Teaching Elephants To Tell Time

Adventures in Distributing Snapshot Isolation

A complicated view...









PostgreSQL creates snapshots as of <u>NOW</u> with serialized local state

Now doesn't work in distributed systems



Now doesn't work in distributed systems



Now doesn't work in distributed systems



NOW doesn't work in distributed systems



PostgreSQL creates snapshots as of NOW with serialized local state

PostgreSQL creates snapshots as of NOW with serialized local state

Time-based MVCC

1.Establish snapshot on current time...GetSnapshotData() irrelevant

1.Establish snapshot on current time...GetSnapshotData() irrelevant

2.Pass snapshot time along with query fragment to shards.

1.Establish snapshot on current time...GetSnapshotData() irrelevant

2.Pass snapshot time along with query fragment to shards.

3.Record commit time in clog

1.Establish snapshot on current time...GetSnapshotData() irrelevant

2.Pass snapshot time along with query fragment to shards.

3.Record commit time in clog

4. Visibility is commit time vs snapshot time

Distributed clocks????



True Time (UTC), t



Time

Pseudo commit algorithm

doCommit()

{

}

// returns immediately with lsn and HLC time of commit log record
pair{lsn,TimeOfCommit} = flushLogToStorageAsync();

// returns when the 3AZ durability point >= the commit record
// Typically returns 1.5ms after commit time @ p50
waitForLsnDurable(lsn);

// returns when ClockBound.earliest > TimeOfCommit
// Typically returns immediately as (commit time - earliest) < 1ms @ p90
waitForEarliest(TimeOfCommit);</pre>

Pseudo commit algorithm

doCommit()

{

}

// returns immediately with lsn and HLC time of commit log record
pair{lsn,TimeOfCommit} = flushLogToStorageAsync();

// returns when the 3AZ durability point >= the commit record
// Typically returns 1.5ms after commit time @ p50
waitForLsnDurable(lsn);

// returns when ClockBound.earliest > TimeOfCommit
// Typically returns immediately as (commit time - earliest) < 1ms @ p90
waitForEarliest(TimeOfCommit);</pre>

No blocking in waitForEarliest observed in cluster sustaining 5MM NOPM (HammerDb)

No blocking in 60 shard cluster performing 2MM commits / sec

That's just the intro

Clog compaction

Snapshot horizon gossip for vacuum control

Distributed commit coordinator (2PC) Non-blocking 2PC algorithm

Hooks that make this an extension

Committing status for long-fork anomaly prevention