

How to build a *Self-Driving* database

An overview

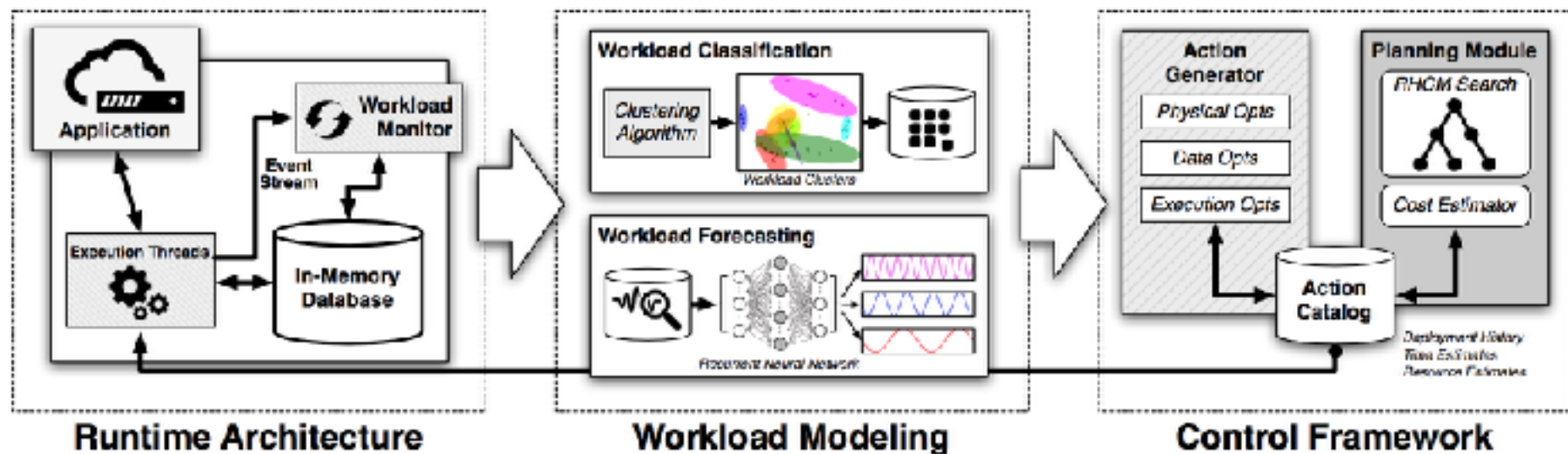
Dongxu

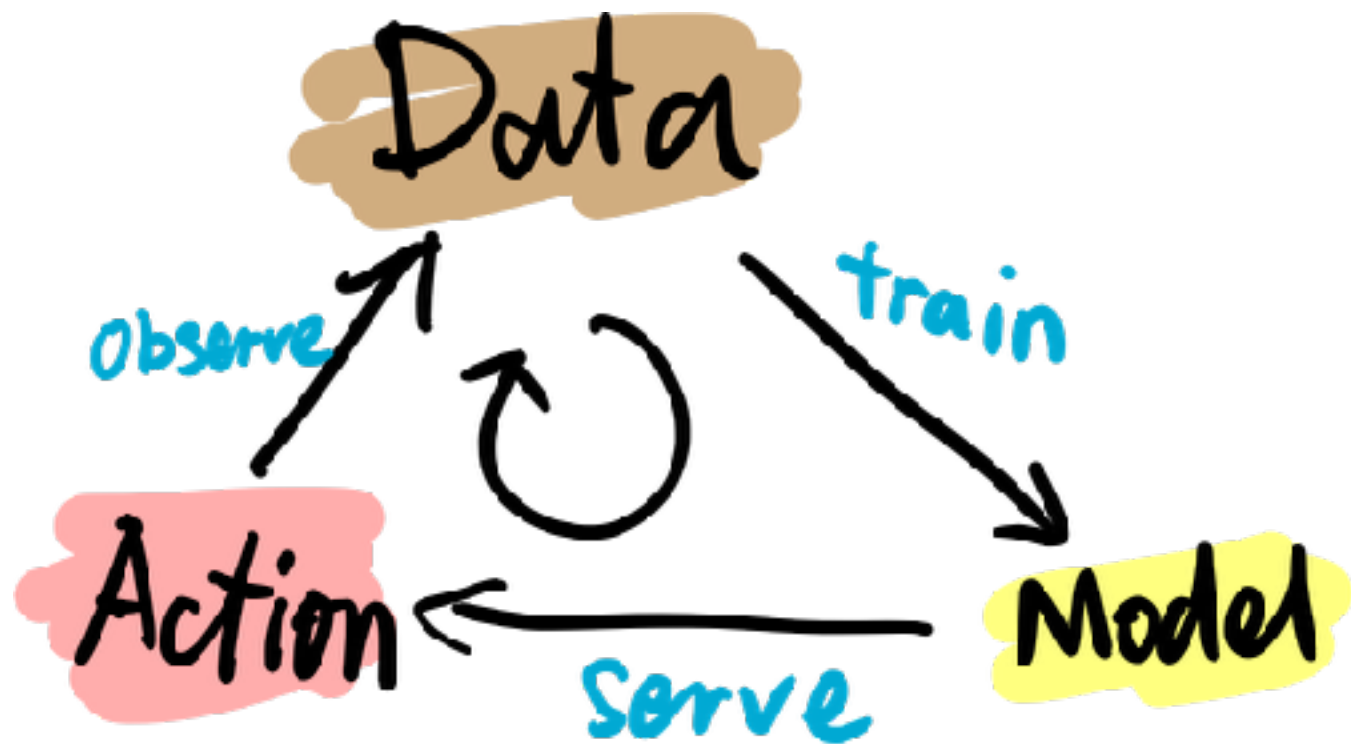
Who am I

- Dongxu Huang
- CTO & Cofounder, PingCAP
- Distributed system engineer / Open source advocator
- h@pingcap.com
- TiDB / TiKV / TiSpark

Self-Driving?





