

A Brief Introduction of TiDB

Dongxu (Edward) Huang
CTO, PingCAP

About me

- Dongxu (Edward) Huang, Cofounder & CTO of PingCAP
- PingCAP, based in Beijing, China.
- Infrastructure software engineer, open source hacker
- Codis / TiDB / TiKV
- Golang / Python / Rust

What would you do when...

- RDBMS is becoming the performance bottleneck of your backend service
- The amount of data stored in RDBMS is overwhelming
- You want to do some complex queries on a sharding cluster
 - e.g. simple JOIN or GROUP BY
- Your application needs ACID transaction on a sharding cluster

TiDB Project - Goal

- SQL is necessary
- Transparent sharding and data movement
- 100% OLTP + 80% OLAP
 - Transaction + Complex query
- Compatible with MySQL, at most cases
- 24/7 availability, even in case of datacenter outages
 - Thanks to Raft consensus algorithm
- Open source, of course.



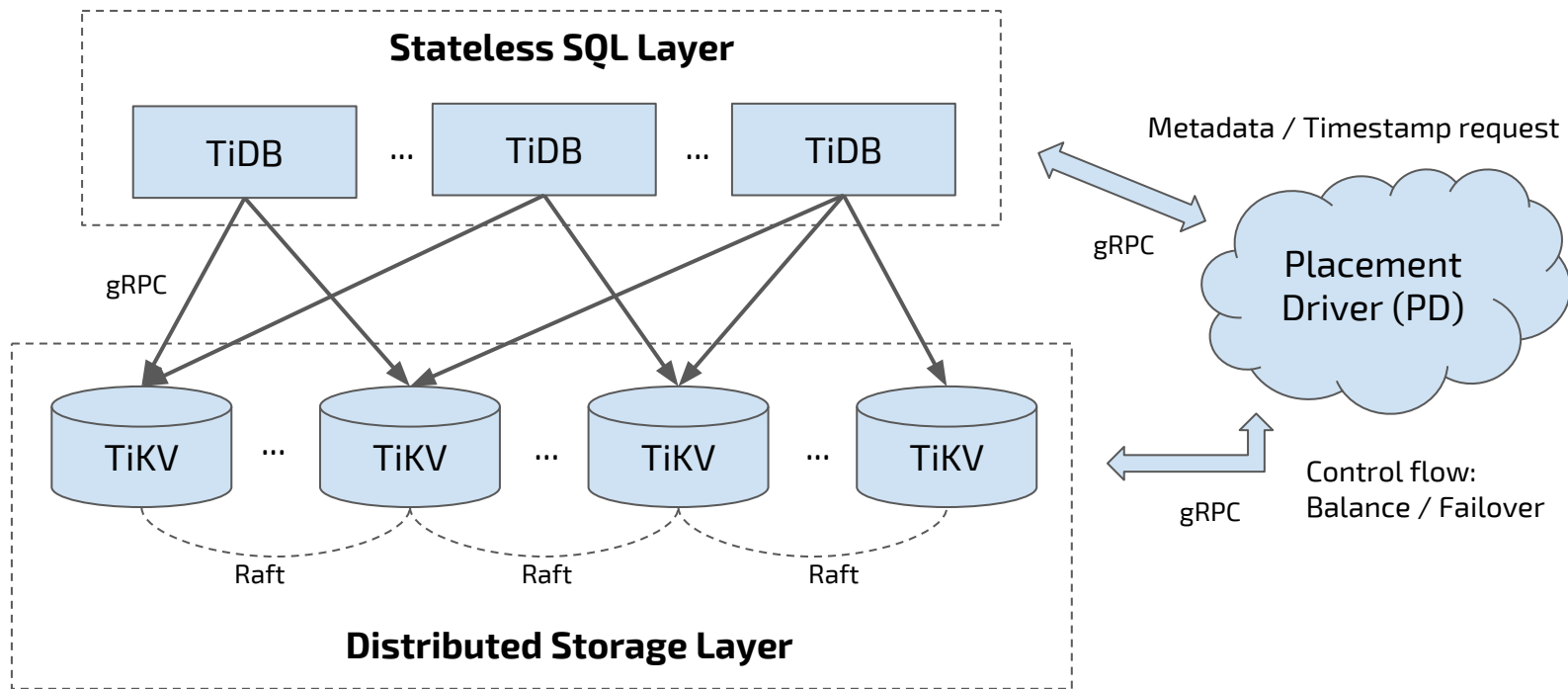
TiDB

A Distributed SQL Database

Agenda

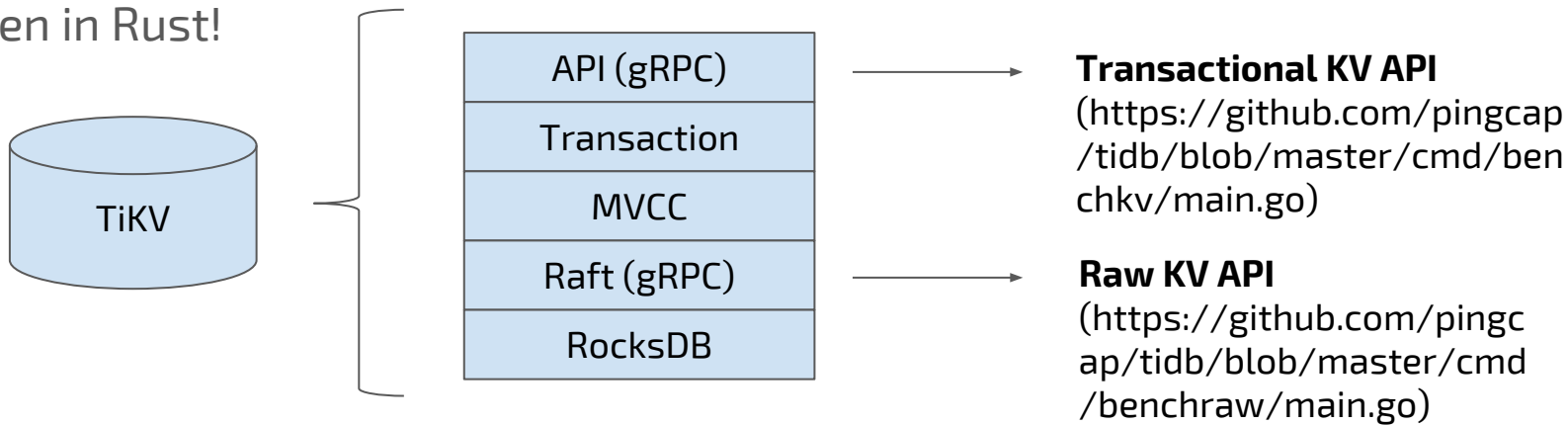
- Technical overview of TiDB / TiKV
 - Storage
 - Distributed SQL
 - Tools
- Real-world cases and benchmarks
- Demo

