



# 小米结构化存储系统 及融合云平台的 设计与实践

林尚泉

linshangquan@xiaomi.com



# 大纲



云计算开源产业联盟  
China Cloud Open Source Alliance

全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 小米结构化数据存储
  - 融合云平台



# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全球新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## 分布式NoSQL数据库服务 对标AWS DynamoDB

弹性可扩展

高可用

低延时

稳定可靠



# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## SDS

### (Structured Data Store)

基

无

简

多

于  
H  
B  
A  
S

状  
态  
公  
网  
访

化  
认  
证  
和  
配

语  
言  
S  
D

多租户

功能扩展

ACL

流  
量  
控  
制

数  
据  
类  
型

二  
级  
索  
引

S  
T  
R  
E  
A  
M

软  
删  
除

数  
据  
冷  
备

# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## 小米内部业务



## 生态链业务



# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 应用规模

国际化

- 北京，天津，美国，新加坡

机器规模

- 100+

业务数

- 20+

数据量

- 100+TB，数千亿行

单集群QPS

- 10万级

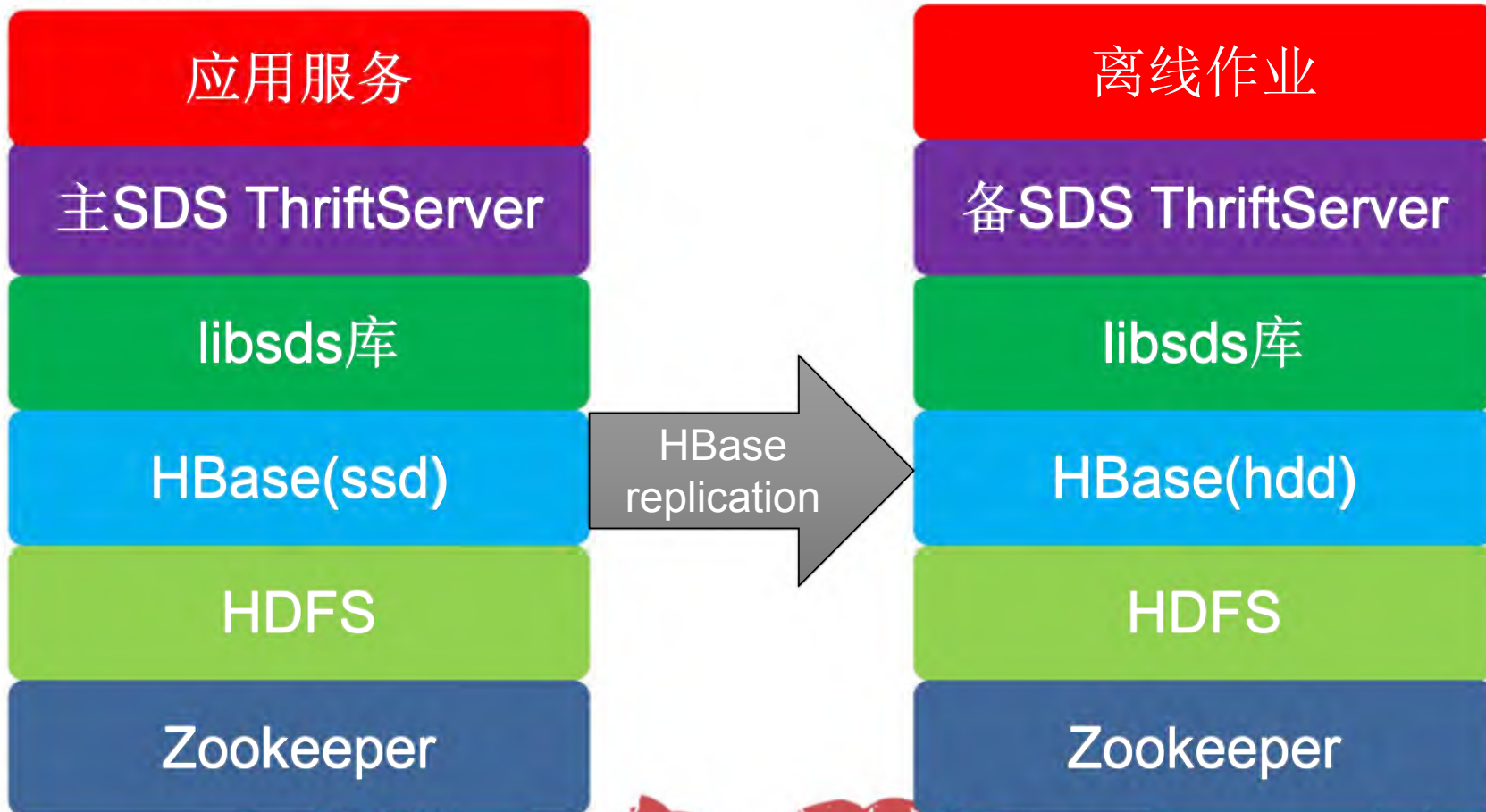


# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全球新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## 典型部署



