

TLA+ specifications of the consistency
guarantees provided by Cosmos DB

Murat Demirbas

[2018-11-01 Thu]

Outline

- 1** Modeling with TLA+
- 2** Specifying consistency
- 3** Single client specifications
- 4** General specifications

Devil is in details

Any large software system is prone to:

- corner cases
- failed assumptions
- race conditions
- cascading faults

Modeling saves you time

Model first and debug your design exhaustively

- catch the errors at design time, before it is too late
- achieve clarity of thinking
- communicate the design precisely
- see alternative ways to implement the design

TLA+ for modeling

TLA+ is a formal language for describing and reasoning about distributed and concurrent systems

It is developed by Dr. Leslie Lamport (Turing Award 2013)

It is based on mathematical logic, set theory, temporal logic, and is supported by tool set (TLC model checker)

Modeling with TLA+

The math-based formal language provides modeling at high-level, yet with precision

PlusCal provides ease-of-use and readability

The integrated model checker exhaustively debugs your model to the face of concurrency and failures, and produces counterexamples for your candidate invariants

Invariant-based reasoning

Instead of happy-path operational thinking, **invariant-based reasoning** focuses on *what needs to go right?*

It avoids the complexities & bugs of operational reasoning for concurrent systems

Invariants are specified as safety and liveness properties

Operationalizing TLA+

- (1) Model your system/protocol; figure out the **invariant**
- (2) Verify your model wrt to requirements
- (3) When you need to add a "feature" to the system:
 - first modify your model
 - model check, and debug design problems
 - implement the correct design

TLA+ is practical

Microsoft Azure Cosmos DB is a globally distributed database

The replication protocol is modeled in TLA+ and is tested for correctness against failures

Other places TLA+ used include

- multi write region modeling
- CRDT modeling
- lease protocol modeling

Specifying consistency

2

Specifying consistency

Consistency defines how the distributed database behaves when a client reads/writes while other clients write at different replicas

The tighter the consistency guarantees provided, the easier it gets to program/develop over the distributed database

Cosmos DB consistency

Cosmos DB provides 5 well-defined consistency properties to the clients:

- strong consistency
- bounded staleness
- session consistency
- consistent prefix
- eventual consistency

Cosmos DB consistency

Strong consistency provides linearizability since the reads are guaranteed to return the most recent version of an item

Bounded staleness consistency relaxes this a bit; the reads lag behind the writes by at most K prefixes or T interval

Session consistency provides monotonic reads, monotonic writes, read-your-own-writes, and write-follows-reads

When using consistent prefix reads, the updates returned are some prefix of all the updates, with no gaps

Eventually-consistent reads may return out of order reads, which stabilize to the most recent version only after all the writes quiesces

Single client specifications

2

Single client writing incremented counter values

```
50 process ( client ∈ Clients )
51 variables
52   m = ⟨ ⟩ ; op = 0 ; chistory = ⟨ 0 ⟩ ; ses = 1 ;
53 {
54   CW: while ( op < MaxNumOp ) {
55     op := op + 1 ;
56     send( Cloud, [type ↦ "Write", dat ↦ op, ses ↦ ses, orig ↦ self] ) ;
57     CWA: receive( m ) ; Ack
58     ses := m.ses ;
59     read
60     CR: send( Cloud, [type ↦ Consistency, ses ↦ ses, orig ↦ self] ) ;
61     CRA: receive( m ) ; Reply
62     chistory := Append( chistory, m.dat ) ;
63     ses := m.ses ;
64   }
65 }
```

Macros for sending and receiving messages

```
1  ┌────────────────────────── MODULE suscop ───────────────────────────┐
2  EXTENDS Naturals, Integers, Sequences, FiniteSets, TLC, Bags
3  CONSTANT NumClients, MaxNumOp, Consistency, K
4  ASSUME Consistency ∈ { "Eventual", "Consistent_Prefix", "Session", "Bounded_Staleness", "Strong" }
5  ASSUME MaxNumOp < 10 ∧ NumClients = 1
6  Cloud ≜ 0
7  Clients ≜ 1 .. NumClients
9  --algorithm suscop{
10 variables
11   chan = [ n ∈ 0 .. NumClients ↦ {} ]; FIFO channels
13   network functions
14   macro send( des, msg ) {
15     chan[des] := Append(chan[des], msg);
16   }
18   macro receive( msg ) {
19     await Len(chan[self]) > 0;
20     msg := Head(chan[self]);
21     chan[self] := Tail(chan[self]);
22   }
```

