSE IF es[1]  $\in$  seen THEN RemDupRec(Tail(es), seen)
91                          ELSE  $\langle es[1] \rangle \circ$  RemDupRec(Tail(es), seen  $\cup \{es[1]\}$ )
92
93     RemoveDuplicates(es)  $\triangleq$  RemDupRec(es,  $\{\}$ )
94
95     SetMax(S)  $\triangleq$  IF S =  $\{\}$  THEN  $-1$ 
96                  ELSE CHOOSE  $i \in S : \forall j \in S : i \geq j$ 
97
98     SeqToSet(s)  $\triangleq$   $\{s[i] : i \in \text{DOMAIN } s\}$ 
99
100    Last(s)  $\triangleq$  s[Len(s)]
101
102    MaxLen(c)  $\triangleq$  LET region  $\triangleq$  CHOOSE  $i \in Regions : \forall j \in Regions : \text{Len}(c[i]) \geq \text{Len}(c[j])$ 
103                IN   Len(c[region])
104
105    MinLen(c)  $\triangleq$  LET region  $\triangleq$  CHOOSE  $i \in Regions : \forall j \in Regions : \text{Len}(c[i]) \leq \text{Len}(c[j])$ 
106                IN   Len(c[region])
107   }
108
109   CLIENT ACTIONS
110
111

```

```

113 Regular write at local region
114 macro write( v )
115 {
116   if ( self[1] ∈ WriteRegions )
117   {
118     when  $\forall i, j \in \text{Regions} : \text{Data}[i] - \text{Data}[j] < \text{Bound}$  ;
119     Database[self[1]] := Append(@, v) ;
120     Data[self[1]] := v ;
121     History := Append(History, [type ↦ "write",
122                               data ↦ v,
123                               region ↦ self[1],
124                               client ↦ self]) ;
125     session_token := v ;
126   }
127 }

129 Reads with consistency checks
130 macro read( )
131 {
132   We check session token for session consistency
133   when Consistency ≠ "session" ∨ Data[self[1]] ≥ session_token ;
134   We check global value for strong consistency
135   when Consistency ≠ "strong" ∨  $\forall i, j \in \text{Regions} : \text{Data}[i] = \text{Data}[j]$  ;
136   History := Append(History, [type ↦ "read",
137                               data ↦ Data[self[1]],
138                               region ↦ self[1],
139                               client ↦ self]) ;
140   session_token := Data[self[1]] ;
141 }

143 REGION ACTIONS
144
145

147 Asynchronously replicates from source region to destination region and merges data history
148 macro replicate( )
149 {
150   with ( s ∈ WriteRegions ; d ∈ Regions )
151   {
152     Database[d] := RemoveDuplicates(SortSeq(Database[d] ∘ Database[s], < )) ;
153     if ( Len(Database[d]) > 0 )
154     {
155       Data[d] := Last(Database[d]) ;
156     }
157   }
158 }

```

```

160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206

```

CLIENT PROCESSES

```

fair process ( client ∈ Clients )
variable session_token = 0 ;
numOp = 0 ;
{
  client_actions:
  while ( numOp < MaxNumOp )
  {
    numOp := numOp + 1 ;
    either
    {
      write:
      value := value + 1 ;
      write(value) ;
    }
    or read: read() ;
  }
}

SERVER PROCESSES

fair process ( CosmosDB = ⟨0, 0⟩ )
{
  database_action:
  while ( TRUE )
  {
    replicate() ;
  }
}

}

BEGIN TRANSLATION
VARIABLES Bound, History, Data, Database, value, pc

define statement
RECURSIVE RemDupRec(-, -)
RemDupRec(es, seen)  $\triangleq$  IF es = ⟨⟩ THEN ⟨⟩
                        ELSE IF es[1] ∈ seen THEN RemDupRec(Tail(es), seen)
                        ELSE ⟨es[1]⟩ ∘ RemDupRec(Tail(es), seen ∪ {es[1]})

RemoveDuplicates(es)  $\triangleq$  RemDupRec(es, {})