# 结构化存储



全球云计算开源峰会 2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- libsdgs库

libsdgs

规范化的  
数据模型

内建数据类型支持

局部二级索引

全局二级索引

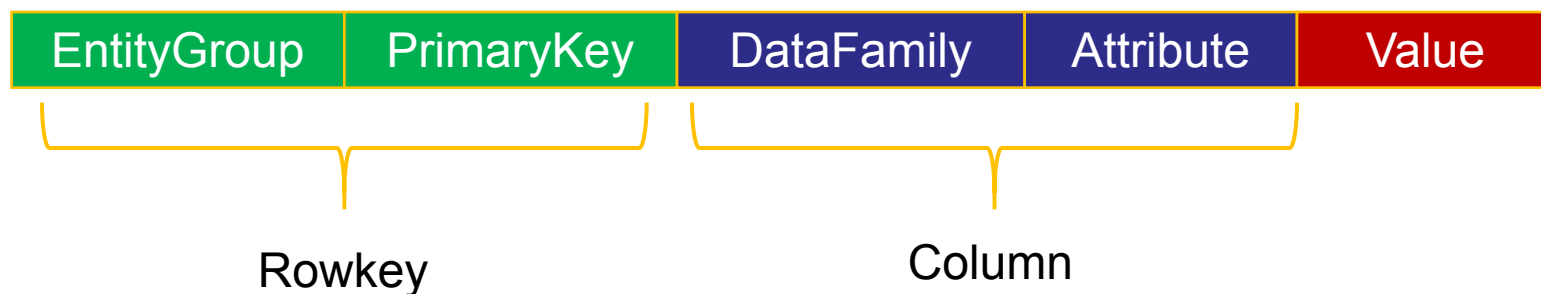
stream



# 结构化存储



- 规范化数据模型
  - Rowkey前缀：实体组键，支持哈希分布
  - Rowkey后缀：主键
  - Column Qualifier：数据列



# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## 数据类型支持

- bool/int8/int16/int32/int64/float/double/string/binary/set

## HBASE-8201和SQLite编码方案<sup>1</sup>

- 顺序保持编码

## 支持逆序

- 按位取反

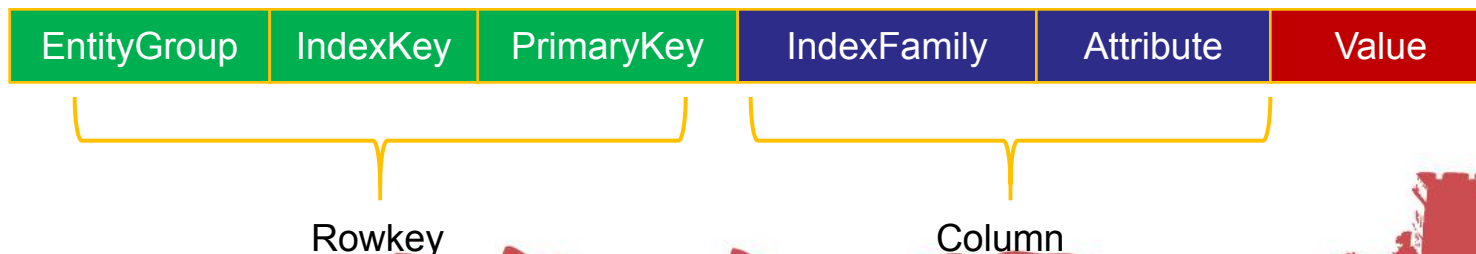
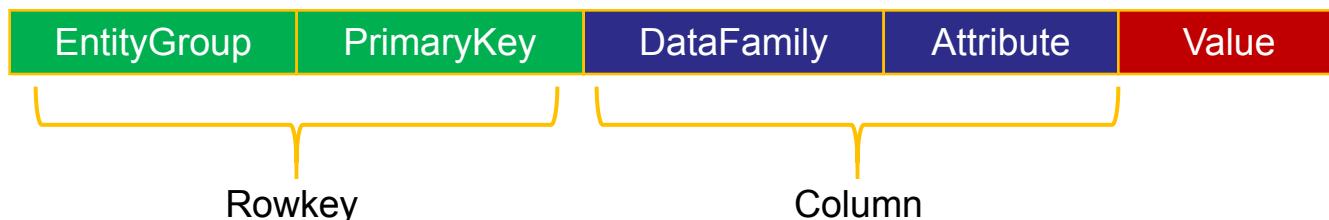
<sup>1</sup>[http://sqlite.org/src4/doc/trunk/www/key\\_encoding.wiki](http://sqlite.org/src4/doc/trunk/www/key_encoding.wiki)



# 结构化存储



- 局部二级索引：实体组内部
  - 同一张表，不同Column Family
  - 基于正则表达式的前缀分割策略(RegexPrefixRegionSplitPolicy):  
保证同一个实体组键数据不能split

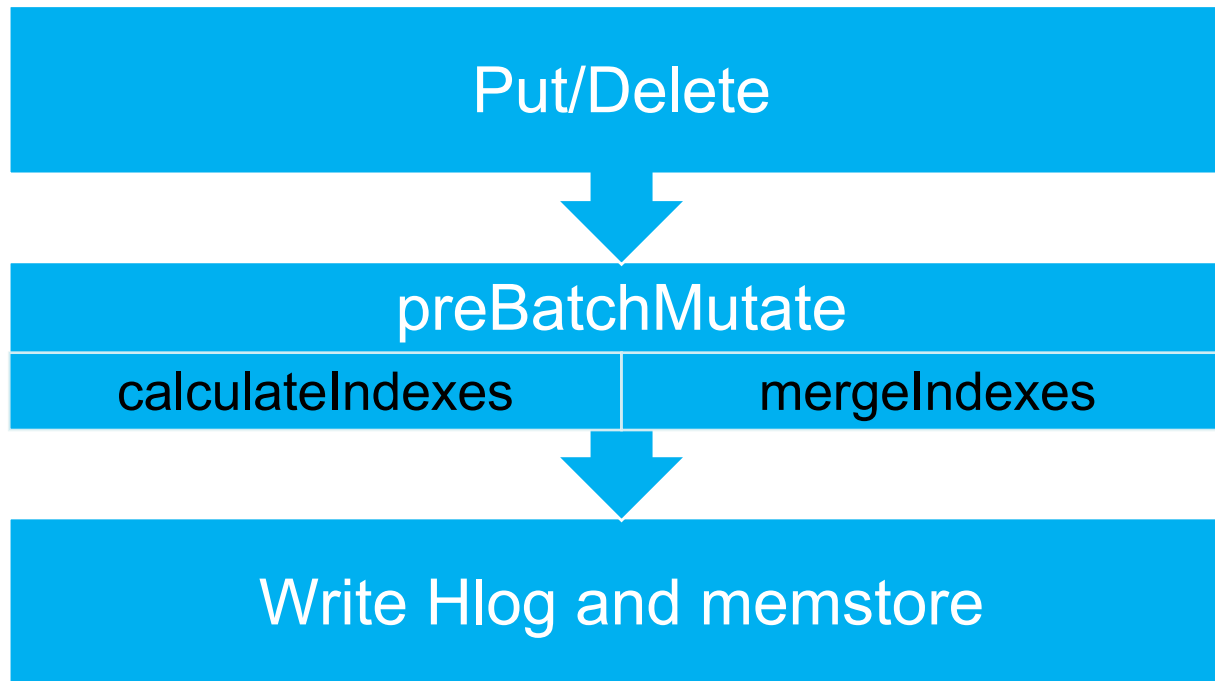


# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 局部二级索引：HBase Coprocessor——IndexObserver



# 结构化存储



- 多种类型局部二级索引

## EAGER

更新删除时同时删除失效索引，适合写少读多

## LAZY

读取时判断索引有效性，更新时不做额外操作，适合写多读少

## IMMUTABLE

适合只读数据一次性写入，读写均无需额外判断

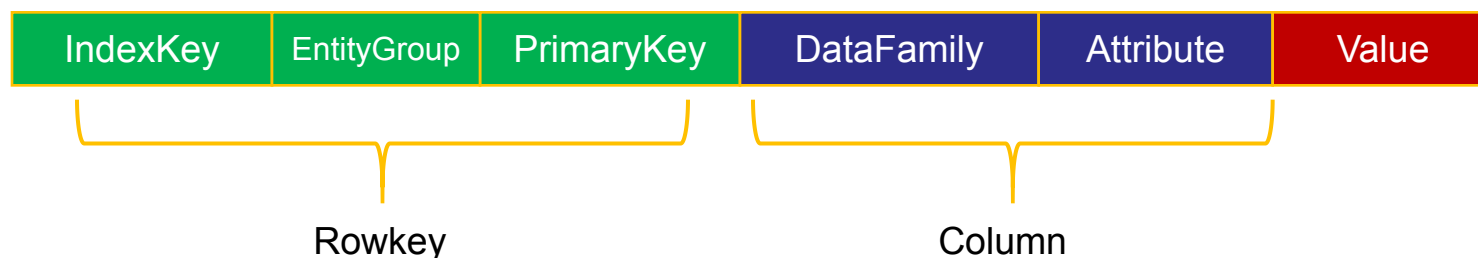
# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 全局二级索引