Cosmos DB as a single process

```
24 process ( cosmosdb ∈ {Cloud} )
25 variables
26   Database = ⟨0⟩; msg = ⟨⟩;
27 { D: while ( TRUE ) {
28   receive(msg);
29   if ( msg.type = "Write" ) {
30     Database := Append(Database, msg.dat);
31   DW: send(msg.orig, [type ↦ "Ack", dat ↦ Database[Len(Database)], ses ↦ Len(Database)]); }
32   else if ( msg.type = "Eventual" )
33   DE: with ( k ∈ 1 .. Len(Database) )
34     send(msg.orig, [type ↦ "Reply", dat ↦ Database[k], ses ↦ k]);
35   else if ( msg.type = "Consistent_Prefix" )
36   DP: with ( k ∈ 1 .. Len(Database) )
37     send(msg.orig, [type ↦ "Reply", dat ↦ Database[k], ses ↦ k]);
38   else if ( msg.type = "Session" )
39   DS: with ( k ∈ msg.ses .. Len(Database) )
40     send(msg.orig, [type ↦ "Reply", dat ↦ Database[k], ses ↦ k]);
41   else if ( msg.type = "Bounded_Staleness" )
42   DB: with ( k ∈ (1F Len(Database) > K THEN Len(Database) - K ELSE 1) .. Len(Database) )
43     send(msg.orig, [type ↦ "Reply", dat ↦ Database[k], ses ↦ k]);
44   else if ( msg.type = "Strong" )
45   DG: with ( k = Len(Database) )
46     send(msg.orig, [type ↦ "Reply", dat ↦ Database[k], ses ↦ k]);
47 }
```

Model checking

```
192 Invariants for single client (ID = 1) writing with op ++  
193 Eventual  $\triangleq$   $chistory[1][Len(chistory[1])] \in \{Database[Cloud][i] : i \in 1 .. Len(Database[Cloud])\}$   
195 Consistent_Prefix  $\triangleq$   $chistory[1][Len(chistory[1])] \in \{Database[Cloud][i] : i \in 1 .. Len(Database[Cloud])\}$   
197 Session  $\triangleq$   $pc[1] = "CW" \Rightarrow chistory[1][Len(chistory[1])] \in \{Database[Cloud][i] :$   
198  $i \in ses[1] .. Len(Database[Cloud])\}$   
200 Bounded_Staleness  $\triangleq$   $pc[1] = "CW" \Rightarrow chistory[1][Len(chistory[1])] \in \{Database[Cloud][i] :$   
201  $i \in (IF Len(Database[Cloud]) > K THEN Len(Database[Cloud]) - K ELSE 1) .. Len(Database[Cloud])\}$   
203 Strong  $\triangleq$   $pc[1] = "CW" \Rightarrow chistory[1][Len(chistory[1])] = Database[Cloud][Len(Database[Cloud])]$   
205
```

Macros for sending and receiving messages

```
1  ┌────────────────────────── MODULE swscop ───────────────────────────┐
2  EXTENDS Naturals, Integers, Sequences, FiniteSets, TLC, Bags
3  CONSTANT NumClients, MaxNumOp, Consistency, K
4  ASSUME Consistency ∈ { "Eventual", "Consistent_Prefix", "Session", "Bounded_Staleness", "Strong" }
5  ASSUME MaxNumOp < 10 ∧ NumClients = 1
6  Cloud ≜ 0
7  Clients ≜ 1 .. NumClients
9  --algorithm swscop{
10 variables
11   chan = [ n ∈ 0 .. NumClients ↦ {} ]; FIFO channels
13   network functions
14   macro send( des, msg ) {
15     chan[des] := Append(chan[des], msg);
16   }
18   macro receive( msg ) {
19     await Len(chan[self]) > 0;
20     msg := Head(chan[self]);
21     chan[self] := Tail(chan[self]);
22   }
```