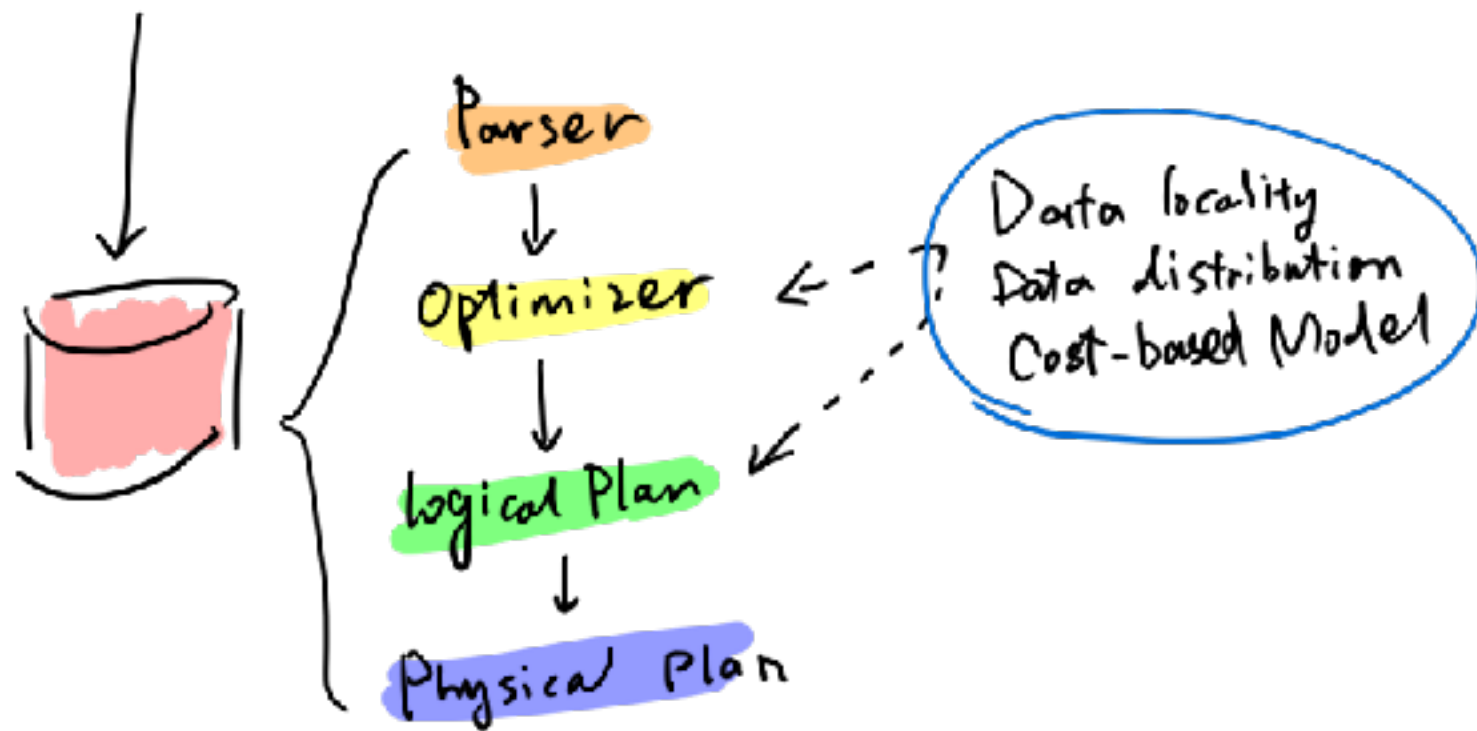


Why now?

- Better hardware
- Better tools
- Workload is different

SQL Tuning

select * from t where id > 10 and id < 100 and name = 'tom';

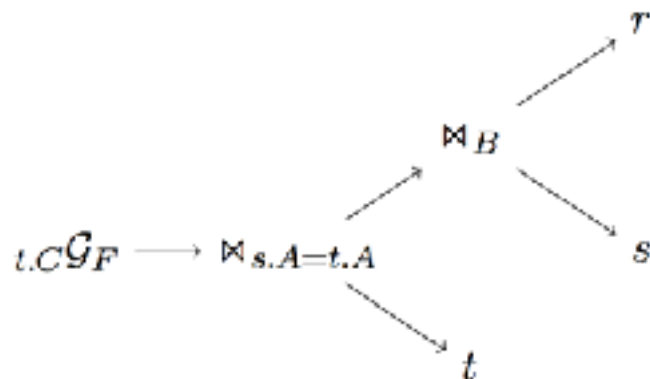


Data : Rows / Index distribution

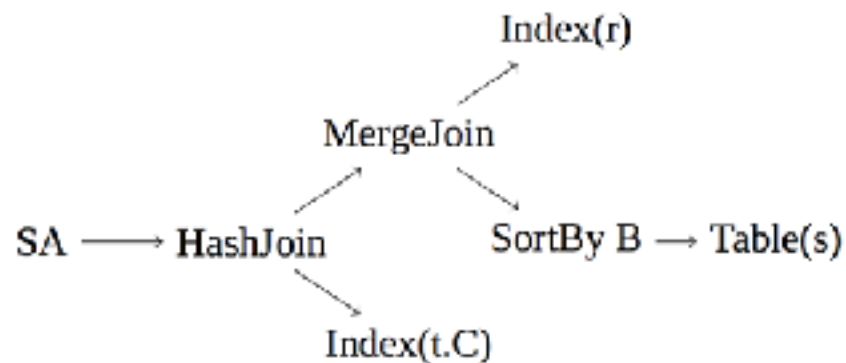
Model : Cost model

Action : Choose a better Physical Plan

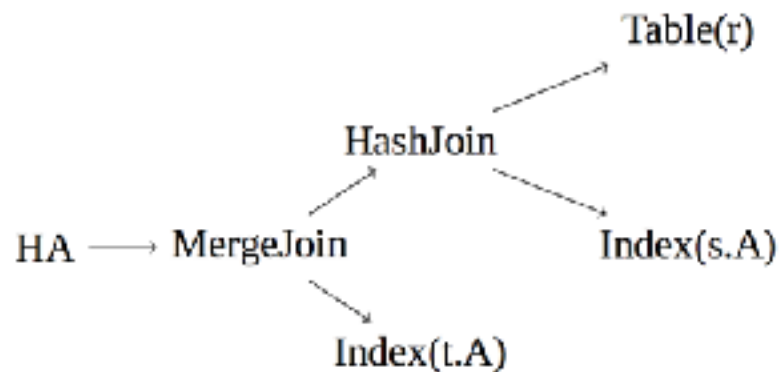
Imagine we got a logical plan:



Its physical plan could be:



or:



Cost estimation

$$Cost(p) = N(p) * F_N + M(p) * F_M + C(p) * F_C$$



Network cost



Memory cost



CPU cost

In TiDB, the default memory factor is 5 and CPU factor is 0.8.

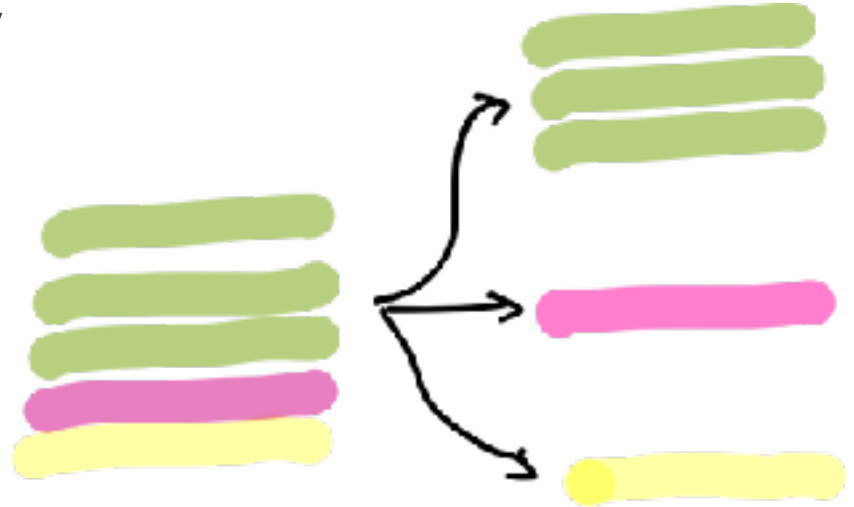
For example: Operator Sort(r), its cost would be:

$$0 + n_r * 5.0 + n_r * \log_2(n_r) * 0.8$$

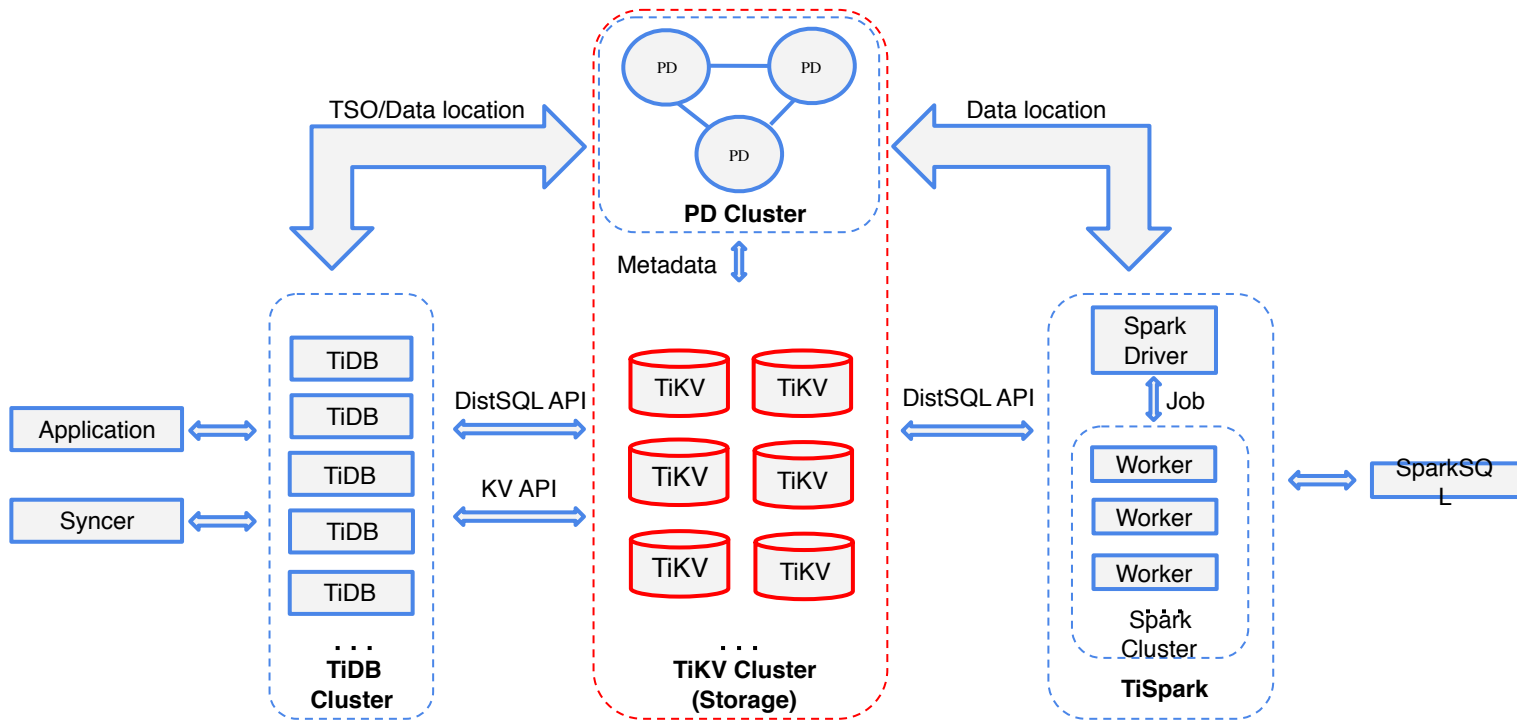
Data placement

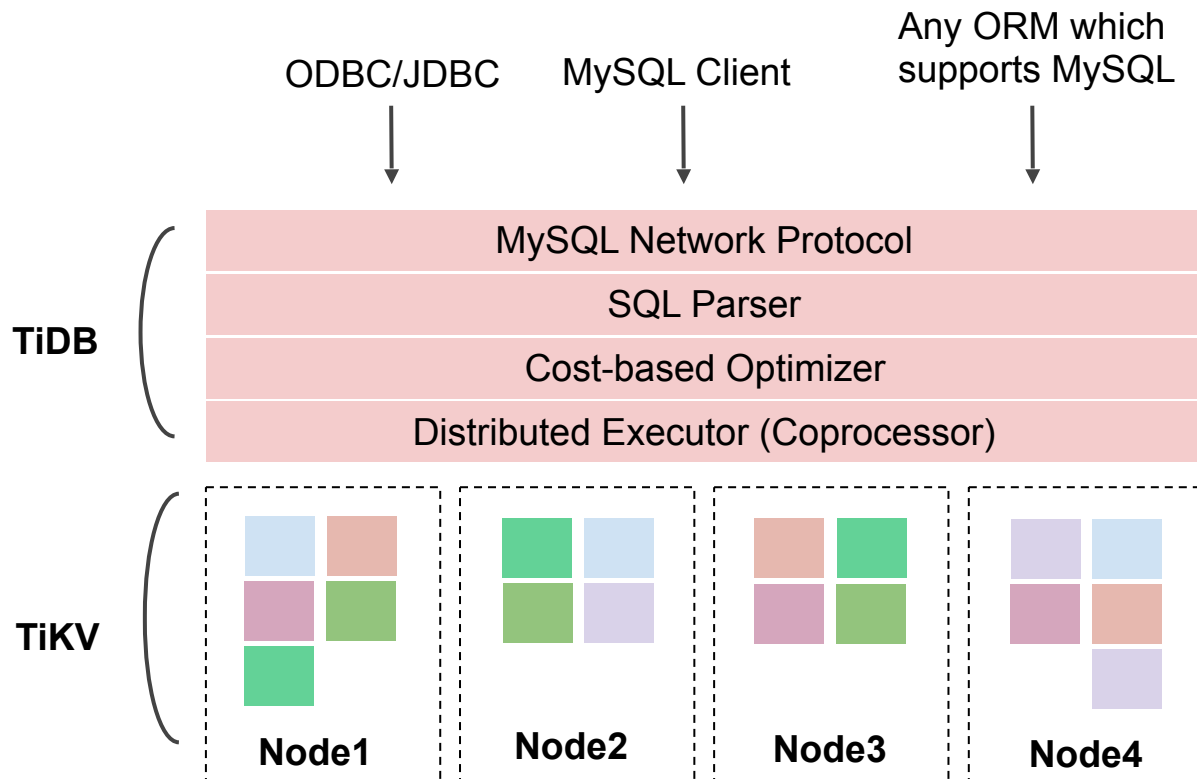
Problem: Uneven Data Distribution

- Dealing with hotspot
 - Choose a wrong sharding key
- Inefficient usage
 - Some are busy
 - Some are idle
- Caused by the nature of RDBMS



A little bit about how data is organized in TiDB





Why Raft?

