



yugabyte**DB**

How YugaByte DB Works as a Distributed PostgreSQL at Scale

Sid Choudhury

Bryn Llewellyn

NoCOUG Summer 2019 Conference

YugaByte DB



Distributed SQL

PostgreSQL Compatible, 100% Open Source (Apache 2.0)



Massive Scale

Millions of IOPS in Throughput, TBs per Node



High Performance

Low Latency Queries



Cloud Native

Fault Tolerant, Multi-Cloud & Kubernetes Ready

Functional Architecture

YugaByte SQL (YSQL)

PostgreSQL-Compatible Distributed SQL API

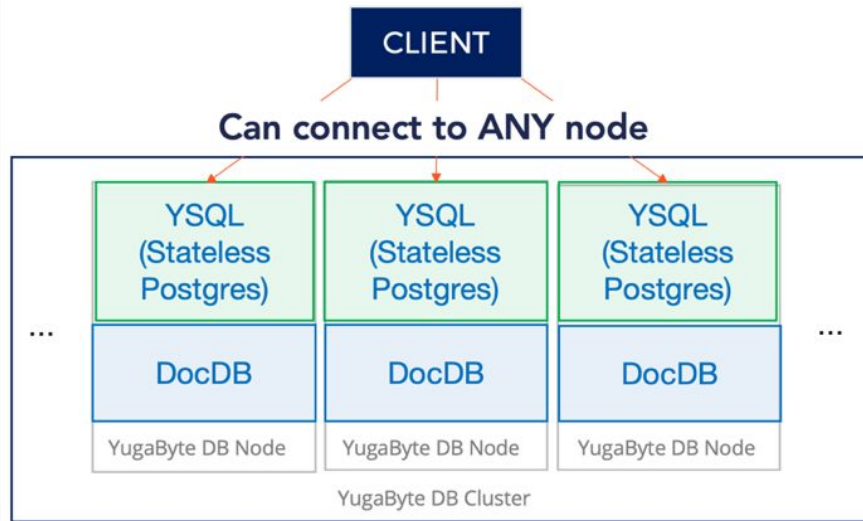
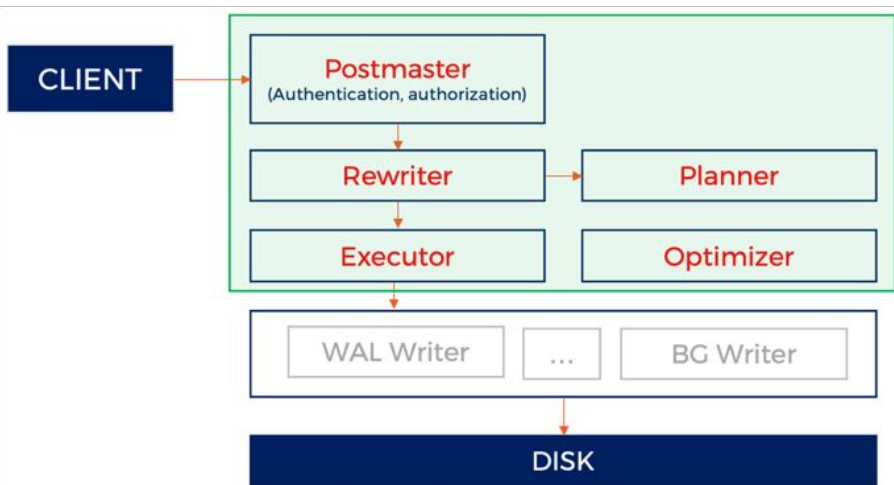
DOCDB

Spanner-Inspired Distributed Document Store
Cloud Neutral: No Specialized Hardware Needed

PostgreSQL Transformed into Distributed SQL



yugabyteDB



Design Goals

- **PostgreSQL compatible**

- Re-uses PostgreSQL query layer
- New changes do not break existing PostgreSQL functionality

- **Enable migrating to newer PostgreSQL versions**

- New features are implemented in a modular fashion
- Integrate with new PostgreSQL features as they are available
- E.g. Moved from PostgreSQL 10.4 → 11.2 in 2 weeks!

- **Cloud native architecture**

- Fully decentralized to enable scaling to 1000s of nodes
- Tolerate rack/zone and datacenter/region failures automatically
- Run natively in containers and Kubernetes
- Zero-downtime rolling software upgrades and machine reconfig

SQL Feature Depth

- **Traditional SQL**

- Data types
- Relational integrity (Foreign keys)
- Built-in functions
- Expressions
- JSON column type
- Secondary indexes
- JOINS
- Transactions
- Views

- **Advanced SQL**

- Partial indexes
 - Stored procedures
 - Triggers
 - Extensions
- And more ...

Create Table & Insert Data

YSQL Tables

- **Tables**

- Each table maps to one DocDB table
- Each DocDB table is sharded into multiple tablets