- 使用单独的HBase 表存储索引
- Google percolator<sup>2</sup>的 HBase实现：Themis<sup>3</sup>，保证跨表更新的原子性
- Chronos<sup>4</sup>全局单调递增时间戳



<sup>2</sup>Large-scale Incremental Processing Using Distributed Transactions and Notifications

<sup>3</sup><https://github.com/xiaomi/themis>

<sup>4</sup><https://github.com/xiaomi/chronos>

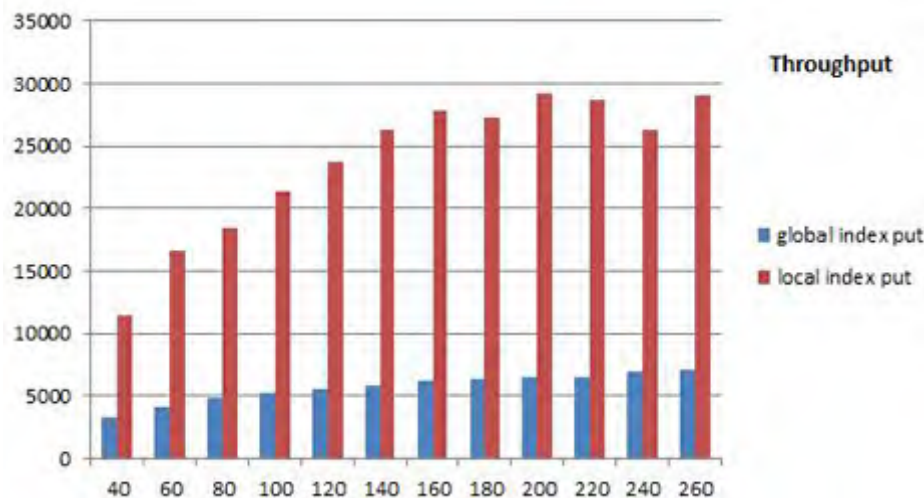


# 结构化存储

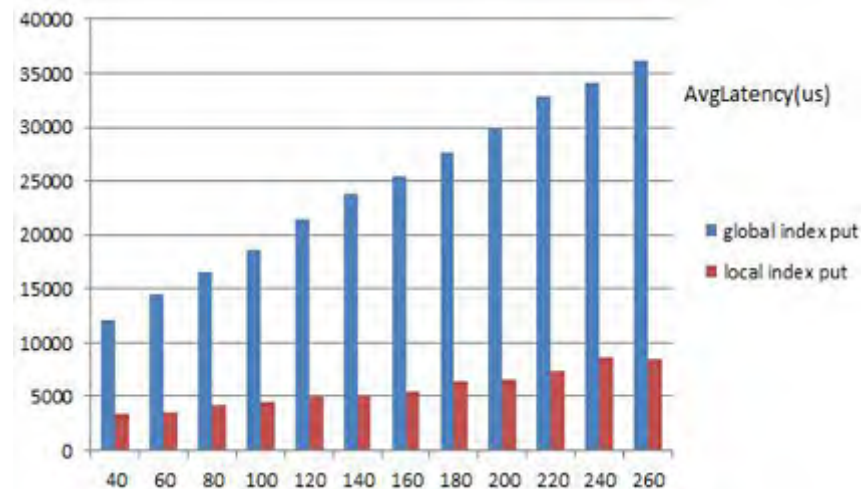


全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 局部二级索引 VS. 全局二级索引(写)



线程数



线程数

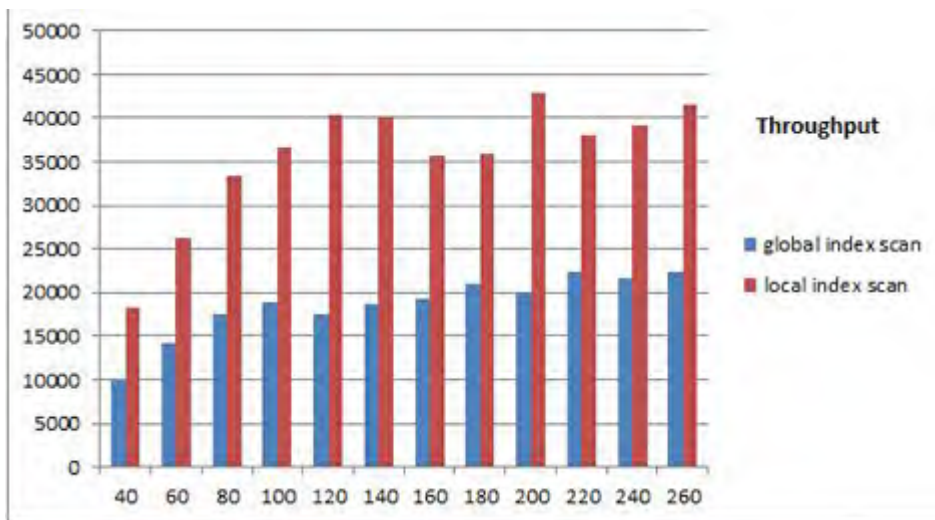


# 结构化存储

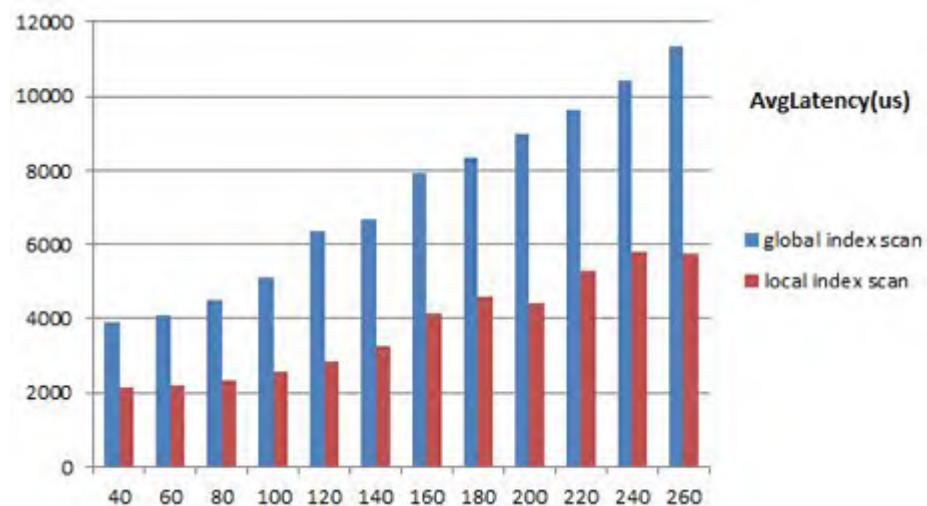


全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 局部二级索引 VS. 全局二级索引(读)



