

Parameterized and Runtime-tunable Snapshot Isolation in Distributed Transactional Key-value Stores

Hengfeng Wei, Yu Huang, Jian Lu

Nanjing University, China

September 19, 2017



Parameterized and Runtime-tunable Snapshot Isolation

RVSI: Relaxed Version Snapshot Isolation

1 CHAMELEON Prototype and RVSI Protocol

Parameterized and Runtime-tunable Snapshot Isolation

RVSI: Relaxed Version Snapshot Isolation

1 CHAMELEON Prototype and RVSI Protocol

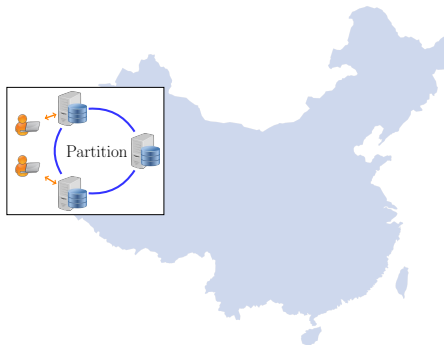
CHAMELEON prototype:

A prototype **partitioned replicated**
distributed transactional **key-value** store

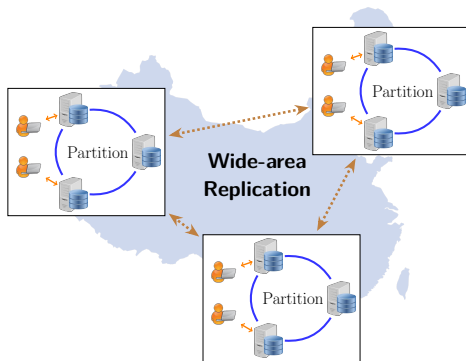
Classic **key-value** data model

Key: (row key, column key)

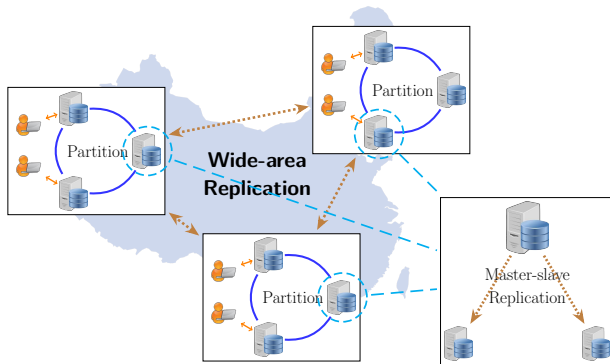




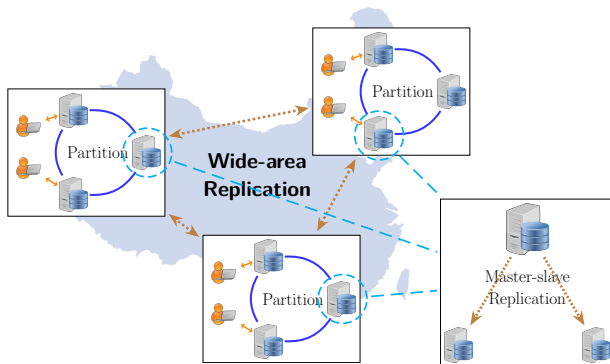
Keys are **partitioned** within a single datacenter.



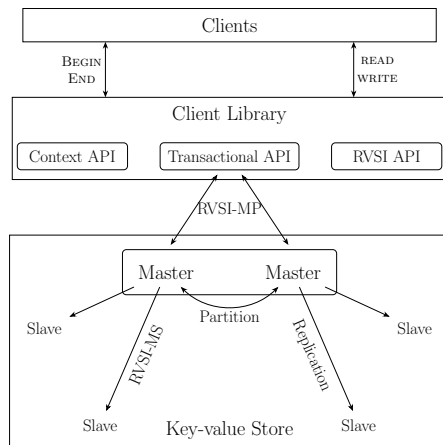
Each key is **replicated** across datacenters

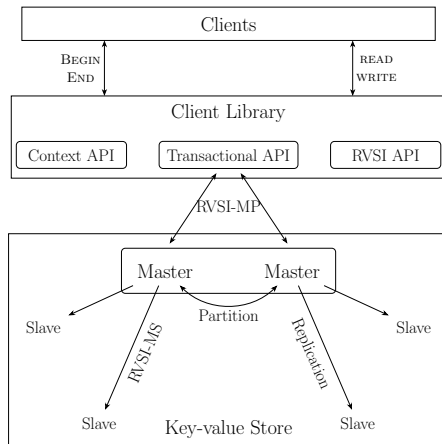


Each key is **replicated** across datacenters in a **master-slave** manner.