Architecture



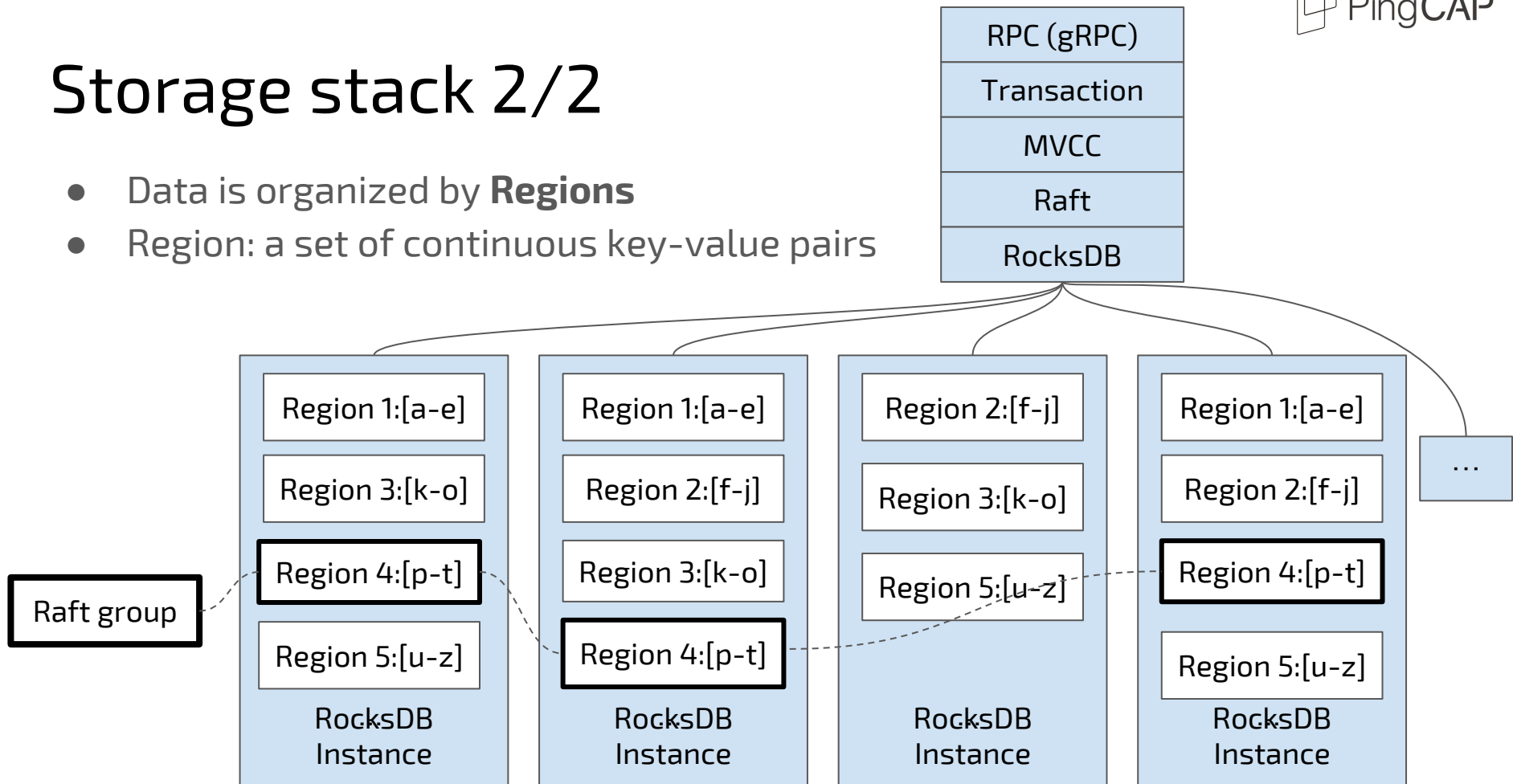
Storage stack 1/2

- TiKV is the underlying storage layer
- Physically, data is stored in RocksDB
- We build a Raft layer on top of RocksDB
 - What is Raft?
- Written in Rust!