线程数



线程数

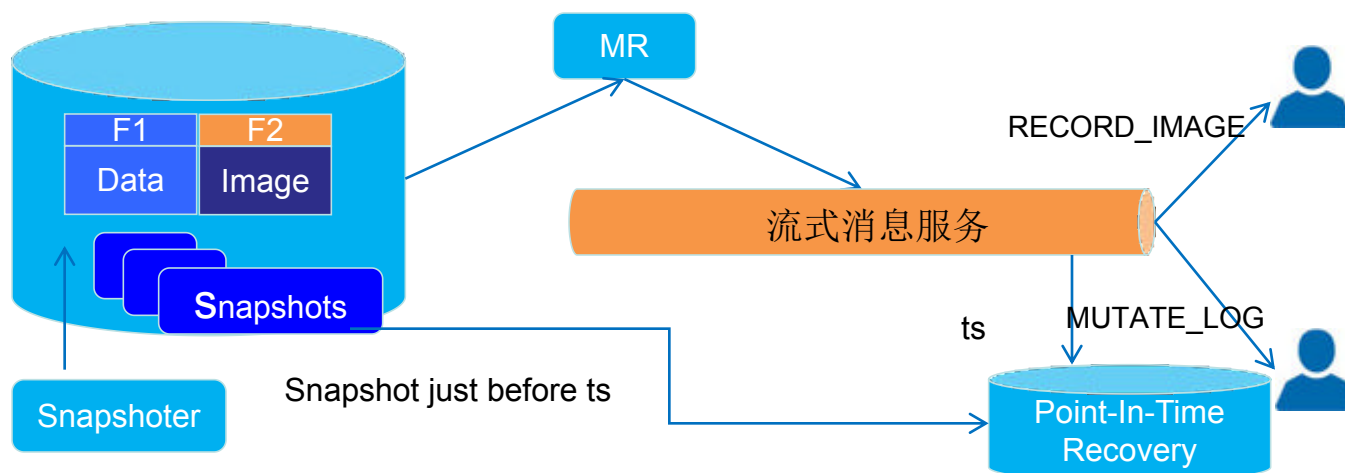


# 结构化存储



- stream

- RECORD\_IMAGE: 增量备份
- MUTATE\_LOG: Point-In-Time Recovery

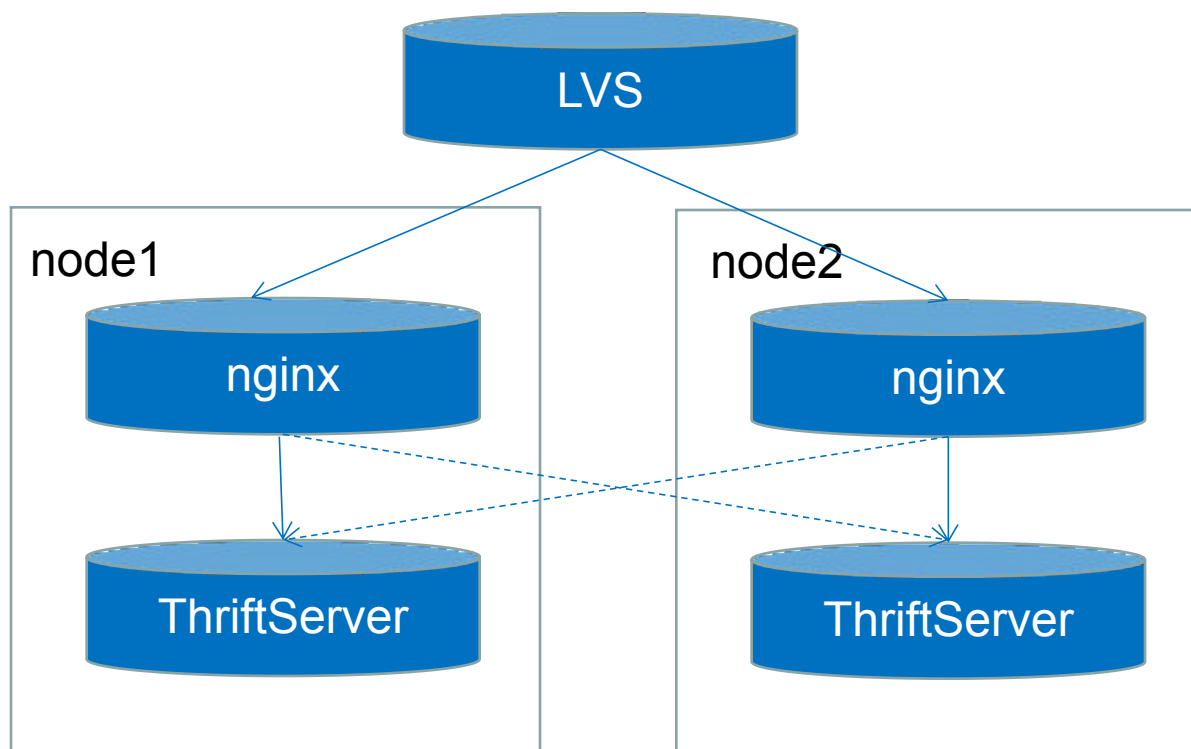


# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- SDS ThriftServer: 无状态 Http Server



# 结构化存储



全球云计算开源峰会 2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## SDS ThriftServer

