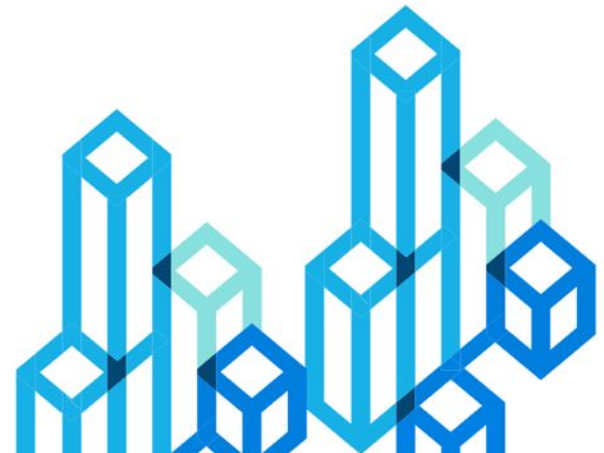
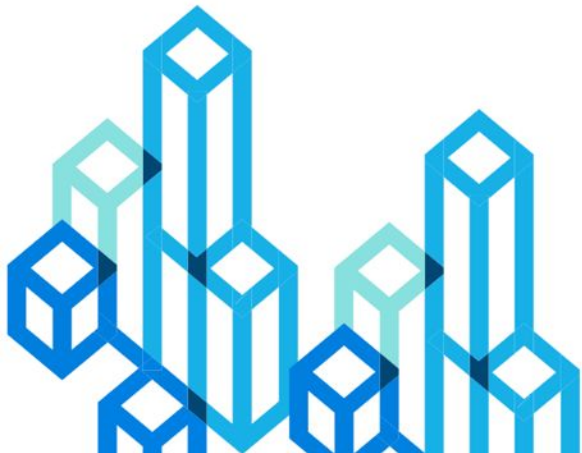


The Evolution of TiDB Architecture



- 张建, 开源爱好者, TiDB Maintainer
- TiKV 产品负责人@PingCAP
- 前 MaxCompute 执行引擎研发@Alibaba
- 关注数据库、分布式系统、开源软件与社区等领域
- Email: zhangjian@pingcap.com

01 What's TiDB

02 The Way to TiDB 5.0

03 Future Work

What's TiDB

TiDB: 全新一栈式实时 HTAP 数据库

- 一键水平扩容或者缩容
- 金融级高可用
- 实时 HTAP
- 云原生的分布式数据库
- 兼容 MySQL 5.7 协议和 MySQL 生态





面向零售高增长交易场景

全球领先餐饮巨头

- **实时高并发**, 高效支撑短时高峰交易及数据访问, **1.5亿用户**
- **弹性高扩展**, TiDB on K8s 随时随地在线扩展
- **安全高可用**, 可靠放心的高可用容灾保障能力



在线核心联机支付交易平台

日本头部金融支付企业

- 在线支付联机交易, 钱包等**高增长交易支撑**, **3000万用户**
- **水平弹性扩展**和几乎线性的扩展能力
- **友好的开发界面**, 应用几乎无需改动



面向金融敏态交易场景

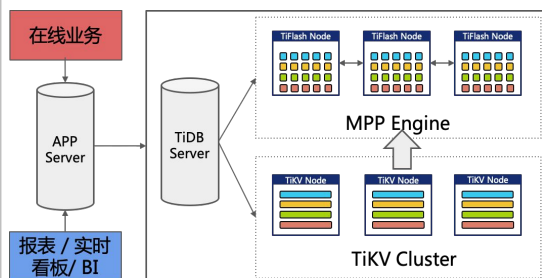
平安人寿在线金融联机服务

- 支撑平安人寿“金管家”单日交易额破**1000亿**, **超1亿用户**
- **高可靠、低延迟、可快速扩展**
- **大大降低敏态应用开发复杂度**, 加快应用上线速度
- **在线弹性扩缩容**满足不确定业务需求