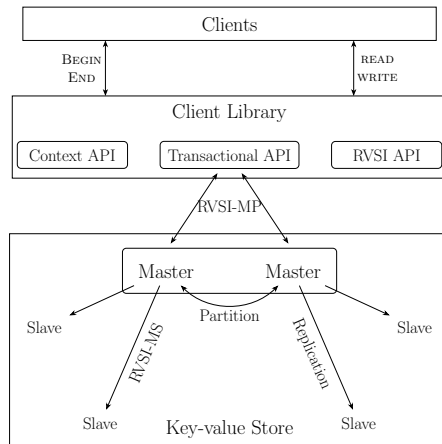


Transactions are first executed and committed on the **masters**, and are then asynchronously propagated to **slaves**.





1. Partitioned replicated transactional key-value store



2. Client library

Code snippet for writing RVSI transactions:

```
// Initialize keys (ck1 and ck2) here

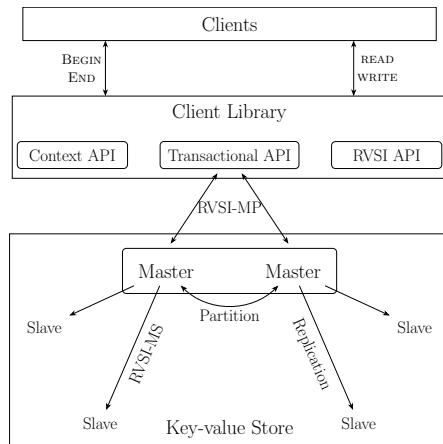
ITx tx = new RVSITx(** context **);

tx.begin();

// Read and write here

// Specify RVSI specs. (e.g., SVSpec)
RVSISpec sv = new SVSpec();
sv.addSpec({ck1, ck2}, 2);
tx.collectRVSISpec(sv);

boolean committed = tx.end();
```



3. RVSI protocol: RVSI-MS + RVSI-MP

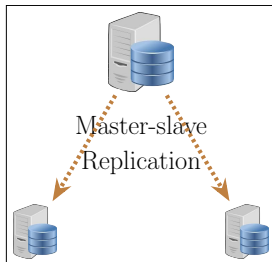
RVSI-MS: RVSI protocol for master-slave replication

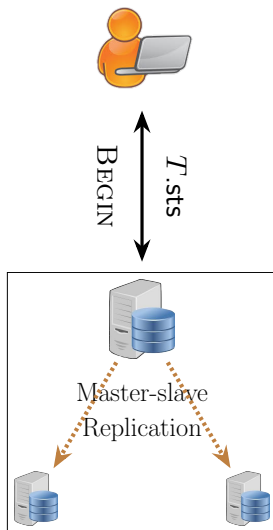
In terms of *event* generation and handling:

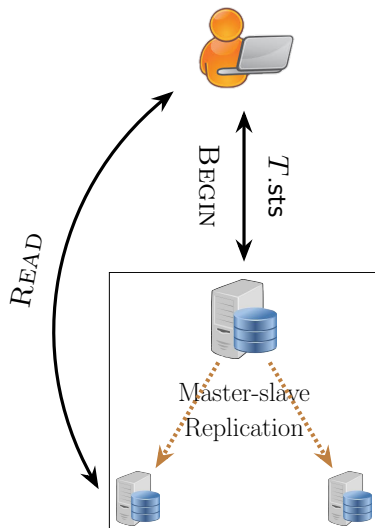
Clients: BEGIN, READ, WRITE, END

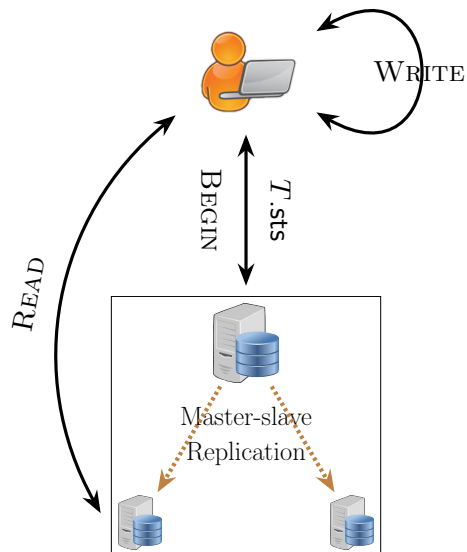
Master: START, COMMIT, SEND

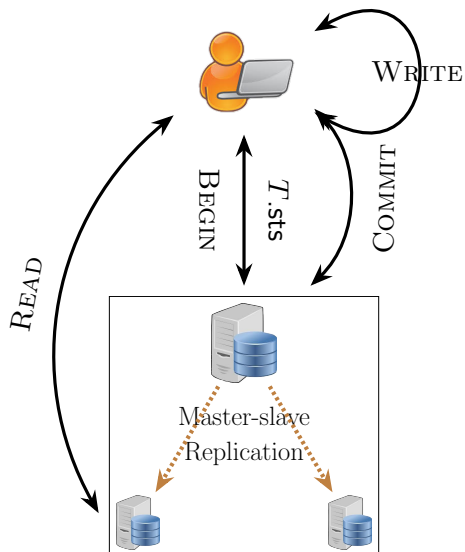
Slaves: RECEIVE

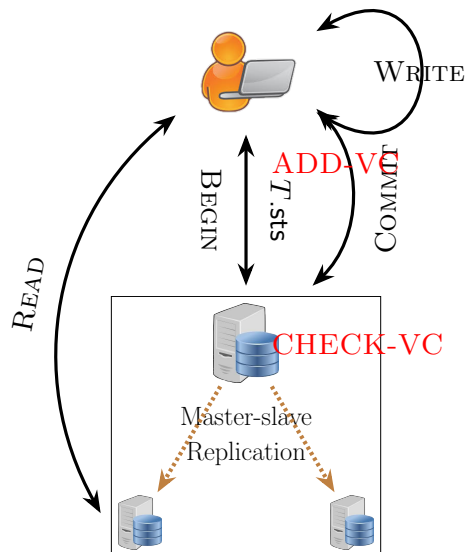












Calculating version constraints for RVSI:

$\mathcal{O}_x(t) = \#$ of versions of x before time t

$$r_i(x_j) \in T_i$$

k_1 -BV:

$$\mathcal{O}_x(T_i.sts) - \mathcal{O}_x(T_j.cts) < k_1$$

k_2 -FV:

$$\mathcal{O}_x(T_j.cts) - \mathcal{O}_x(T_i.sts) \leq k_2$$

Calculating version constraints for RVSI:

$\mathcal{O}_x(t) = \#$ of versions of x before time t

$$r_i(x_j) \in T_i$$

k_1 -BV:

$$\mathcal{O}_x(T_i.sts) - \mathcal{O}_x(T_j.cts) < k_1$$

k_2 -FV:

$$\mathcal{O}_x(T_j.cts) - \mathcal{O}_x(T_i.sts) \leq k_2$$

k_3 -SV:

$$r_i(x_j), r_i(y_l) \in T_i$$

$$\mathcal{O}_{\textcolor{red}{x}}(T_l.cts) - \mathcal{O}_{\textcolor{red}{x}}(T_j.cts) \leq k_3$$

Algorithm 1 RVSI-MS Protocol for Executing Transaction T (Client).

```

1: procedure BEGIN()
2:    $T.sts \leftarrow$  rpc-call START() at master  $\mathcal{M}$ 
3: procedure READ( $x$ )
4:    $x.ver \leftarrow$  rpc-call READ( $x$ ) at any site
5: procedure WRITE( $x, v$ )
6:   add  $(x, v)$  to  $T.writes$ 
7: procedure END( $T$ )
8:    $T.vc \leftarrow$  ADD-VC()
9:    $c/a \leftarrow$  rpc-call COMMIT( $T.writes, T.vc$ ) at  $\mathcal{M}$ 
  
```

Algorithm 1 RVSI-MS Protocol for Executing Transaction T (Master).