SetMax(S)  $\triangleq$  IF S = {} THEN -1

```

```

207         ELSE CHOOSE  $i \in S : \forall j \in S : i \geq j$ 
209  $SeqToSet(s) \triangleq \{s[i] : i \in \text{DOMAIN } s\}$ 
211  $Last(s) \triangleq s[Len(s)]$ 
213  $MaxLen(c) \triangleq \text{LET } region \triangleq \text{CHOOSE } i \in Regions : \forall j \in Regions : Len(c[i]) \geq Len(c[j])$ 
214           IN  $Len(c[region])$ 
216  $MinLen(c) \triangleq \text{LET } region \triangleq \text{CHOOSE } i \in Regions : \forall j \in Regions : Len(c[i]) \leq Len(c[j])$ 
217           IN  $Len(c[region])$ 
219 VARIABLES  $session\_token, numOp$ 
221  $vars \triangleq \langle Bound, History, Data, Database, value, pc, session\_token, numOp \rangle$ 
223  $ProcSet \triangleq (Clients) \cup \{\langle 0, 0 \rangle\}$ 
225  $Init \triangleq$  Global variables
226    $\wedge Bound = (\text{CASE } Consistency = \text{"strong"} \rightarrow 1$ 
227      $\square Consistency = \text{"bounded\_staleness"} \rightarrow K$ 
228      $\square Consistency = \text{"session"} \rightarrow MaxNumOp$ 
229      $\square Consistency = \text{"consistent\_prefix"} \rightarrow MaxNumOp$ 
230      $\square Consistency = \text{"eventual"} \rightarrow MaxNumOp)$ 
231    $\wedge History = \langle \rangle$ 
232    $\wedge Data = [r \in Regions \mapsto 0]$ 
233    $\wedge Database = [r \in Regions \mapsto \langle \rangle]$ 
234    $\wedge value = 0$ 
235   Process client
236    $\wedge session\_token = [self \in Clients \mapsto 0]$ 
237    $\wedge numOp = [self \in Clients \mapsto 0]$ 
238    $\wedge pc = [self \in ProcSet \mapsto \text{CASE } self \in Clients \rightarrow \text{"client\_actions"}$ 
239      $\square self = \langle 0, 0 \rangle \rightarrow \text{"database\_action"}]$ 
241  $client\_actions(self) \triangleq \wedge pc[self] = \text{"client\_actions"}$ 
242    $\wedge \text{IF } numOp[self] < MaxNumOp$ 
243     THEN  $\wedge numOp' = [numOp \text{ EXCEPT } ![self] = numOp[self] + 1]$ 
244        $\wedge \vee \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"write"}]$ 
245        $\vee \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"read"}]$ 
246     ELSE  $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"Done"}]$ 
247        $\wedge numOp' = numOp$ 
248    $\wedge \text{UNCHANGED } \langle Bound, History, Data, Database, value,$ 
249      $session\_token \rangle$ 
251  $write(self) \triangleq \wedge pc[self] = \text{"write"}$ 
252    $\wedge value' = value + 1$ 
253    $\wedge \text{IF } self[1] \in WriteRegions$ 
254     THEN  $\wedge \forall i, j \in Regions : Data[i] - Data[j] < Bound$ 

```

```

255       $\wedge Database' = [Database \text{ EXCEPT } ![self[1]] = Append(@, value')]$ 
256       $\wedge Data' = [Data \text{ EXCEPT } ![self[1]] = value']$ 
257       $\wedge History' = Append(History, [type \mapsto \text{"write"},$ 
258                           $data \mapsto value',$ 
259                           $region \mapsto self[1],$ 
260                           $client \mapsto self])$ 
261       $\wedge session\_token' = [session\_token \text{ EXCEPT } ![self] = value']$ 
262      ELSE  $\wedge \text{TRUE}$ 
263       $\wedge \text{UNCHANGED } \langle History, Data, Database,$ 
264                           $session\_token \rangle$ 
265       $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"client\_actions"}]$ 
266       $\wedge \text{UNCHANGED } \langle Bound, numOp \rangle$ 

268   $read(self) \triangleq \wedge pc[self] = \text{"read"}$ 
269       $\wedge Consistency \neq \text{"session"} \vee Data[self[1]] \geq session\_token[self]$ 
270       $\wedge Consistency \neq \text{"strong"} \vee \forall i, j \in Regions : Data[i] = Data[j]$ 
271       $\wedge History' = Append(History, [type \mapsto \text{"read"},$ 
272                           $data \mapsto Data[self[1]],$ 
273                           $region \mapsto self[1],$ 
274                           $client \mapsto self])$ 
275       $\wedge session\_token' = [session\_token \text{ EXCEPT } ![self] = Data[self[1]]]$ 
276       $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"client\_actions"}]$ 
277       $\wedge \text{UNCHANGED } \langle Bound, Data, Database, value, numOp \rangle$ 