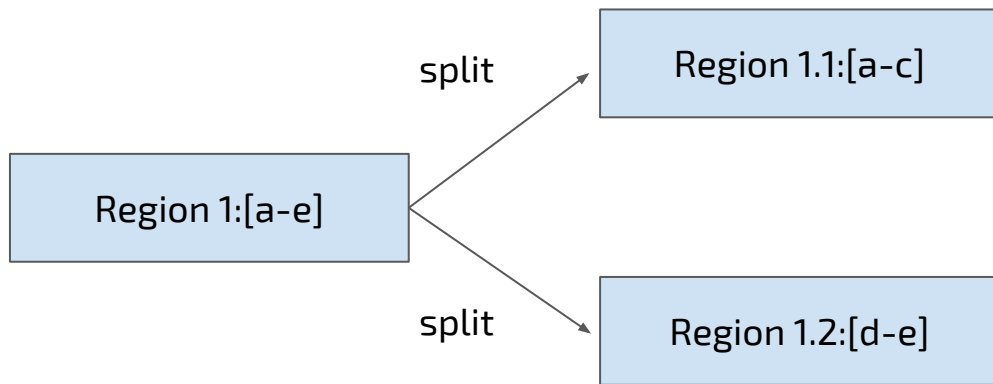
Storage stack 2/2

- Data is organized by **Regions**
- Region: a set of continuous key-value pairs



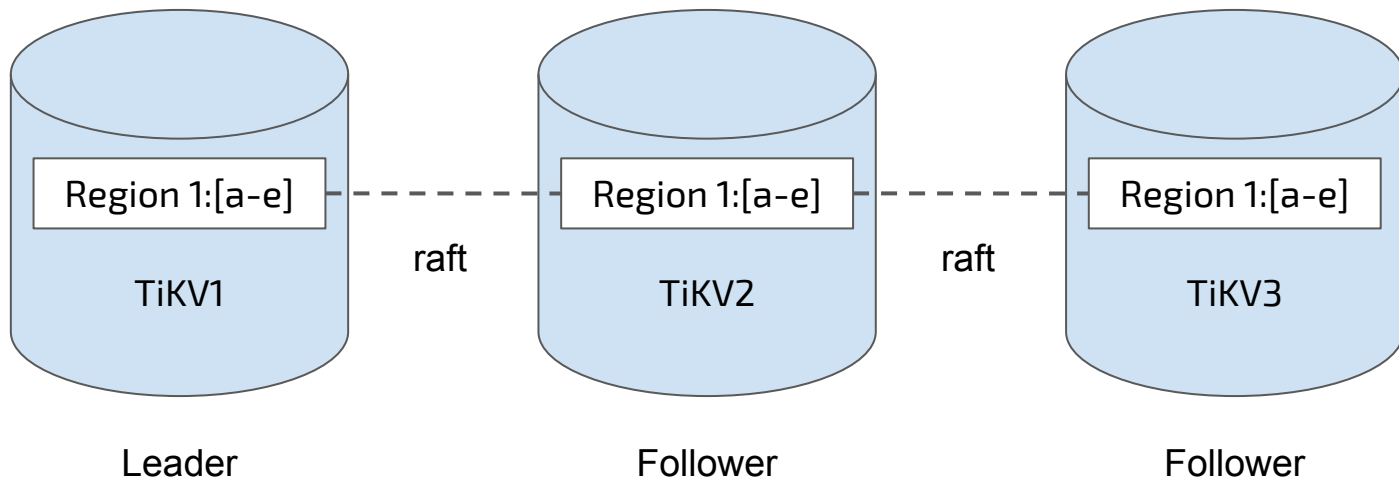
Dynamic Multi-Raft

- What's Dynamic Multi-Raft?
 - Dynamic split / merge
- Safe split / merge

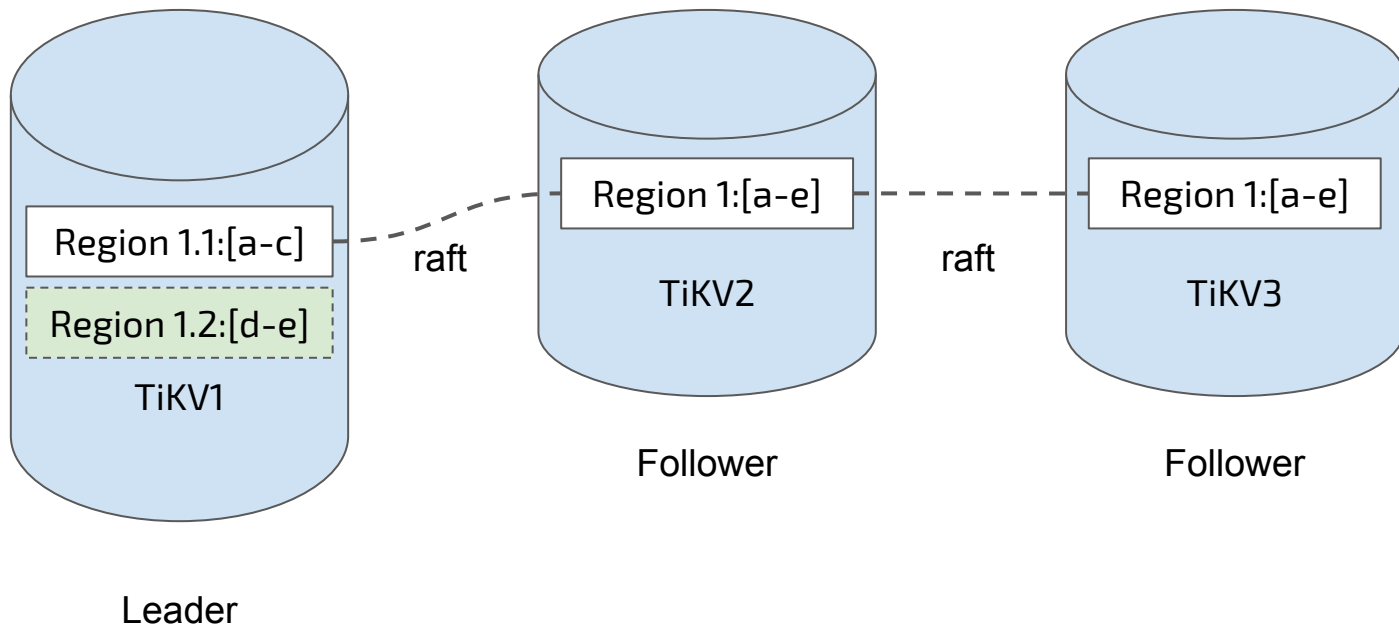


Safe Split: 1/4

