



构建生态化 分布式数据库架构体系



About 自我介绍

陶勇

brave.taoy@alibaba-inc.com

Twitter & Sina: @bravetao



Index 内容概要

- 分布式数据库解决方案 @ Alibaba B2B
 - 分布式数据存储与访问
 - Cobar
 - 准实时增量数据获取与消费
 - Erosa/Eromanga
 - 多维度数据同步与网站镜像
 - Otter
- 构建分布式数据库生态架构 @ Alibaba B2B
 - 全站数据架构
 - 思考与展望



Keywords 关键词

MySQL protocol

Schema垂直拆分

Table水平拆分

Global Failover

实时日志解析

事务顺序

Global ID

实时镜像

双向同步

同步事务支持

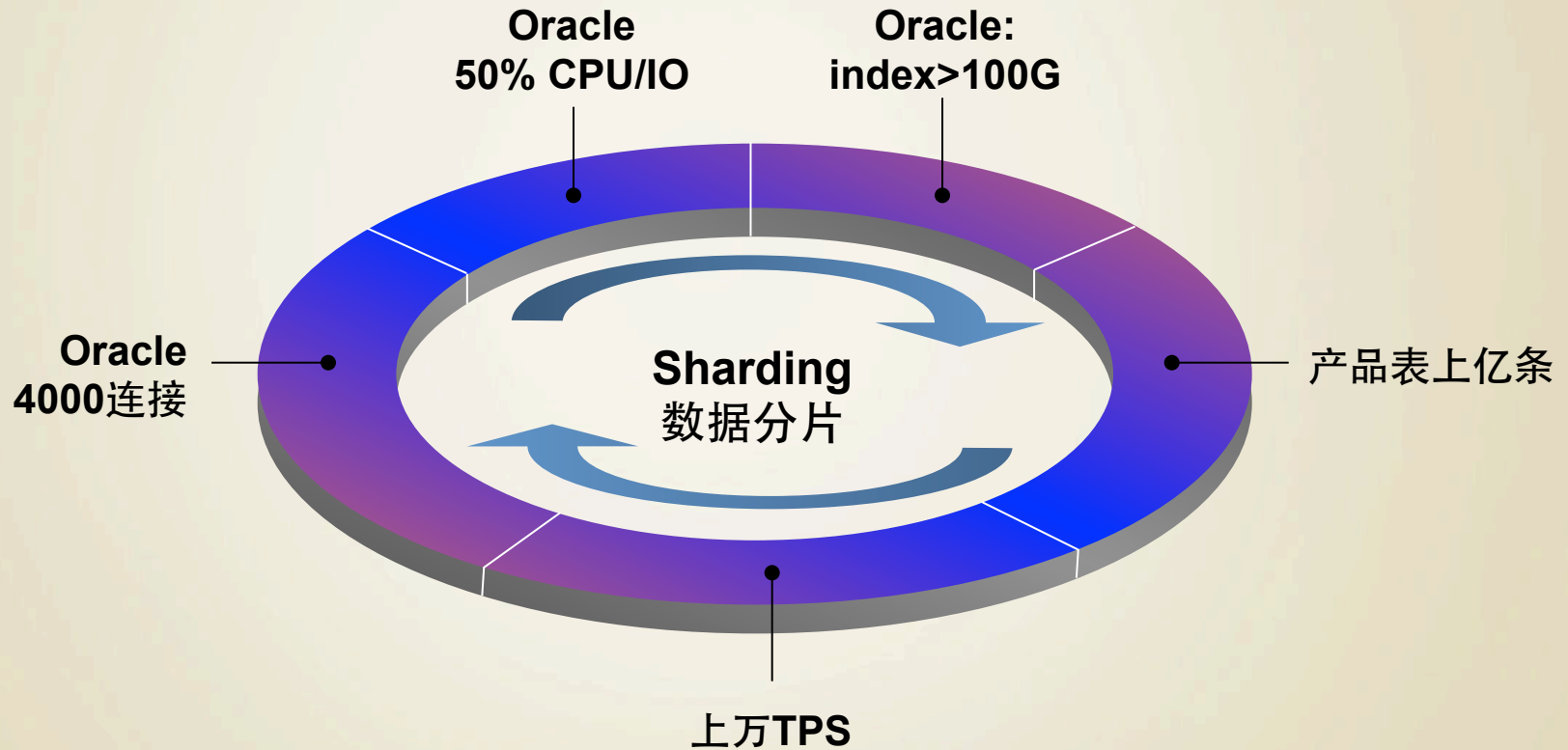


Cobar

分布式数据存储与访问

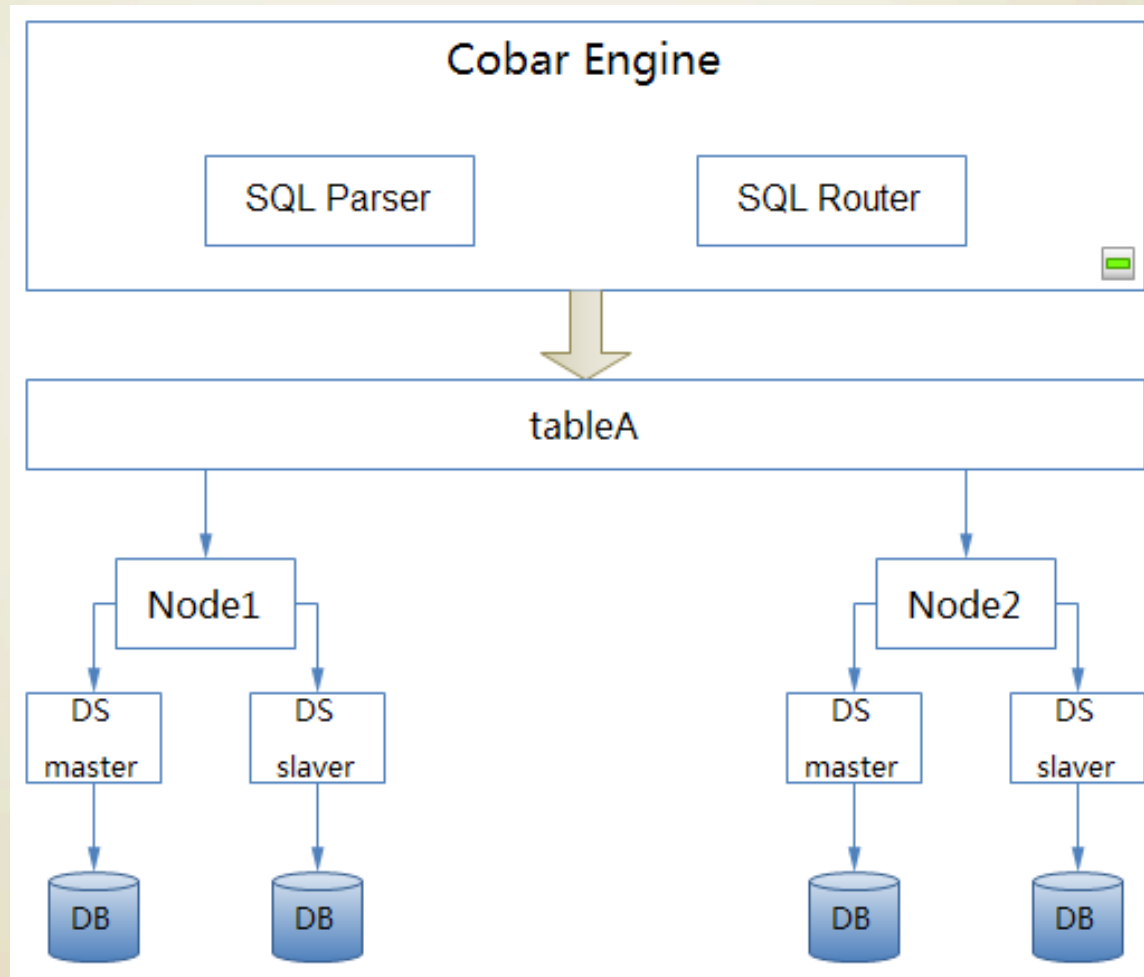


遇到的问题 - 2008



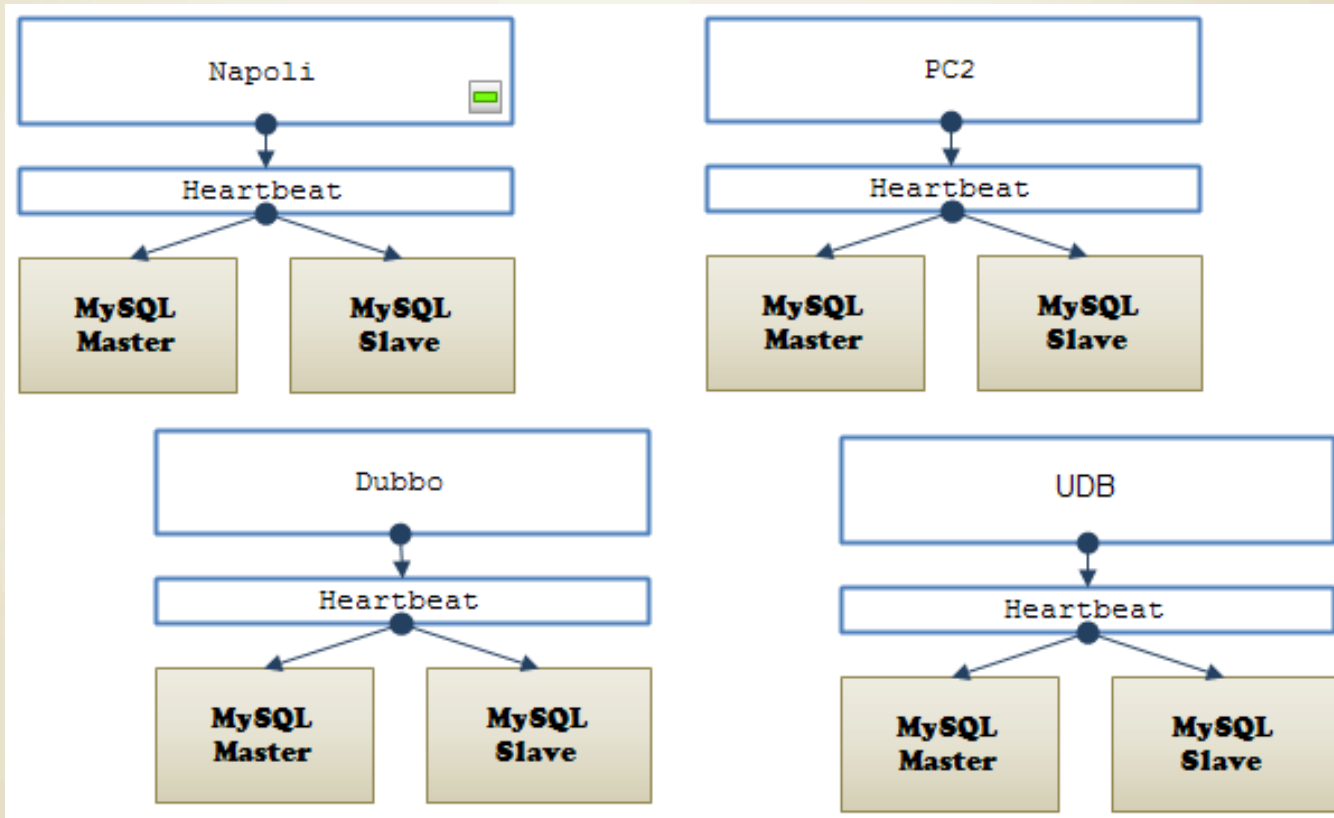


v0.6-1.0.x (08-10)