Single client writing incremented counter values

```
50 process ( client ∈ Clients )
51 variables
52   m = ⟨ ⟩ ; op = 0 ; chistory = ⟨ 0 ⟩ ; ses = 1 ;
53 {
54   CW: while ( op < MaxNumOp ) {
55     op := op + 1 ;
56     send( Cloud, [type ↦ "Write", dat ↦ op, ses ↦ ses, orig ↦ self] ) ;
57     CWA: receive( m ) ; Ack
58     ses := m.ses ;
59     read
60     CR: send( Cloud, [type ↦ Consistency, ses ↦ ses, orig ↦ self] ) ;
61     CRA: receive( m ) ; Reply
62     chistory := Append( chistory, m.dat ) ;
63     ses := m.ses ;
64   }
65 }
```

Single client reading, incrementing, writing-back counter values

In this scenario, the client reads **value** from the database (with one of the 5 consistency levels configured), increments it, and writes it back

The **value** stored in the database become consecutive only using the session and strong consistency reads

The client model

I

The database model is unchanged; the client is modified

```
50 process ( client ∈ Clients )
51 variables
52   m = ⟨ ⟩ ; op = 0 ; v = 0 ; chistory = ⟨ 0 ⟩ ; ses = 1 ;
53 {
54   CR: while ( op < MaxNumOp ) {
55     send(Cloud, [type ↦ Consistency, ses ↦ ses, orig ↦ self]); read
56   CRA: receive(m); Reply
57     chistory := Append(chistory, m.dat);
58     v := m.dat ;
59     ses := m.ses ;
60     write v + 1
61   CW: send(Cloud, [type ↦ "Write", dat ↦ v + 1, ses ↦ ses, orig ↦ self]);
62   CWA: receive(m) ; Ack
63     ses := m.ses ;
64     op := op + 1 ;
65   }
66 }
```

Strong consistency trace

I

```
> <Action line 150, col 13 to line 156, col 69 of... | State (num = 58)
/\ Database = {0 :> <<0, 1, 2, 3, 4, 5, 6, 7>>}
/\ chan = {0 :> <<>> @@ 1 :> <<>>}
/\ chistory = <<<<0, 0, 1, 2, 3, 4, 5, 6>>>>
/\ m = <<[ses |-> 8, type |-> "Ack", dat |-> 7]>>
/\ msg = {0 :> [ses |-> 7, type |-> "Write", dat |-> 7, orig |-> 1]}
/\ op = <<7>>
/\ pc = {0 :> "D" @@ 1 :> "Done"}
/\ ses = <<8>>
/\ v = <<6>>
```

Bounded consistency trace

I

▶ <Action line 150, col 13 to line 156, col 69 of... State (num = 58)

```
/\ Database = (0 :> <<0, 1, 1, 1, 2, 2, 2, 3>>)
/\ chan = (0 :> <<>> @@ 1 :> <<>>)
/\ chistory = <<<<0, 0, 0, 0, 1, 1, 1, 2>>>>
/\ m = <<[ses |-> 8, type |-> "Ack", dat |-> 3]>>
/\ msg = (0 :> [ses |-> 5, type |-> "Write", dat |-> 3, orig |-> 1])
/\ op = <<7>>
/\ pc = (0 :> "D" @@ 1 :> "Done")
/\ ses = <<8>>
/\ v = <<2>>
```

General specifications

I

To implement a linearizability checker, the clients share a **backchannel**: the counter **value** (or a real time clock)

Multiple clients across different regions put Cosmos DB to test with respect to the consistency guarantees it provides in the presence of concurrent read and writes

For brevity we use shared memory model instead of message passing

Database model exposing regions

I

The clients perform their writes and reads to their respective local regions

The writes are appended to the Database [r] of the corresponding region r

The Cosmos DB servers perform *anti-entropy* across the regions and merge their Database respecting the happened-before relationship to implement last writer wins

The client

I

```
142  fair process ( client ∈ Clients )
143  variable session_token = 0;
144  numOp = 0;
145  {
146    client_actions:
147    while ( numOp < MaxNumOp )
148    {
149      numOp := numOp + 1;
150      either
151      {
152        write:
153        value := value + 1;
154        write(value);
155      }
156      or read: read();
157    }
```