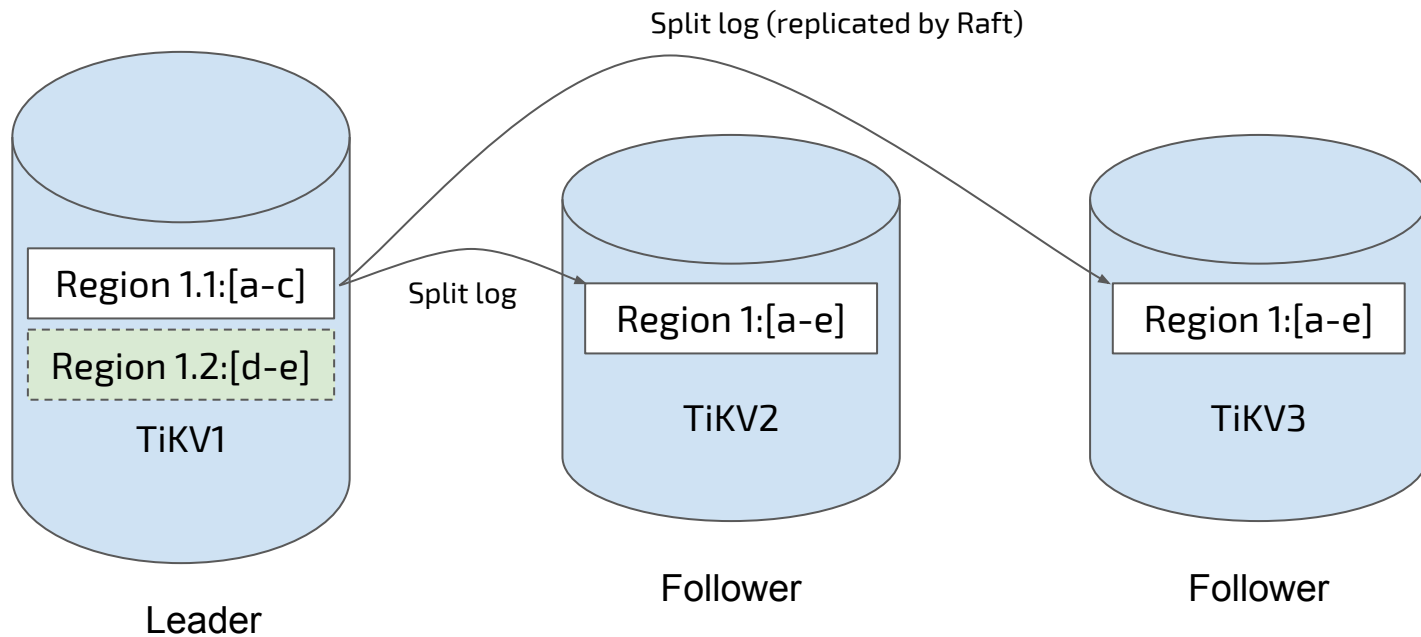
Raft group



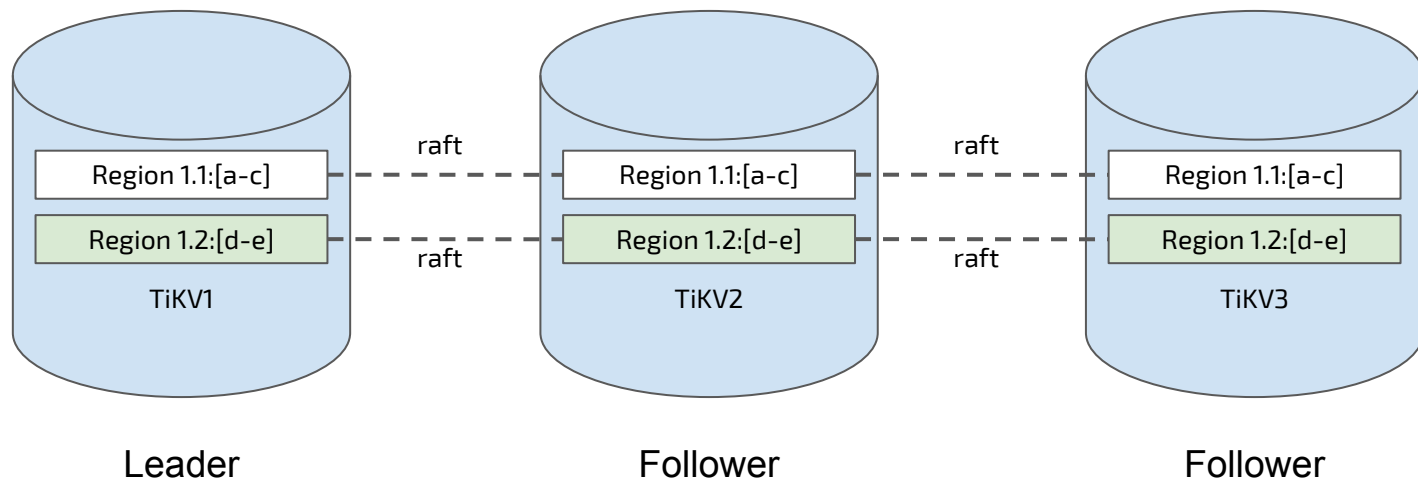
Safe Split: 2/4



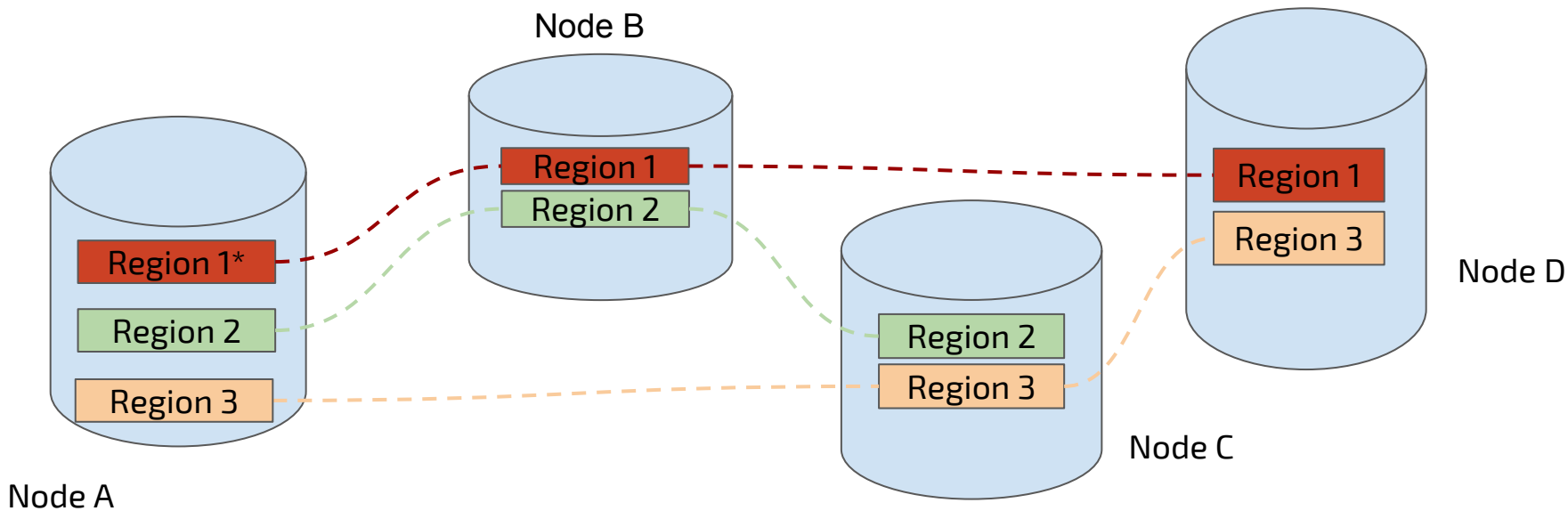
Safe Split: 3/4



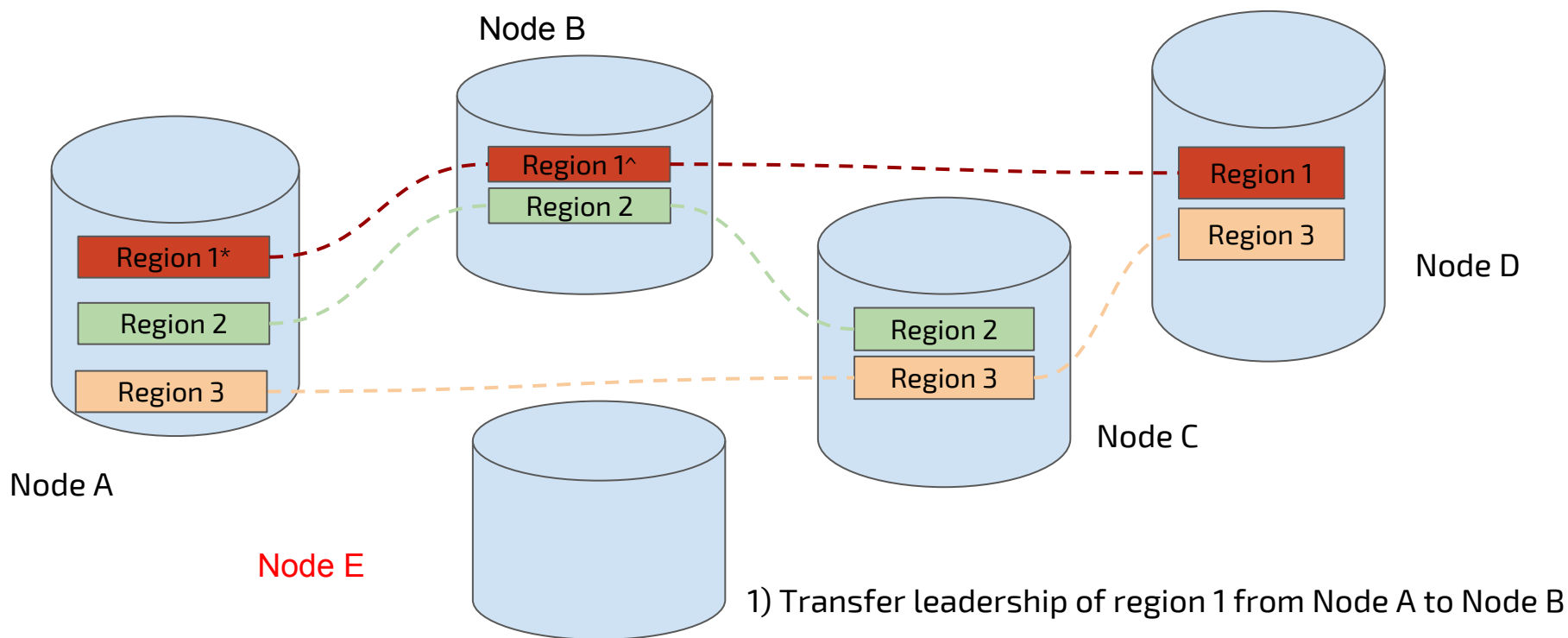
Safe Split: 4/4



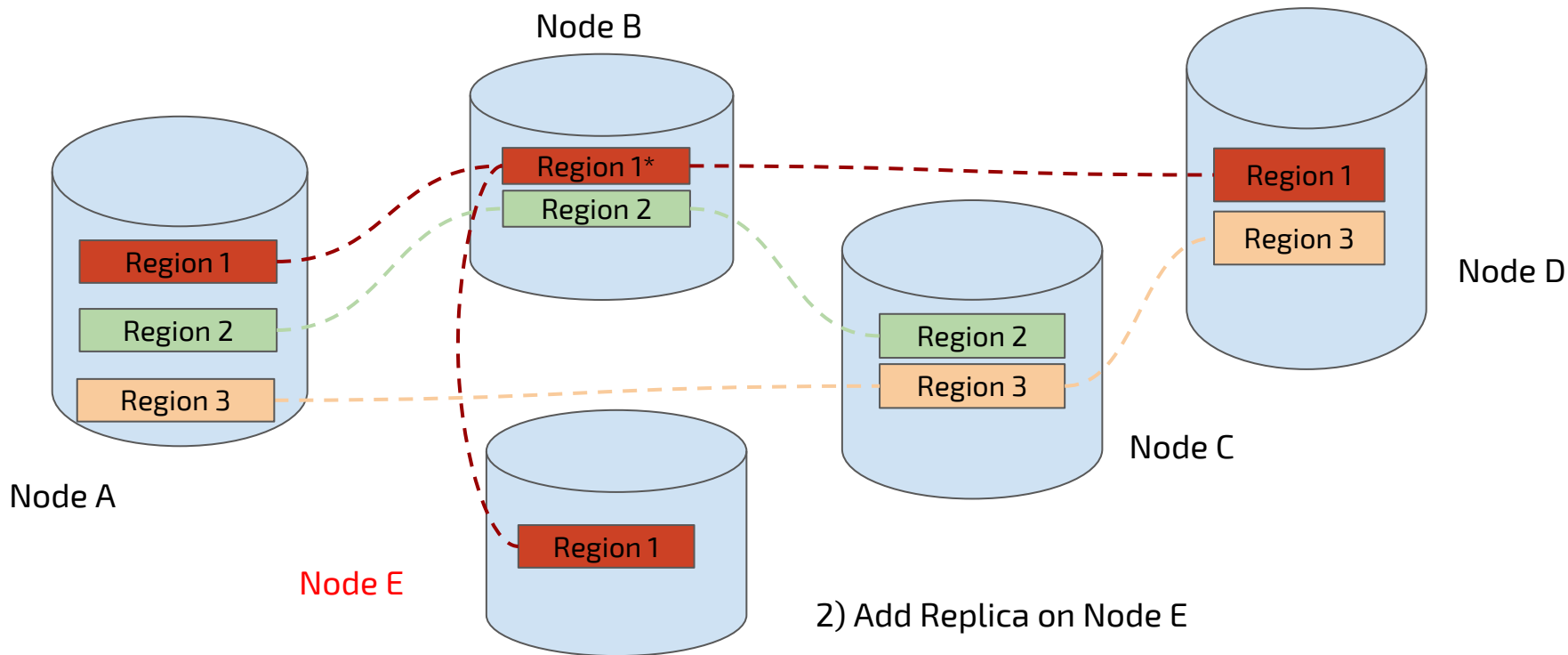
Scale-out (initial state)



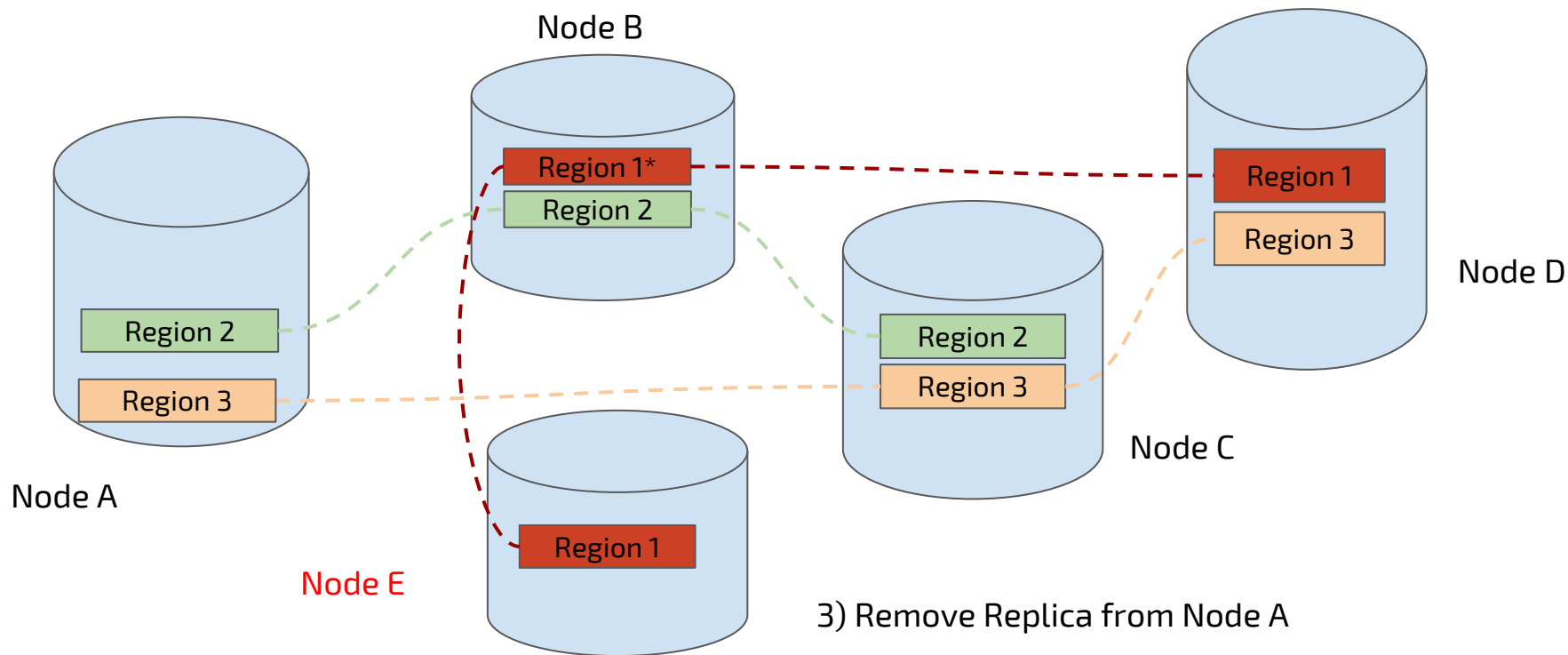
Scale-out (add new node)