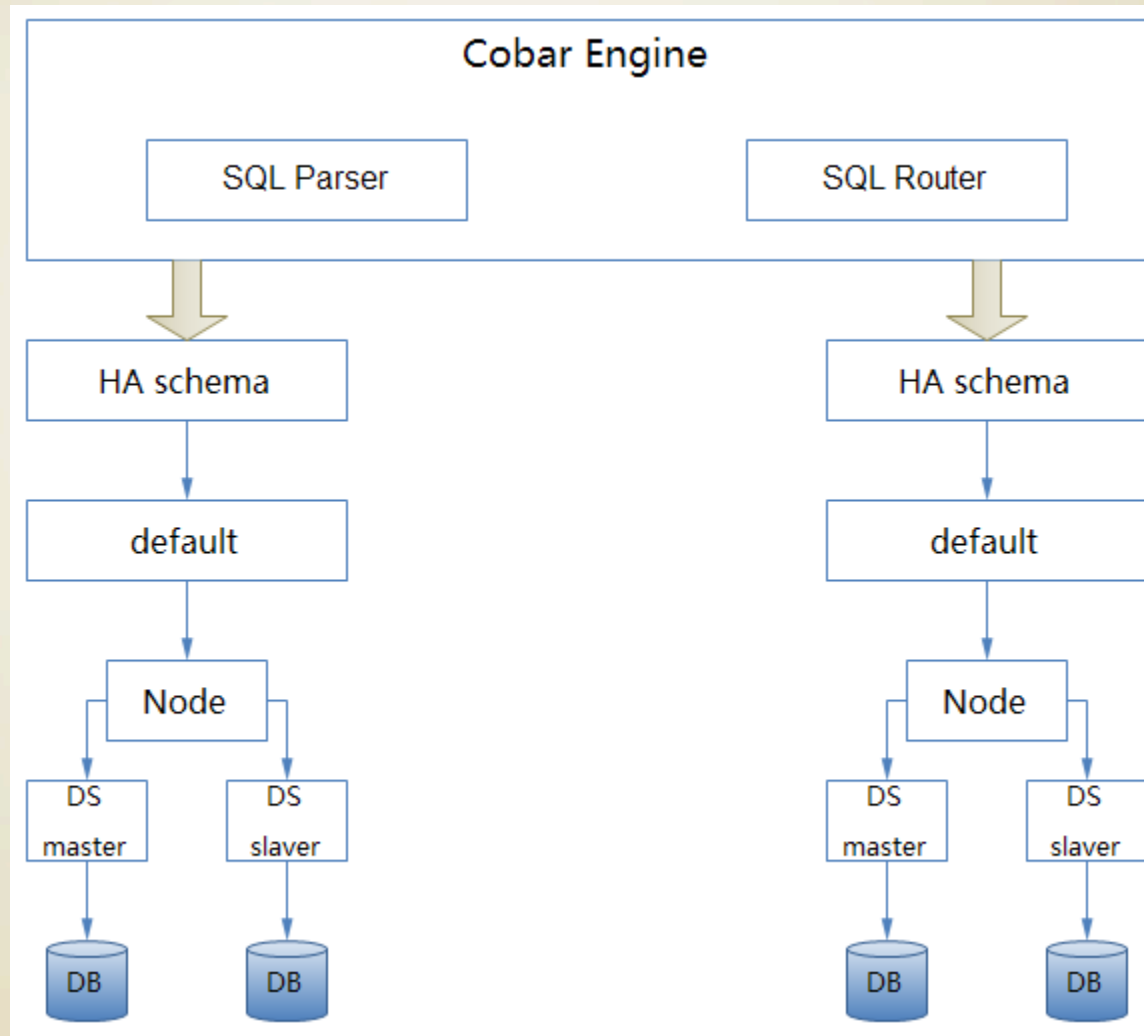


遇到的问题 - 2010



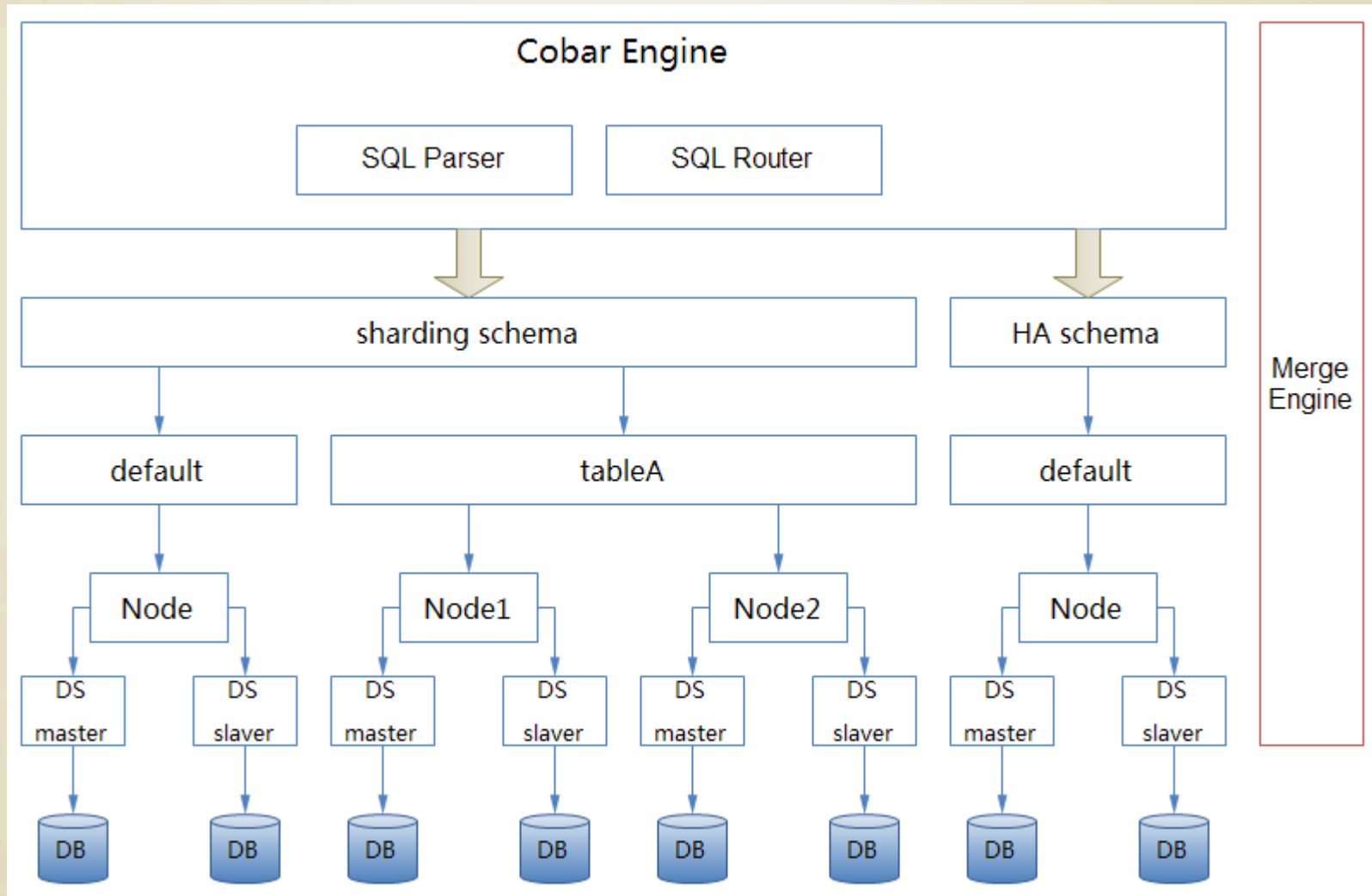


v1.1.x (10-11)