Http 无  
状态公  
网访问

简化认  
证和配  
置

多语言  
SDK

多租户

ACL

流量  
控制

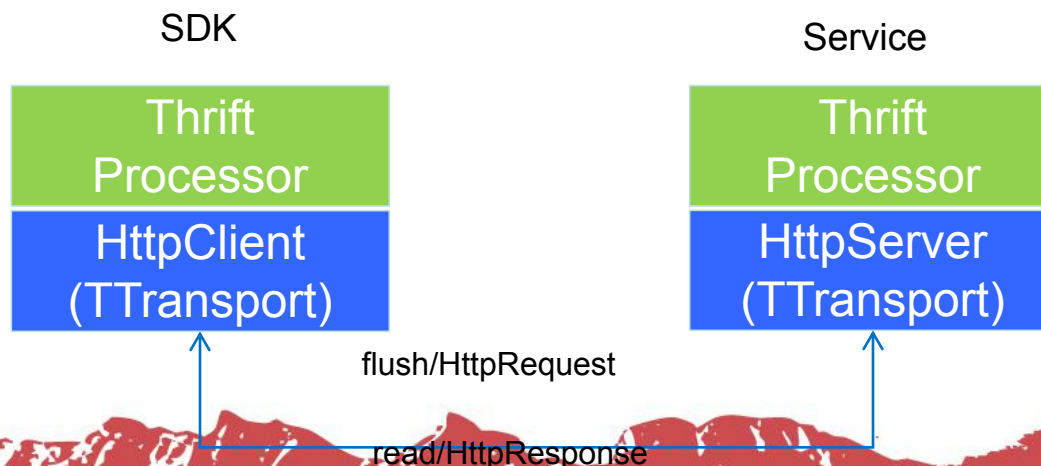
# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## • SDS ThriftServer

- Http 无状态公网访问
- 简化认证和配置：使用小米融合云或开放平台认证，对外屏蔽了Zookeeper配置和Kerberos认证
- 多语言SDK(java/python/php/node.js/c++/go/javascript)

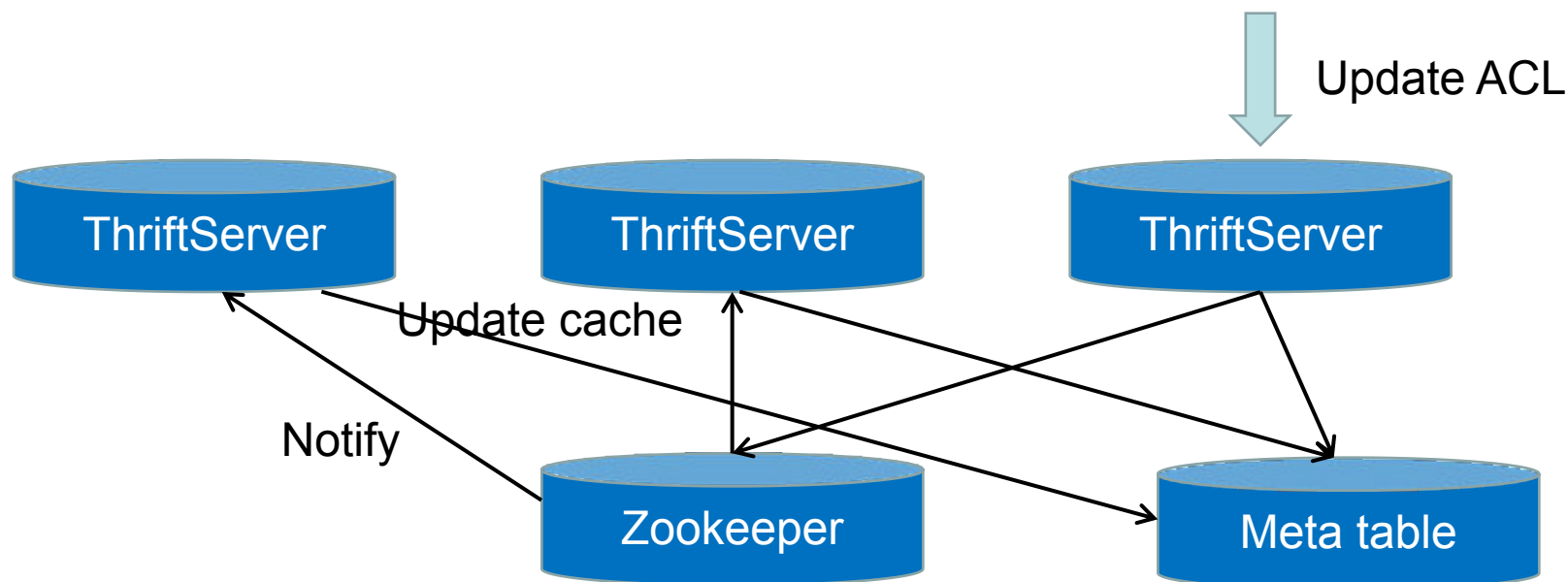




# 结构化存储



- SDS ThriftServer (ACL)
  - 元数据表: tableId => {ACL信息}
  - 本地cache元数据, zookeeper watcher 同步



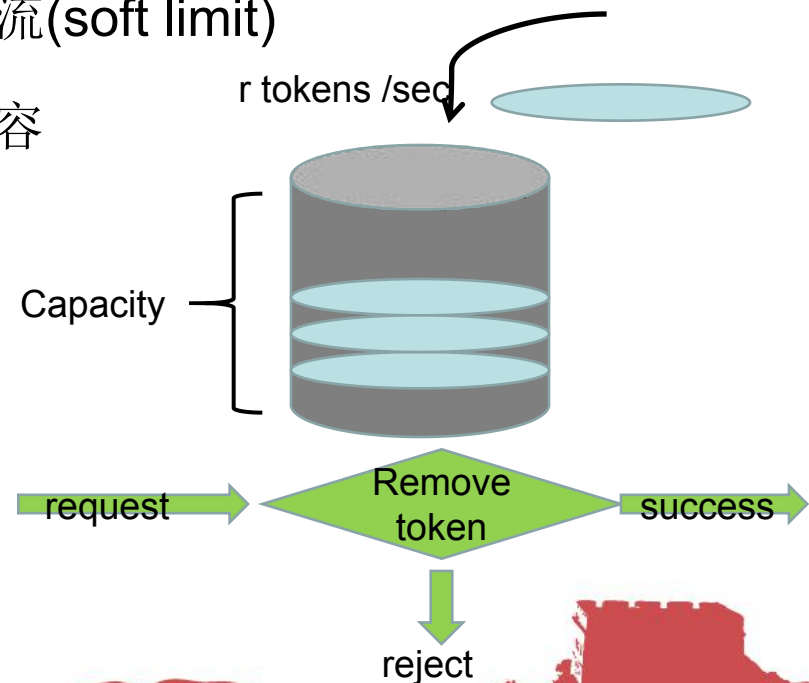
# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## • SDS ThriftServer (流量控制)

- 用户设置表的读写quota，设置表quota时检查集群能力
- 基于token bucket算法<sup>5</sup>进行限流(soft limit)
- 集群能力使用超80%时提醒扩容