混合负载场景

汽车之家
看车·买车·用车

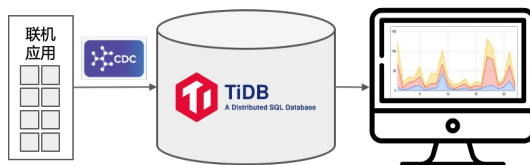
交易处理和在线分析同步规模化
报表业务不影响在线交易



流式计算场景

360
www.360.cn

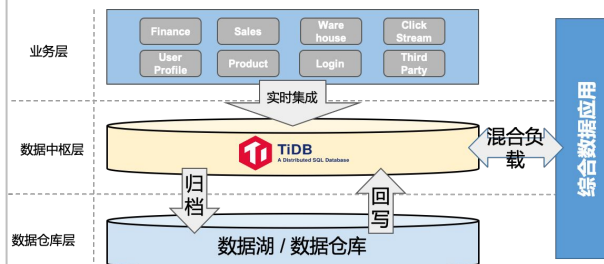
实时到账、实时统计
缩短业务等待周期，提升用户体验



数据中枢场景

贝壳
链家大平台

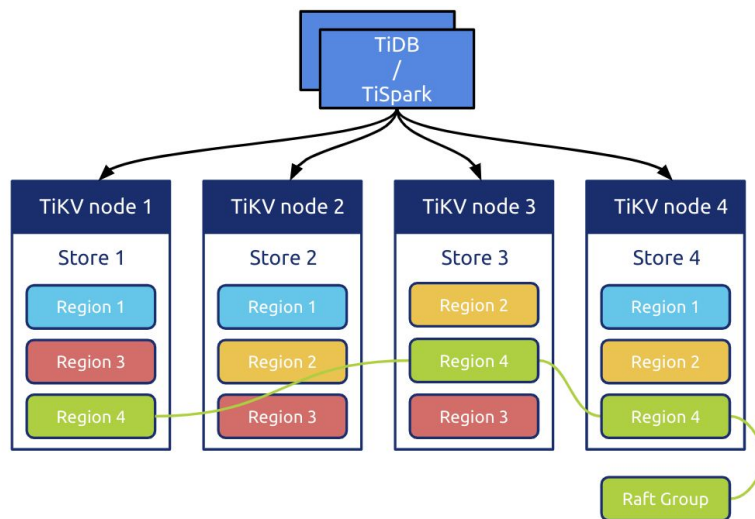
多源汇聚，统一数据服务
节省系统和数据聚合代码开发时间



写 SQL 就行了

一个TiDB 系统，一个访问入口，一份数据

The Way to TiDB 5.0

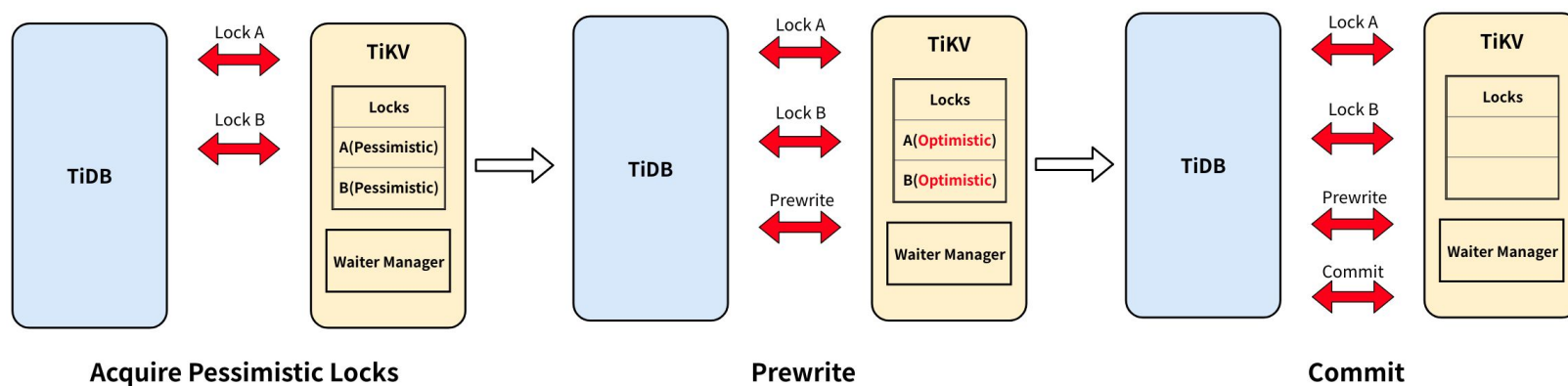


- 从 DBA 熬夜的日子出发
- 设计像 iPhone 一样简单
- 所有能够分离的都分离
- 回归经典数据库体验

乐观事务:中间的 DML 语句不检查锁, Commit 时再检查冲突

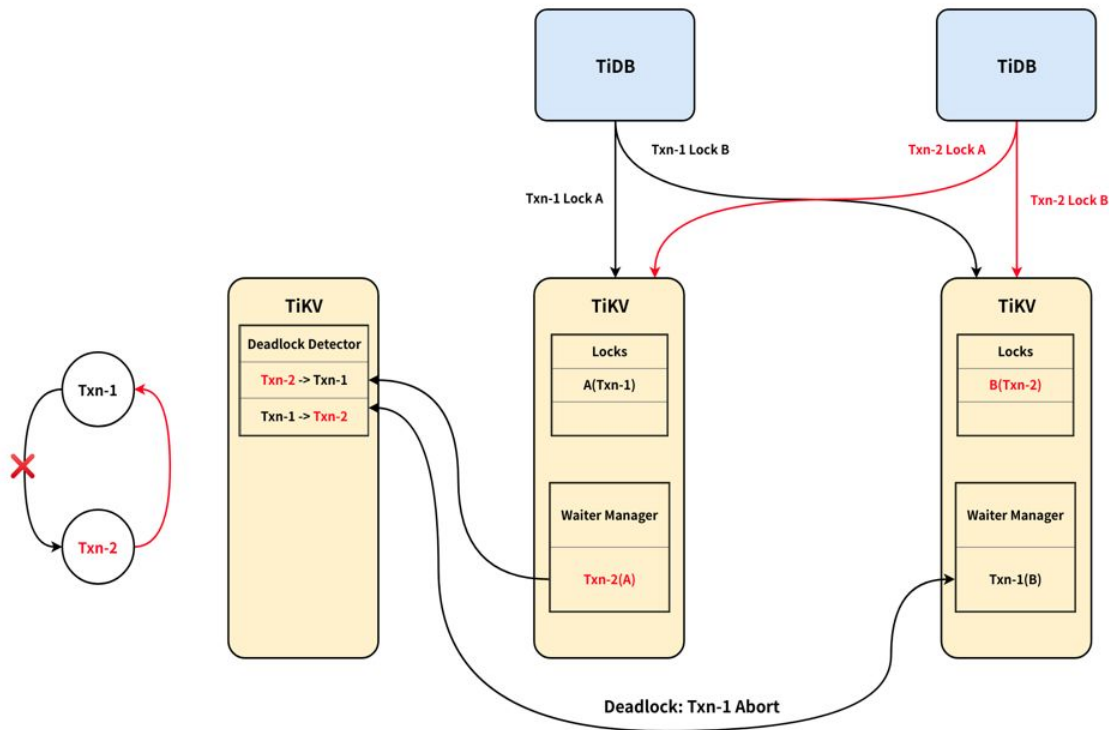
不足:

1. 高并发冲突的情况下,性能严重下降
2. 大量应用基于悲观锁设计,学习和改造成本高



TiDB 新特性漫谈: 悲观事务

- 动态选举出一个 TiKV 负责死锁检测
- 悲观事务在等锁的时候, 异步进行死锁检测

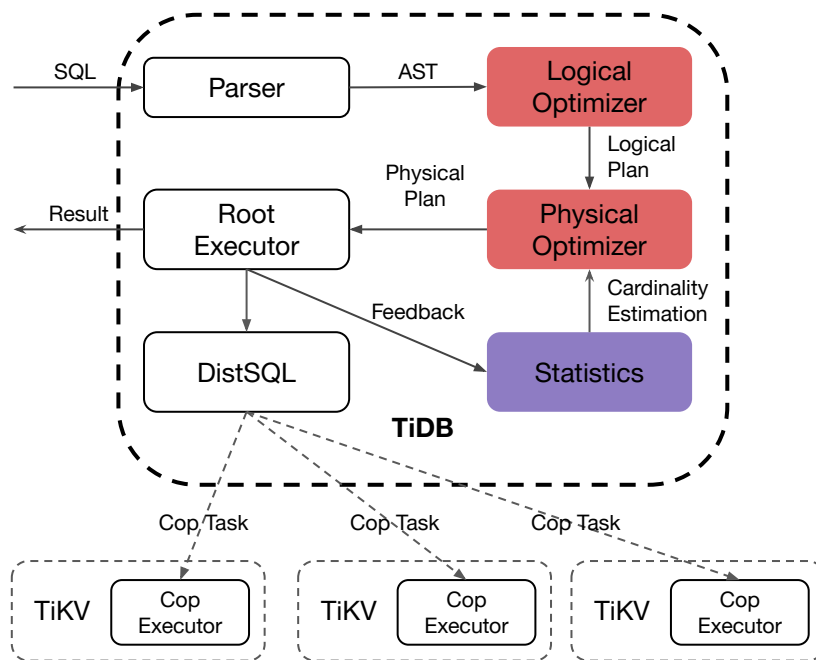


挑战二：执行计划的稳定性

场景一：线上 SQL 忽然执行计划变了，选错了索引导致慢查询，业务运行异常

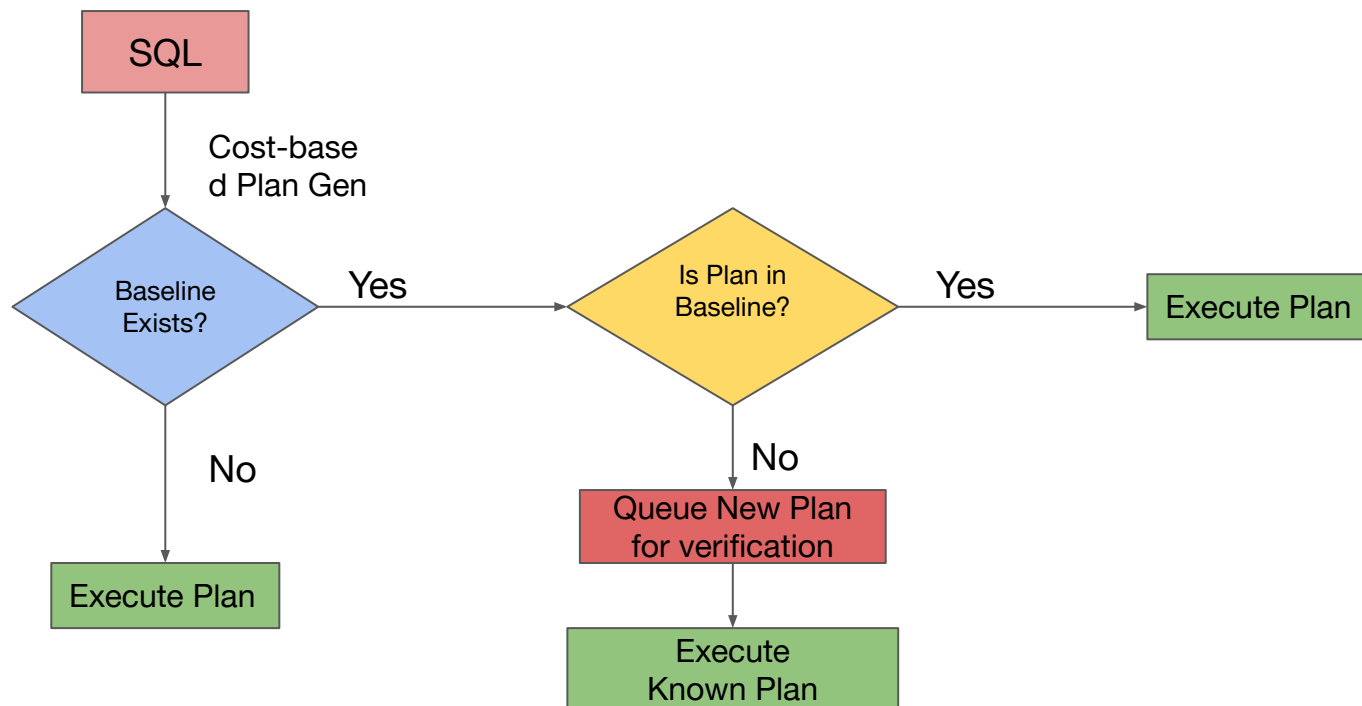
场景二：升级后，因为优化器的改动导致执行计划变化，升级后业务运行异常

...



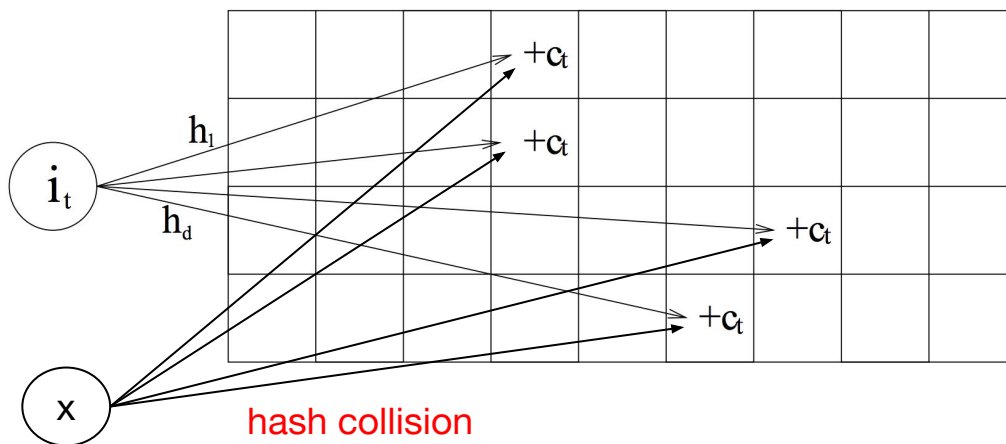
先提供规避和手工纠正的能力：

- 需要改业务 SQL: 增加 SQL Hint, 从 TiDB v1.0 加到 v5.0, 形成了一套体系
- 不需要修改业务: 从 SQL Binding (TiDB v3.0)
- 不想一条条加: SPM (TiDB v4.0)