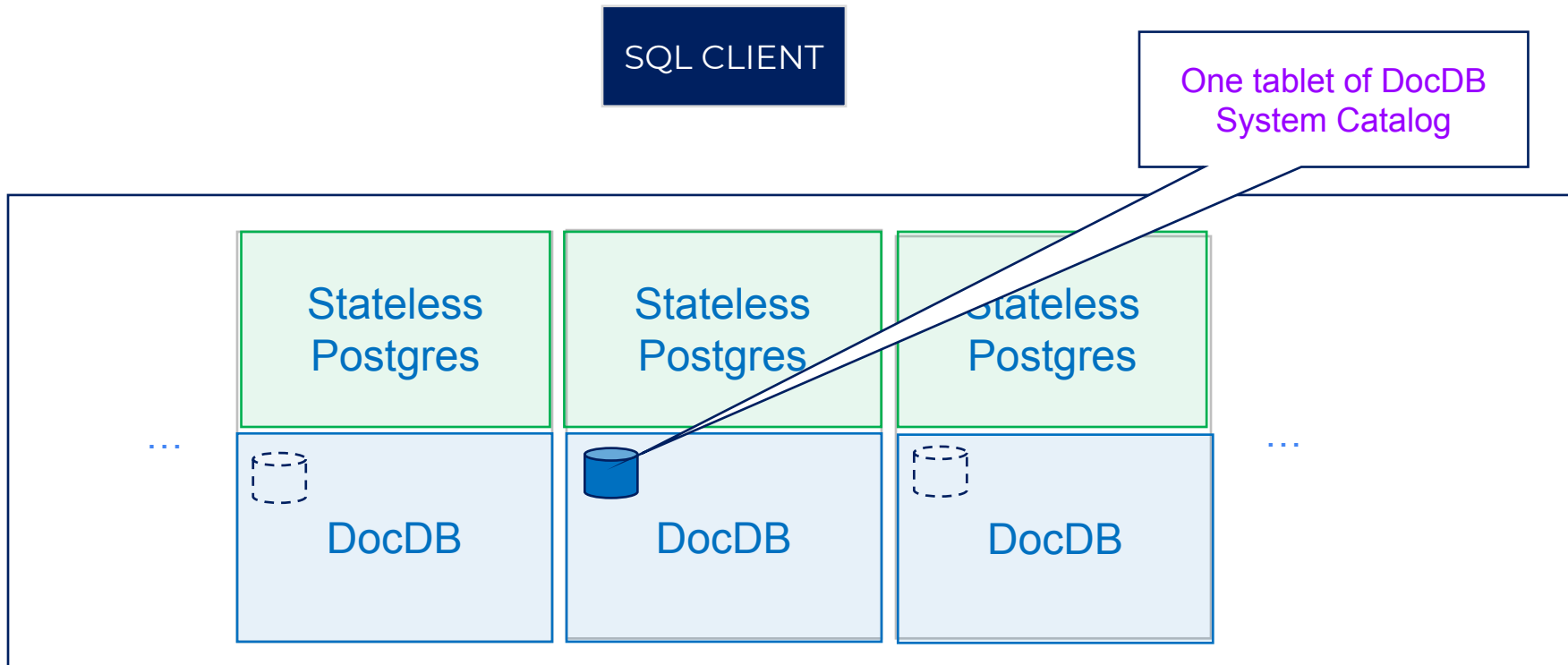
- **System tables**

- PostgreSQL system catalog tables map to special DocDB tables
- All such special DocDB tables use a single tablet

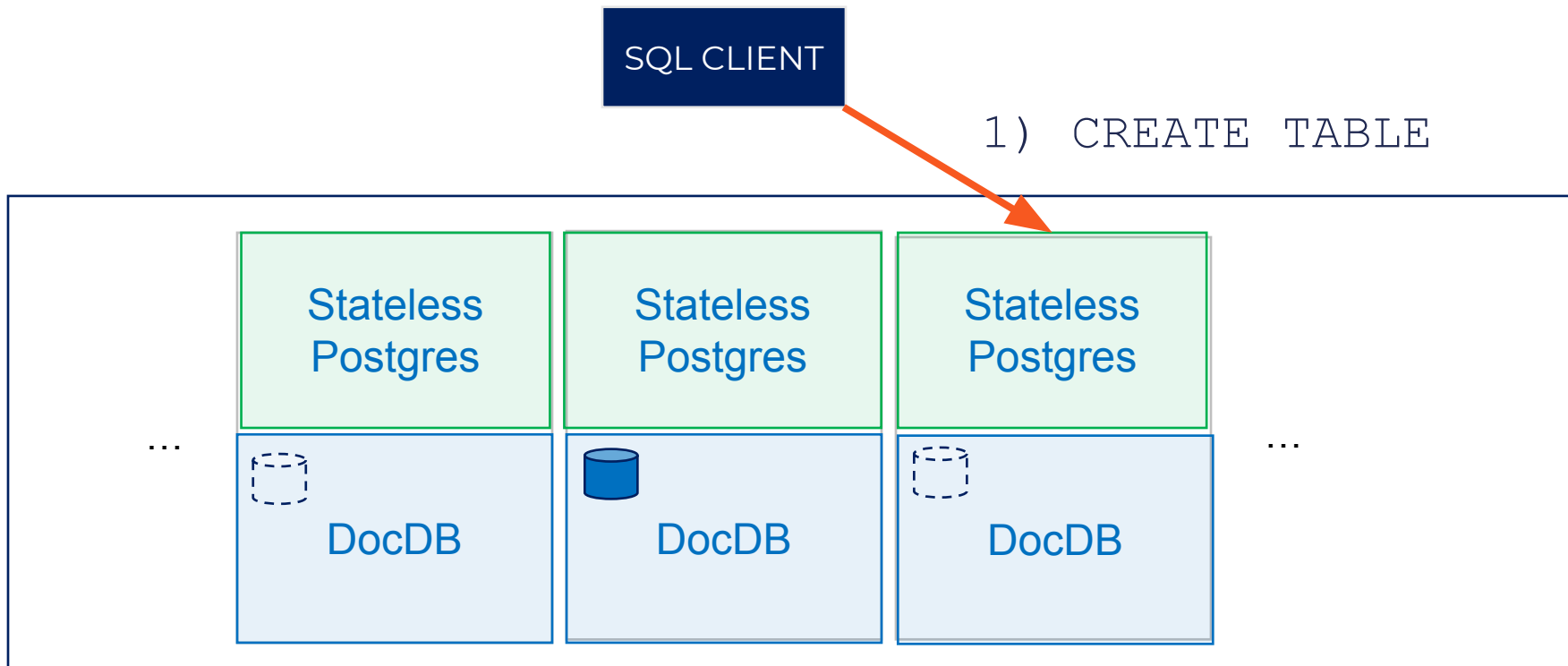
- **(Internal) DocDB tables**

- Have same key → document format
- Schema enforcement using the table schema metadata