The client write operation

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

The client read operation

```
128 macro read( )
129 {
130     We check session token for session consistency
131     when Consistency  $\neq$  "session"  $\vee$  Last(Database[self[1]])  $\geq$  session_token ;
132     We check global value for strong consistency
133     when Consistency  $\neq$  "strong"  $\vee \forall i, j \in$  Regions : Last(Database[i]) = Last(Database[j]) ;
134     History := Append(History, [type  $\mapsto$  "read",
135                             data  $\mapsto$  Last(Database[self[1]]),
136                             region  $\mapsto$  self[1],
137                             client  $\mapsto$  self]) ;
138     session_token := Last(Database[self[1]]) ;
139 }
```

The database

```
164  fair process ( CosmosDB ∈ Servers )
165  {
166    database_action:
167    while ( TRUE )
168    {
169      with ( s ∈ WriteRegions ) {
170        Database[self[1]] := RemoveDuplicates(SortSeq(Database[self[1]] ◦ Database[s], < ));
171      }
172    }
173  }
```

General specifications

To implement a linearizability checker, the clients share a **backchannel**: the counter **value** (or a real time clock)

Multiple clients across different regions put Cosmos DB to test with respect to the consistency guarantees it provides in the presence of concurrent read and writes

For brevity we use shared memory model instead of message passing

Invariant

```
309 Operation in history  $h$  is monotonic  
310  $Monotonic(h) \triangleq \forall i, j \in \text{DOMAIN } h : i \leq j \Rightarrow h[i].data \leq h[j].data$   
  
312 Reads in region  $r$  are monotonic  
313  $MonotonicReadPerRegion(r) \triangleq \text{LET reads} \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"read"}$   
314  $\wedge History[j].region = r\}$   
315  $\mapsto History[i]]$   
316  $\text{IN } Monotonic(reads)$   
  
318 Reads from client  $c$  are monotonic  
319  $MonotonicReadPerClient(c) \triangleq \text{LET reads} \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"read"}$   
320  $\wedge History[j].client = c\}$   
321  $\mapsto History[i]]$   
322  $\text{IN } Monotonic(reads)$   
  
324  $MonotonicWritePerRegion(r) \triangleq \text{LET writes} \triangleq [i \in \{j \in \text{DOMAIN } History : \wedge History[j].type = \text{"write"}$   
325  $\wedge History[j].region = r\}$   
326  $\mapsto History[i]]$   
327  $\text{IN } Monotonic(writes)$   
  
329 Clients read their own writes  
330  $ReadYourWrite \triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$   
331  $\wedge History[i].type = \text{"write"}$   
332  $\wedge History[j].type = \text{"read"}$   
333  $\wedge History[i].client = History[j].client$   
334  $\Rightarrow History[j].data \geq History[i].data$   
  
336 Read the latest writes  
337  $ReadAfterWrite \triangleq \forall i, j \in \text{DOMAIN } History : \wedge i < j$   
338  $\wedge History[i].type = \text{"write"}$   
339  $\wedge History[j].type = \text{"read"}$   
340  $\Rightarrow History[j].data \geq History[i].data$ 
```

Invariant

366 *Linearizability* $\triangleq \forall i, j \in \text{DOMAIN History} : \wedge i < j$
367 $\Rightarrow \text{History}[j].\text{data} \geq \text{History}[i].\text{data}$

369 *BoundedStaleness* $\triangleq \wedge \forall i, j \in \text{Regions} : \text{Last}(\text{Database}[i]) - \text{Last}(\text{Database}[j]) \leq K$
370 $\wedge \forall r \in \text{Regions} : \text{MonotonicReadPerRegion}(r)$
371 $\wedge \text{ReadYourWrite}$

373 *ConsistentPrefix* $\triangleq \forall r \in \text{Regions} : \wedge \text{MonotonicWritePerRegion}(r)$
374 $\wedge \text{AnyReadPerRegion}(r)$

376 *Strong* $\triangleq \wedge \text{Linearizability}$
377 $\wedge \text{Monotonic}(\text{History})$
378 $\wedge \text{ReadAfterWrite}$

380 *Session* $\triangleq \wedge \forall c \in \text{Clients} : \text{MonotonicReadPerClient}(c)$
381 $\wedge \text{ReadYourWrite}$

383 *Eventual* $\triangleq \forall i \in \text{DOMAIN History} :$
384 $\text{LET } r \triangleq \text{History}[i].\text{region}$
385 $\text{IN } \text{History}[i].\text{data} \in \{\text{Database}[r][j] : j \in \text{DOMAIN Database}[r]\} \cup \{0\}$