- Saft split/merge
- Self-healing
- Easy to implement

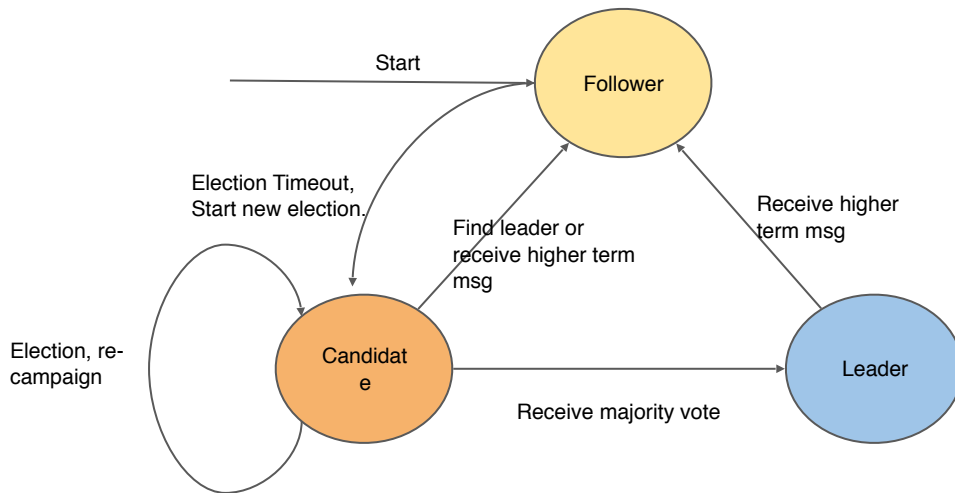
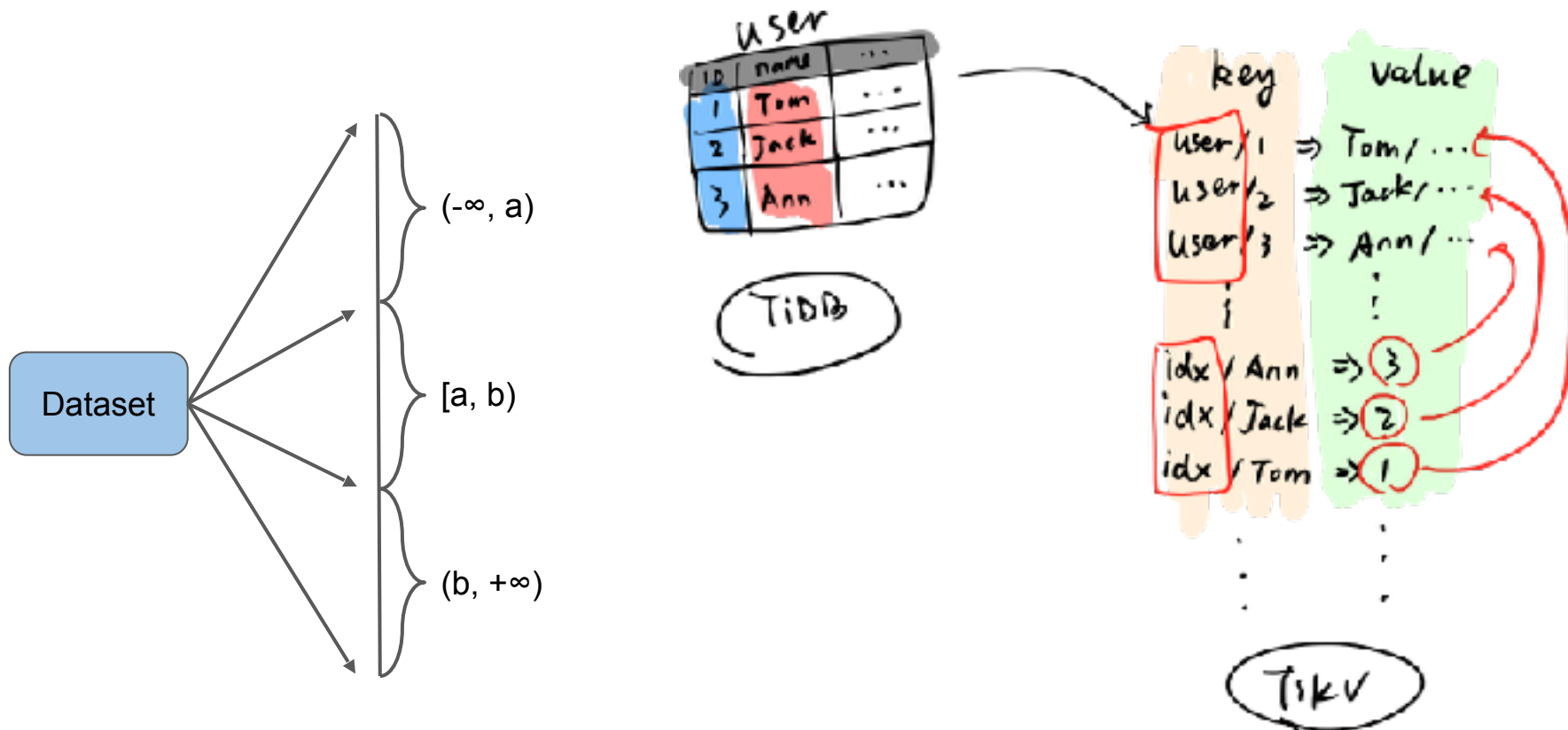
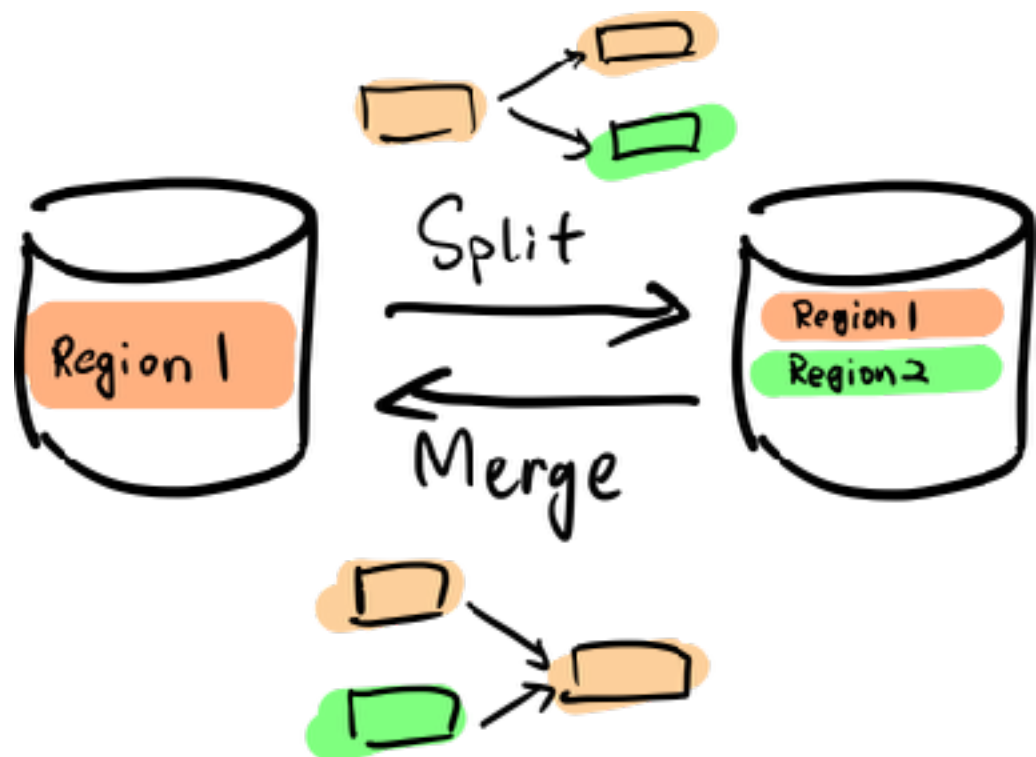
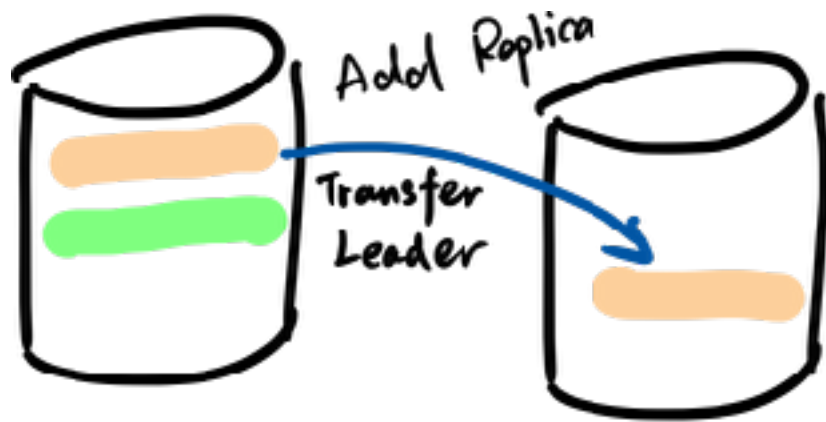


Table mapping: Rows \Rightarrow Key-Value pairs

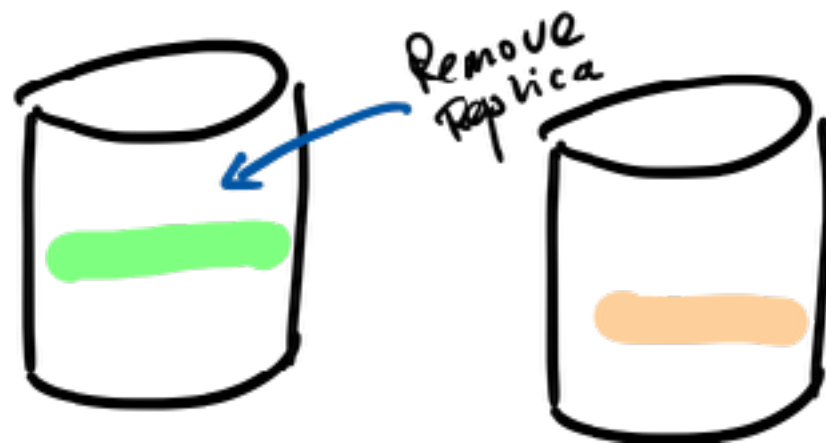


Data movement
Step 1:

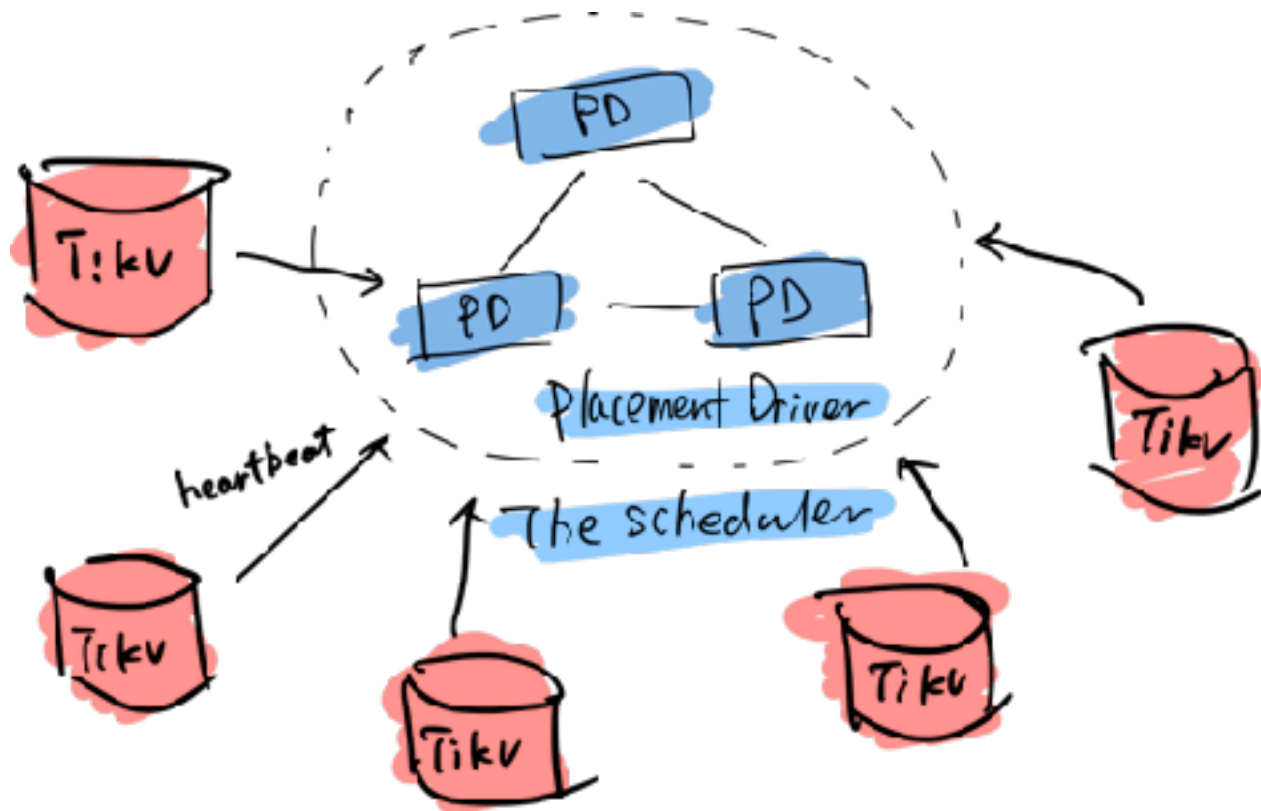




Data movement
Step 2:



Placement Driver: The scheduler

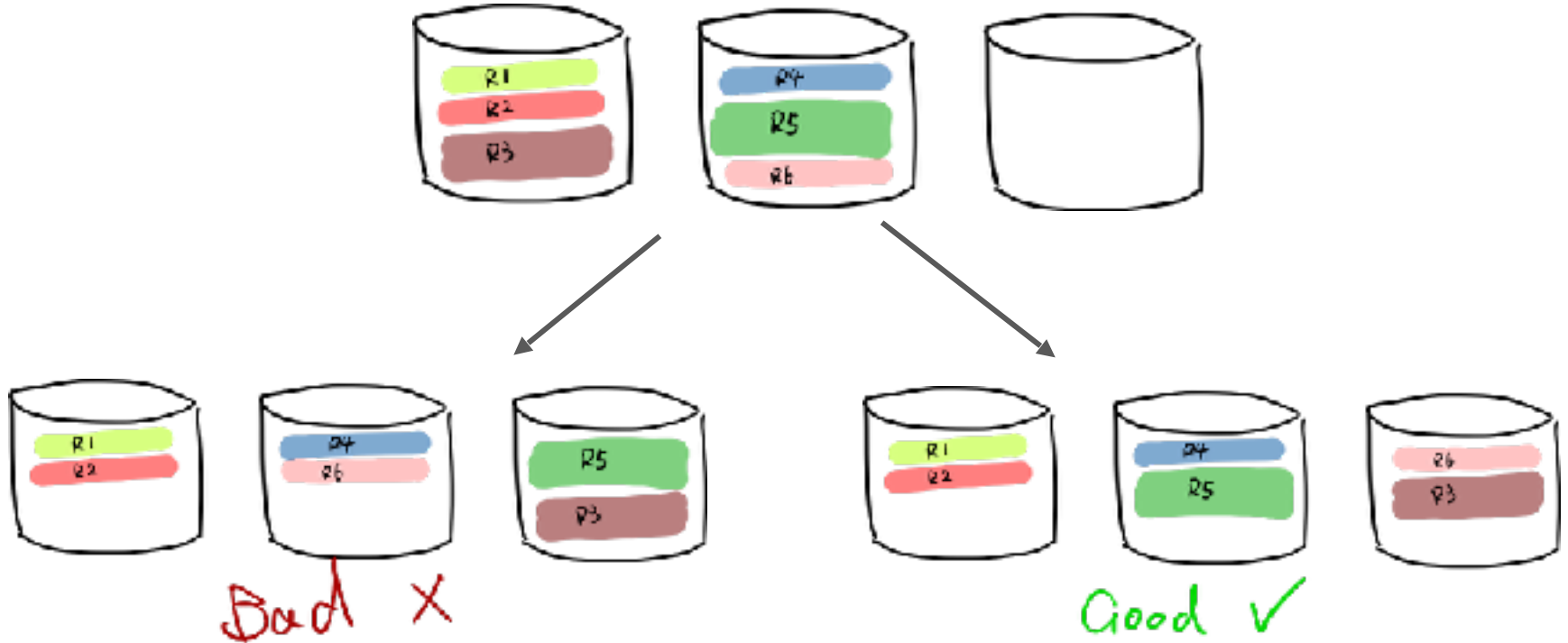


Data : Machine load Info / Disk capacity
Region Info / load on Region
... of each node via heartbeat