System Catalog Tables are Special Tables



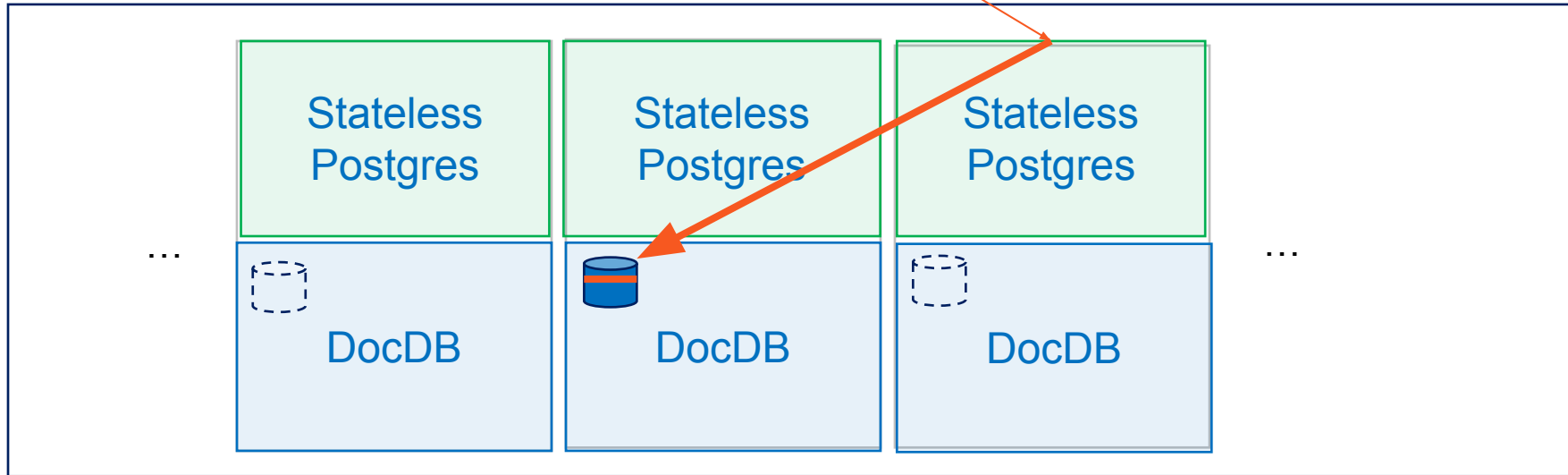
Create a Table



Create a Table

SQL CLIENT

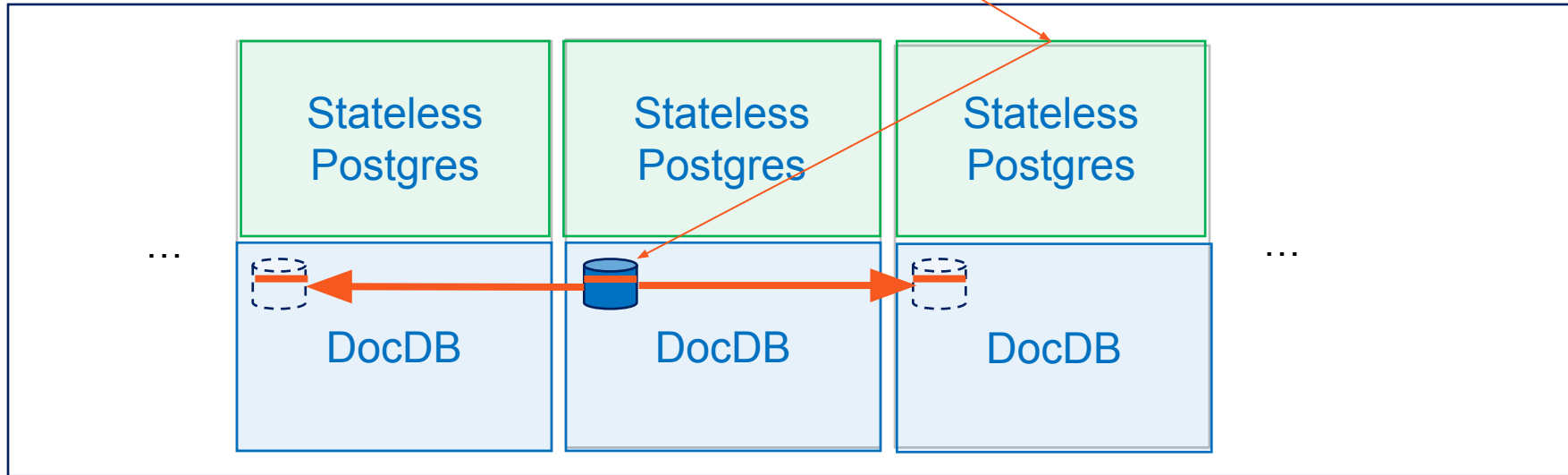
2) RECORD SCHEMA



Create a Table

SQL CLIENT

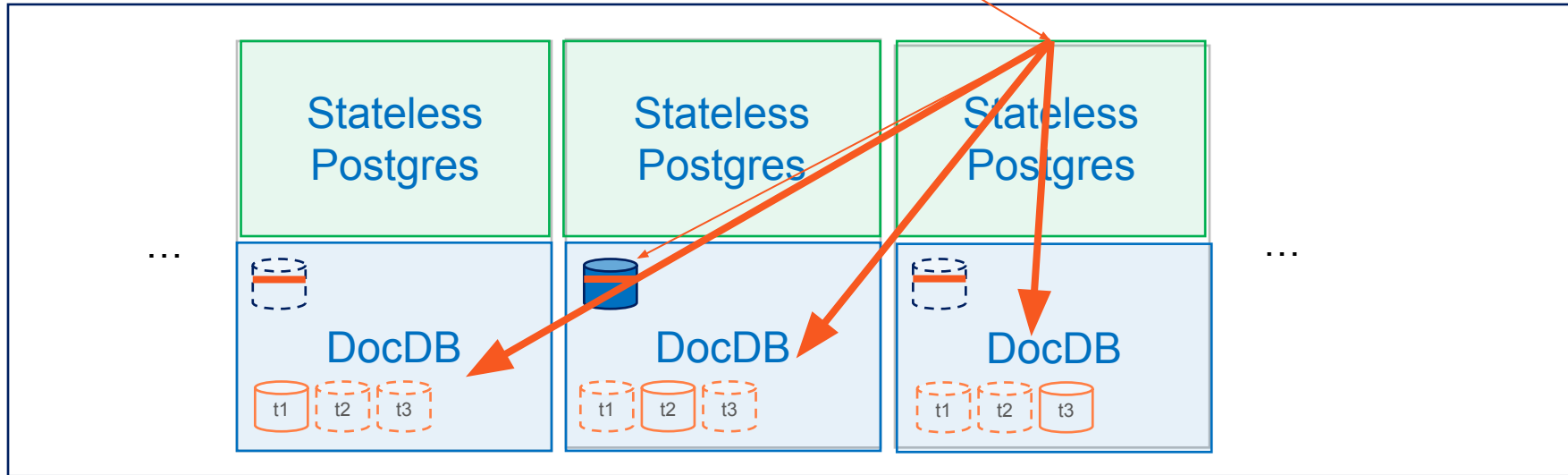
3) RAFT REPLICATE



Create a Table

CLIENT

4) CREATE TABLETS



Insert Data into Tables

- **Primary keys**

- The primary key column(s) map to a single document key
- Each row maps to one document in DocDB