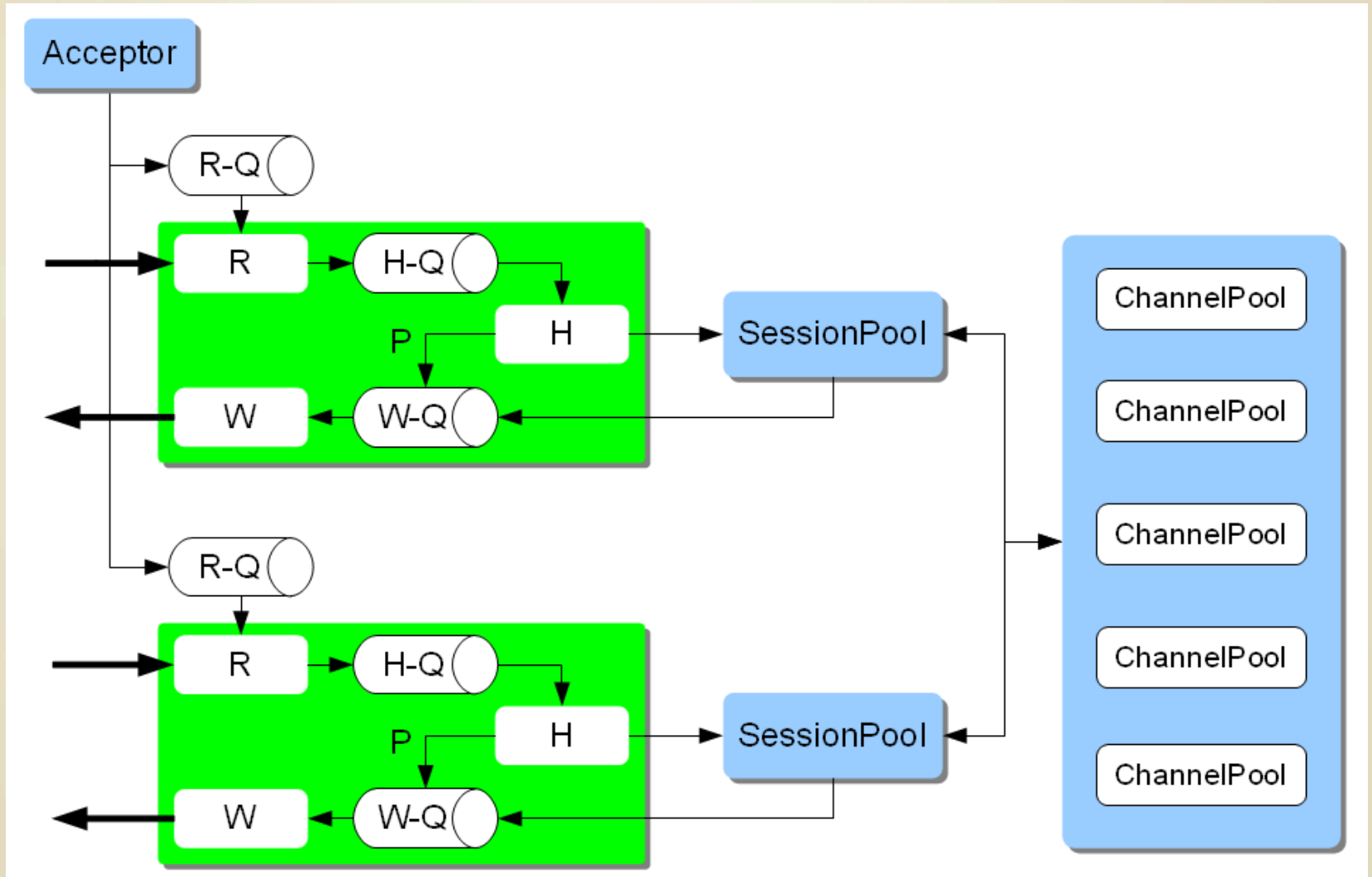


现在 (10-2012.12.23)