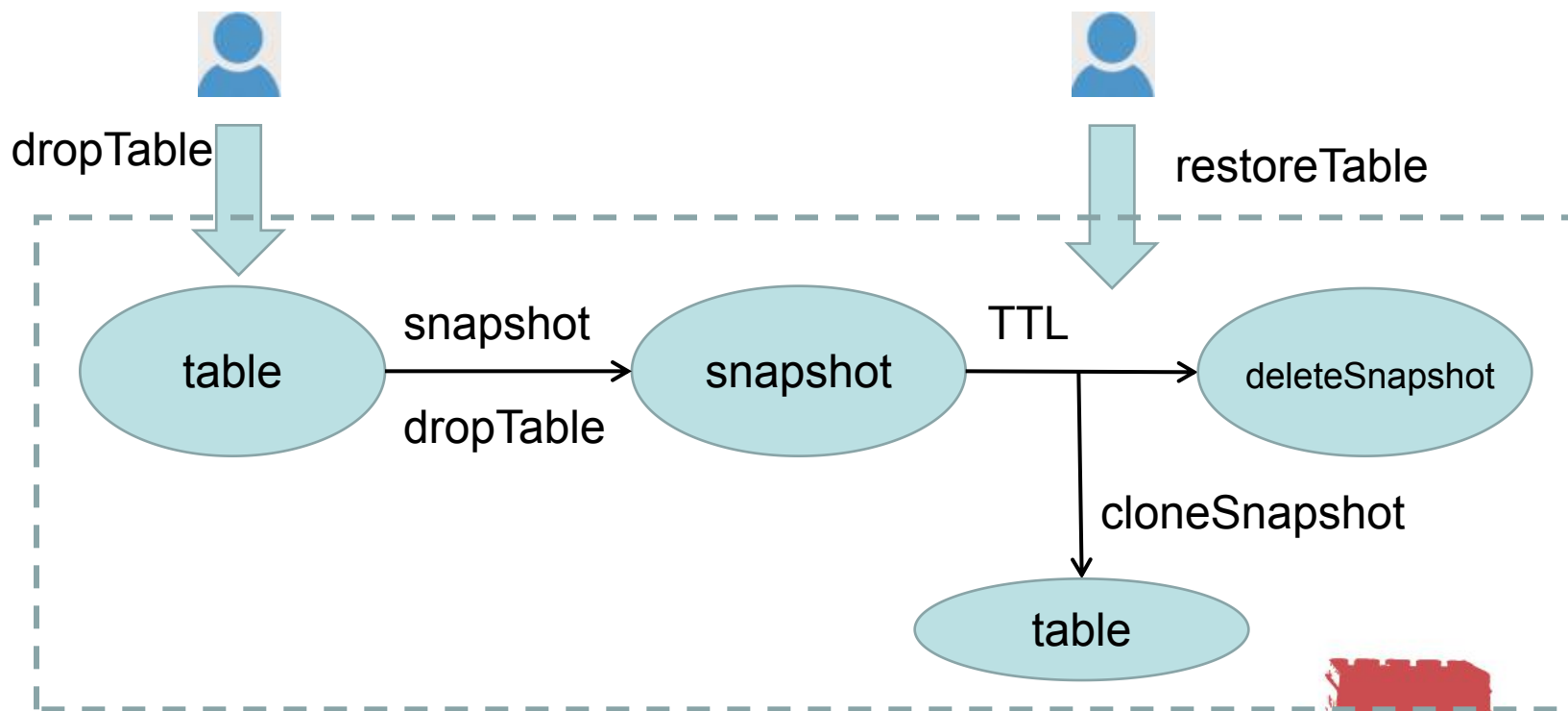


<sup>5</sup>[https://en.wikipedia.org/wiki/Token\\_bucket](https://en.wikipedia.org/wiki/Token_bucket)