Scale-out (balancing)



Scale-out (balancing)



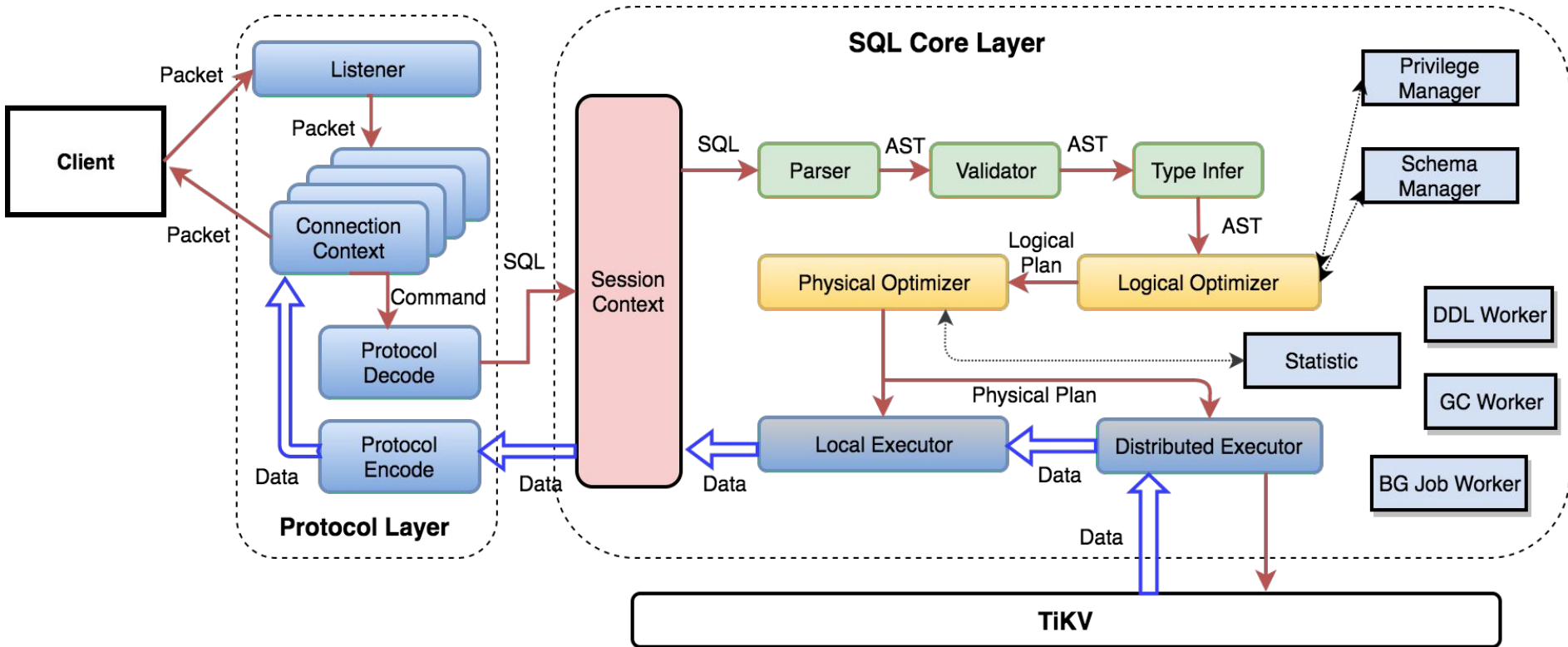
ACID Transaction

- Based on Google Percolator
- 'Almost' decentralized 2-phase commit
 - Timestamp Allocator
- Optimistic transaction model
- Default isolation level: Repeatable Read
- External consistency: Snapshot Isolation + Lock
 - `SELECT ... FOR UPDATE`

Distributed SQL

- Full-featured SQL layer
- Predicate pushdown
- Distributed join
- Distributed cost-based optimizer (Distributed CBO)

TiDB SQL Layer overview

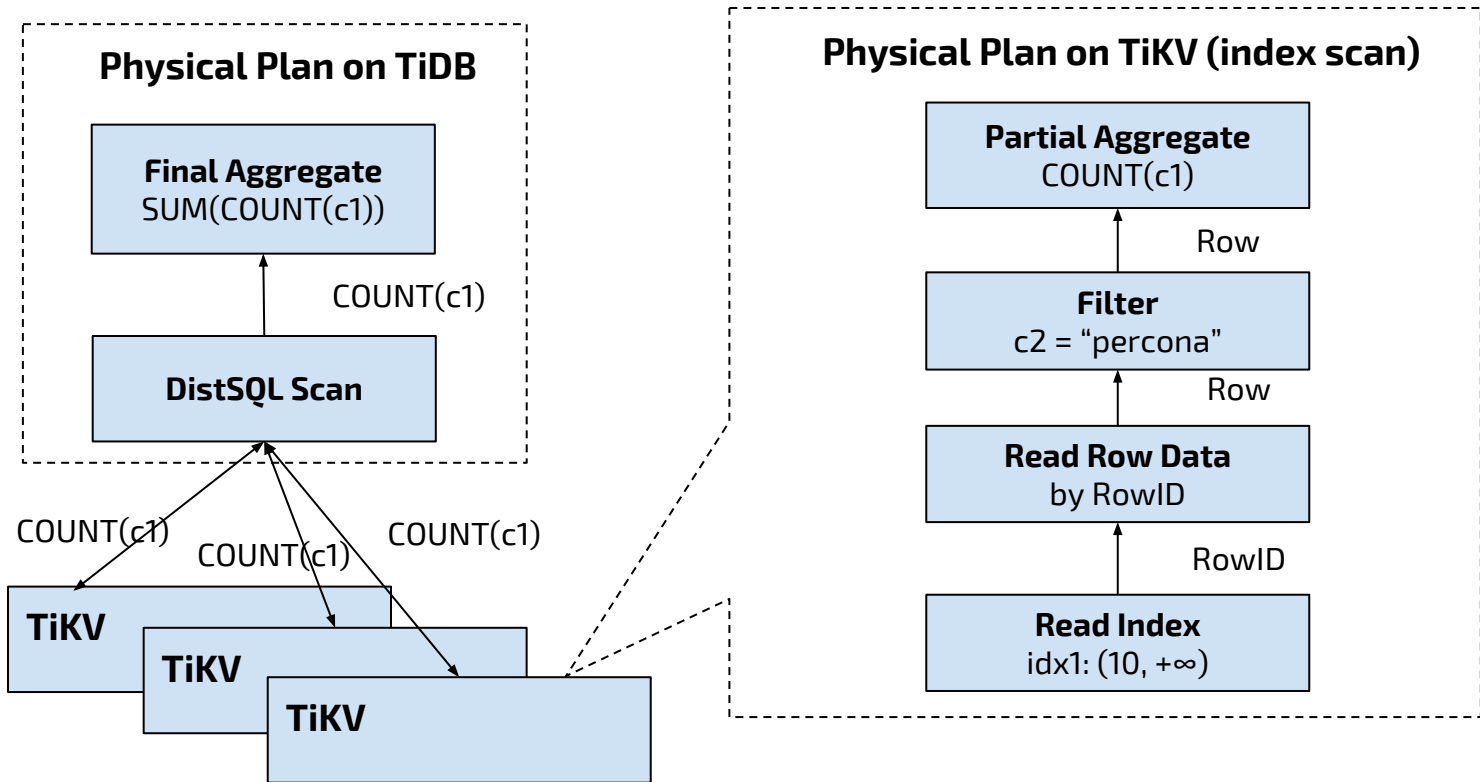


What happens behind a query

```
CREATE TABLE t (c1 INT, c2 TEXT, KEY idx_c1(c1));
```

```
SELECT COUNT(c1) FROM t WHERE c1 > 10 AND c2 = 'percona';
```

