



The Accord Consensus Protocol

Fast General Purpose Transactions

Elliott Smith, Benedict
Miller, Alex

Zhang, Tony
Capwell, David

Eggleson, Blake Andreas, Scott
Yeschenko, Aleksey

Dr. Howard, Heidi
Weisberg, Ariel

True ACID Compliance
Transactions
Resilience & Linear scalability
Performance
Accessible

Novel

- Global Strict Serialisability
- Multi-partition, Multi-table
- Leaderless, Unique Recovery (No rollbacks)
- Single Network Round-trip
- Commodity Clocks
- Standalone library (not tied to C*)
- Flexible Fast-Path Electorates
- Timestamp Reorder Buffer



Algorithm 1 Consensus Protocol

```

receive t on coordinator C from client:
1:  $t_c \leftarrow (\text{now}, 0, C)$ 
2: send PreAccept(t, tc) to  $\forall p \in \mathbb{P}^t$ 
receive PreAccept(t, tc) on p:
3: if  $t_c > \max(t_p \mid \gamma \sim t)$  then
4:    $t_c \leftarrow t$ 
5: else
6:    $t_c \leftarrow \max(T_p \mid \gamma \sim t)$ 
7:    $t_c(\text{seq}, id) \leftarrow (t_c, \text{seq} + 1, p)$ 
8: end if
9:  $T_c \leftarrow \{t_c\}$ 
10:  $Accepted_c \leftarrow \text{true}$ 
11: reply PreAcceptOK(tc, deps) :  $\{\gamma \mid \gamma \sim t \wedge t_{\gamma} < t_c\}$ 
receive PreAcceptOK(tc, deps) from  $\mathbb{Q}^t$ :
12:  $deps = \bigcup \{p.deps \mid p \in \mathbb{Q}^t\}$ 
13: if  $\exists \mathbb{P}^t \subseteq \mathbb{Q}^t \mid (p \in \mathbb{P}^t \wedge p = t_c)$  then
14:   send Commit(tc, tc, deps) to  $\forall p \in \mathbb{P}^t$ 
15:   go to Execution Protocol
16: else
17:    $t = \max(p.t \mid p \in \mathbb{Q}^t)$ 
18:   send Accept(t, tc, t, deps) to  $\forall p \in \mathbb{P}^t$ 
19: end if
20: reply PreAcceptOK(tc, deps) :  $\{\gamma \mid \gamma \sim t \wedge t_{\gamma} < t_c\}$ 
receive PreAcceptOK(tc, deps) from  $\mathbb{Q}^t$ :
21:  $deps = \bigcup \{p.deps \mid p \in \mathbb{Q}^t\}$ 
22: if  $\exists \mathbb{P}^t \subseteq \mathbb{Q}^t \mid (p \in \mathbb{P}^t \wedge p = t_c)$  then
23:   send Commit(tc, tc, deps) to  $\forall p \in \mathbb{P}^t$ 
24:   go to Execution Protocol
25: end if
26:  $t_c \leftarrow \max(t, t_c) \leq t$ 
27: if  $\neg Accepted_c \wedge \neg Applied_c$  then
28:   reply NACK
29: else if  $Committed_c \vee Applied_c$  then
30:   return
31: else
32:    $MaxBallot_c \leftarrow b$ 
33:    $Accepted_c \leftarrow b$ 
34:    $T_c \leftarrow \max(t_c, T_c)$ 
35:    $Accepted_c \leftarrow \text{true}$ 
36:   reply AcceptOK(deps) :  $\{\gamma \mid \gamma \sim t \wedge t_{\gamma} < t_c\}$ 
37: end if
receive NACK on C:
38: yield to competing coordinator
receive AcceptOK(deps) on C from  $\mathbb{Q}^t$ :
39:  $deps = \bigcup \{p.deps \mid p \in \mathbb{Q}^t\}$ 
40: send Commit(tc, tc, deps) to  $\forall p \in \mathbb{P}^t \cup \mathbb{Q}^t$ 
41: go to Execution Protocol