- Tables without primary key use an internal ID (logically a row-id)

- **Secondary indexes**

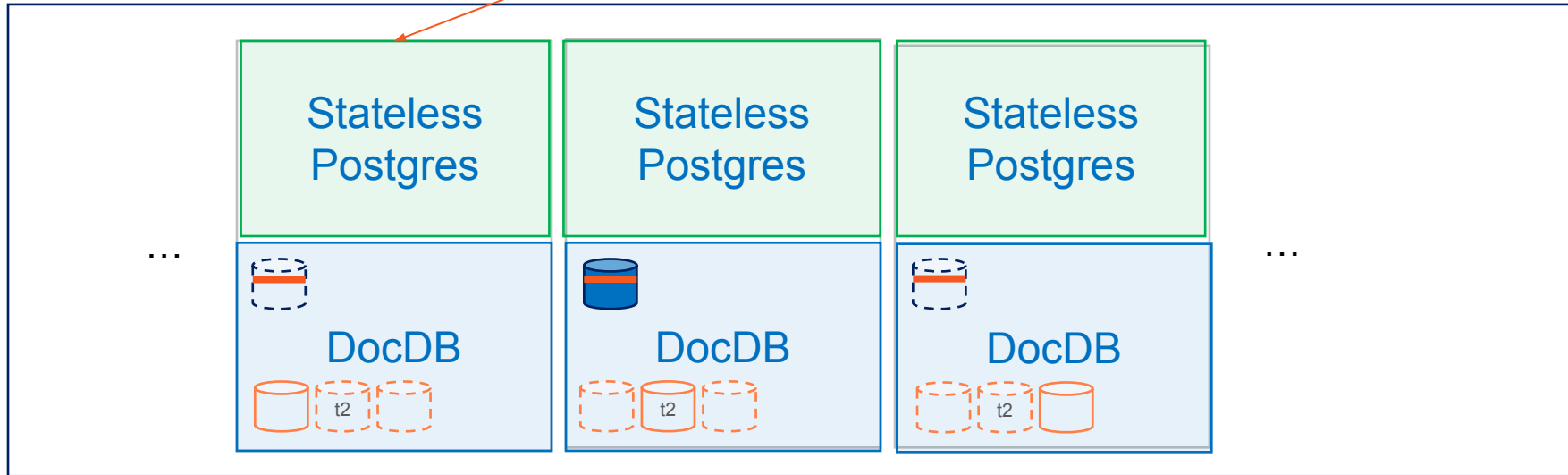
- Each index maps to a separate distributed DocDB table
- DML implemented using **DocDB distributed transactions**
- E.g: insert into table with one index will perform the following:

```
BEGIN DOCDB DISTRIBUTED TRANSACTION
    insert into index values (...)
    insert into table values (...)
COMMIT
```