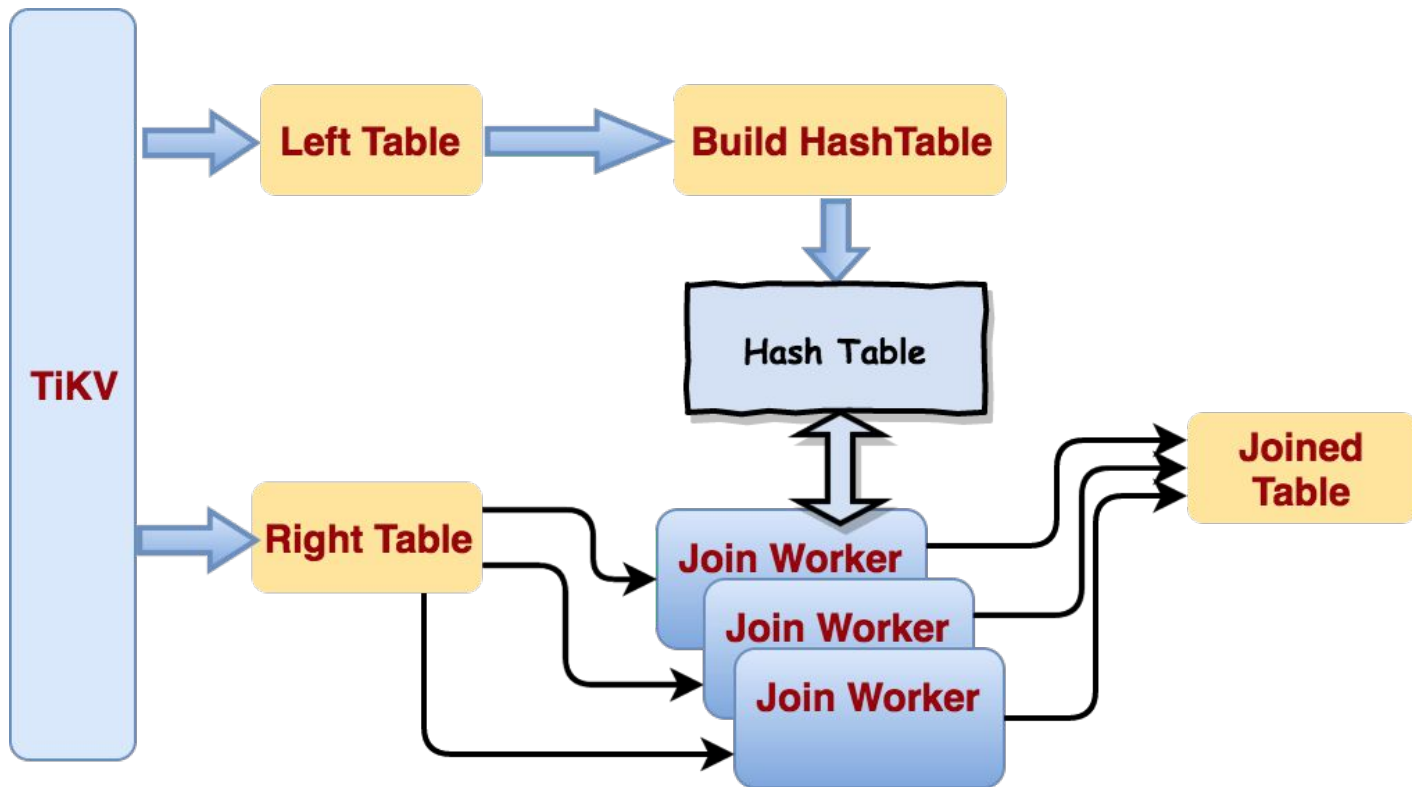
Query Plan



What happens behind a query

```
CREATE TABLE left (id INT, email TEXT, KEY idx_id(id));  
CREATE TABLE right (id INT, email TEXT, KEY idx_id(id));  
SELECT * FROM left join right WHERE left.id = right.id;
```

Distributed Join (HashJoin)



Supported Distributed Join Type

- Hash Join
- Sort merge Join
- Index-lookup Join

No silver bullet (anti-patterns for TiDB SQL)

- Join between large tables without index or any hints
- Get distinct values from large tables without index
- Sort without index
- Result set is too large (forget LIMIT N?)

Best practices

- **Random**, massive, read / write workload
- No hot small table
- Use transaction, but not much conflicts

Tools matter

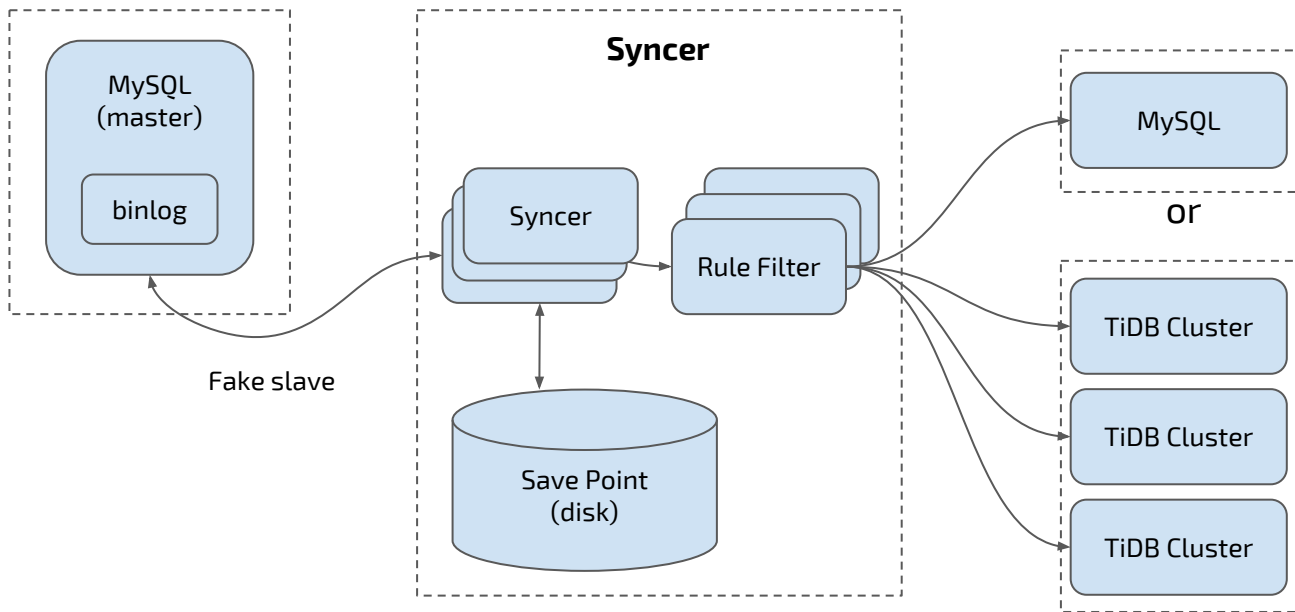
- Syncer
- TiDB-Binlog
- Mydumper/MyLoader(loader)

Open sourced, too.

<https://github.com/pingcap/tidb-tools>

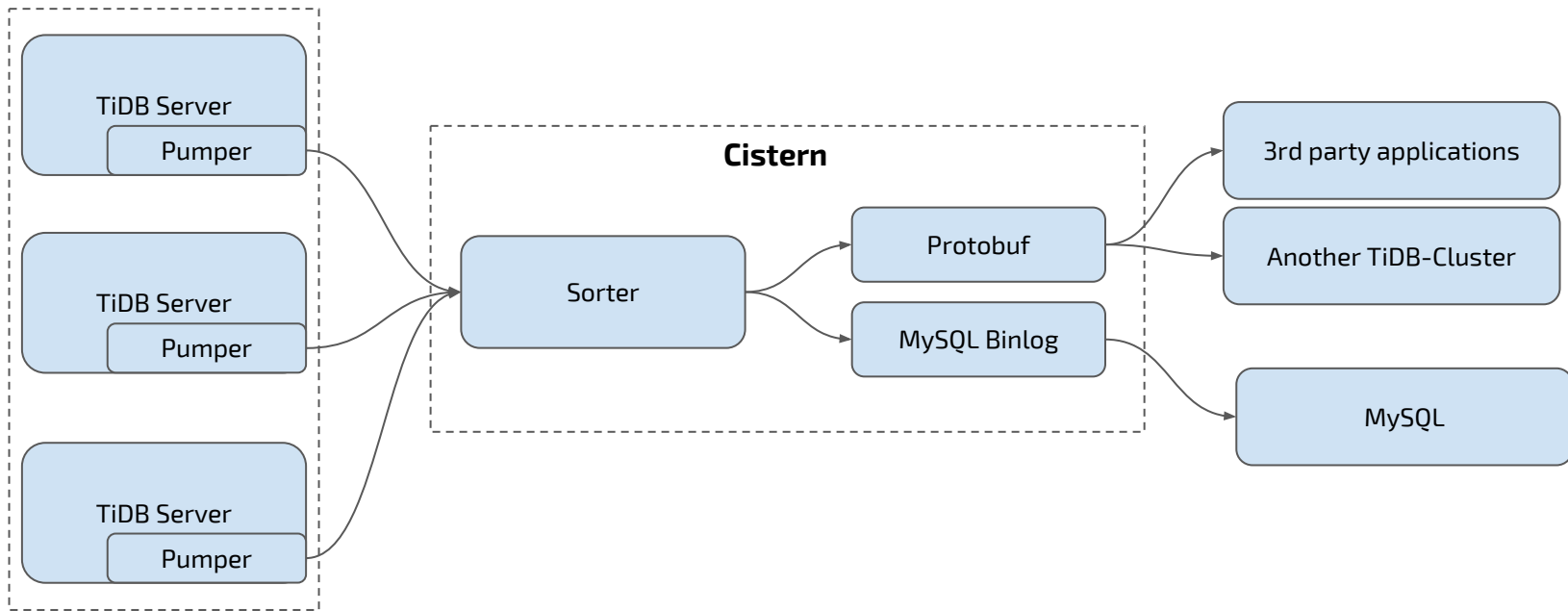
Syncer

- Synchronize data from MySQL in real-time
- Hook up as a MySQL replica



TiDB-Binlog

- Subscribe the incremental data from TiDB
- Output Protobuf formatted data or MySQL Binlog format(WIP)