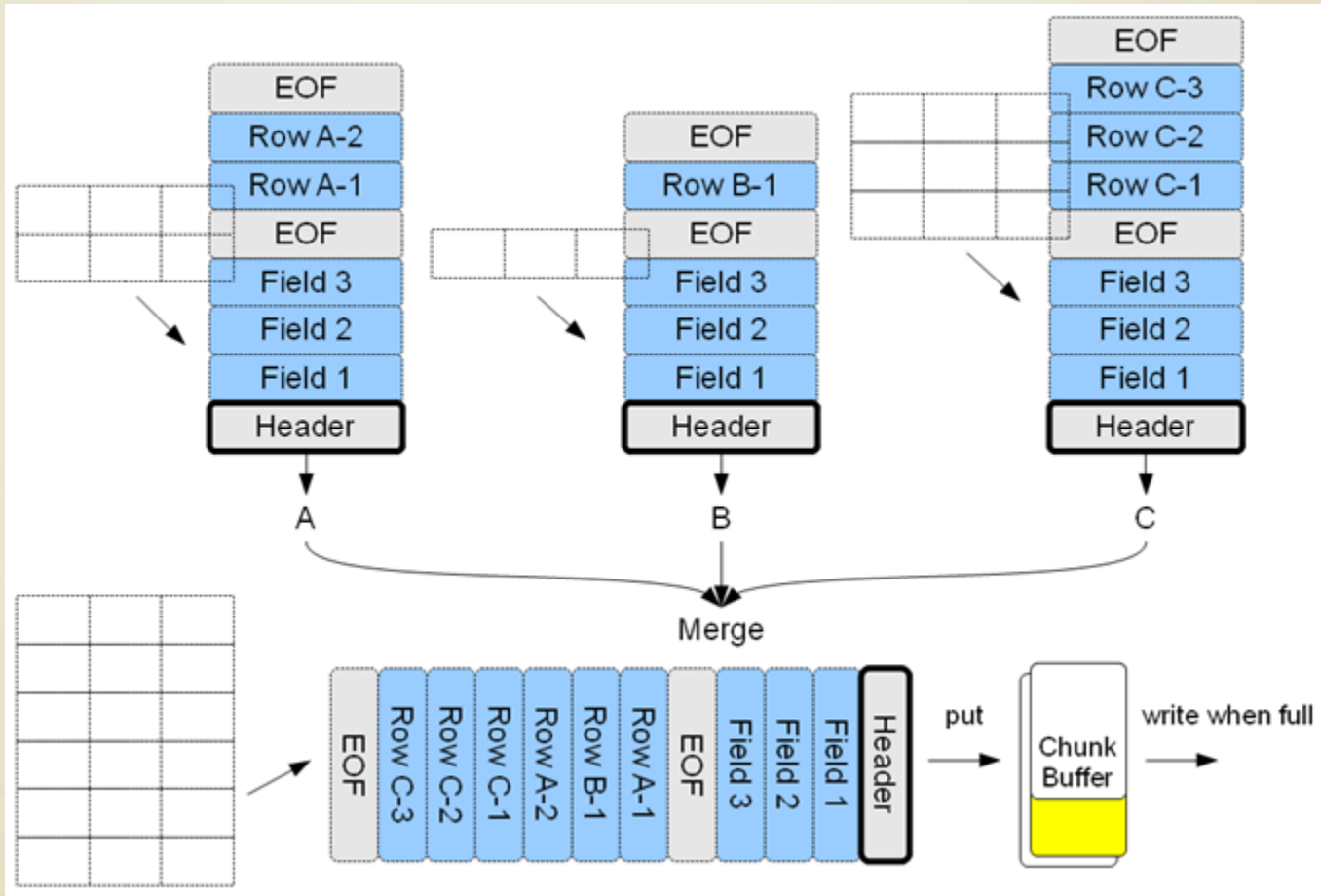


细节：线程复用模型



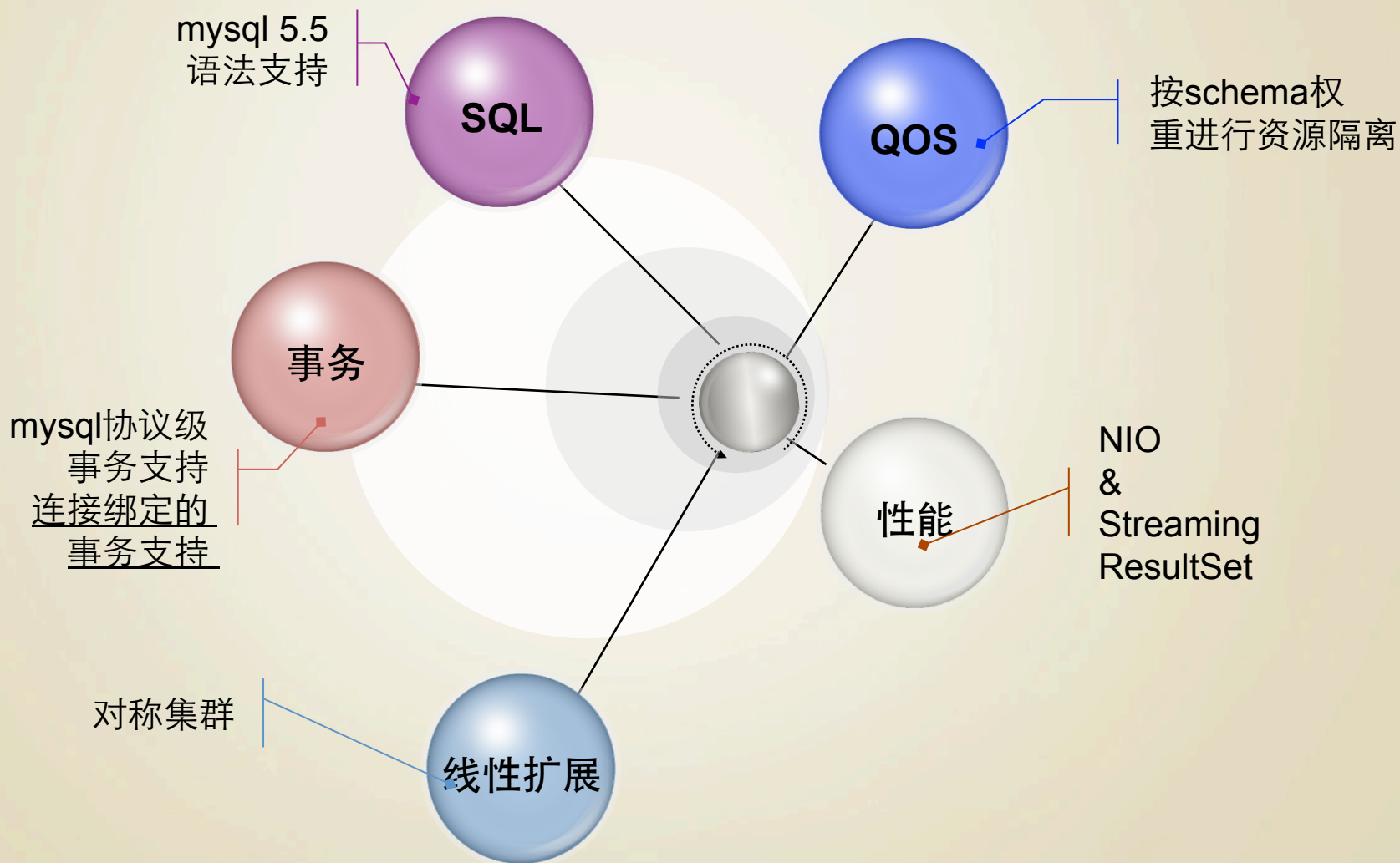


细节：事务、ResultSet





特性





一些数据

- 20+ Cobar
- 线上高峰期 4-5w TPS，单机性能 10w TPS
- 15并发以上，整体性能超过直接访问mysql
- 支持mysql 5.5 全部DML和部分DDL语法