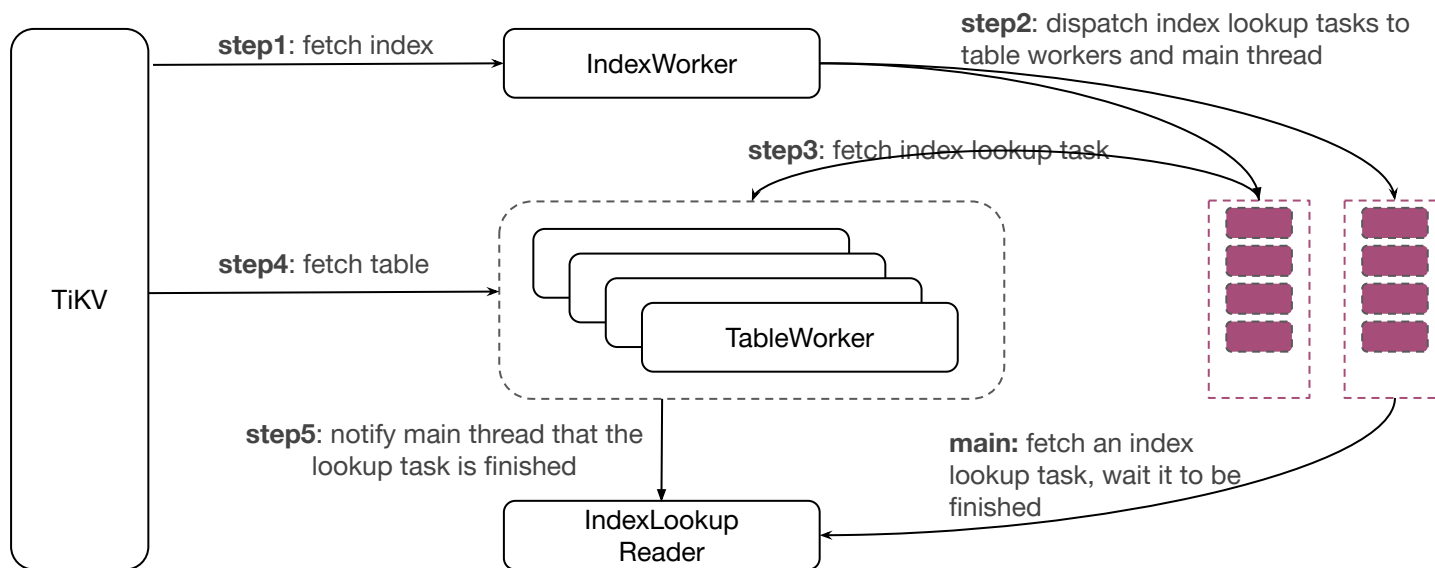
优化内功:

- 降低犯错的概率:引入 Skyline Pruning 裁剪不优索引, 引入更多启发式规则
- 提升做对的能力:重构统计信息, 降低基数估算误差
- 建立优化器测试评估体系



情况一：单个大查询

情况二：高并发中等查询

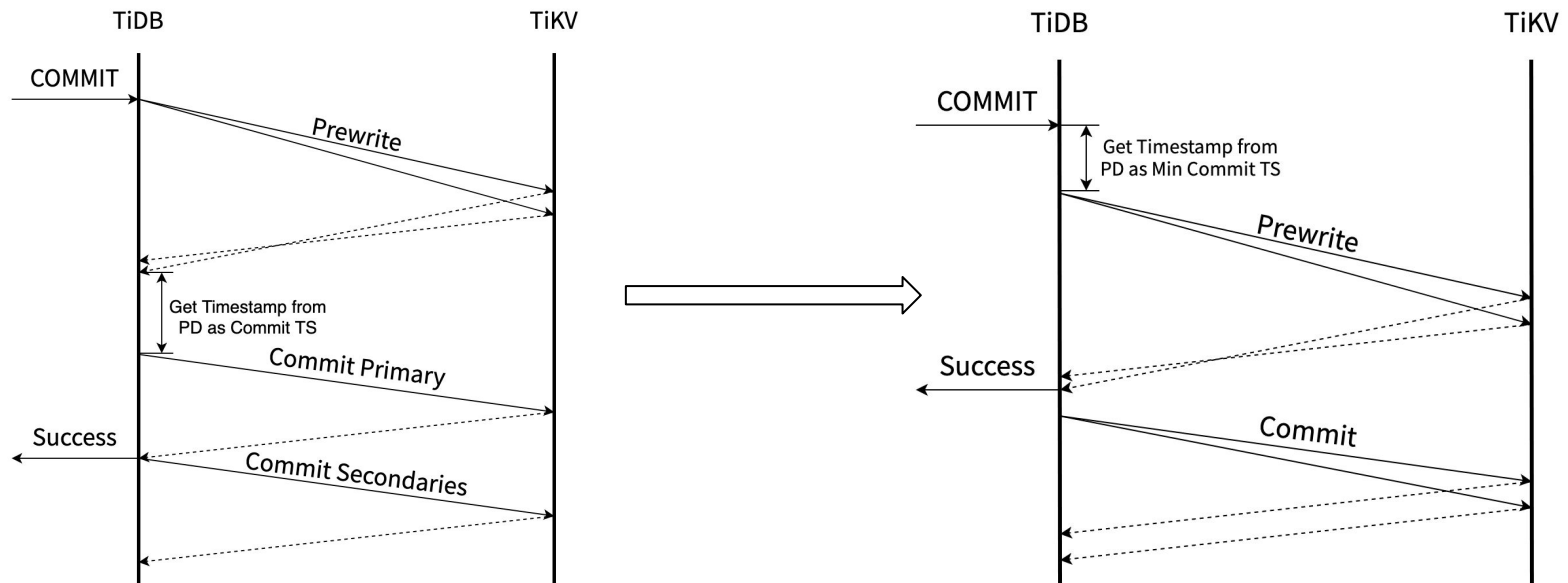


情况一：单个大查询

情况二：高并发中等查询

- 隔离火势：追踪内存使用，当超过阈值时，能够将该查询 kill 掉，不影响其他查询
- 灭火：通过算子落盘，让内存占用过多的大查询不被 kill 顺利执行完

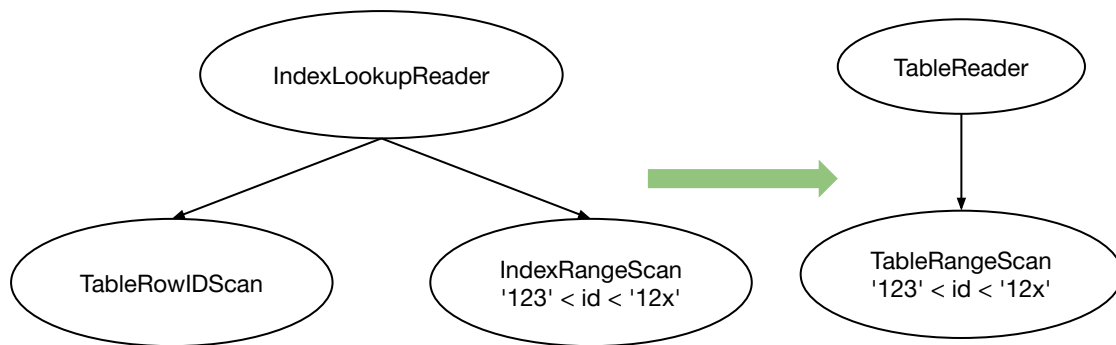
Async Commit (TiDB v5.0) : 把 2PC 变为 1PC



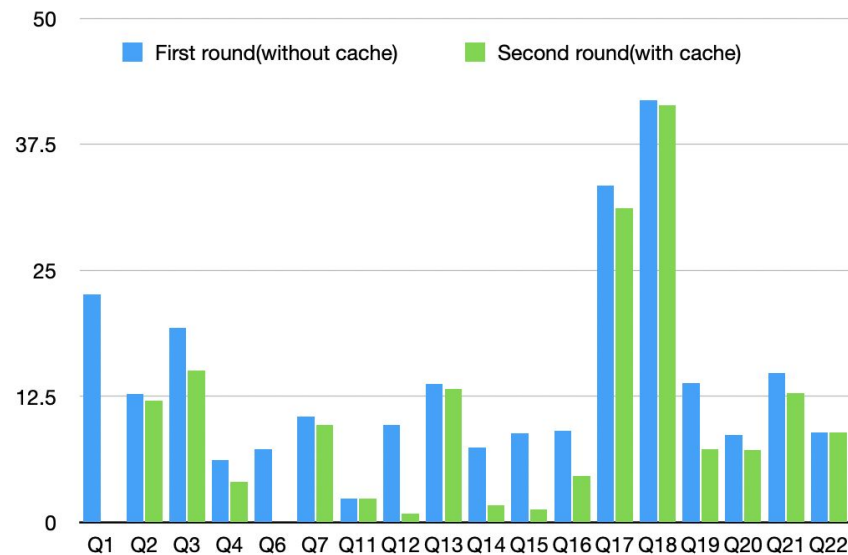
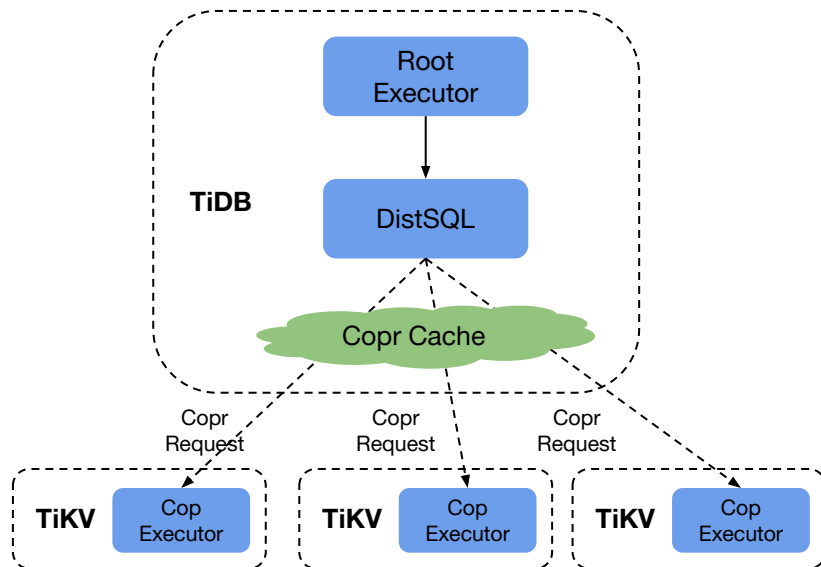
聚簇索引(TiDB v5.0): 减少回表次数, 让查询时间减半

```
CREATE TABLE `user` (  
  `id` varchar(20) NOT NULL,  
  `name` varchar(20) DEFAULT NULL,  
  PRIMARY KEY (`id`) CLUSTERED  
)
```

```
SELECT * FROM user  
WHERE id > '123'  
AND id < '12x';
```