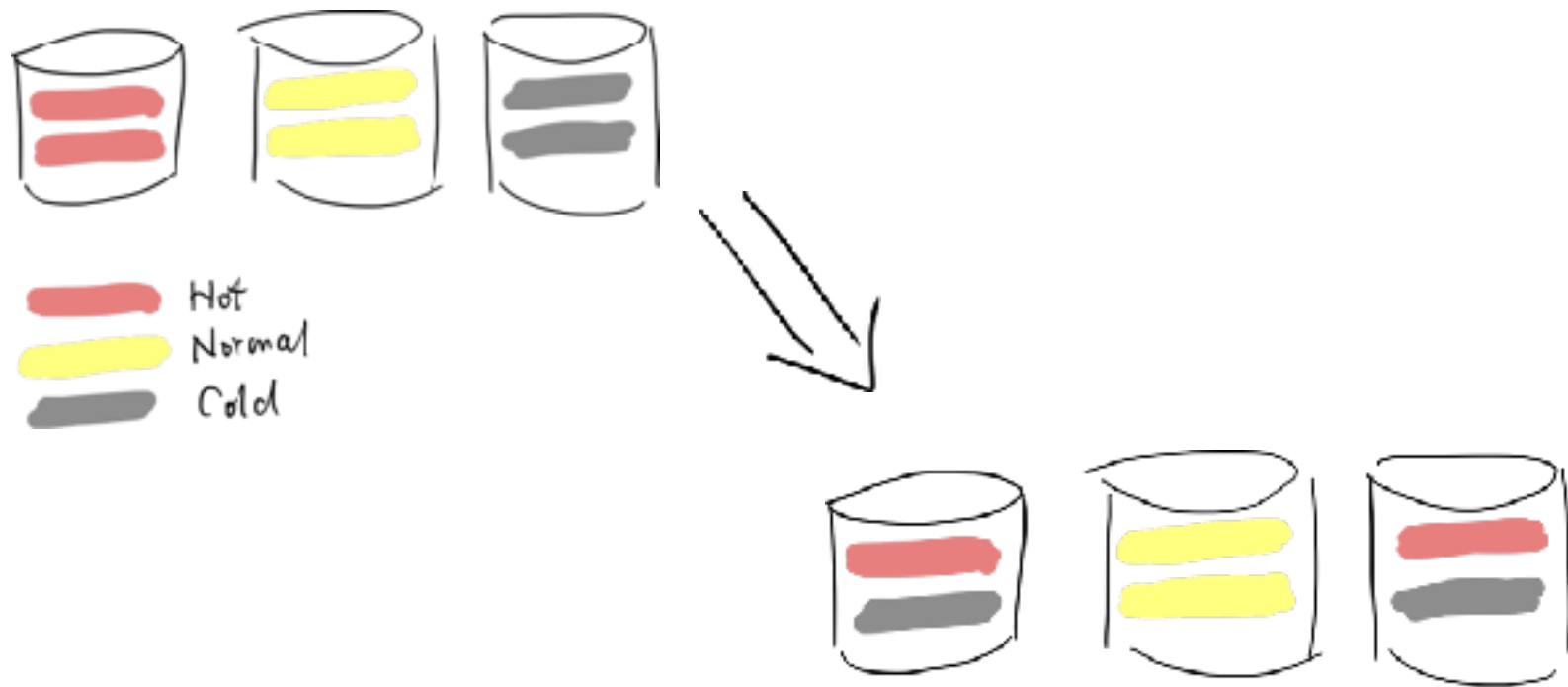
Model : Rules

Action : Add Replica
Remove Replica
Transfer leader
Force Split
⋮

Only by Region count? Size also matters.



Hot Balance



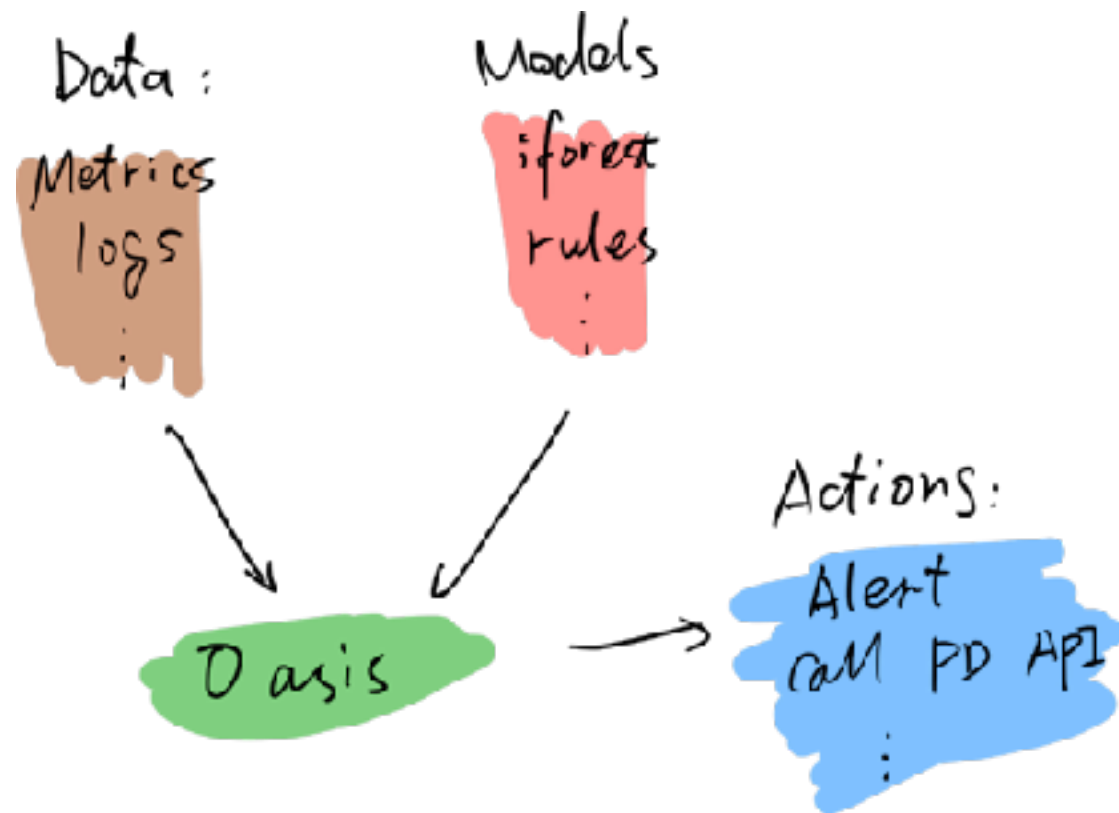
Scheduler - More

- More...
 - Weight Balance - High-weight TiKV will save more data
 - Evict Leader Balance - Some TiKV node can't have any Raft leader
- OpInfluence - Avoid over frequent balancing

Autonomy

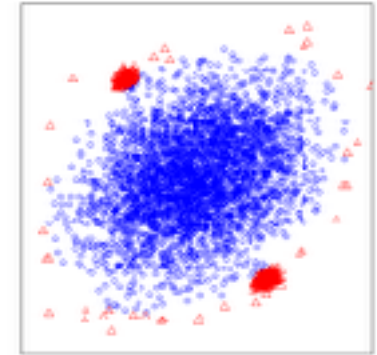
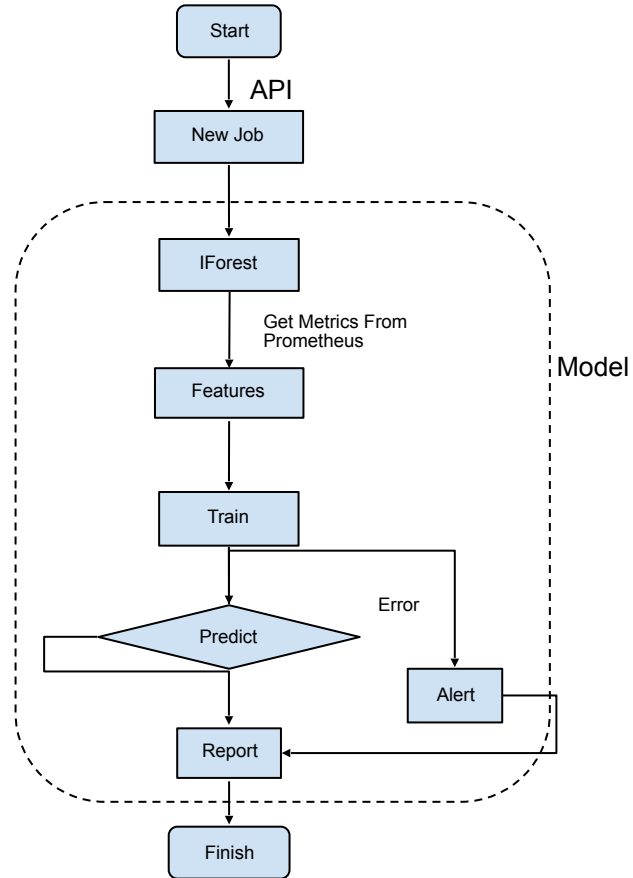


Our work: Oasis

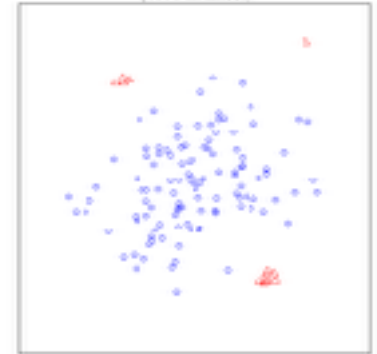


Isolation Forest

- Anomaly detection
- Easy to implement
- [Paper](#)