279   $client(self) \triangleq client\_actions(self) \vee write(self) \vee read(self)$ 

281   $database\_action \triangleq \wedge pc[\langle 0, 0 \rangle] = \text{"database\_action"}$ 
282       $\wedge \exists s \in WriteRegions :$ 
283       $\exists d \in Regions :$ 
284       $\wedge Database' = [Database \text{ EXCEPT } ![d] = RemoveDuplicates(SortSeq(Database[d] \circ$ 
285       $\wedge \text{IF } Len(Database'[d]) > 0$ 
286       $\text{ THEN } \wedge Data' = [Data \text{ EXCEPT } ![d] = Last(Database'[d])]$ 
287       $\text{ ELSE } \wedge \text{TRUE}$ 
288       $\wedge Data' = Data$ 
289       $\wedge pc' = [pc \text{ EXCEPT } ![\langle 0, 0 \rangle] = \text{"database\_action"}]$ 
290       $\wedge \text{UNCHANGED } \langle Bound, History, value, session\_token, numOp \rangle$ 

292   $CosmosDB \triangleq database\_action$ 

294   $Next \triangleq CosmosDB$ 
295       $\vee (\exists self \in Clients : client(self))$ 

297   $Spec \triangleq \wedge Init \wedge \Box [Next]_{vars}$ 
298       $\wedge \forall self \in Clients : WF_{vars}(client(self))$ 
299       $\wedge WF_{vars}(CosmosDB)$ 

301  END TRANSLATION

```

```

306 enable these invariants in model checker

308 Check elements in History are type of Opertion
309  $TypeOK \triangleq \{History[i] : i \in \text{DOMAIN } History\} \subseteq Operations$ 

311 Read value in any regional database history
312  $AnyReadPerRegion(r) \triangleq \forall i \in \text{DOMAIN } History : \wedge History[i].type = \text{"read"}$ 
313  $\wedge History[i].region = r$ 
314  $\Rightarrow History[i].data \in SeqToSet(Database[r]) \cup \{0\}$ 

316 Operation in history h is monotonic
317  $Monotonic(h) \triangleq \forall i, j \in \text{DOMAIN } h : i \leq j \Rightarrow h[i].data \leq h[j].data$ 

319 Reads in region r are monotonic
320  $MonotonicReadPerRegion(r) \triangleq \text{LET } reads \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"read"}$ 
321  $\wedge History[j].region = r\}$ 
322  $\mapsto History[i]]$ 
323 IN  $Monotonic(reads)$ 

325 Reads from client c are monotonic
326  $MonotonicReadPerClient(c) \triangleq \text{LET } reads \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"read"}$ 
327  $\wedge History[j].client = c\}$ 
328  $\mapsto History[i]]$ 
329 IN  $Monotonic(reads)$ 

331  $MonotonicWritePerRegion(r) \triangleq \text{LET } writes \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"write"}$ 
332  $\wedge History[j].region = r\}$ 
333  $\mapsto History[i]]$ 
334 IN  $Monotonic(writes)$ 

336 Clients read their own writes
337  $ReadYourWrite \triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$ 
338  $\wedge History[i].type = \text{"write"}$ 
339  $\wedge History[j].type = \text{"read"}$ 
340  $\wedge History[i].client = History[j].client$ 
341  $\Rightarrow History[j].data \geq History[i].data$ 

343 Read the latest writes
344  $ReadAfterWrite \triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$ 
345  $\wedge History[i].type = \text{"write"}$ 
346  $\wedge History[j].type = \text{"read"}$ 
347  $\Rightarrow History[j].data \geq History[i].data$ 

349  $Linearizability \triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$ 
350  $\Rightarrow History[j].data \geq History[i].data$ 

```