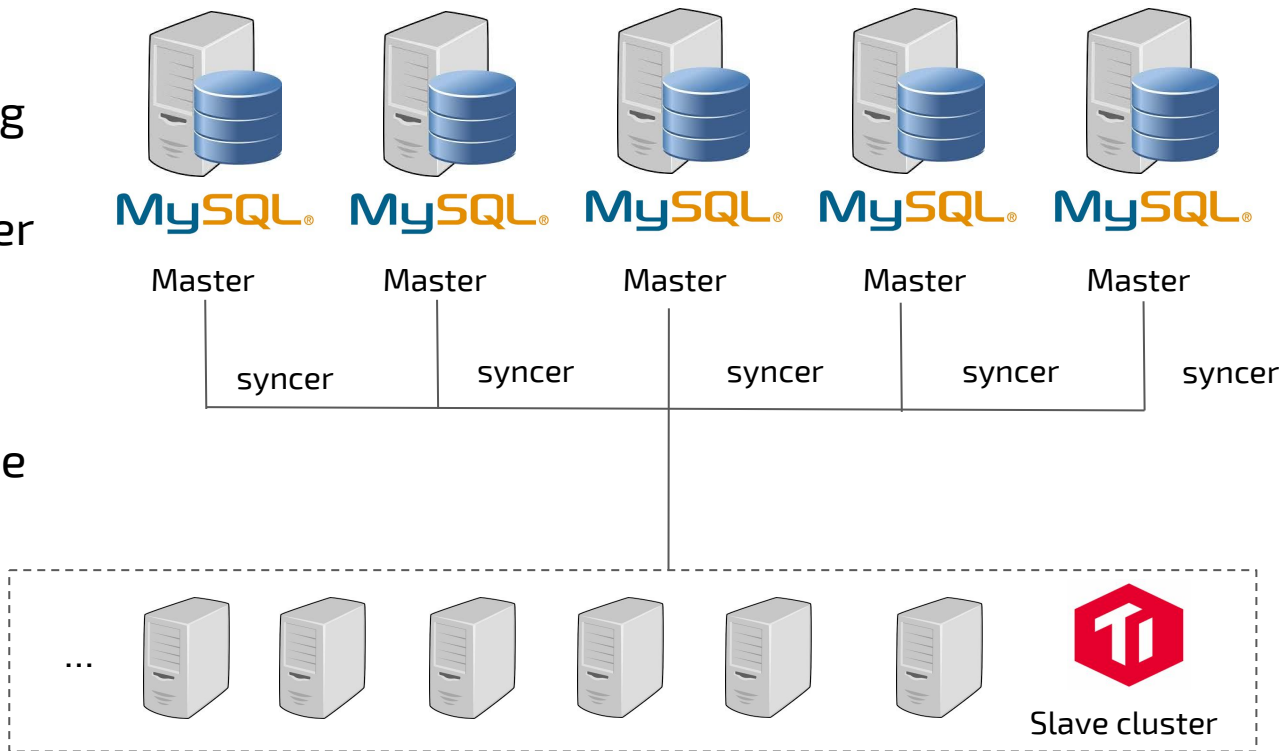


MyDumper / Loader

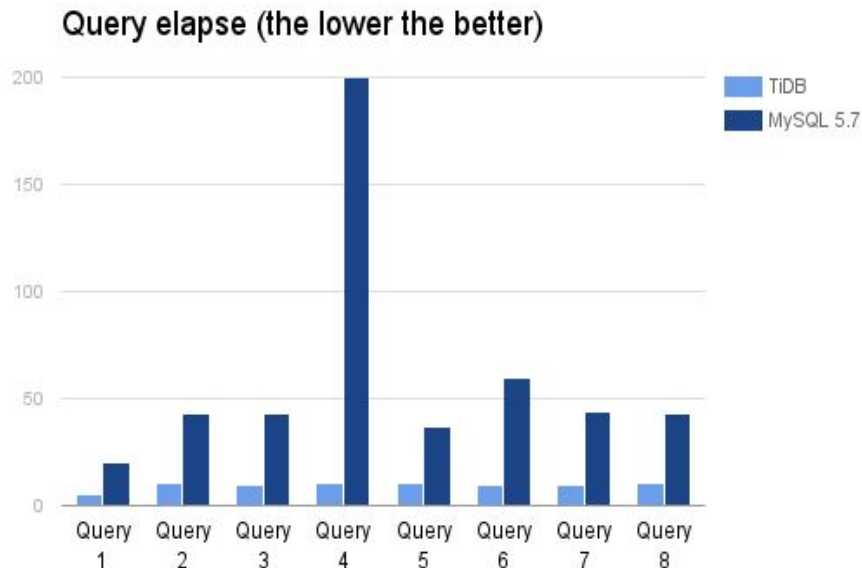
- Backup/restore in parallel
- Works for TiDB too
- Actually, we don't have our own data migration tool for now

Use case 1: OLTP + OLAP

- One of the most popular bike sharing companies in China
- 7-nodes TiDB cluster for order storage (OLTP).
- Hook up as MySQL Replica, synchronize data to a 10-nodes TiDB cluster for Ad-hoc OLAP.