Coprocessor Cache (TiDB v5.0) : 减少下推的计算量, 降低重复查询的延迟



TiDB 有哪些后台任务：

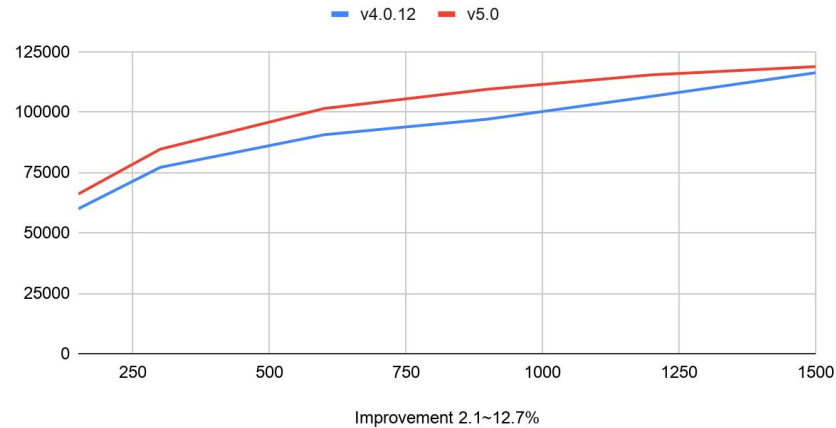
1. MVCC GC
2. RocksDB Compaction
3. Region Scheduling
4. Analyze Table
5. Update Global Variables
6. ...

- MVCC GC in Compaction Filter
- RocksDB's IO Rate Limiter
- Avoid redundant region rebalance caused by empty regions

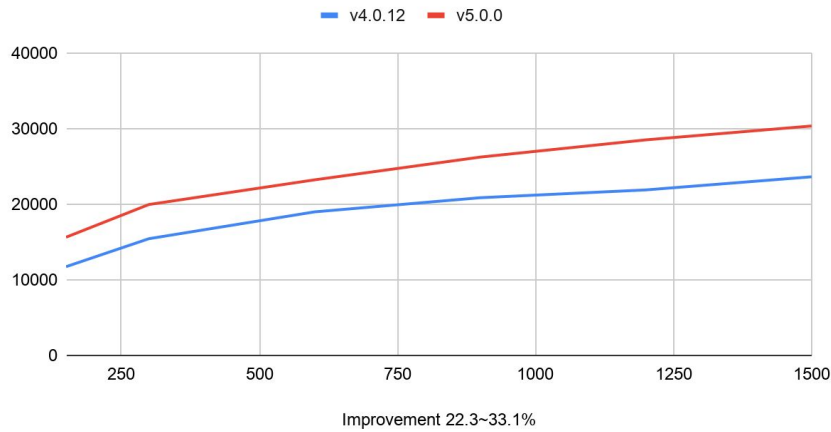
Standard deviation value for measuring TPC-C tpm-C jitter \leq 2% (8 hours)