$\mathcal{M}.ts$: for start-timestamps and commit-timestamps

$\{x.ver = (x.ts, x.ord, x.val)\}$: set of versions of x

```

1: procedure START()
2:   return ++ $\mathcal{M}.ts$ 
3: procedure READ( $x$ )
4:   return the latest  $x.ver$  installed
5: procedure COMMIT( $T.writes, T.vc$ )
6:   if CHECK-VC( $T.vc$ ) && write-conflict freedom then
7:      $T.cts \leftarrow ++\mathcal{M}.ts$ 
8:     ▷ apply  $T.writes$  locally and propagate it
9:      $T.upvers = \emptyset$                                 ▷ collect updated versions
10:    for  $(x, v) \in T.writes$  do
11:       $x.new-ver \leftarrow (T.cts, ++x.ord, v)$ 
12:      add  $x.new-ver$  to  $\{x.ver\}$  and  $T.upvers$ 
13:    broadcast  $\langle PROP, T.upvers \rangle$  to slaves
14:    return  $c$  denoting “committed”
15:  return  $a$  denoting “aborted”

```

Algorithm 1 RVSI-MS Protocol for Executing Transaction T (Slave).

$x.ver = (x.ts, x.ord, x.val)$: the latest version of x

```

1: procedure READ( $x$ )
2:   return  $x.ver$ 
3: upon RECEIVED( $\langle \text{PROP}, T.upvers \rangle$ )
4:   for  $(x.ver' = (x.ts', x.ord', x.val')) \in T.upvers$  do
5:     if  $x.ord' > x.ord$  then
6:        $x.ver \leftarrow x.ver'$ 
  
```

Distributed transactions spanning multiple masters
need to be committed atomically.

Distributed transactions spanning multiple masters
need to be committed atomically.

Using the two-phase commit (2PC) protocol.

Distributed transactions spanning multiple masters
need to be committed atomically.

Using the two-phase commit (2PC) protocol.

Two issues to address

Assumes a timestamp oracle [\[Peng@OSDI'10\]](#):

Client: asks for the start-timestamp in BEGIN

Coordinator: asks for the commit-timestamp in COMMIT

Split the RVSI version constraints according to partitions:

$$r_i(x_j) \in T_i$$

k_1 -BV:

$$\mathcal{O}_x(T_i.sts) - \mathcal{O}_x(T_j.cts) < k_1$$

k_2 -FV:

$$\mathcal{O}_x(T_j.cts) - \mathcal{O}_x(T_i.sts) \leq k_2$$

k_3 -SV:

$$r_i(x_j), r_i(y_l) \in T_i$$

$$\mathcal{O}_x(T_l.cts) - \mathcal{O}_x(T_j.cts) \leq k_3$$

All version constraints involve only one data item.

Algorithm 2 RVSI-MP for Executing Transaction T (Client).

```
1: procedure BEGIN()  
2:   return rpc-call GETTS() at  $\mathcal{T}$   
3: procedure END()  
4:    $T.vc \leftarrow \text{ADD-VC}()$   
5:    $c/a \leftarrow \text{rpc-call C-COMMIT}(T.writes, T.vc)$  at  $\mathcal{C}$ 
```

Algorithm 2 RVSI-MP for Executing Transaction T (Timestamp Oracle).

$\mathcal{T}.ts$: for start-timestamps and commit-timestamps

```
1: procedure GETTS()  
2:   return ++ $\mathcal{T}.ts$ 
```

Algorithm 2 RVSI-MP for Executing Transaction T (Coordinator).

```

1: procedure C-COMMIT( $T.writes, T.vc$ )
2:   split  $T.writes$  and  $T.vc$  with the data partitioning strategy
3:    $\triangleright$  the prepare phase:
4:   rpc-call PREPARE( $T.writes, T.vc$ ) at each  $\mathcal{M}$ 
5:    $\triangleright$  the commit phase:
6:   if all PREPARE( $T.writes, T.vc$ ) return true then
7:      $T.cts \leftarrow$  rpc-call GETTS() at  $\mathcal{T}$ 
8:     rpc-call COMMIT( $T.cts, T.writes$ ) at each  $\mathcal{M}$ 
9:   else
10:    rpc-call ABORT() at each  $\mathcal{M}$ 
11:    return  $a$  denoting “aborted”
12:   if all COMMIT( $T.cts, T.writes$ ) return true then
13:     return  $c$  denoting “committed”
14:   else
15:     return  $a$  denoting “aborted”
  
```

Algorithm 2 RVSI-MP for Executing Transaction T (Master).

```

1: procedure PREPARE( $T.writes, T.vc$ )
2:   return CHECK-VC( $T.vc$ ) && write-conflict freedom
3: procedure COMMIT( $T.cts, T.writes$ )
4:   ▷ apply  $T.writes$  locally and propagate it
5: procedure ABORT()
6:   ▷ abort
  
```