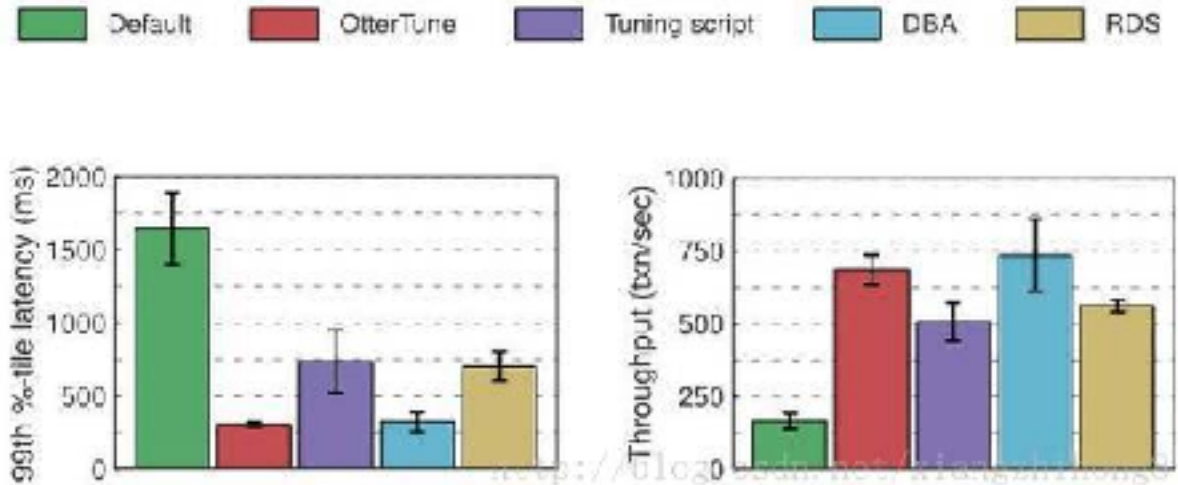
(a) Original sample
(4096 instances)



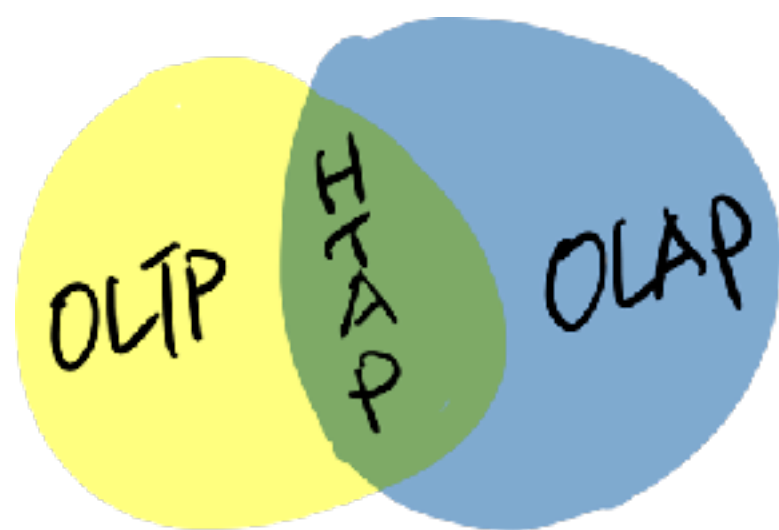
(b) Sub-sample
(128 instances)

Related work - OtterTune

- Database Tuning-as-a-Service
 - Automatically generates DBMS knob configurations
 - Reuse data from previous tuning sessions
- Supported systems
 - PostgreSQL
 - MySQL
 - Greenplum
 - Vectorwise



Serving different workloads at the same time



10:46

< My Posts Details

黄东旭 | PingCAP

我个人是很看好 spark 及 sparksql 一统 clap 的天下。而且 tikv 和 sparksql 的深度融合已经开始，作为进一步完善 tidb 生态的一部分。另外，正榜的时间做正榜的事情重要。

大数据那些事(29)从Spark到Spark

5 March 2017 16:28 Delen



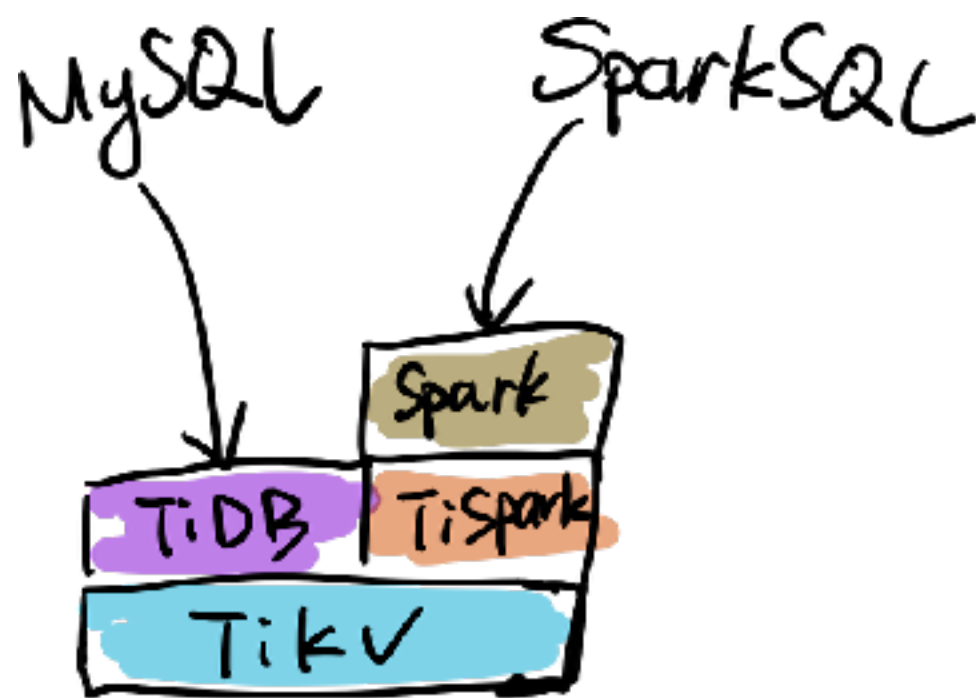
5 March 2017 16:35

tikv作为spark的data source，有点意思。

黄东旭 | PingCAP 5 March 2017 16:38

而且作为一个还算有怀疑精神的人，我一直怀疑 hadoop 在大数据

Comment



[Features](#)[Business](#)[Explore](#)[Marketplace](#)[Pricing](#)[This repository](#)[Sign in](#) or [Sign up](#)[pingcap](#) / [tispark](#)[Watch](#)

26

[★ Star](#)

223

[Fork](#)

51

[Code](#)[Issues](#) 37[Pull requests](#) 3[Projects](#) 0[Wiki](#)[Insights](#)

TiSpark is built for running Apache Spark on top of TiDB/TiKV

[176 commits](#)[10 branches](#)[6 releases](#)[11 contributors](#)[Apache-2.0](#)Branch: [master](#)[New pull request](#)[Find file](#)[Clone or download](#)[birdstorm and lovesoup](#) Make tiDBMapTable replace existing view if any (#350)

Latest commit 61586ed an hour ago

[R](#)

1.1 snapshot prepare (#342)

27 days ago

[config](#)

add tikv, pd, and tidb's config files (#216)

4 months ago

[core](#)

Make tiDBMapTable replace existing view if any (#350)

an hour ago

[docs](#)

1.1 snapshot prepare (#342)

27 days ago

[python](#)

1.1 snapshot prepare (#342)

