

Two Phase Commit Protocol

Murat Demirbas

[2019-08-26 Mon]

2 phase commit

A transaction is performed over **resource managers (RMs)**

The **transaction manager (TM)** finalizes the transaction

- For the transaction to be committed, each participating RM must be prepared to commit it
- Otherwise, the transaction must be aborted

Definitions

```
8 --fair algorithm 2PC{
9   variable  $rmState = [rm \in RM \mapsto \text{"working"}]$ ,
10           $tmState = \text{"init"} ;$ 

12  define {
13     $canCommit \triangleq \forall rm \in RM : rmState[rm] \in \{\text{"prepared"}\}$ 
14     $canAbort \triangleq \exists rm \in RM : rmState[rm] \in \{\text{"aborted"}\}$ 
15  }
```

TM modeling

```
37  fair process ( TManager = 0 ) {
38  TM: either
39      { await canCommit ;
40        tmState := "commit" ;
41  F1:   if ( TMMAYFAIL ) tmState := "unavailable" ; }
42      or
43      { await canAbort ;
44        tmState := "abort" ;
45  F2:   if ( TMMAYFAIL ) tmState := "unavailable" ; }
46      }
```

TM modeling

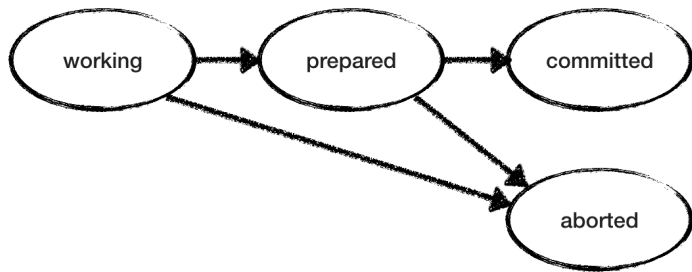
TM checks if it `canCommit` or `canAbort` and updates `tmState` accordingly.

TM can also fail making `tmState` "unavailable"

To keep things simple yet interesting, TM fails only after it makes a decision. These two updates are nonatomic:

- `tmState` is available for RMs to read for a duration
- labels at fail actions provide nonatomicity

RM modeling



RM modeling

```
20  fair process ( RManager ∈ RM ) {
21  RS: while ( rmState[self] ∉ { "committed", "aborted" } ) {
22      either {
23          await rmState[self] = "working" ;
24          with ( x ∈ { "prepared", "aborted" } ) rmState[self] := x ; }
25      or {
26          await rmState[self] = "prepared" ∧ tmState = "commit" ;
27  RC:   rmState[self] := "committed" ; }
28      or {
29          await rmState[self] = "prepared" ∧ tmState = "abort" ;
30  RA:   rmState[self] := "aborted" ; }
31      }
32  }
```

Invariants

175 *Consistency* \triangleq

A state predicate asserting that two *RM*s have not arrived at conflicting decisions.

180 $\forall rm1, rm2 \in RM : \neg \wedge rmState[rm1] = \text{"aborted"}$

181 $\wedge rmState[rm2] = \text{"committed"}$

183 *Completed* $\triangleq \diamond (\forall rm \in RM : rmState[rm] \in \{ \text{"committed"}, \text{"aborted"} \})$

Consistency checks that there are no 2 *RM*s such that one says "committed" and other says "aborted"

Model checking

If TM does not fail, the 2-phase commit algorithm is correct

When TM fails, **termination** can be violated

We add a Backup TM, to take over, and achieve **termination**

BTM modeling

```
50  fair process ( BTManager = 10 ) {  
51  BTM: either  
52      { await canCommit  $\wedge$  tmState = "unavailable" ;  
53  BTC:   tmState := "commit" ; }  
54      or  
55      { await canAbort  $\wedge$  tmState = "unavailable" ;  
56  BTA:   tmState := "abort" ; }  
57  }
```