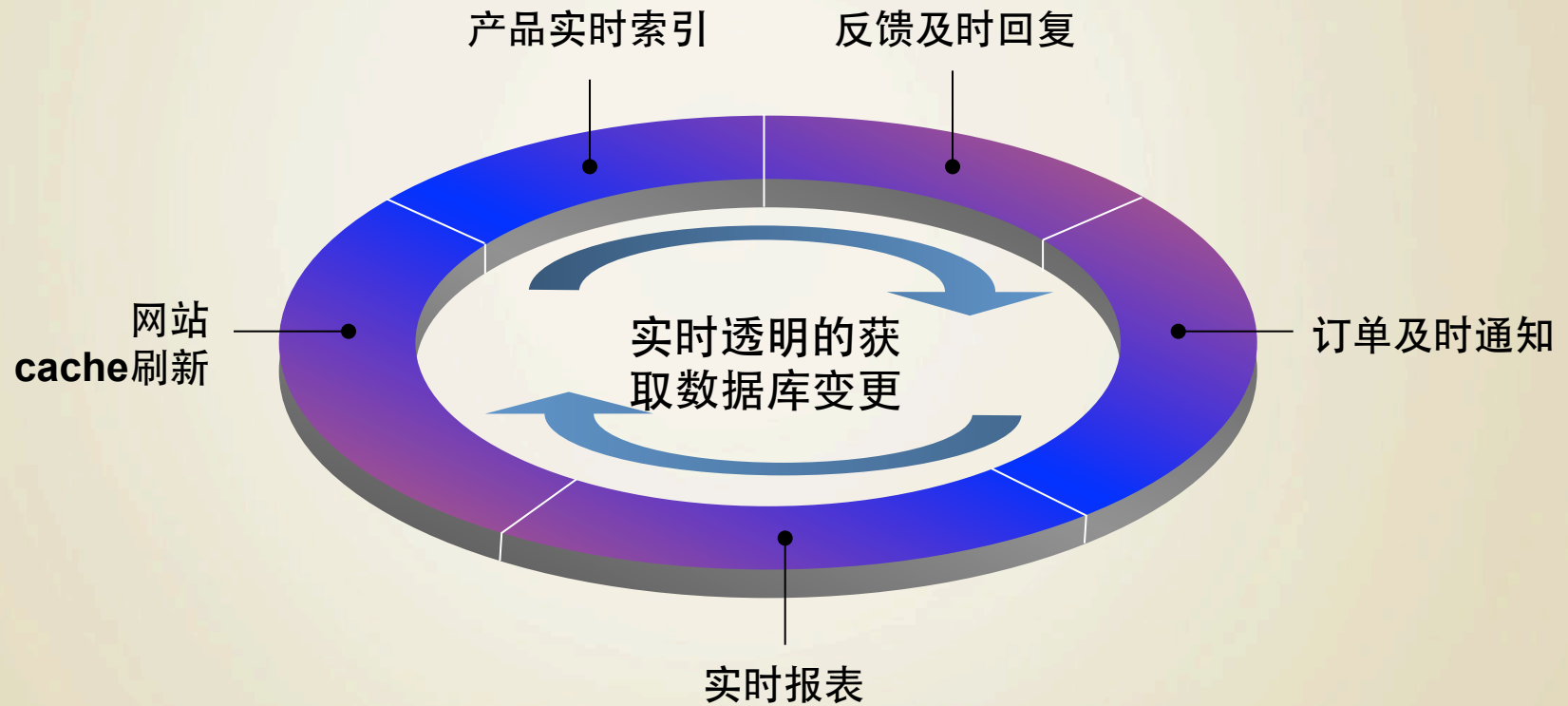
分布式数据库解决方案

E&E

准实时增量数据获取与消费



遇到的问题



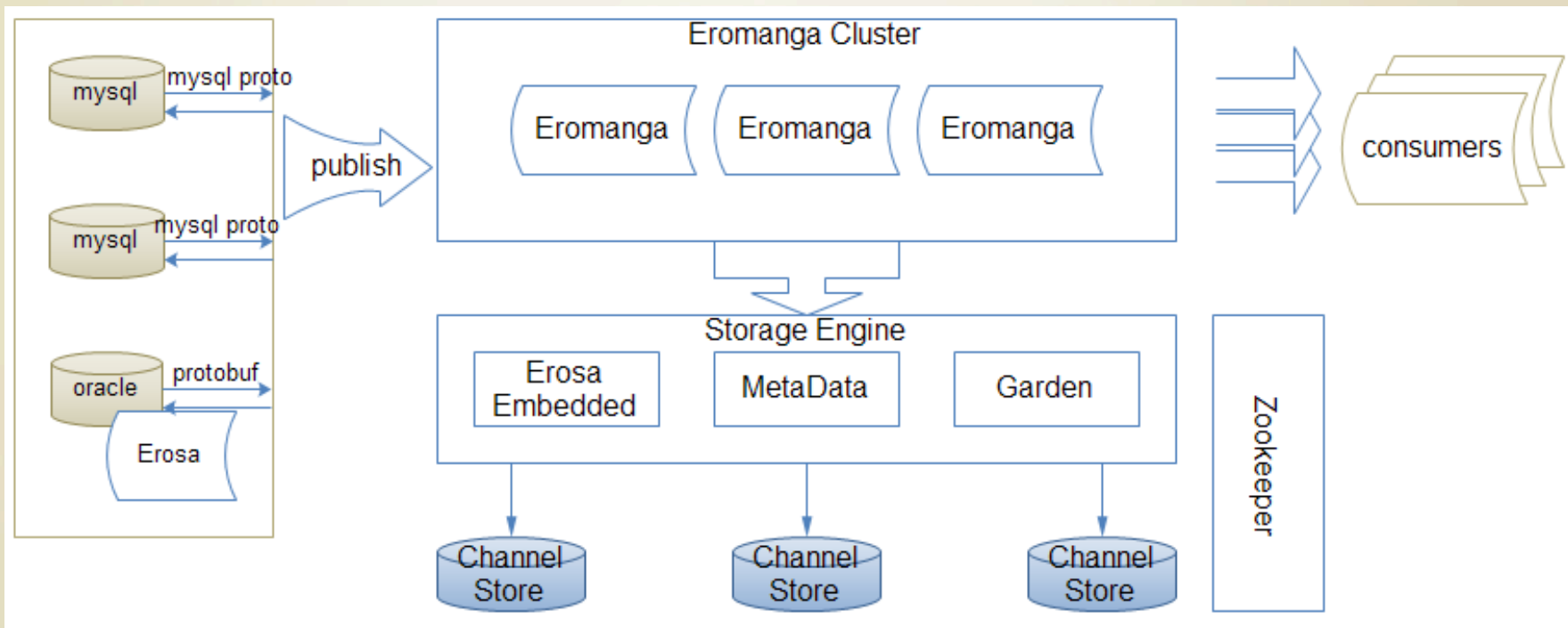


准实时增量数据获取与消费

- 以前的做法
 - DB Trigger
 - Dump table
 - Application MQ
- 问题
 - 运维困难
 - 数据库、网络瞬时压力大
 - 业务侵入性强