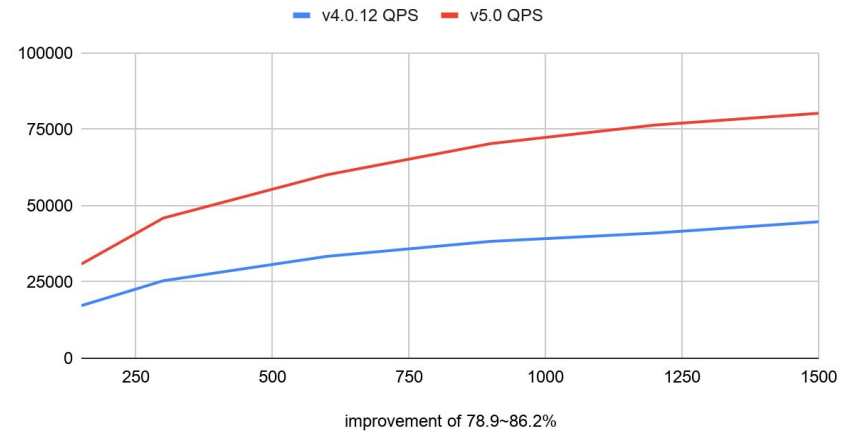
Sysbench Read Write QPS



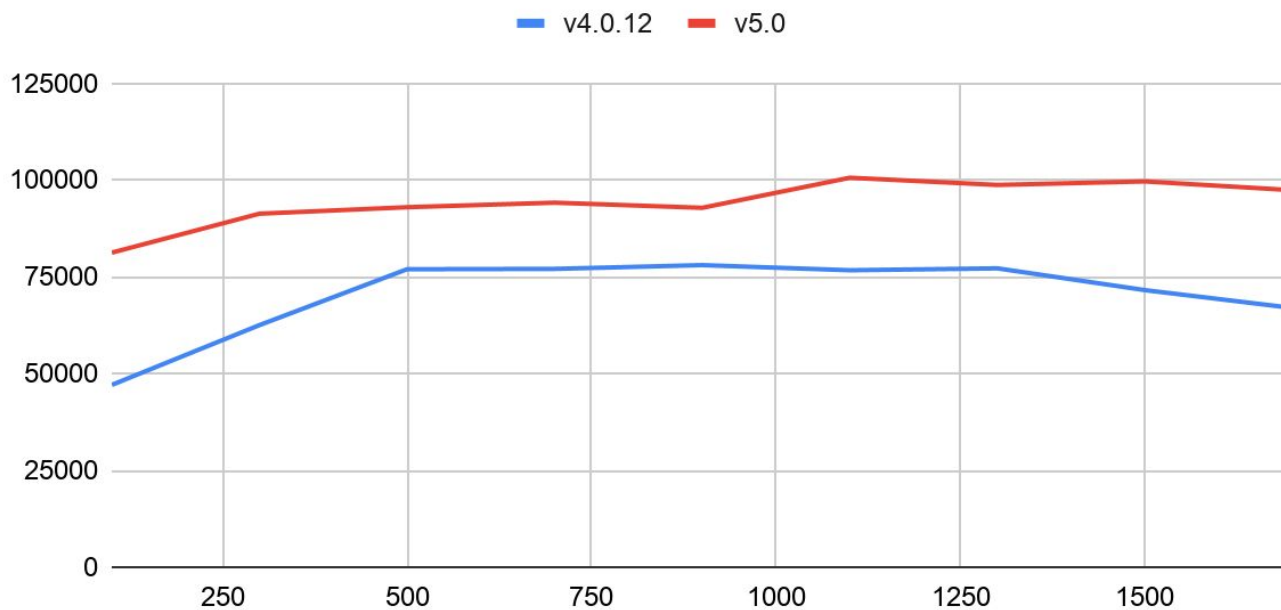
Sysbench Update Index QPS



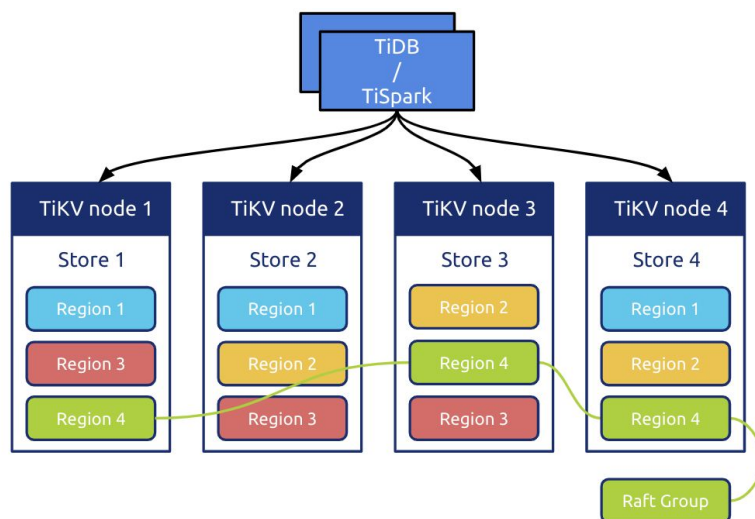
Sysbench Update Non-index QPS



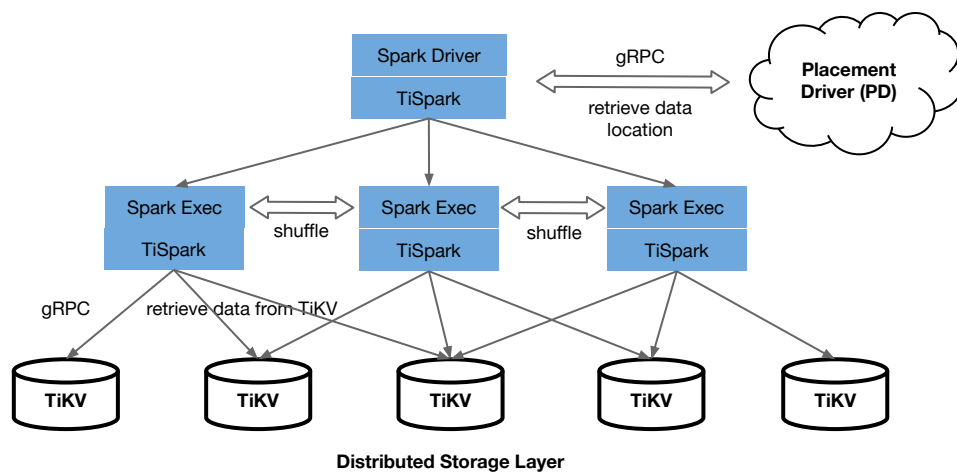
TPC-C



Improvement 18.9~72.5%



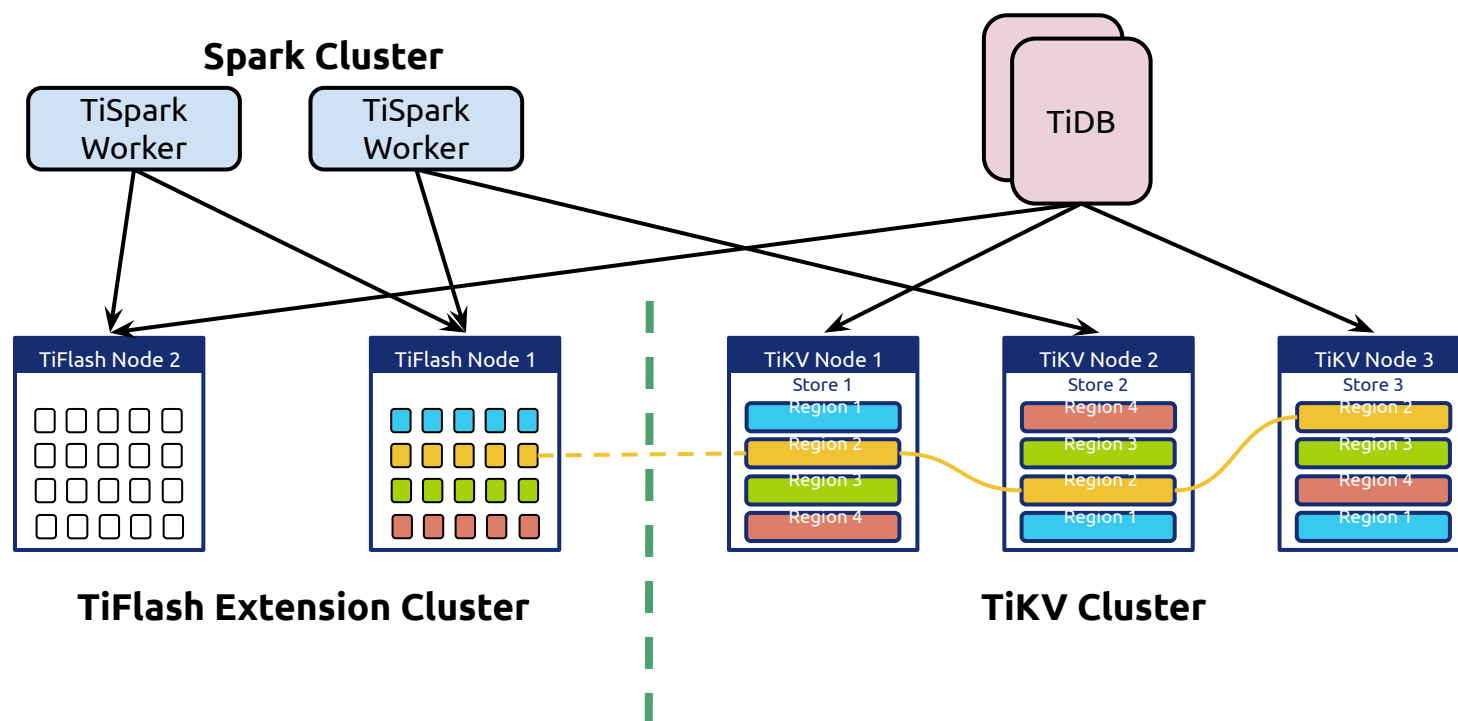
让我们来看看 HTAP 的发展

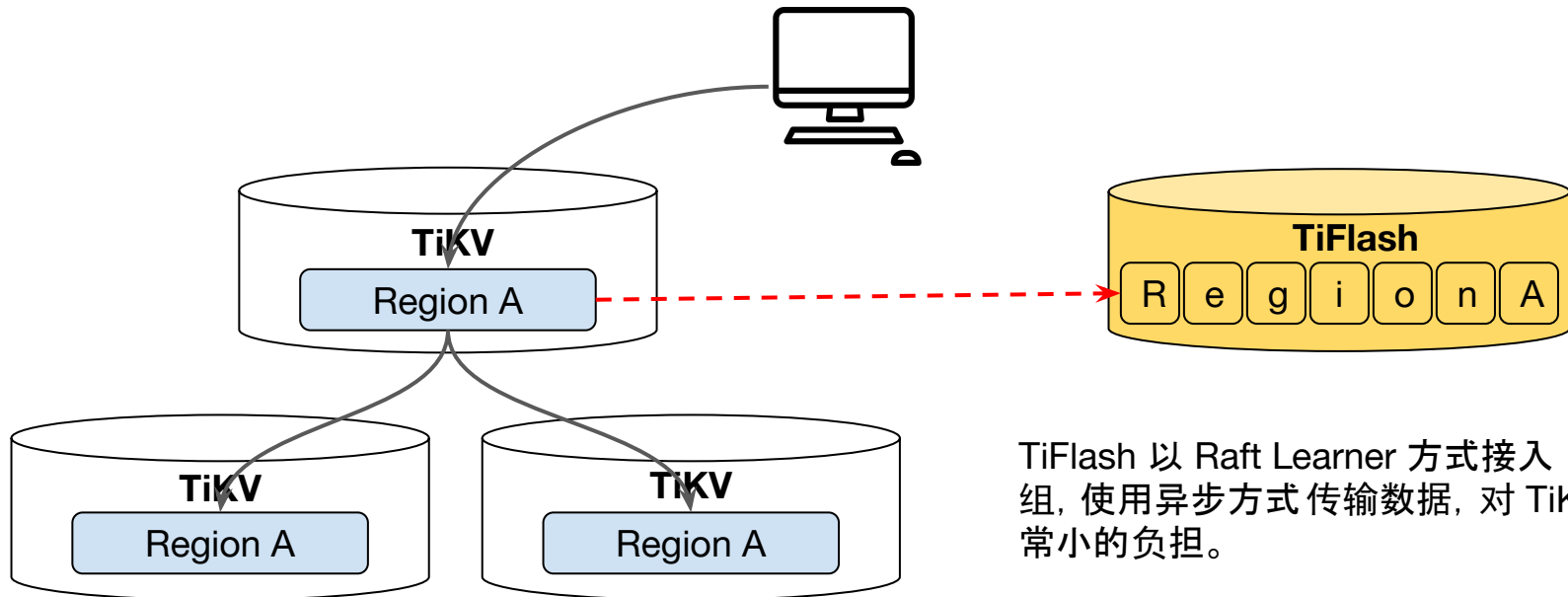


- 使用存储在 TiKV 的同一份数据
- 用 Spark 做分析查询

- 用户在很多场合下需要高并发中小规模 AP 能力
- TiKV 这样的 **行存** 无法满足分析场景需要的 **性能**
- TiKV 上跑在线业务时跑 OLAP 负载, 没有 **资源隔离** 简直是个灾难