```

Algorithm 2 Execution Protocol (with reconfiguration and ballots)

```

Coordinator C:
42: for  $p \in \mathbb{P}^t$  do
43:    $deps_p \leftarrow \{\gamma \mid \gamma \in deps \wedge p \in \mathbb{P}^t\}$ 
44:   await Read(tc, tc, depsp) to some nearby correct  $p \in \mathbb{P}^t$ 
45: end for
receive Commit(tc, tc, tc, depsp):
46:  $Committed_c \leftarrow \text{true}$ 
receive Read(tc, tc, depsp) on p:
47: await Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
48: reads  $\leftarrow \text{read}(\gamma)$ 
49: reply ReadOK(reads)
receive ReadOK(reads) from each shard:
50:  $read \leftarrow \text{read}(\gamma, \text{read})$ 
51:  $read \leftarrow \text{read}(t_c, \text{deps}_p, \text{read})$  to  $\forall p \in \mathbb{P}^t$ 
52: send read to client
receive Apply(tc, tc, depsp, result) :
53: if  $\neg Applied_c$  then
54:   await Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
55:   apply writes(tc)
56:   wait Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
57:   apply writes(tc)
58:   Applied_c  $\leftarrow \text{true}$ 
59: end if

```

Algorithm 3 Recovery Protocol

```

Coordinator C:
1:  $b \leftarrow \text{fresh ballot}$ 
2: send Recover(tc, tc) to  $\forall p \in \mathbb{P}^t$ 
receive Recover(tc, tc) on p:
3: if  $\neg MaxBallot_c$  then
4:   reply NACK(b)
5: else
6:   MaxBallot_c  $\leftarrow \max(t_c \mid t_c \leq t \wedge t_c \leq t_{\gamma} \wedge t_{\gamma} \geq t_c)$ 
7:   Commit_c  $\leftarrow \{\gamma \mid \gamma \sim t_c \wedge t_c \leq t_{\gamma} \wedge t_{\gamma} \geq t_c\}$ 
8:   Wait_c  $\leftarrow \{\gamma \in \text{Accepts} \mid t_{\gamma} \leq t_c \wedge t_{\gamma} \geq t_c\}$ 
9:   Superseding_c  $\leftarrow \{\gamma \in \text{Accepts} \mid t_{\gamma} > t_c\}$ 
10:  end if
11: if  $\neg PreAccepted_c$  then
12:   run Consensus Protocol 1|10
13:  end if
14: if  $\neg Accepted_c \wedge \neg Committed_c \wedge \neg Applied_c$  then
15:   reply RecoverOK(tc, Superseding_c, Wait_c)
16: end if
receive NACK on C:
17: yield to competing coordinator
receive RecoverOK(*, Superseding, Wait) from  $p \in \mathcal{R}^t$ :
21: if  $\exists p \in \mathcal{R}^t \mid (p.Applied_c \wedge \neg p.Committed_c)$  then
22:   send response(result) to  $\forall p \in \mathcal{R}^t$ 
23:   send response(result) to  $\forall p \in \mathcal{R}^t$ 
24: else if  $\exists p \in \mathcal{R}^t \mid p.Committed_c$  then
25:   send Commit(tc, tc, p.deps) to  $\forall p \in \mathcal{R}^t$ 
26: else if  $\exists p \in \mathcal{R}^t \mid \neg p.Committed_c$  then
27:   elect  $p$  as the new coordinator
28:   select  $p$  with highest accepted ballot
29:    $t = p.t$ ,  $deps = p.deps$ 
30:   go to Consensus Protocol 1|10
31: else if  $\exists p \in \mathcal{R}^t \mid p.Applied_c \wedge \neg p.Committed_c$  then
32:    $t = t_c$ ,  $deps = \bigcup \{p.deps \mid p \in \mathcal{R}^t\}$ 
33:   if  $\exists \{p \in \mathcal{R}^t \mid p.t > p.t_c\} > |\mathcal{R}^t| - |\mathcal{R}|$  then
34:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
35:   else
36:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
37:   end if
38:   if  $\exists \{p \in \mathcal{R}^t \mid p \in \mathcal{R}^t\} \neq \emptyset$  then
39:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
40:   end if
41:   go to Consensus Protocol 1|10
42: end if