整体架构

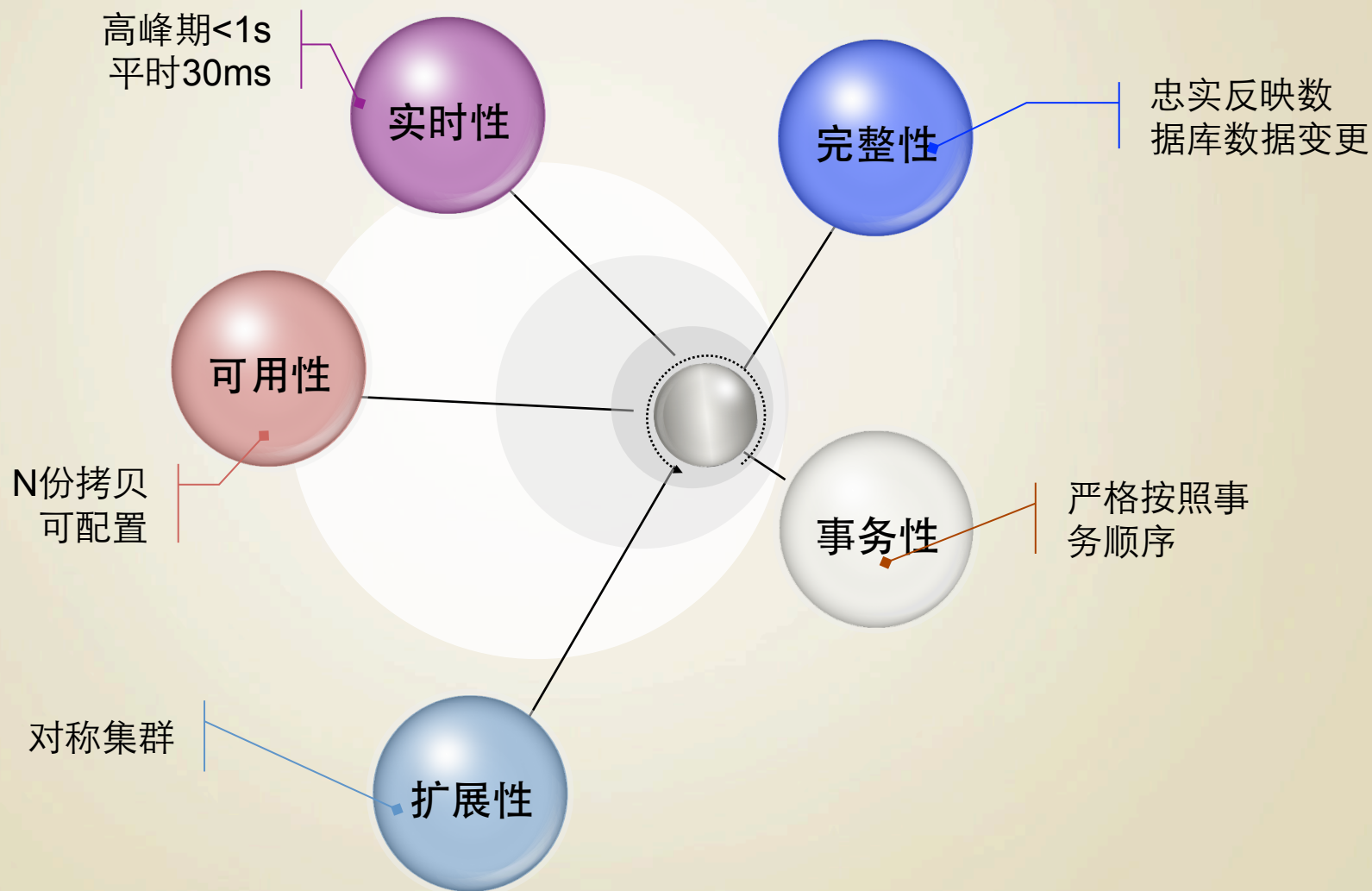




- 数据解析
 - oracle : redolog parser
 - Archive log ~ 2min
 - Online log < 10ms
 - Open column supplemental log
 - mysql : binlog parser
 - COM_BINLOG_DUMP
- 数据消费
 - 对称集群
 - Data cursor : ZooKeeper