# 结构化存储



- 软删除

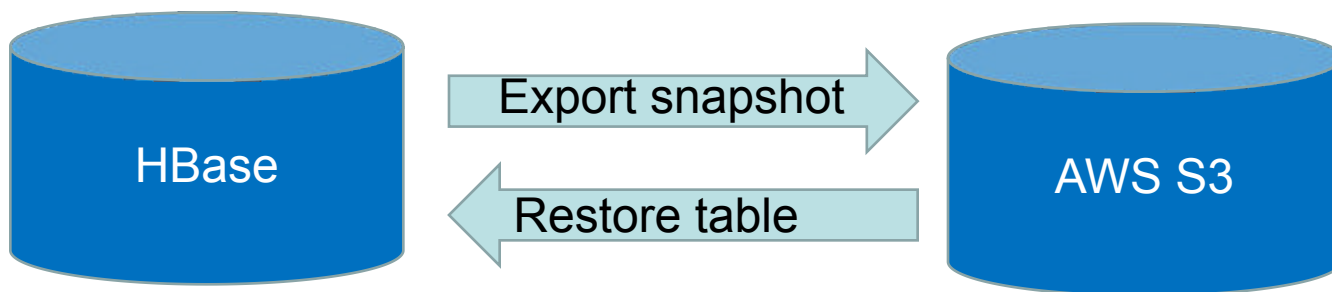


# 结构化存储



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 数据冷备



# 大纲



云计算开源产业联盟  
China Cloud Open Source Alliance

全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 小米结构化数据存储
  - 融合云平台



# 融合云



云计算开源产业联盟

全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT



小米融合云

产品

文档



团队管理



Docker 镜像



结构化存储



文件存储



HBASE



HDFS



消息队列



流式消息队列

LCS



OLAP



深度学习

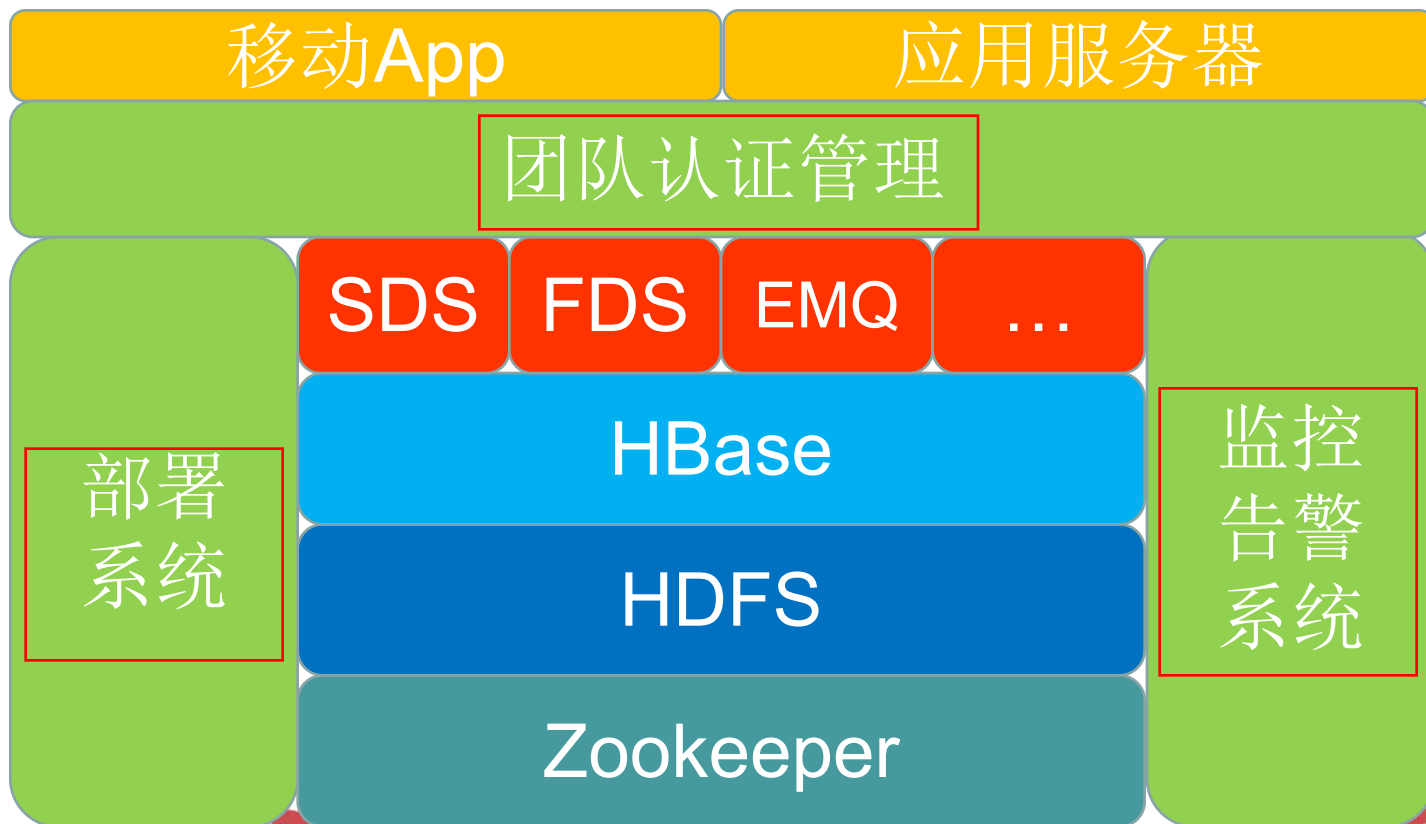


# 融合云



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 架构

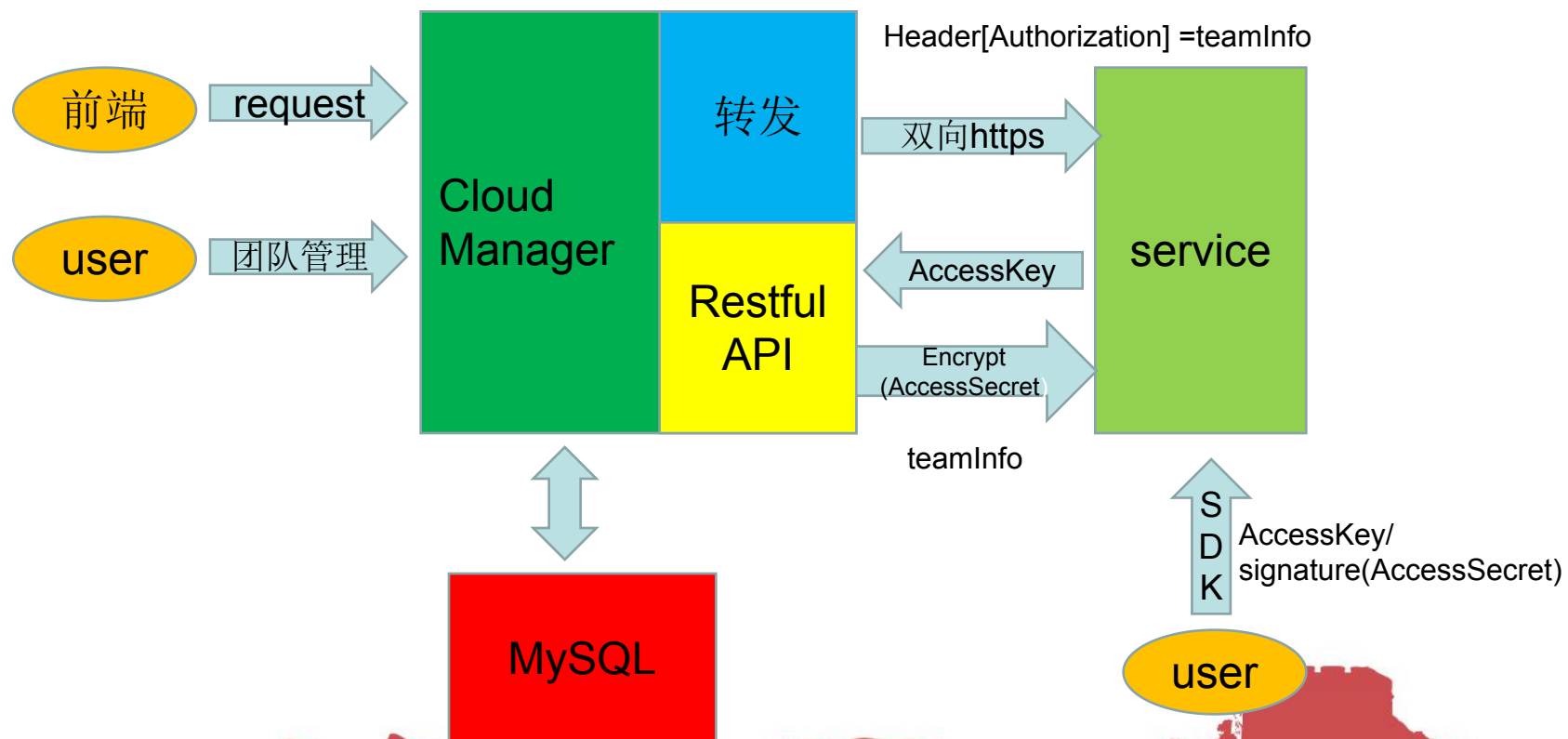