```

Algorithm 4 Reconfiguration Protocol

```

receive new configuration Conf(tc+1) on process p:
1: await ReadyEpoch = c
2: if  $p \in \mathbb{E}_c$  then
3:   await ReadyElectoratec
4:   send JoinElectorate(tc+1, tc+1, c+1) to  $\forall p \in (\mathbb{E}_{c+1} \setminus \mathbb{E}_c)$ 
5: end if
6: if  $p \in \mathbb{E}_c \wedge \neg p.Applied_c$  then
7:   ReadyElectoratec+1  $\leftarrow \text{false}$ 
8: else
9:   ReadyElectoratec+1  $\leftarrow \text{true}$ 
10: end if
11: Epoch  $\leftarrow c+1$ 
12: ReadyElectoratec+1  $\leftarrow \text{true}$ 
13: if  $p \in \mathbb{P} \setminus (\mathbb{P}_c \cup \mathbb{E}_c)$  then
14:   await ReadyShard
15:   await ReadyShard
16:   if  $\exists \gamma \in \mathcal{R} \mid \gamma \in [t_c | PreAccepted_c \wedge t_c \leq \gamma \leq c]$  then
17:     incoming_shard  $\leftarrow \text{true}$ 
18:   end if
19:   if  $\neg Applied_c$  then
20:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
21:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
22:   end if
23:   if  $\neg Applied_c$  then
24:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
25:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
26:   end if
27:   if  $\neg Applied_c$  then
28:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
29:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
30:   end if
31:   if  $\neg Applied_c$  then
32:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
33:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
34:   end if
35:   if  $\neg Applied_c$  then
36:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
37:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
38:   end if
39:   if  $\neg Applied_c$  then
40:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
41:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
42:   end if

```

Algorithm 5 Execution Protocol (with reconfiguration and ballots)

```

Coordinator C:
42: for  $p \in \mathbb{P}^t$  do
43:    $deps_p \leftarrow \{\gamma \mid \gamma \in deps \wedge p \in \mathbb{P}^t\}$ 
44:   await Read(tc, tc, depsp) to some nearby correct  $p \in \mathbb{P}^t$ 
45: end for
receive Commit(tc, tc, tc, depsp):
46:  $Committed_c \leftarrow \text{true}$ 
receive Read(tc, tc, depsp) on p:
47: await Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
48: reads  $\leftarrow \text{read}(\gamma)$ 
49: reply ReadOK(reads)
receive ReadOK(reads) from each shard:
50:  $read \leftarrow \text{read}(\gamma, \text{read})$ 
51:  $read \leftarrow \text{read}(t_c, \text{deps}_p, \text{read})$  to  $\forall p \in \mathbb{P}^t$ 
52: send read to client
receive Apply(tc, tc, depsp, result) :
53: if  $\neg Accepted_c$  then
54:   await Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
55:   apply writes(tc)
56:   wait Applied(tc, tc, depsp, tc) to  $\forall p \in \mathbb{P}^t$ 
57:   apply writes(tc)
58:   Accepted_c  $\leftarrow \text{true}$ 
59: end if

```

Algorithm 6 Reconfiguration Protocol

```

receive new configuration Conf(tc+1) on process p:
1: await ReadyEpoch = c
2: if  $p \in \mathbb{E}_c$  then
3:   await ReadyElectoratec
4:   send JoinElectorate(tc+1, tc+1, c+1) to  $\forall p \in (\mathbb{E}_{c+1} \setminus \mathbb{E}_c)$ 
5: end if
6: if  $p \in \mathbb{E}_c \wedge \neg p.Applied_c$  then
7:   ReadyElectoratec+1  $\leftarrow \text{false}$ 
8: else
9:   ReadyElectoratec+1  $\leftarrow \text{true}$ 
10: end if
11: Epoch  $\leftarrow c+1$ 
12: ReadyElectoratec+1  $\leftarrow \text{true}$ 
13: if  $p \in \mathbb{P} \setminus (\mathbb{P}_c \cup \mathbb{E}_c)$  then
14:   await ReadyShard
15:   await ReadyShard
16:   if  $\exists \gamma \in \mathcal{R} \mid \gamma \in [t_c | PreAccepted_c \wedge t_c \leq \gamma \leq c]$  then
17:     incoming_shard  $\leftarrow \text{true}$ 
18:   end if
19:   if  $\neg Applied_c$  then
20:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
21:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
22:   end if
23:   if  $\neg Applied_c$  then
24:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
25:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
26:   end if
27:   if  $\neg Applied_c$  then
28:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
29:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
30:   end if
31:   if  $\neg Applied_c$  then
32:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
33:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
34:   end if
35:   if  $\neg Applied_c$  then
36:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
37:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
38:   end if
39:   if  $\neg Applied_c$  then
40:      $\gamma = \text{Commit}(t_{c+1}, t_{c+1}, t_c, epoch=c+1)$ 
41:     send JoinShard(tc+1, tc+1, epoch=c+1) to  $\forall p \in (\mathbb{P}_{c+1} \setminus \mathbb{P}_c)$ 
42:   end if