特性





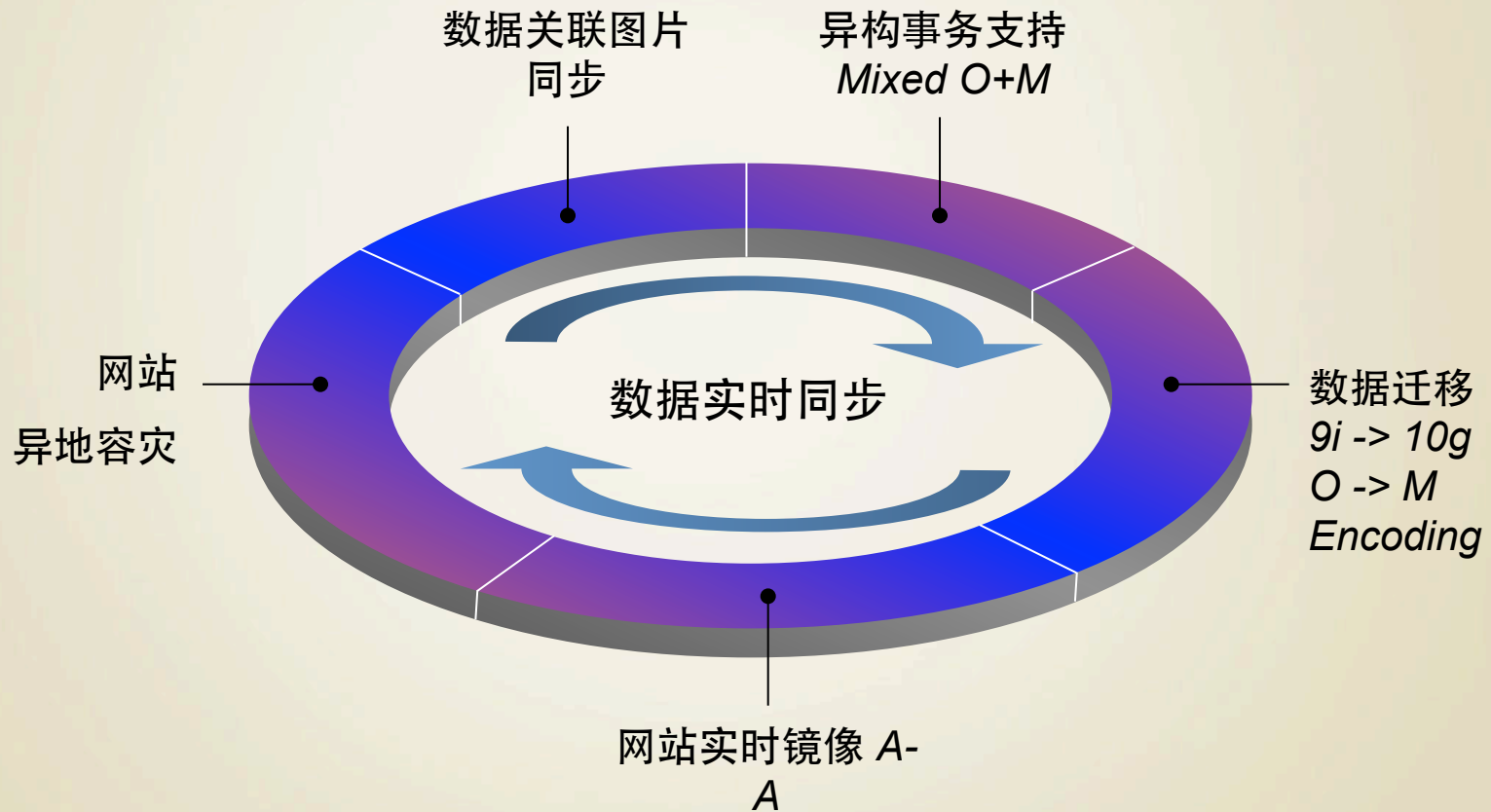
分布式数据库解决方案

Otter

多维度数据同步与网站镜像

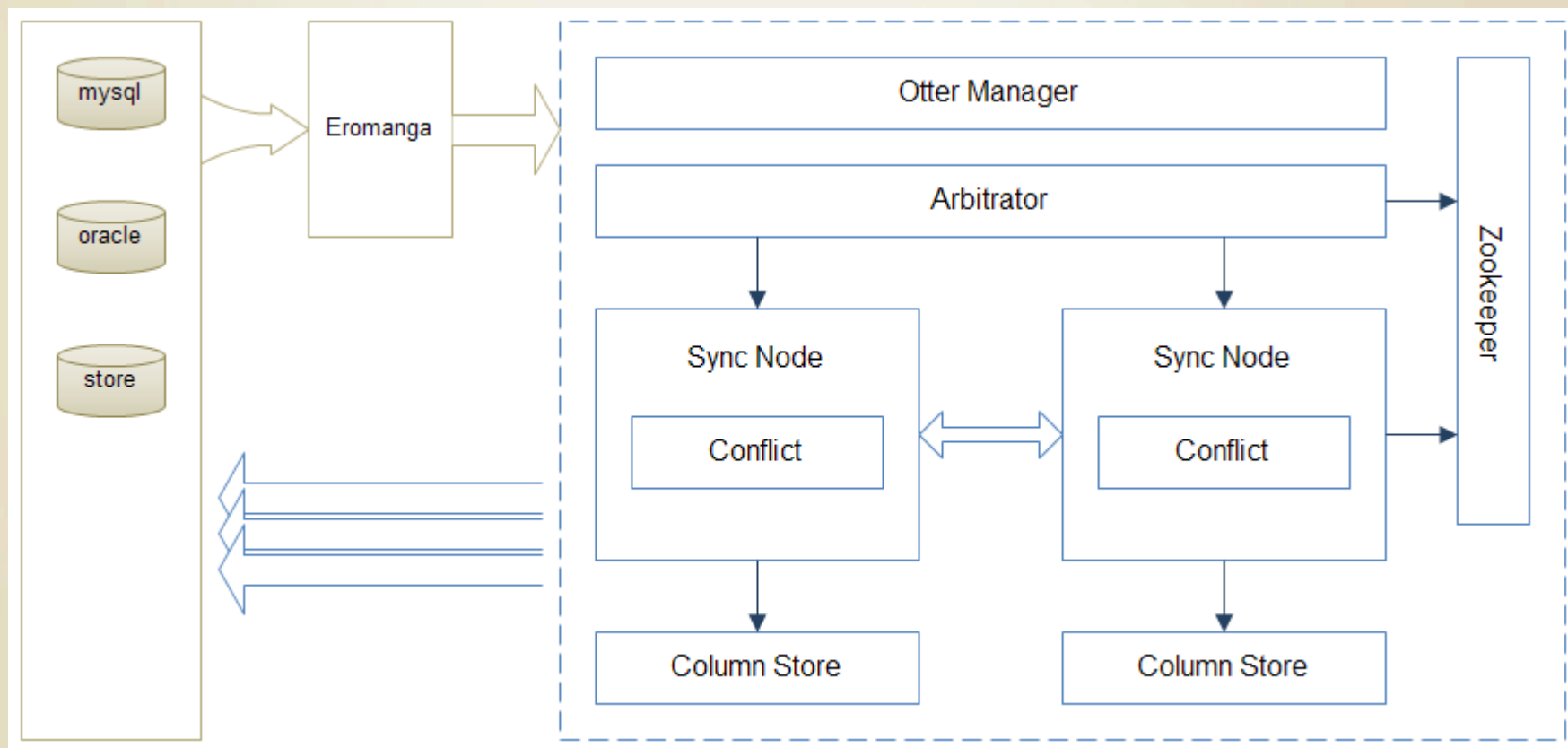


遇到的问题





整体架构





- 数据同步
 - 自定义字段过滤
 - 自定义文件同步逻辑
 - 按事务表并发加载
 - 按表PK hash并发加载
- 数据冲突
 - 实时字段级冲突合并