Strengthening canCommit

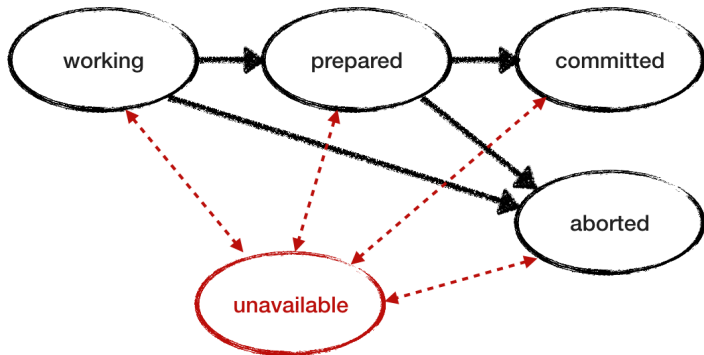
```
8 --fair algorithm 2PC{
9   variable rmState = [rm ∈ RM ↦ "working"],
10      tmState = "init" ;
11
12   define {
13     canCommit ≜    ∀ rm ∈ RM : rmState[rm] ∈ {"prepared"}
14                  ∨   ∃ rmc ∈ RM : rmState[rmc] ∈ {"committed"} for BTM
15     canAbort ≜    ∃ rm ∈ RM : rmState[rm] ∈ {"aborted"}
16   }
```

BTM modeling

BTM takes over when TM is unavailable and uses the same logic to make decisions

For simplicity we assume the BTM does not fail

What if RMs could also fail?



What if RMs could also fail?

```
20  fair process ( RManager ∈ RM )
21  variables pre = "" ; {
22  RS: while ( rmState[self] ∉ { "committed", "aborted" } ) {
23      either {
24          await rmState[self] = "working" ;
25          with ( x ∈ { "prepared", "aborted" } ) rmState[self] := x ; }
26      or {
27          await rmState[self] = "prepared" ∧ tmState = "commit" ;
28  RC:   rmState[self] := "committed" ; }
29      or {
30          await rmState[self] = "prepared" ∧ tmState = "abort" ;
31  RA:   rmState[self] := "aborted" ; }
32      or {
33          await RM MAY FAIL ∧ pre ≠ rmState[self] ;
34          pre := rmState[self] ;
35          rmState[self] := "unavailable" ;
36  RR:   rmState[self] := pre ;
37          }
38      }
39  }
```

Strengthening canAbort

```
12 define {  
13   canCommit  $\triangleq$   $\forall rmc \in RM : rmState[rmc] \in \{\text{"prepared"}\}$   
14    $\vee \exists rm \in RM : rmState[rm] \in \{\text{"committed"}\}$  for BTM  
15   canAbort  $\triangleq$   $\exists rm \in RM : rmState[rm] \in \{\text{"aborted"}, \text{"unavailable"}\}$   
16    $\wedge \neg \exists rmc \in RM : rmState[rmc] = \text{"committed"}$  for BTM  
17 }
```

Inconsistency problem!