```

Algorithm 7 Recovery Protocol

```

Coordinator C:
1:  $b \leftarrow \text{fresh ballot}$ 
2: send Recover(tc, tc) to  $\forall p \in \mathbb{P}^t$ 
receive Recover(tc, tc) on p:
3: if  $\neg MaxBallot_c$  then
4:   reply NACK(b)
5: else
6:   MaxBallot_c  $\leftarrow \max(t_c \mid t_c \leq t \wedge t_c \leq t_{\gamma} \wedge t_{\gamma} \geq t_c)$ 
7:   Commit_c  $\leftarrow \{\gamma \mid \gamma \sim t_c \wedge t_c \leq t_{\gamma} \wedge t_{\gamma} \geq t_c\}$ 
8:   Wait_c  $\leftarrow \{\gamma \in \text{Accepts} \mid t_{\gamma} \leq t_c \wedge t_{\gamma} \geq t_c\}$ 
9:   Superseding_c  $\leftarrow \{\gamma \in \text{Accepts} \mid t_{\gamma} > t_c\}$ 
10:  end if
11: if  $\neg PreAccepted_c$  then
12:   run Consensus Protocol 1|10
13:  end if
14: if  $\neg Accepted_c \wedge \neg Committed_c \wedge \neg Applied_c$  then
15:   reply RecoverOK(tc, Superseding_c, Wait_c)
16: end if
receive NACK on C:
17: yield to competing coordinator
receive RecoverOK(*, Superseding, Wait) from  $p \in \mathcal{R}^t$ :
21: if  $\exists p \in \mathcal{R}^t \mid (p.Applied_c \wedge \neg p.Committed_c)$  then
22:   send response(result) to  $\forall p \in \mathcal{R}^t$ 
23:   send response(result) to  $\forall p \in \mathcal{R}^t$ 
24: else if  $\exists p \in \mathcal{R}^t \mid p.Committed_c$  then
25:   send Commit(tc, tc, p.deps) to  $\forall p \in \mathcal{R}^t$ 
26: else if  $\exists p \in \mathcal{R}^t \mid \neg p.Committed_c$  then
27:   elect  $p$  as the new coordinator
28:   select  $p$  with highest accepted ballot
29:    $t = p.t$ ,  $deps = p.deps$ 
30:   go to Consensus Protocol 1|10
31: else if  $\exists p \in \mathcal{R}^t \mid p.Applied_c \wedge \neg p.Committed_c$  then
32:    $t = t_c$ ,  $deps = \bigcup \{p.deps \mid p \in \mathcal{R}^t\}$ 
33:   if  $\exists \{p \in \mathcal{R}^t \mid p.t > p.t_c\} > |\mathcal{R}^t| - |\mathcal{R}|$  then
34:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
35:   else
36:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
37:   end if
38:   if  $\exists \{p \in \mathcal{R}^t \mid p \in \mathcal{R}^t\} \neq \emptyset$  then
39:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
40:   end if
41:   if  $\exists \{p \in \mathcal{R}^t \mid p \in \mathcal{R}^t\} \neq \emptyset$  then
42:      $t = \max(p.t \mid p \in \mathcal{R}^t)$ 
43:   end if
44:   if  $\exists \{p \in \mathcal{R}^t \mid p \in \mathcal{R}^t\} \neq$ 
```