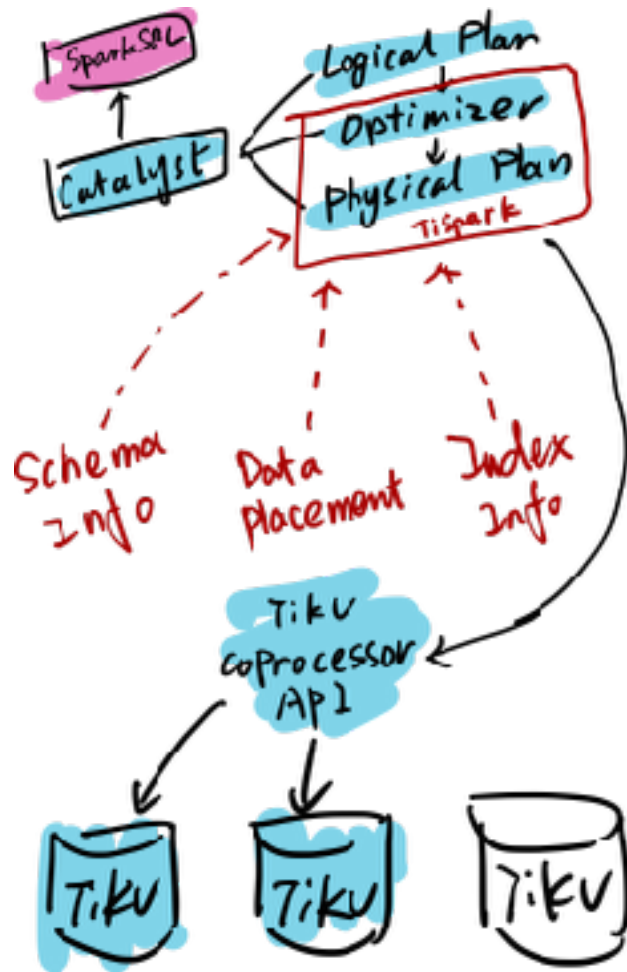
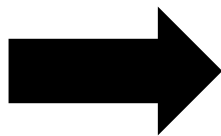
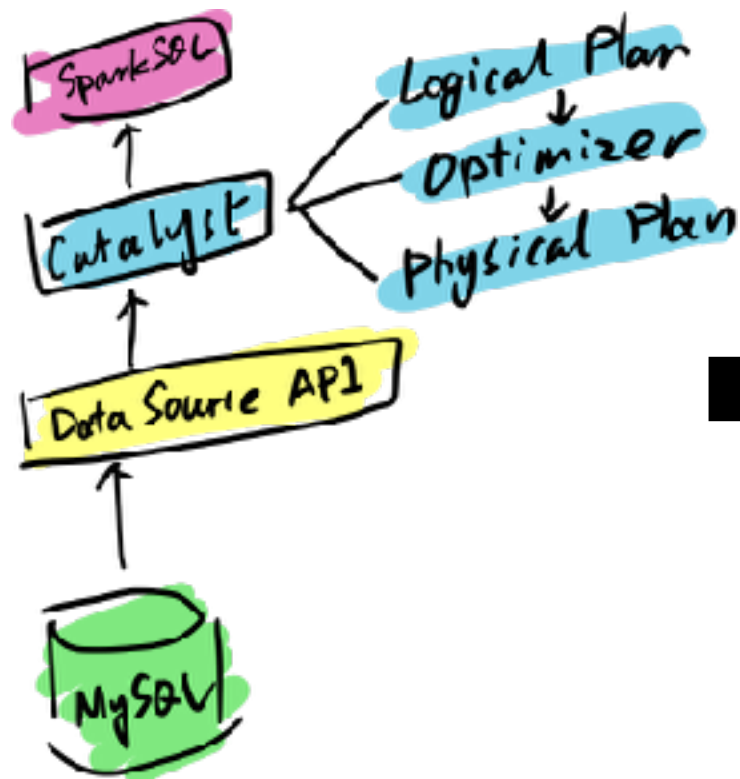
27 days ago

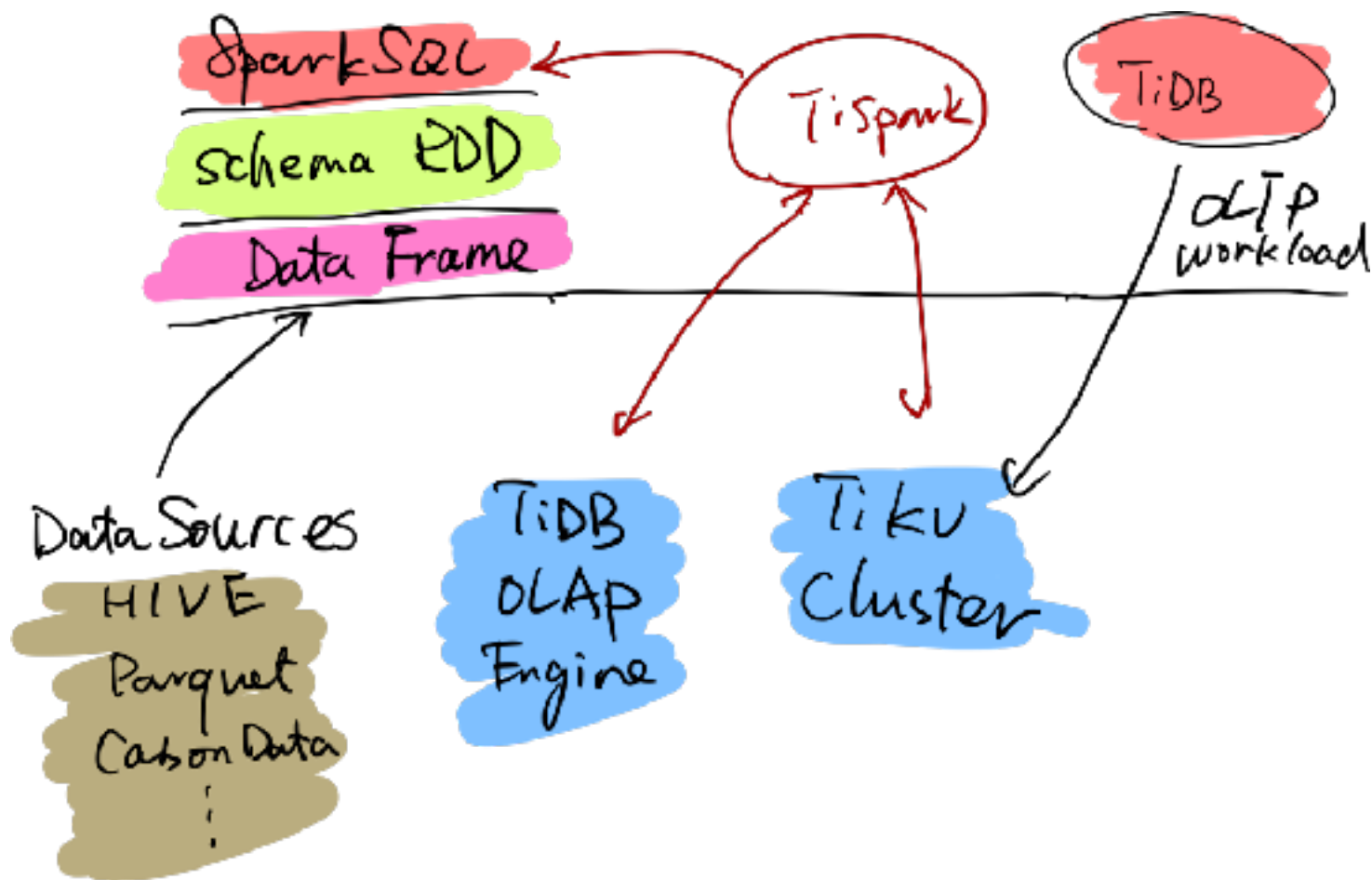
[tikv-client](#)

Fix incorrect totalRowCount calculation (#353)

an hour ago

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





How to test

- Unit Tests
- Simulator
- Jepsen Test
- 24/7 Chaos Test
- TLA+ ([open sourced](#))

...

**You Don't Choose Chaos Monkey...
Chaos Monkey Chooses You**



Thanks

Q&A