Links

`https://docs.microsoft.com/en-us/azure/cosmos-db/consistency-levels`

`https://github.com/Azure/azure-cosmos-tla`

`https://lamport.azurewebsites.net/tla/tla.html`

`http://learntla.com`

```

cosmos_client.tla  Model_1
/Users/mad/TLA-DEV/ignite/cosmos_client.tla

TLA Module
308         i-> History[i]]
309         IN Monotonic(reads)
310
311 MonotonicWritePerRegion(r) == LET writes == [i \in {j \in DOMAIN History : ^ History[j].type = "write"
312                                     ^ History[j].region = r}
313         i-> History[i]]
314         IN Monotonic(writes)
315
316 (* Clients read their own writes *)
317 ReadYourWrite == \A i, j \in DOMAIN History : ^ i < j
318             ^ History[i].type = "write"
319             ^ History[j].type = "read"
320             ^ History[i].client = History[j].client
321             => History[j].data >= History[i].data
322
323 (* Read the latest writes *)
324 ReadAfterWrite == \A i, j \in DOMAIN History : ^ i < j
325             ^ History[i].type = "write"
326             ^ History[j].type = "read"
327             => History[j].data >= History[i].data
328
329
330 (*          HIGHER LEVEL INVARIANTS          *)
331
332 Linearizability == \A i, j \in DOMAIN History : ^ i < j
333             => History[j].data >= History[i].data
334
335 BoundedStaleness == ^ \A i, j \in Regions : Last(Database[i]) - Last(Database[j]) <= K
336             ^ \A r \in Regions : MonotonicReadPerRegion(r)
337             ^ ReadYourWrite
338
339 ConsistentPrefix == \A r \in Regions : ^ MonotonicWritePerRegion(r)
340             ^ AnyReadPerRegion(r)
341
342 Strong == ^ Linearizability
343             ^ Monotonic(History)
344             ^ ReadAfterWrite
345
346 Session == ^ \A c \in Clients : MonotonicReadPerClient(c)
347             ^ ReadYourWrite
348
349 Eventual == \A i \in DOMAIN History :
350     LET r == History[i].region
351     IN History[i].data \in {Database[r][j] : j \in DOMAIN Database[r]} \union {0}
352
353 Invariant == ^ TypeOK
354             ^ CASE Consistency = "strong" -> Strong
355                 ^ Consistency = "bounded_staleness" -> BoundedStaleness
356                 ^ Consistency = "session" -> Session
357                 ^ Consistency = "consistent_prefix" -> ConsistentPrefix
358                 ^ Consistency = "eventual" -> Eventual
359
360 Liveness == << \A i, j \in Regions : Database[i] = Database[j]
361
362

```

TLC Errors

Model_1

Invariant Strong is violated.

Error-Trace Exploration

Error-Trace

Name	Value
<Initial predicate>	State (num = 1)
<Action line 222, col 25 to line 230, col 56 of module...>	State (num = 2)
<Action line 222, col 25 to line 230, col 56 of module...>	State (num = 3)
<Action line 232, col 16 to line 245, col 46 of module...>	State (num = 4)
<Action line 247, col 15 to line 256, col 62 of module...>	State (num = 5)

```

/\ Bound = 5
/\ Database = <<<<>>, <<1>>, <<>>>>
/\ History = << [type i-> "write", data i-> 1, region i-> 2, client i-> <<2, 1>>],
  [type i-> "read", data i-> 0, region i-> 1, client i-> <<1, 1>>] >>
/\ numOp = {<<1, 1>> i> 1 @@ <<2, 1>> i> 1 @@ <<3, 1>> i> 0}
/\ pc = { <<1, 0>> i> "database_action" @@
  <<1, 1>> i> "client_actions" @@
  <<2, 0>> i> "database_action" @@
  <<2, 1>> i> "client_actions" @@
  <<3, 0>> i> "database_action" @@
  <<3, 1>> i> "client_actions" }
/\ session_token = {<<1, 1>> i> 0 @@ <<2, 1>> i> 1 @@ <<3, 1>> i> 0}
/\ value = 1

```

360 : 3

135M of 329M

Spec Status: passed