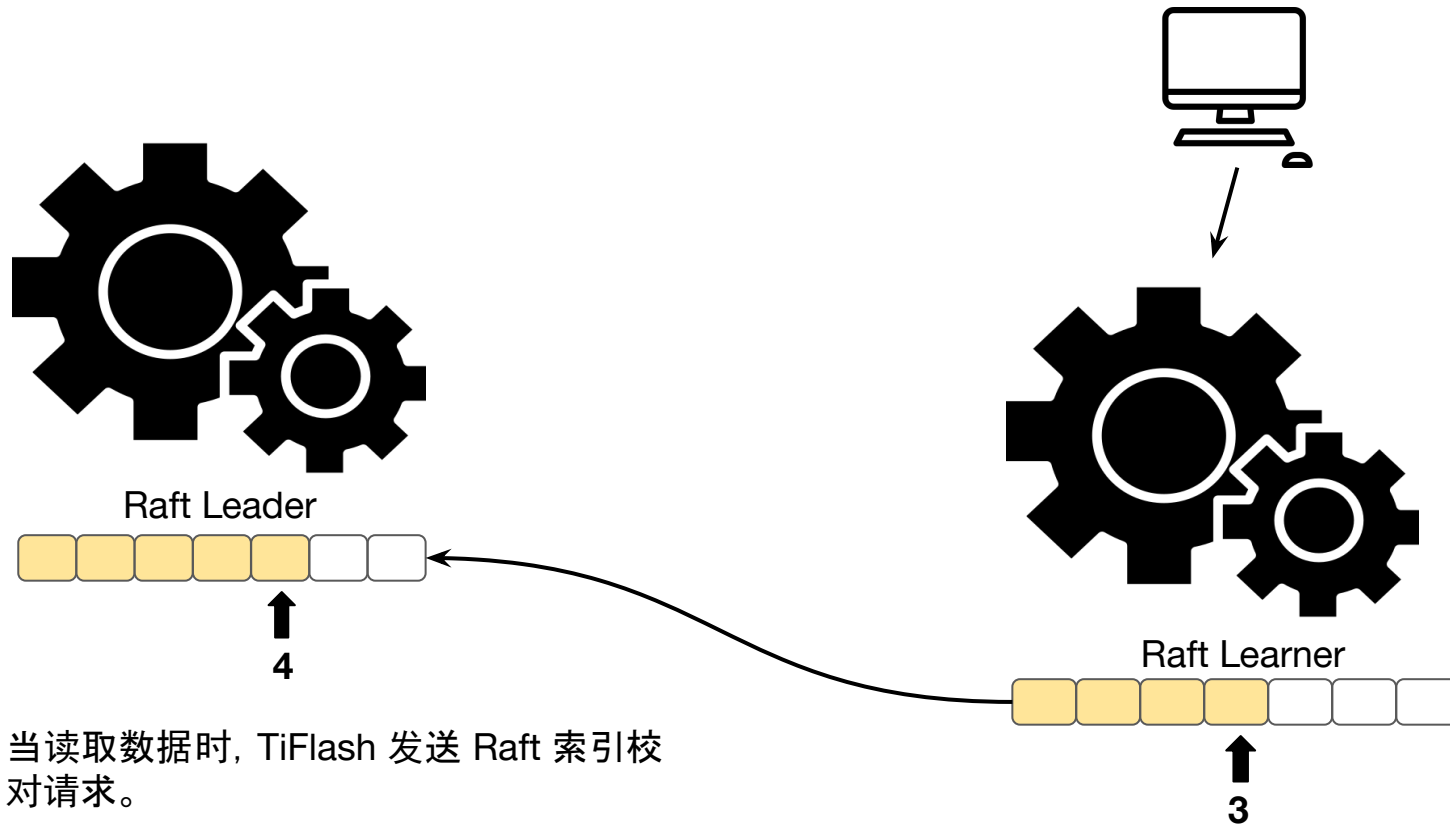


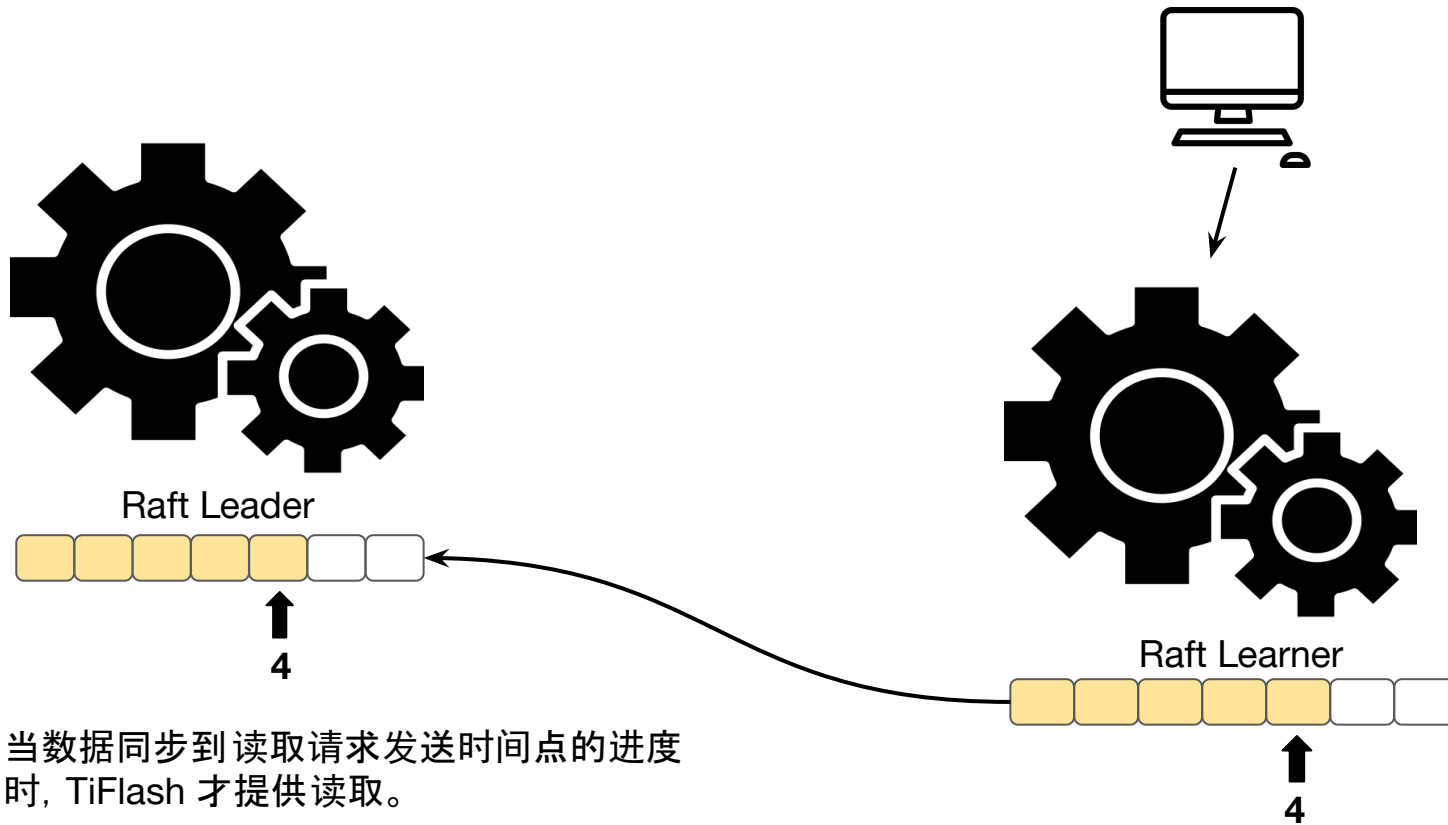


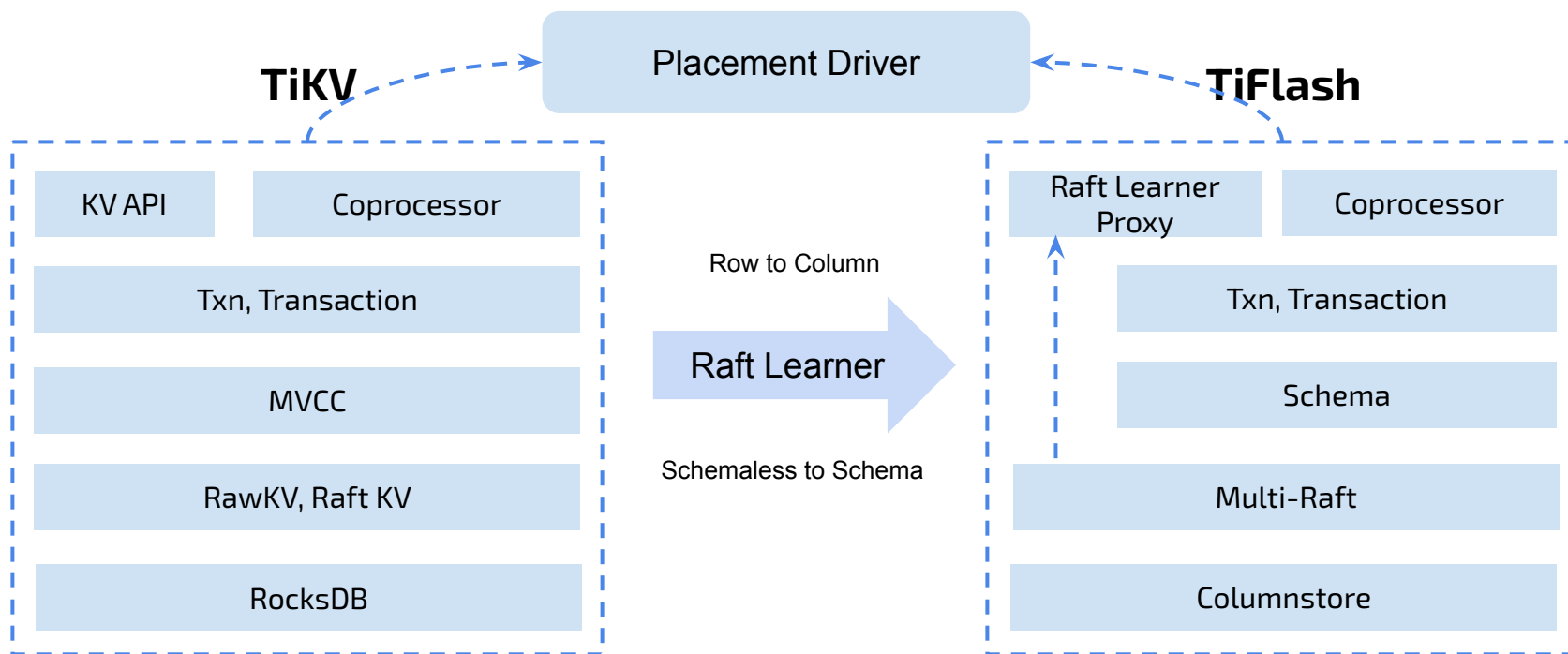


TiFlash 以 Raft Learner 方式接入 Multi-Raft 组, 使用异步方式传输数据, 对 TiKV 产生非常小的负担。

当数据同步到 TiFlash 时, 会被从行格式拆解为列格式。