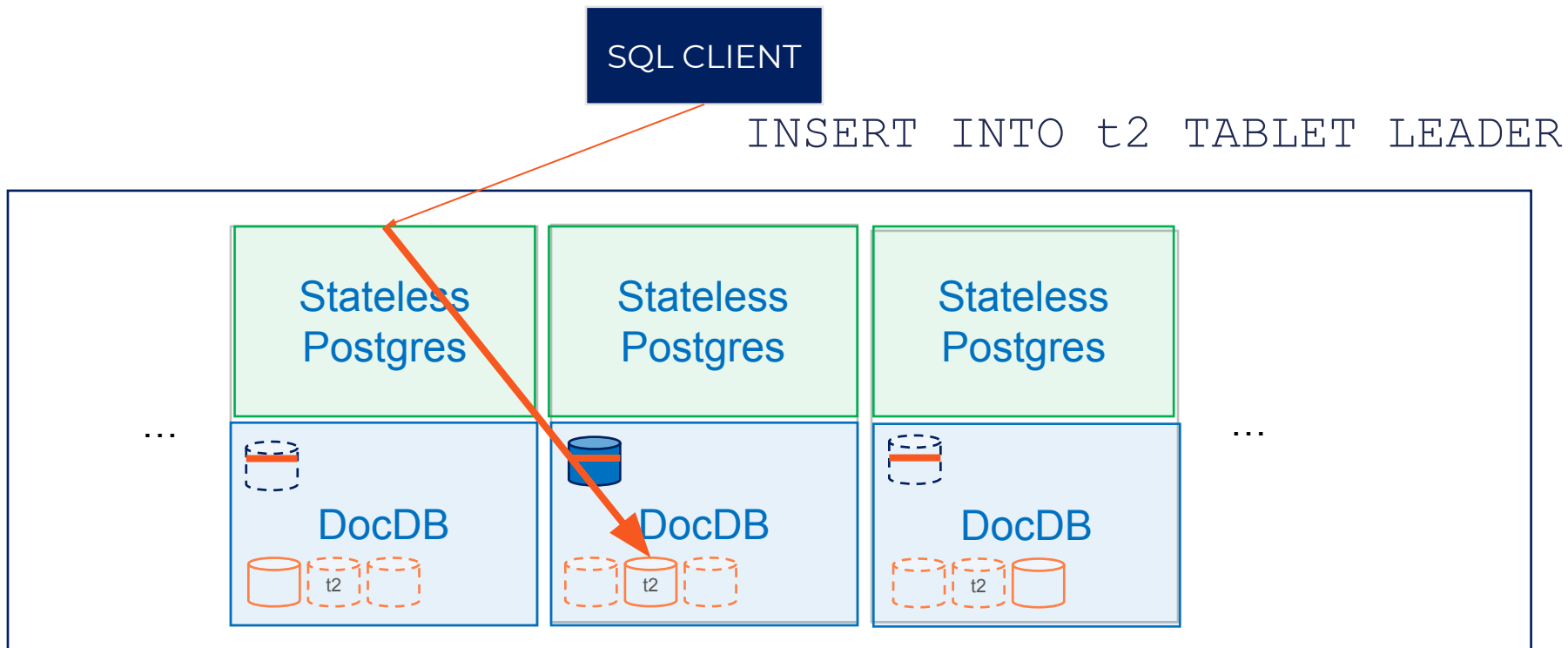
Insert Data

SQL CLIENT

INSERT ROW



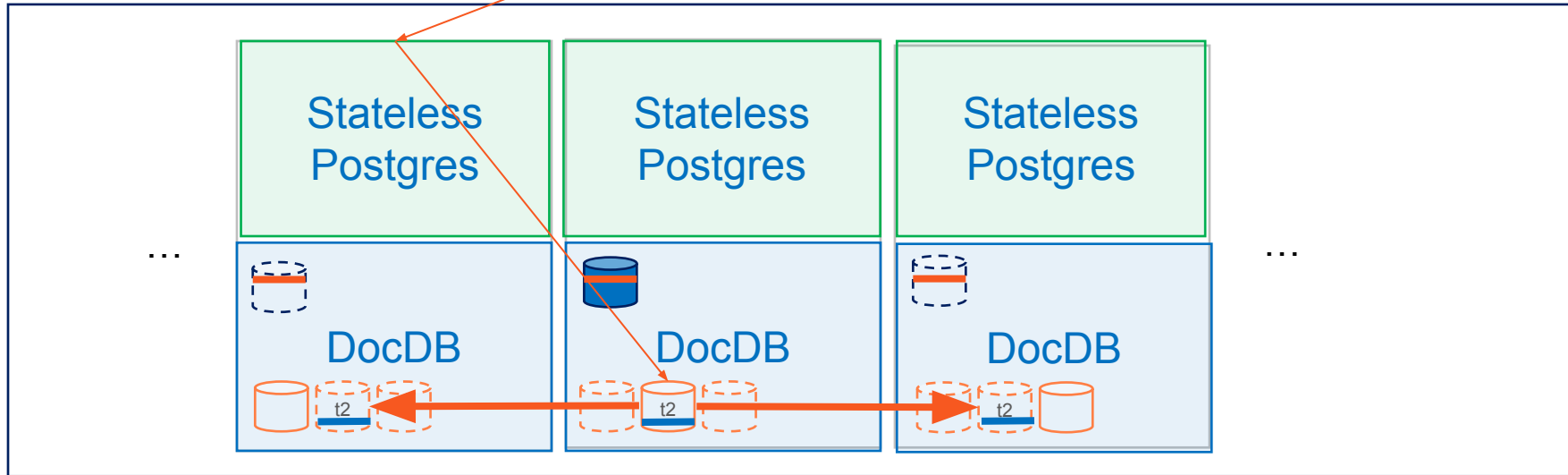
Insert Data



Insert Data

SQL CLIENT

RAFT REPLICATE DATA

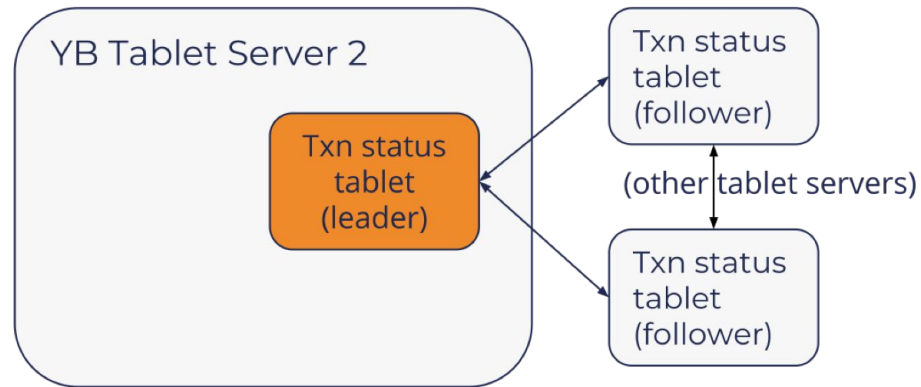
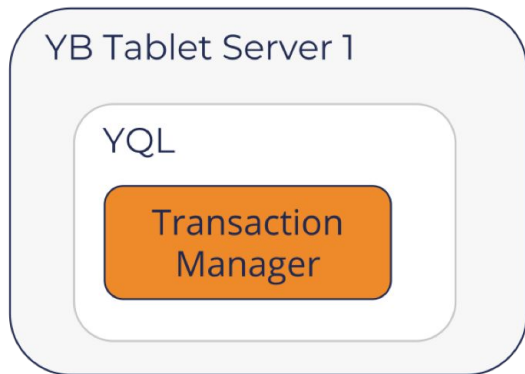


Distributed Transactions

Fully Decentralized Architecture

- **No single point of failure or bottleneck**
 - Any node can act as a Transaction Manager
- **Transaction status table distributed across multiple nodes**