- ▼ ▲ <RS line 90, col 13 to line 111, col 32 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <F1 line 139, col 7 to line 145, col 37 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <RS line 90, col 13 to line 111, col 32 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <BTM line 157, col 8 to line 162, col 47 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <BTA line 169, col 8 to line 172, col 38 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <RC line 113, col 13 to line 116, col 43 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <RS line 90, col 13 to line 111, col 32 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState
 - ▶ tmState
- ▼ ▲ <RA line 118, col 13 to line 121, col 43 of module 2PC2withBTM>
 - ▶ pc
 - ▶ pre
 - ▶ rmState

```
State (num = 6)
(0 :> "F1" @@ 1 :> "RC" @@ 2 :> "RS" @@ 3 :> "RS" @@ 10 :> "BTM")
<<"", "", "">>
<<"prepared", "prepared", "prepared">>
"commit"
State (num = 7)
(0 :> "Done" @@ 1 :> "RC" @@ 2 :> "RS" @@ 3 :> "RS" @@ 10 :> "BTM")
<<"", "", "">>
<<"prepared", "prepared", "prepared">>
"unavailable"
State (num = 8)
(0 :> "Done" @@ 1 :> "RC" @@ 2 :> "RR" @@ 3 :> "RS" @@ 10 :> "BTM")
<<"", "prepared", "">>
<<"prepared", "unavailable", "prepared">>
"unavailable"
State (num = 9)
(0 :> "Done" @@ 1 :> "RC" @@ 2 :> "RR" @@ 3 :> "RS" @@ 10 :> "BTA")
<<"", "prepared", "">>
<<"prepared", "unavailable", "prepared">>
"unavailable"
State (num = 10)
(0 :> "Done" @@ 1 :> "RC" @@ 2 :> "RR" @@ 3 :> "RS" @@ 10 :> "Done")
<<"", "prepared", "">>
<<"prepared", "unavailable", "prepared">>
"abort"
State (num = 11)
(0 :> "Done" @@ 1 :> "RC" @@ 2 :> "RR" @@ 3 :> "RS" @@ 10 :> "Done")
<<"", "prepared", "">>
<<"committed", "unavailable", "prepared">>
"abort"
State (num = 12)
(0 :> "Done" @@ 1 :> "RS" @@ 2 :> "RR" @@ 3 :> "RA" @@ 10 :> "Done")
<<"", "prepared", "">>
<<"committed", "unavailable", "prepared">>
"abort"
State (num = 13)
(0 :> "Done" @@ 1 :> "RS" @@ 2 :> "RR" @@ 3 :> "RS" @@ 10 :> "Done")
<<"", "prepared", "">>
<<"committed", "unavailable", "aborted">>
```


What went on?

- RM1 sees commit from TM
- TM becomes unavailable
- RM2 becomes unavailable
- BTM takes over for TM
- BTM decides on abort seeing <prepared, unavailable, prepared> from RMs. (It may also be that RM1 may also look unavailable due to unreachability)
- RM1 acts on commit from TM
- RM3 sees abort from BTM
- RM3 acts on abort from BTM
- Consistency is violated

Inconsistency problem!

If BTM decides, it may decide incorrectly, violating **consistency**

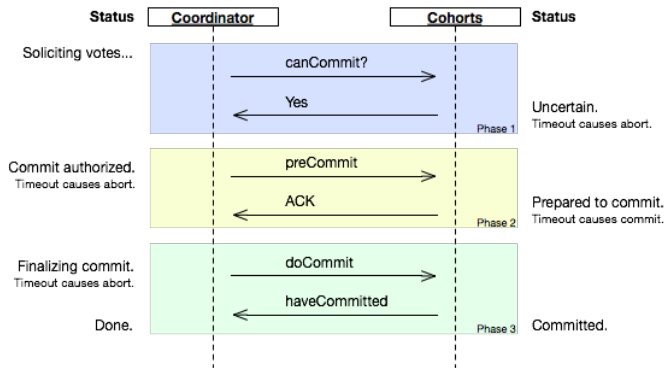
If BTM waits, it may be waiting forever on a crashed node, and violating **termination** (i.e., **progress**)

Timeouts may be incorrect, due to inopportune timing

FLP impossibility

In an **asynchronous** system, it is impossible to solve consensus (both safety and progress) in the presence of crash faults

What about 3PC?



3PC problems

BTM may act as if TM is down

RMs go with whatever TM or BTM says (2 leaders)

This asymmetry of information is the root of all evil in distributed systems

Lamport & Gray on transaction commit

"Several 3-phase protocols have been proposed, and a few have been implemented. They have usually attempted to "fix" the 2-Phase Commit protocol by choosing another TM if the first TM fails. However, we know of none that provides a complete algorithm proven to satisfy a clearly stated correctness condition. For example, the discussion of non-blocking commit in the classic text of Bernstein, Hadzilacos, and Goodman fails to explain what a process should do if it receives messages from two different processes, both claiming to be the current TM. Guaranteeing that this situation cannot arise is a problem that is as difficult as implementing a transaction commit protocol."

Paxos!

Paxos is always **safe** even in presence of inaccurate failure detectors, asynchronous execution, faults, and eventually makes progress when consensus gets in the realm of possibility

Paxos makes **progress** when the system is outside the realm of consensus impossibility

You can emulate TM with a Paxos cluster of 3 nodes and solve the inconsistent TM/BTM problem (Google Spanner approach)