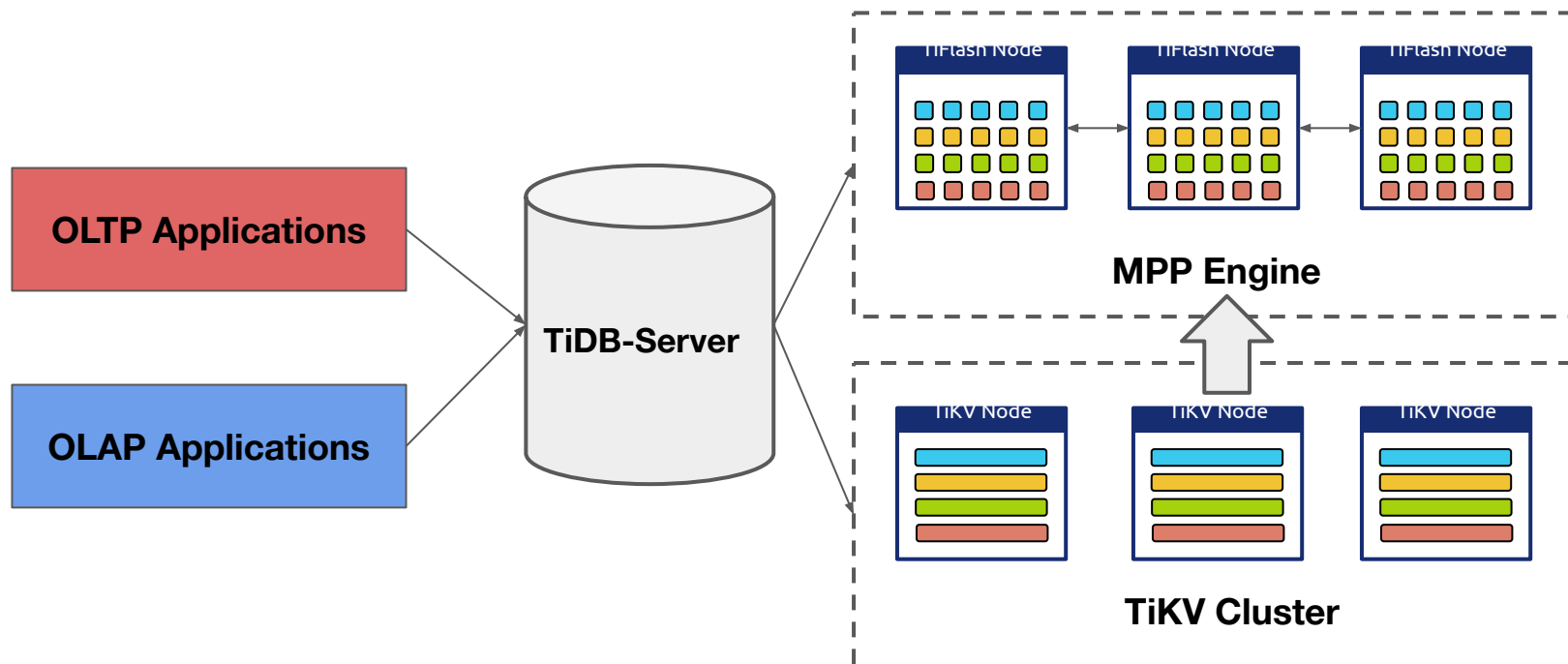




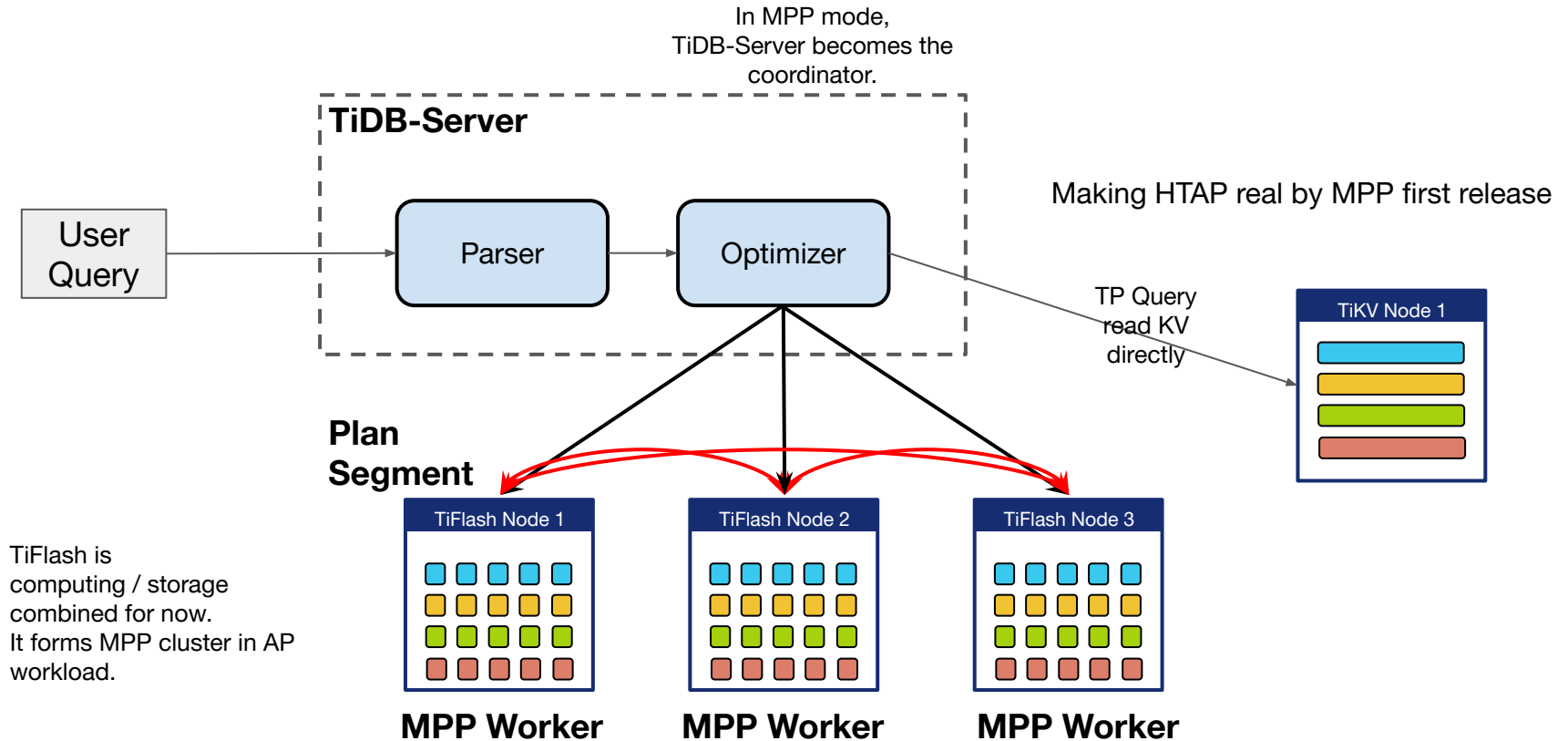


- TiSpark 作为唯一分布式计算引擎
 - 缺少中规模快速查询的解决方案
 - 略重的模型(MR 模型)- 仍需要 MPP 引擎
- 写入需要通过 TiKV
 - 大批量写入速度吞吐不够
 - 副本必须先以行存方式写入再同步为列存