# 融合云



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

## • 团队认证管理——CloudManager

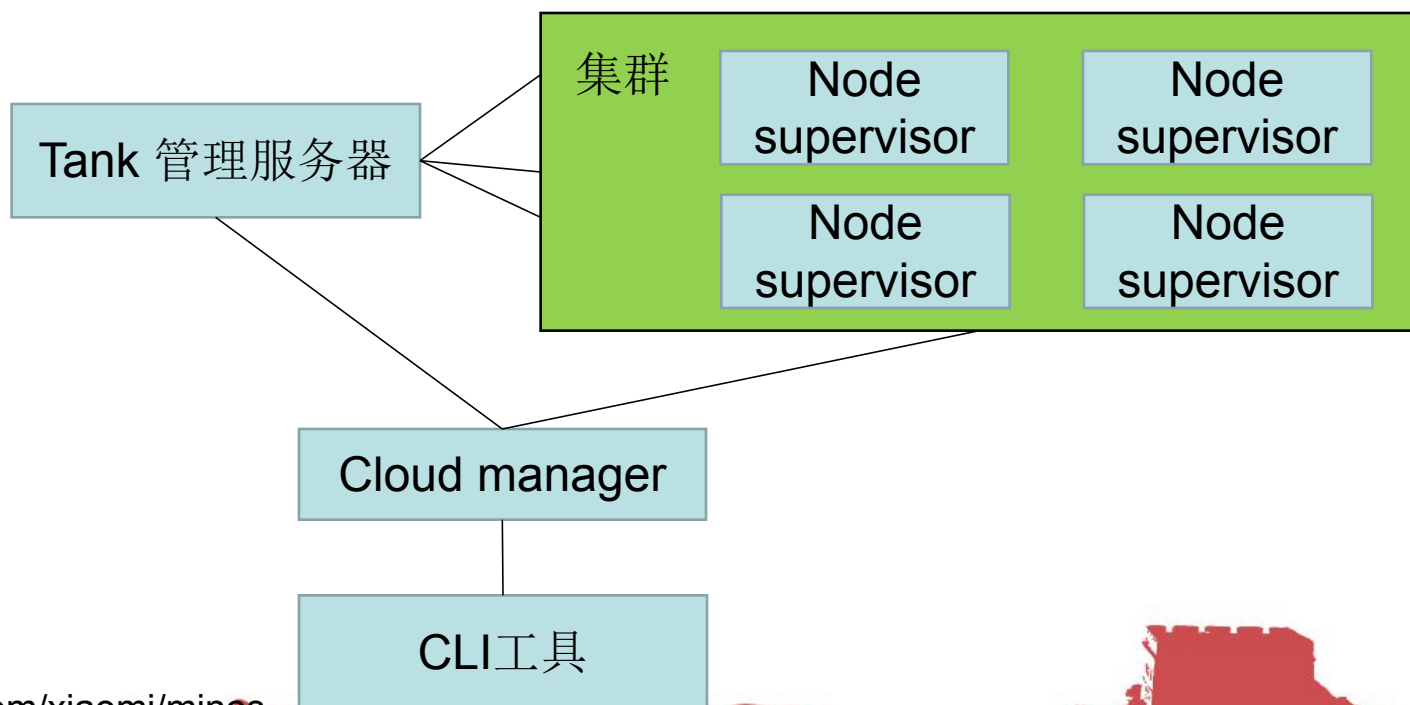


# 融合云



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 部署系统——Minos<sup>6</sup> 2.0
  - 增加了认证授权模块



<sup>6</sup><https://github.com/xiaomi/minos>

# 融合云



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 监控告警系统

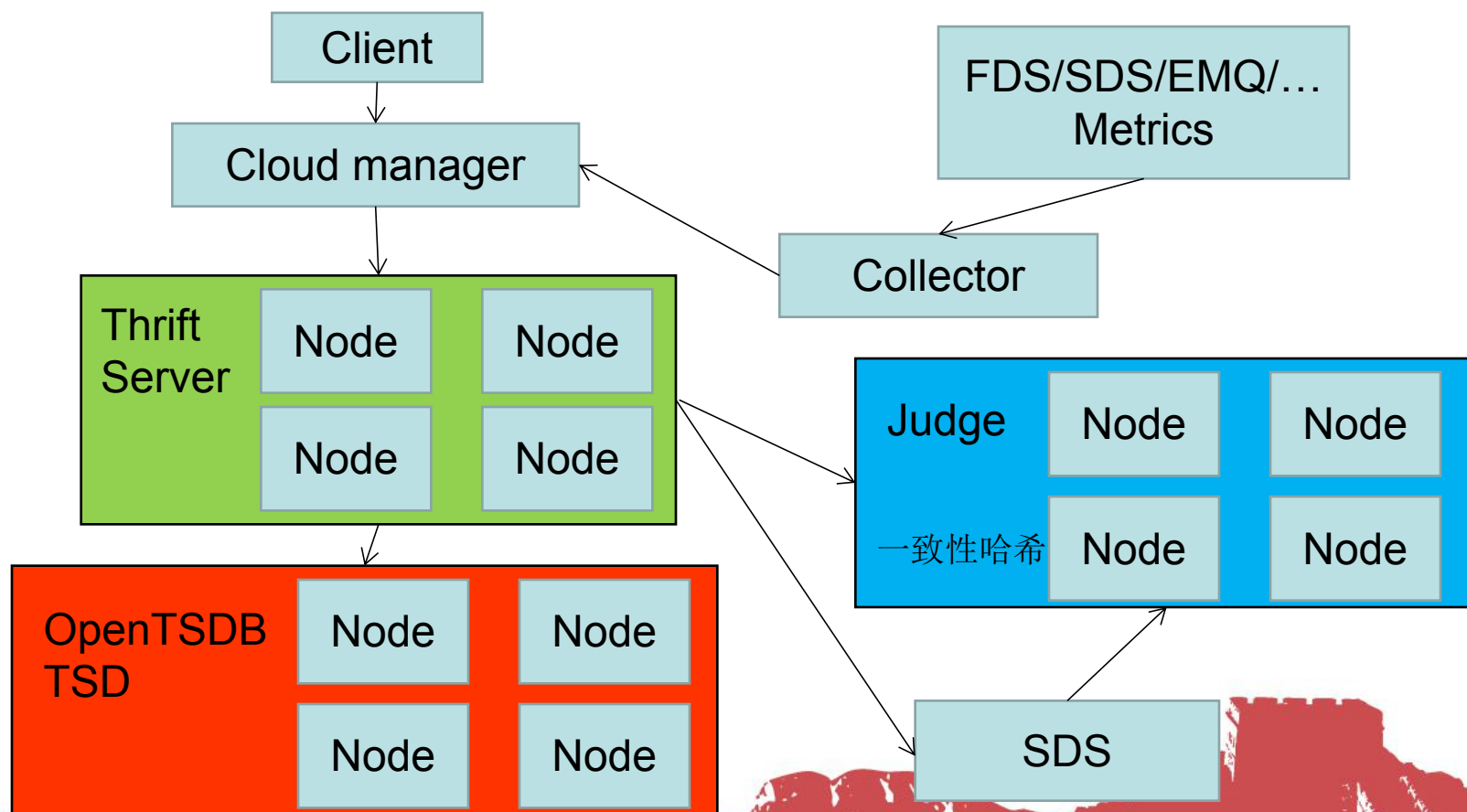


# 融合云



全球云计算开源峰会2017  
聚合云计算新势力，拥抱全世界新开源  
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 监控告警系统——架构





云计算开源产业联盟  
www.cioa.org.cn

全球云计算开源峰会2017

聚合云计算新势力，拥抱全世界新开源

GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT



THANK YOU