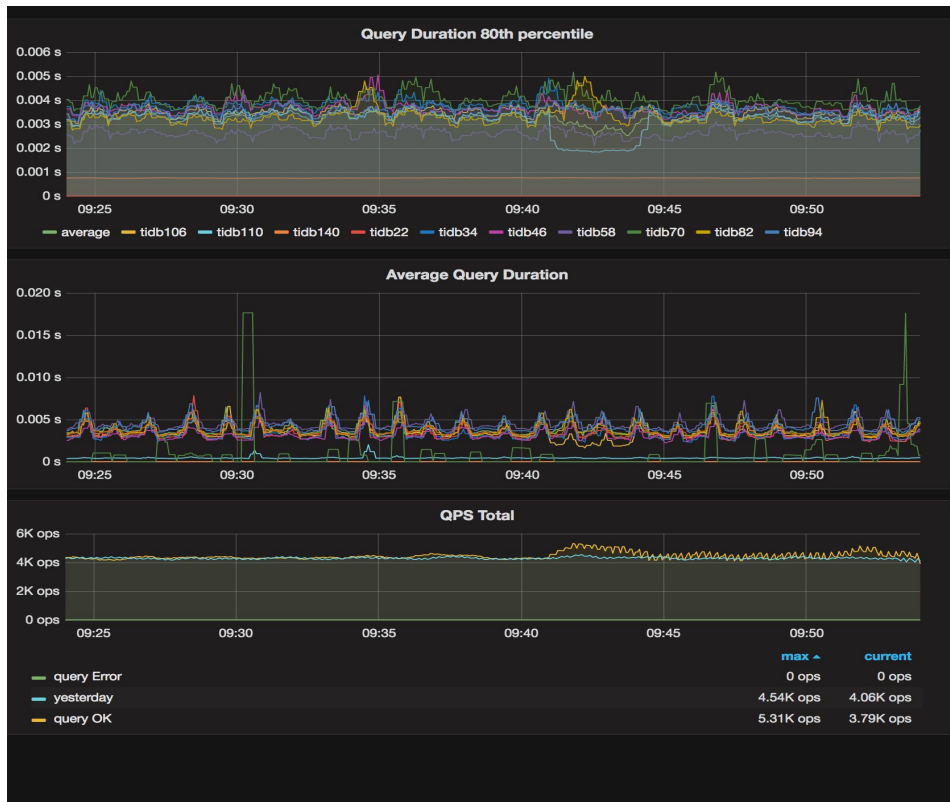
Use case 1: Ad-hoc OLAP

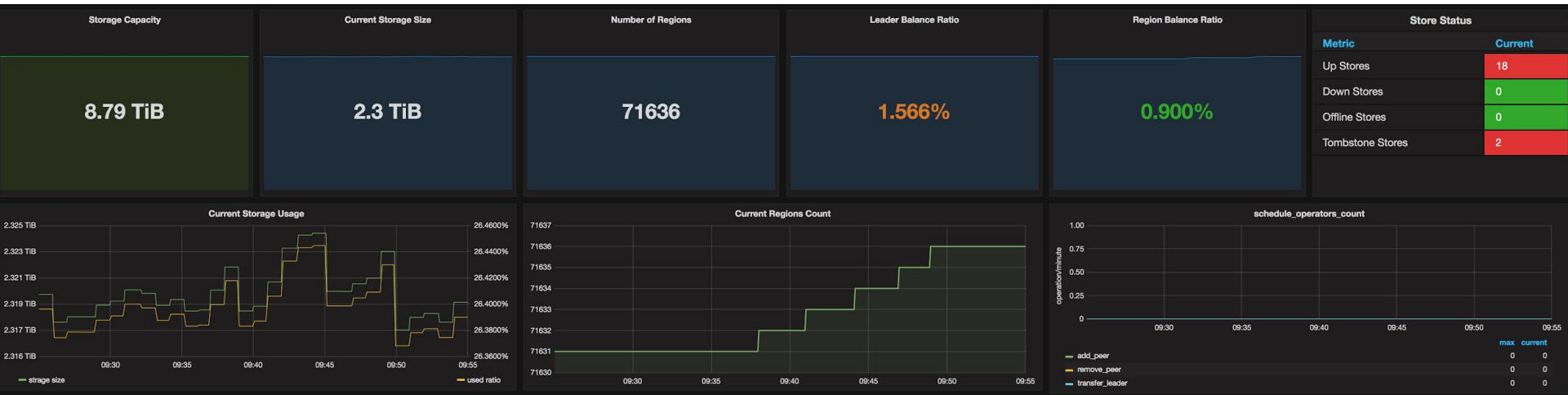


TiDB Elapse (3 nodes)	MySQL Elapse
5.07699437s	19.93s
10.524703077s	43.23s
10.077812714s	43.33s
10.285957629s	>20 mins
10.462306097s	36.81s
9.968078965s	1 min 0.27 sec
9.998030375s	44.05s
10.866549284s	43.18s

Use case 2: Distributed OLTP

- One of the biggest MMORPG game in China.
- 2.2 T, 18 nodes.
- Drop-in replacement for MySQL
- Distributed OLTP





Sysbench

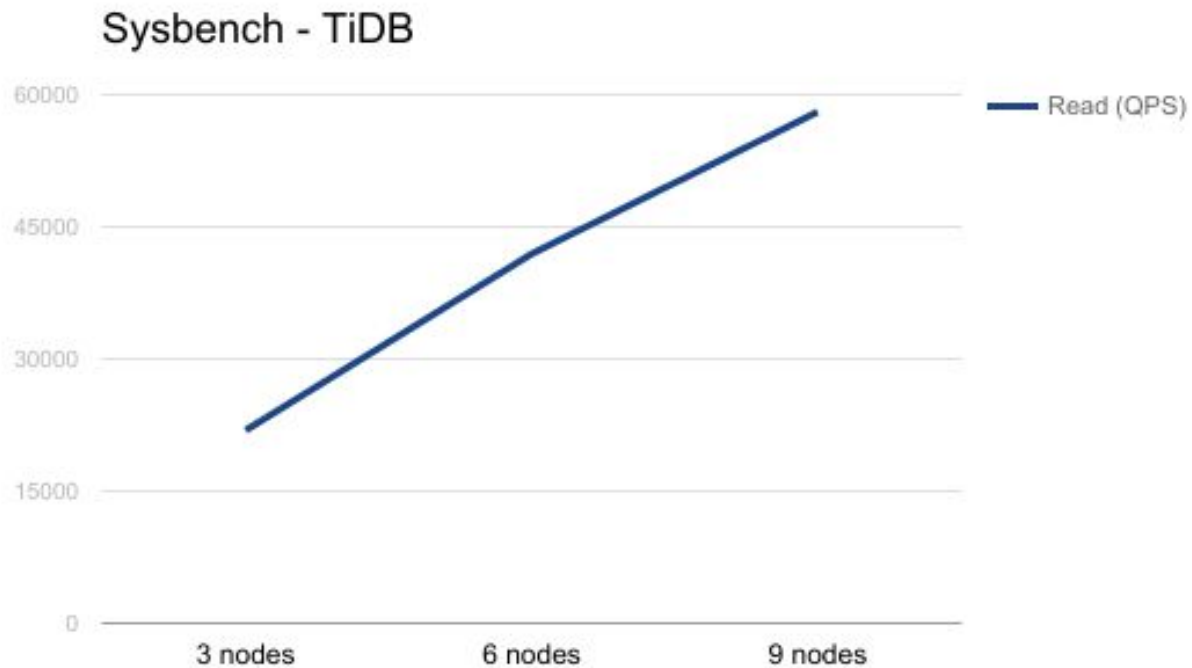
OS	linux (ubuntu 14.04)
CPU	28 ECUs, 8 vCPUs, 2.8 GHz, Intel Xeon E5-2680v2
RAM	16 G
DISK	80 G (SSD)

Notice: 3 replicas

Sysbench (Read)

	table count	table size	sysbench threads	qps	latency(avg/.95)
3 nodes	16	1M rows	256	21899.59	11.69ms / 19.87ms
6 nodes	16	1M rows	256	41928.84	6.10ms / 10.96ms
9 nodes	16	1M rows	256	58044.80	4.41ms / 7.36ms

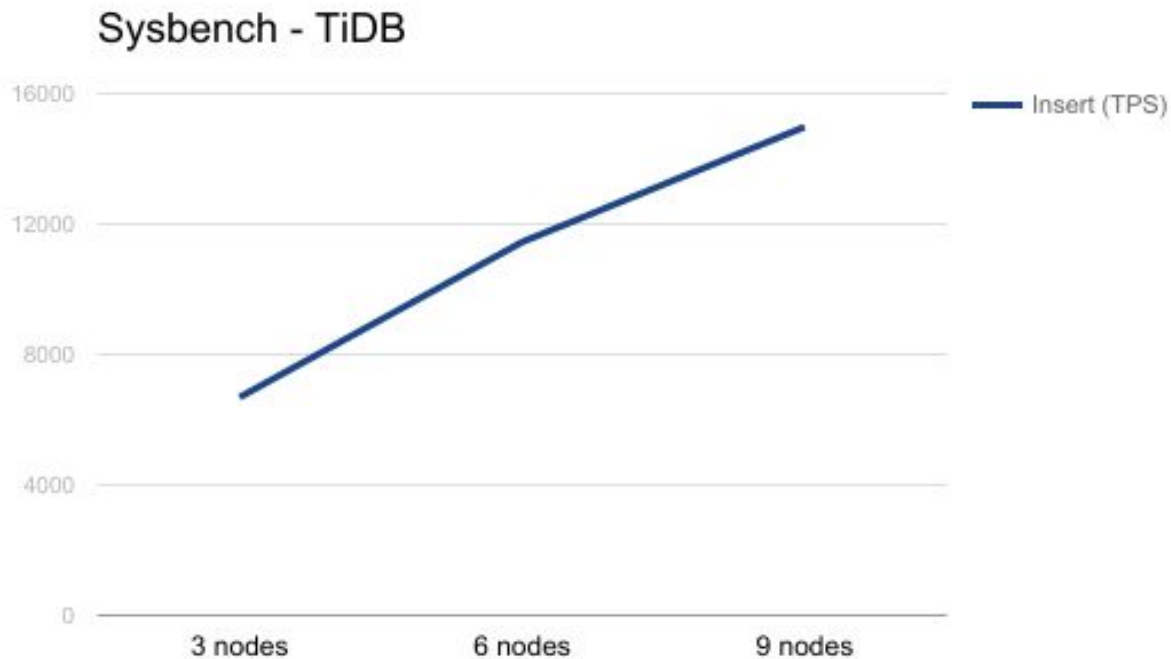
Systembench (Read)



Sysbench (Insert)

	table count	table size	sysbench threads	TPS	latency(avg/.95)
3 nodes	16	1M rows	256	6686.59	38.28ms / 78.21ms
6 nodes	16	1M rows	256	11448.08	22.36ms / 44.61ms
9 nodes	16	1M rows	512	14977.01	34.18ms / 86.85ms

Sysbench (Insert)



Roadmap

- TiSpark: Integrate TiKV with SparkSQL
- Better optimizer (Statistic && CBO)
- Json type and document store for TiDB
 - MySQL 5.7.12+ X-Plugin
- Integrate with Kubernetes
 - Operator by CoreOS

Thanks

<https://github.com/pingcap/tidb>

<https://github.com/pingcap/tikv>

Contact me:

huang@pingcap.com