 - Tracks state of active transactions
- **Transactions have 3 states**
 - Pending
 - Committed
 - Aborted
- **Reads served only for Committed Transactions**
 - Clients never see inconsistent data

Distributed Transactions - Write Path



Distributed Transactions - Write Path

1

Client's request:
set k1=v1, k2=v2

YB Tablet Server 1

YQL

Transaction
Manager

YB Tablet Server 2

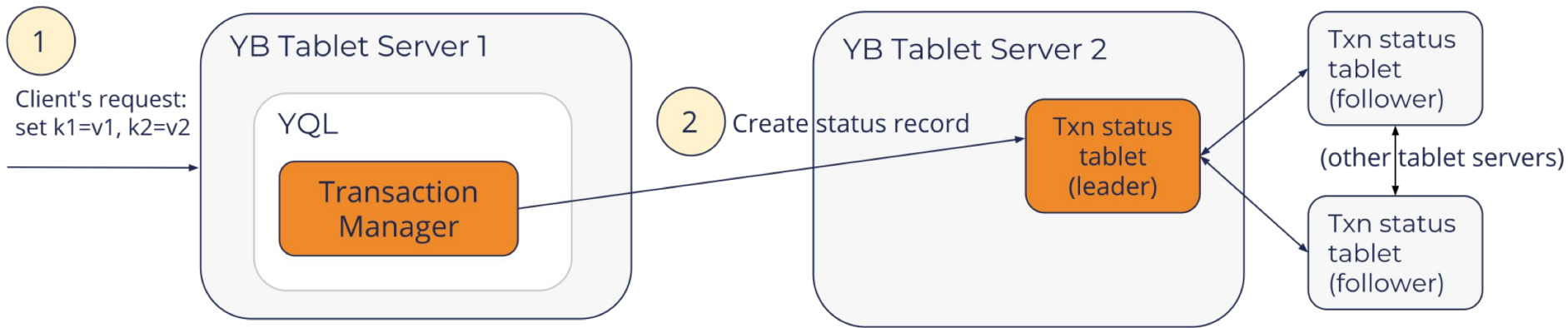
Txn status
tablet
(leader)

Txn status
tablet
(follower)