```
306 enable these invariants in model checker
```

308 Check elements in *History* are type of *Opertion*

309 $TypeOK \stackrel{\Delta}{=} \{History[i] : i \in \text{DOMAIN } History\} \subseteq Operations$

311	Read value in any regional database history
-----	---

312 $AnyReadPerRegion(r) \stackrel{\Delta}{=} \forall i \in \text{DOMAIN } History : \wedge History[i].type = \text{“read”}$

$$313 \quad \wedge \text{History}[i].\text{region} = r$$
$$\Rightarrow History[i].data \in SeqToSet(Database[r]) \cup \{0\}$$

316 Operation in history h is monotonic

$$317 \text{ Monotonic}(h) \triangleq \forall i, j \in \text{DOMAIN } h : i \leq j \Rightarrow h[i].data \leq h[j].data$$

319 Reads in region r are monotonic

320 $MonotonicReadPerRegion(r) \triangleq \text{LET } reads \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"read"}$
321 $\wedge History[j].region = r\}$

322 $\mapsto \text{History}[i]$

323 IN *Monotonic(reads)*

325 Reads from client c are monotonic

$$\text{MonotonicReadPerClient}(c) \stackrel{\Delta}{=} \text{LET } reads \stackrel{\Delta}{=} [i \in \{j \in \text{DOMAIN History} : \wedge \text{History}[j].\text{type} = \text{"read"} \\ \wedge \text{History}[j].\text{client} = c\}$$
328 $\mapsto \text{History}[i]$

329 IN *Monotonic(reads)*

331 $MonotonicWritePerRegion(r) \triangleq \text{LET } writes \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"write"}$
332 $\wedge History[j].region = r\}$

333 $\mapsto History[i]$

334 IN *Monotonic(writes)*

336 Clients read their own writes

$$337 \text{ ReadYourWrite} \triangleq \forall i, j \in \text{DOMAIN History} : \wedge i < j$$
338 $\wedge \text{History}[i].\text{type} = \text{"write"}$ 339 $\wedge History[j].type = \text{"read"}$
$$\wedge \text{History}[i].client = \text{History}[j].client$$
341 $\Rightarrow History[j].data \geq History[i].data$

343 Read the latest writes

$$344 \text{ ReadAfterWrite} \stackrel{\Delta}{=} \forall i, j \in \text{DOMAIN } History : \wedge i < j$$
345 $\wedge History[i].type = \text{"write"}$
$$\wedge \text{History}[j].type = \text{"read"}$$
$$\Rightarrow History[j].data \geq History[i].data$$

349 *Linearizability* $\triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$

$$\Rightarrow History[j].data \geq History[i].data$$

```

352 LastOperation(c)  $\triangleq$  LET i  $\triangleq$  SetMax( $\{j \in \text{DOMAIN } \textit{History} : \textit{History}[j].\textit{client} = c\}$ )
353      IN IF i > 0 THEN History[i] ELSE  $\langle \rangle$ 

356 BoundedStaleness  $\triangleq$   $\wedge \forall i, j \in \textit{Regions} : \textit{Data}[i] - \textit{Data}[j] \leq K$ 
357       $\wedge \forall r \in \textit{Regions} : \textit{MonotonicReadPerRegion}(r)$ 
358       $\wedge \textit{ReadYourWrite}$ 

360 ConsistentPrefix  $\triangleq$   $\forall r \in \textit{Regions} : \wedge \textit{MonotonicWritePerRegion}(r)$ 
361       $\wedge \textit{AnyReadPerRegion}(r)$ 

363 Strong  $\triangleq$   $\wedge \textit{Linearizability}$ 
364       $\wedge \textit{Monotonic}(\textit{History})$ 
365       $\wedge \textit{ReadAfterWrite}$ 

367 Session  $\triangleq$   $\wedge \forall c \in \textit{Clients} : \textit{MonotonicReadPerClient}(c)$ 
368       $\wedge \textit{ReadYourWrite}$ 

370 Eventual  $\triangleq$   $\forall i \in \text{DOMAIN } \textit{History} :$ 
371      LET r  $\triangleq$  History[i].region
372      IN History[i].data  $\in \{\textit{Database}[r][j] : j \in \text{DOMAIN } \textit{Database}[r]\} \cup \{0\}$ 

374 Invariant  $\triangleq$   $\wedge \textit{TypeOK}$ 
375       $\wedge$  CASE Consistency = "strong"  $\rightarrow$  Strong
376       $\square$  Consistency = "bounded_staleness"  $\rightarrow$  BoundedStaleness
377       $\square$  Consistency = "session"  $\rightarrow$  Session
378       $\square$  Consistency = "consistent_prefix"  $\rightarrow$  ConsistentPrefix
379       $\square$  Consistency = "eventual"  $\rightarrow$  Eventual

381 Liveness  $\triangleq$   $\diamond \square (\forall i, j \in \textit{Regions} : \textit{Database}[i] = \textit{Database}[j])$ 
383  $\square$ 

```