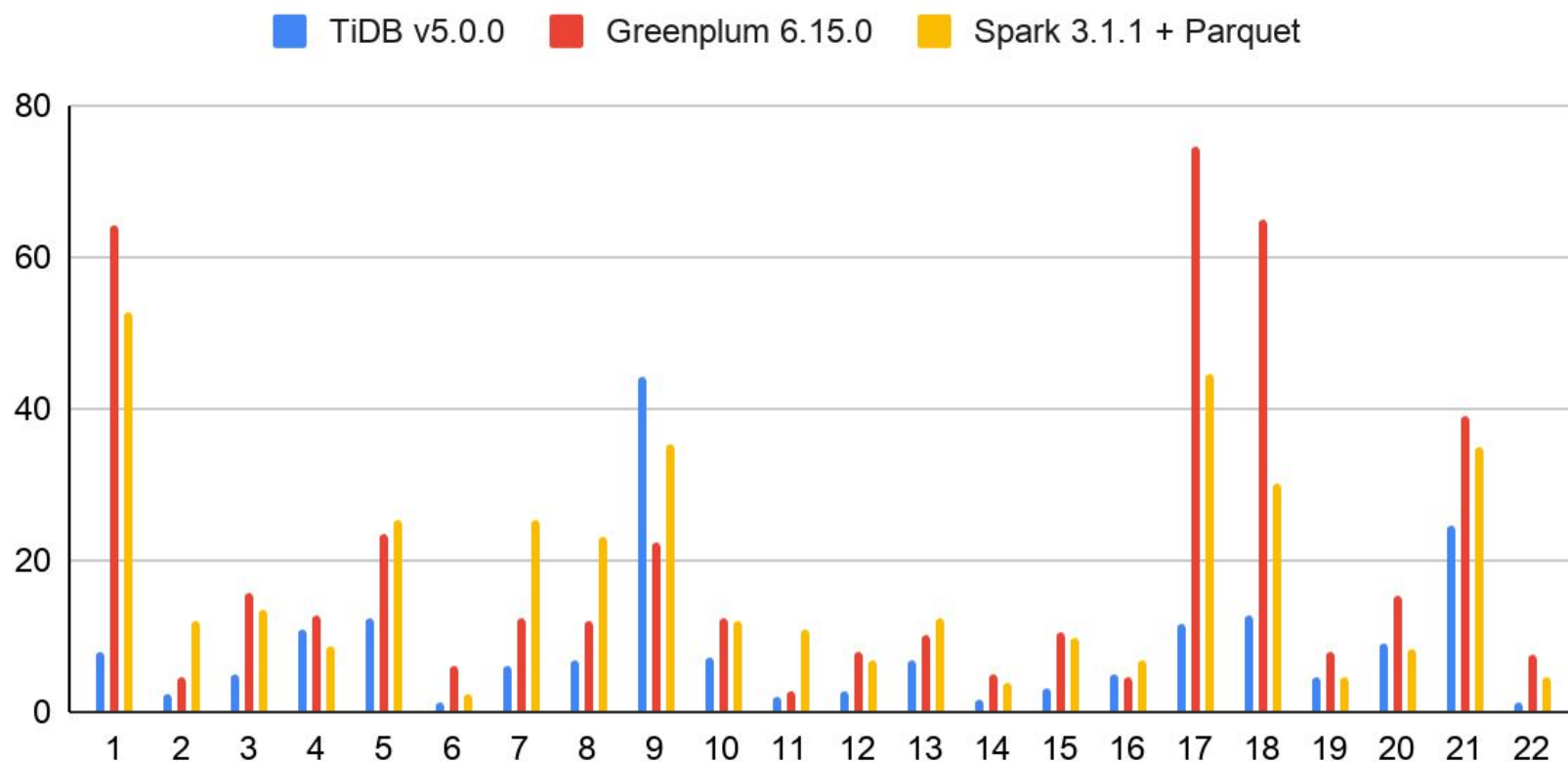
先做 MPP，解决中规模快速查询的问题



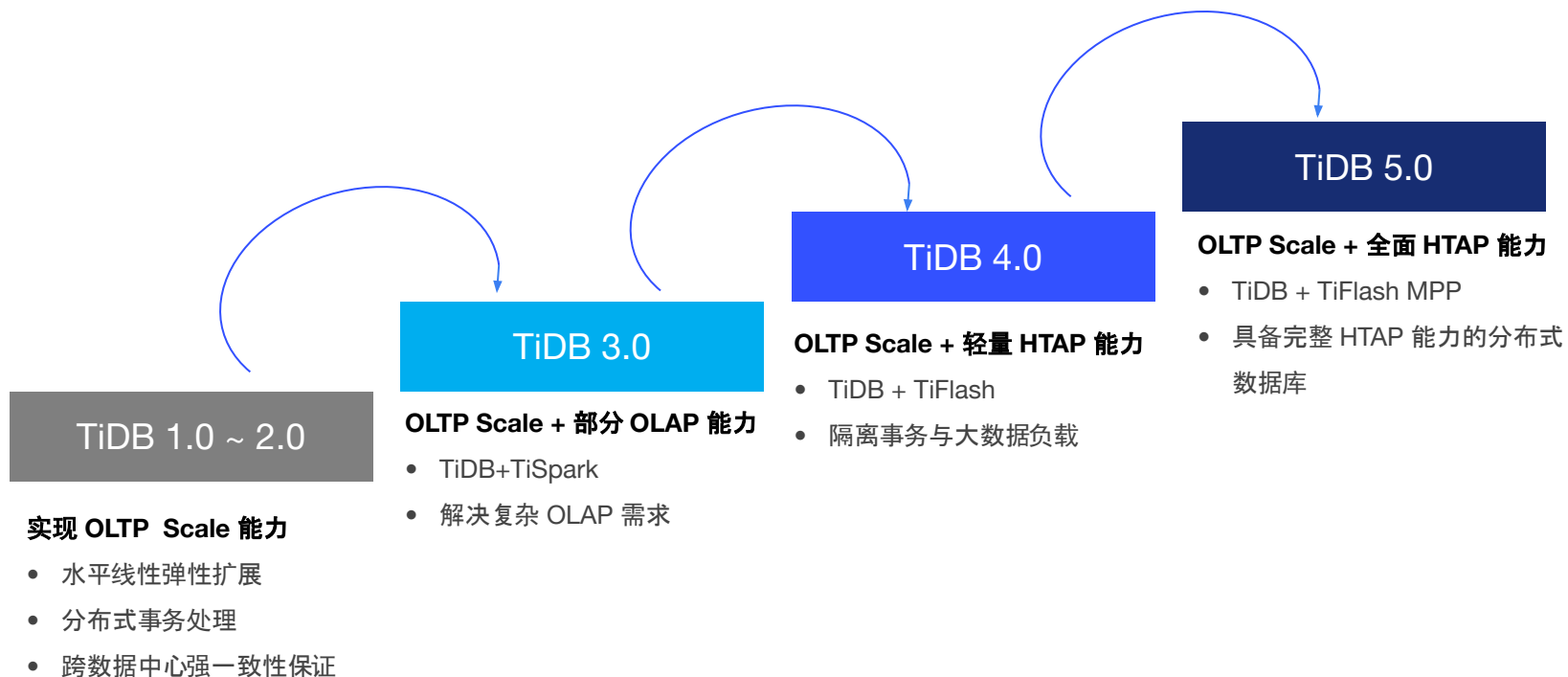
先做 MPP，解决中规模快速查询的问题



TPC-H 100GB on 3 Nodes



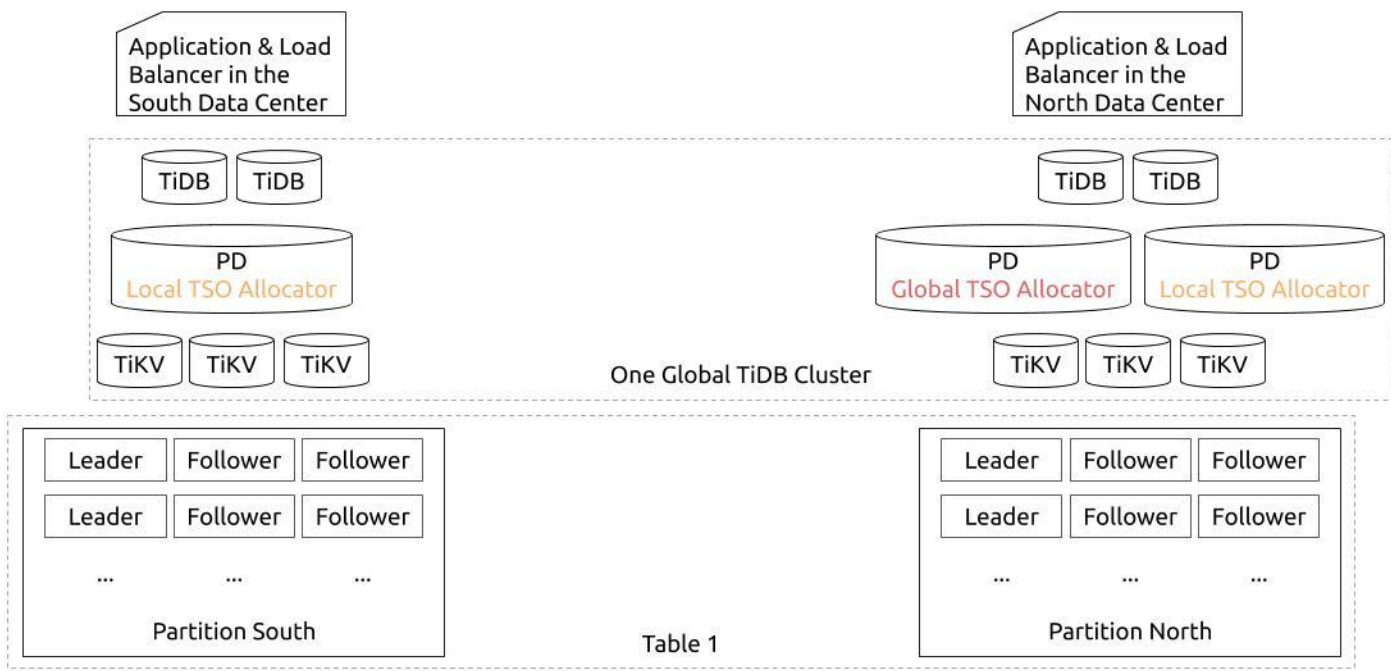
TiDB vs Greenplum Improvement 17.7~793.6%, TiDB vs Spark improvement 34.1~651.4%



Future Work

OLTP Scale

- Big Cluster
- Plan Stability
- Memory Management & Query Performance
- Geo-Partition



OLTP Scale

- Big Cluster
- Plan Stability
- Memory Management & Query Performance
- Geo-Partition

Real-time HTAP

- Direct Write
- Bigdata Ecosystem Integration



麦思博(msup)有限公司是一家面向技术型企业的培训咨询机构, 携手2000余位中外客座导师, 服务于技术团队的能力提升、软件工程效能和产品创新迭代, 超过3000余家企业续约学习, 是科技领域占有率第1的客座导师品牌, msup以整合全球领先经验实践为己任, 为中国产业快速发展提供智库。



高可用架构公众号主要关注互联网架构及高可用、可扩展及高性能领域的知识传播。订阅用户覆盖主流互联网及软件领域系统架构技术从业人员。

高可用架构系列社群是一个社区组织, 其精神是“分享+交流”, 提倡社区的人人参与, 同时从社区获得高质量的内容。