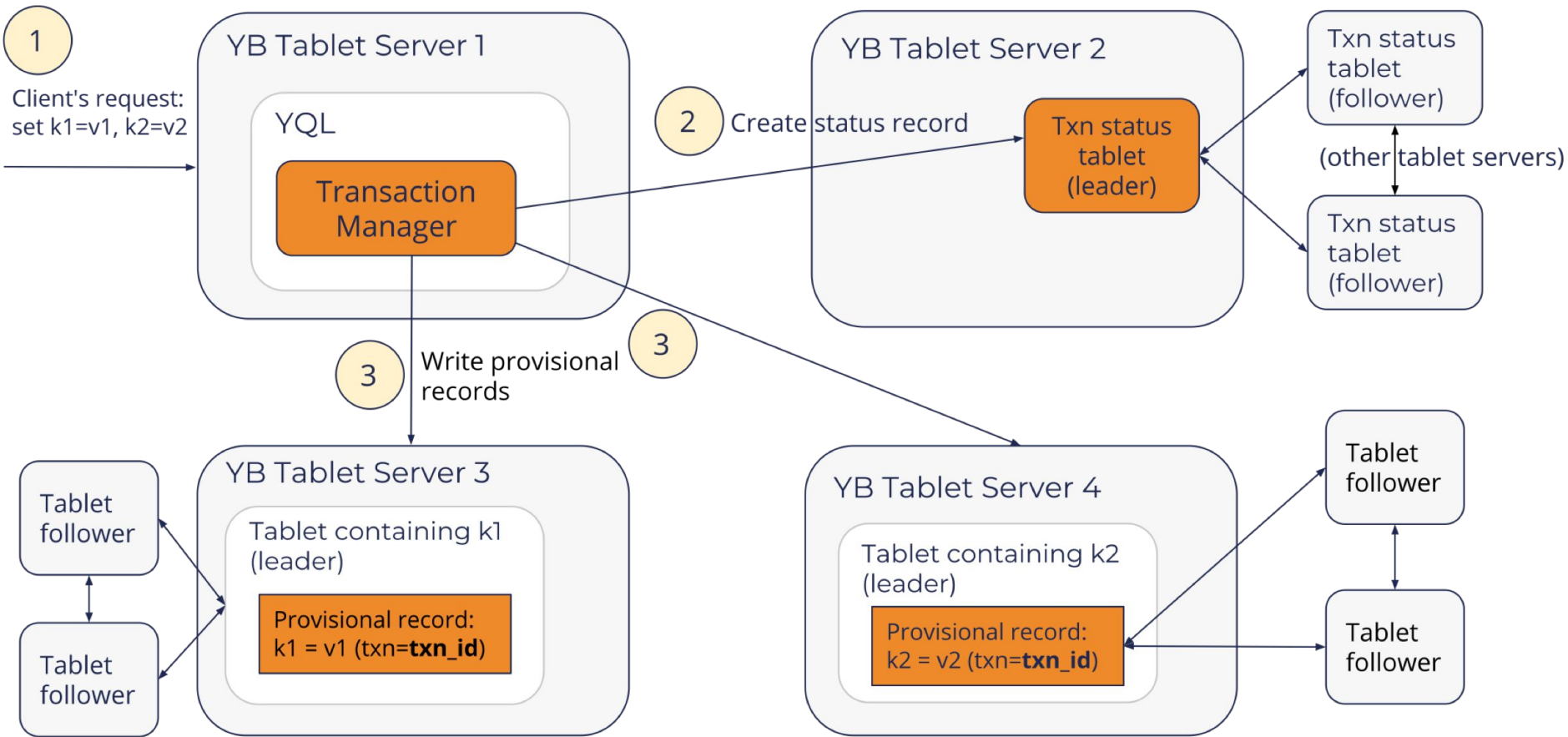
(other tablet servers)

Txn status
tablet
(follower)

