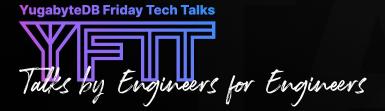
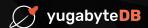
# Don't Let Tombstones Affect Scan Performance

Premkumar Thangamani Feb/24/2023





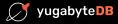
#### **Topics of discussion**

- Effect of Large Range Deletes
- Point Lookups vs Range Scan
- Simple Overview
  - Log Structured Merge Tree (LSM)
  - Sorted Strings Table (SST)
- Possible Workarounds
- Future Work



### **Effect of Large Deletes**

**Events Scenario** 



#### **Events / Job Queue / Ticketing**

```
id int,
  ts timestamp,
  data text,
  PRIMARY KEY (id)
);
```

#### **Earliest event**

SELECT \* FROM events ORDER BY ts ASC LIMIT 1

#### Helpful Index

CREATE INDEX idx\_ts on events(ts ASC)



#### Earliest event - usually very fast (~2 ms)

#### SELECT \* FROM events ORDER BY ts ASC LIMIT 1

Time: 2.257 ms



#### After a lot of deletes ... (~2000 ms)

SELECT \* FROM events ORDER BY ts ASC LIMIT 1

Time: 2271.798 ms



### I see dead records!

#### **On Deletes**

- Records not removed immediately
- Tombstone Markers will be placed
- **LSM**: Log Structured Merge-Tree
  - Multi-Level Tree like storage
- MVCC: Multi Version Concurrency Control
  - Consistent View w/o Locking for reads



#### Simple visual of SST

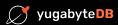


#### Compaction

- Fixes the problem
- Removes the Tombstones

- But,
  - Resource Intensive
  - To support MVCC, some records may stay on

## Simple Trick



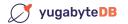
#### Last processed timestamp

- The problem: not knowing where to start.

- Let's help the Query Executor
  - ts > :last\_processed\_ts

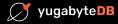
#### Use last processed timestamp as a hint

Time: 2.850 ms



14

### **DEMO**



#### **Storing last timestamp**

- Start with a base timestamp of zero
  - Store it as a local variable

- Store the last\_processed\_ts in a separate table
  - Background thread updates it every ~30s
  - Useful in a multiple worker scenario



#### Follow up

#### **Future Work**

- Identifying scenarios of dense deletes
- Automatically triggering compaction

#### **Further Watching**

- Check out John Meehan's Tech Talk on Compaction
- Full Compactions in YugabyteDB



### **Thank You**

Join us on Slack: yugabyte.com/slack (#yftt channel)

Star us on Github: github.com/yugabyte/yugabyte-db