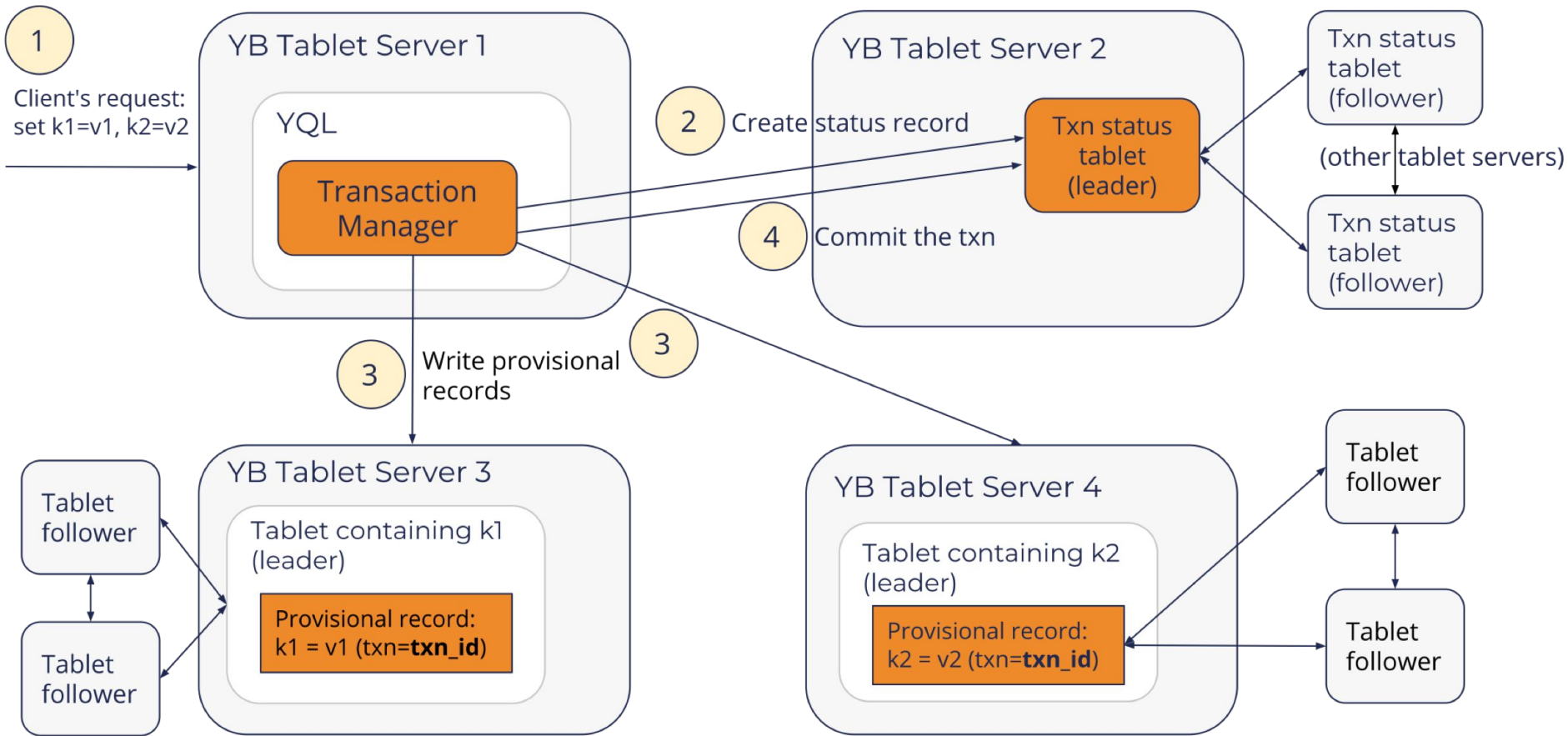
Distributed Transactions - Write Path



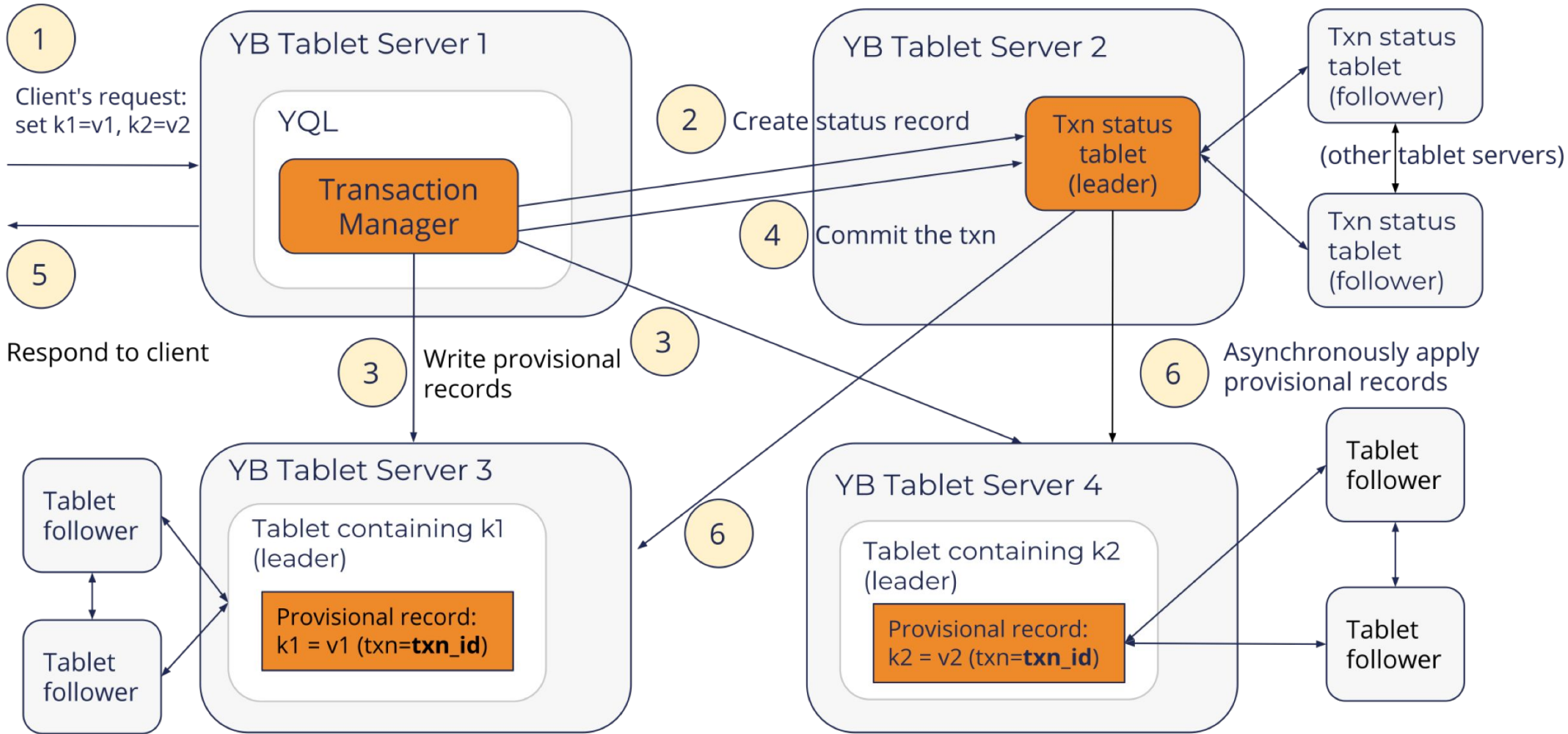
Distributed Transactions - Write Path



Distributed Transactions - Write Path



Distributed Transactions - Write Path



Isolation Levels

- **Serializable Isolation**

- Read-write conflicts get auto-detected
- Both reads and writes in read-write txns need provisional records
- Maps to SERIALIZABLE in PostgreSQL

- **Snapshot Isolation**

- Write-write conflicts get auto-detected
- Only writes in read-write txns need provisional records
- Maps to REPEATABLE READ, READ COMMITTED & READ UNCOMMITTED in PostgreSQL

- **Read-only Transactions**

- Lock free

Summary

Most Advanced Open Source Distributed SQL



PostgreSQL
Query Layer

World's Most Advanced
Open Source SQL Engine



Google Spanner
Storage Layer

World's Most Advanced
Distributed OLTP Architecture



yugabyteDB



Distributed SQL Summit

📅 September 20, 2019 📍 Hilton, San Jose, CA

A full day of talks from experts on what it takes to build, deploy and scale distributed SQL databases in the cloud and on Kubernetes

Hear from the creators of Google Spanner, Amazon Aurora, Facebook DBs & YugaByte DB

Register today at
distributedsql.org

Read more at
blog.yugabyte.com

Storage Layer

blog.yugabyte.com/distributed-postgresql-on-a-google-spanner-architecture-storage-layer

Query Layer

blog.yugabyte.com/distributed-postgresql-on-a-google-spanner-architecture-query-layer



Questions?

Download

download.yugabyte.com

Join Slack Discussions

yugabyte.com/slack

Star on GitHub

github.com/YugaByte/yugabyte-db