 - 冲突判断及解决



现有的应用场景

- 数据库
 - 备份：多master -> 单slaver
 - 异构迁移、跨版本迁移
 - *Oracle Active-Active*
- 网站容灾
 - 容灾备份
 - 读写分离
- 网站镜像
 - 双向读写
 - 按字段同步
 - 按事务并发同步



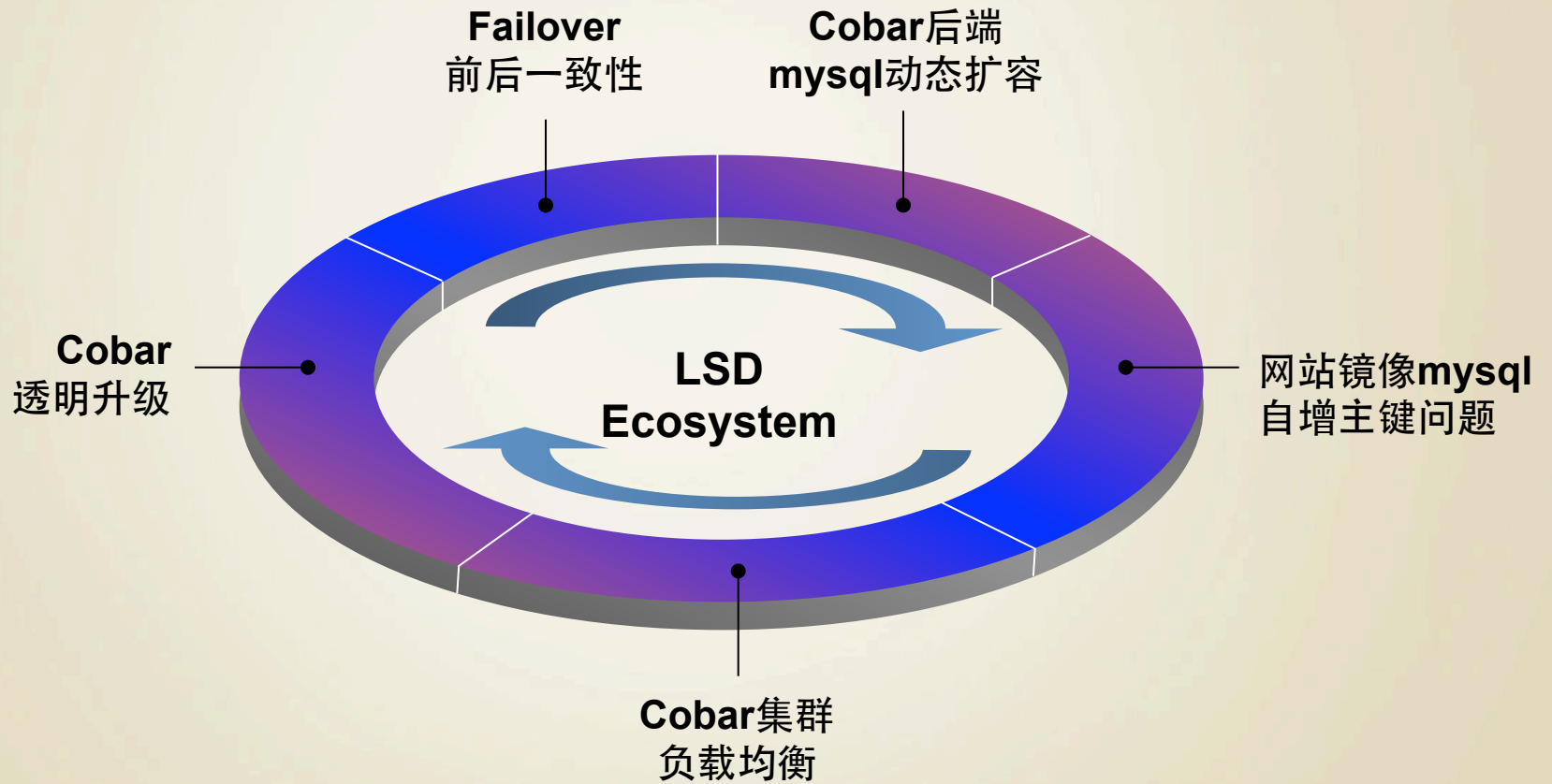
分布式数据库解决方案

Ecosystem

构建分布式数据库生态架构

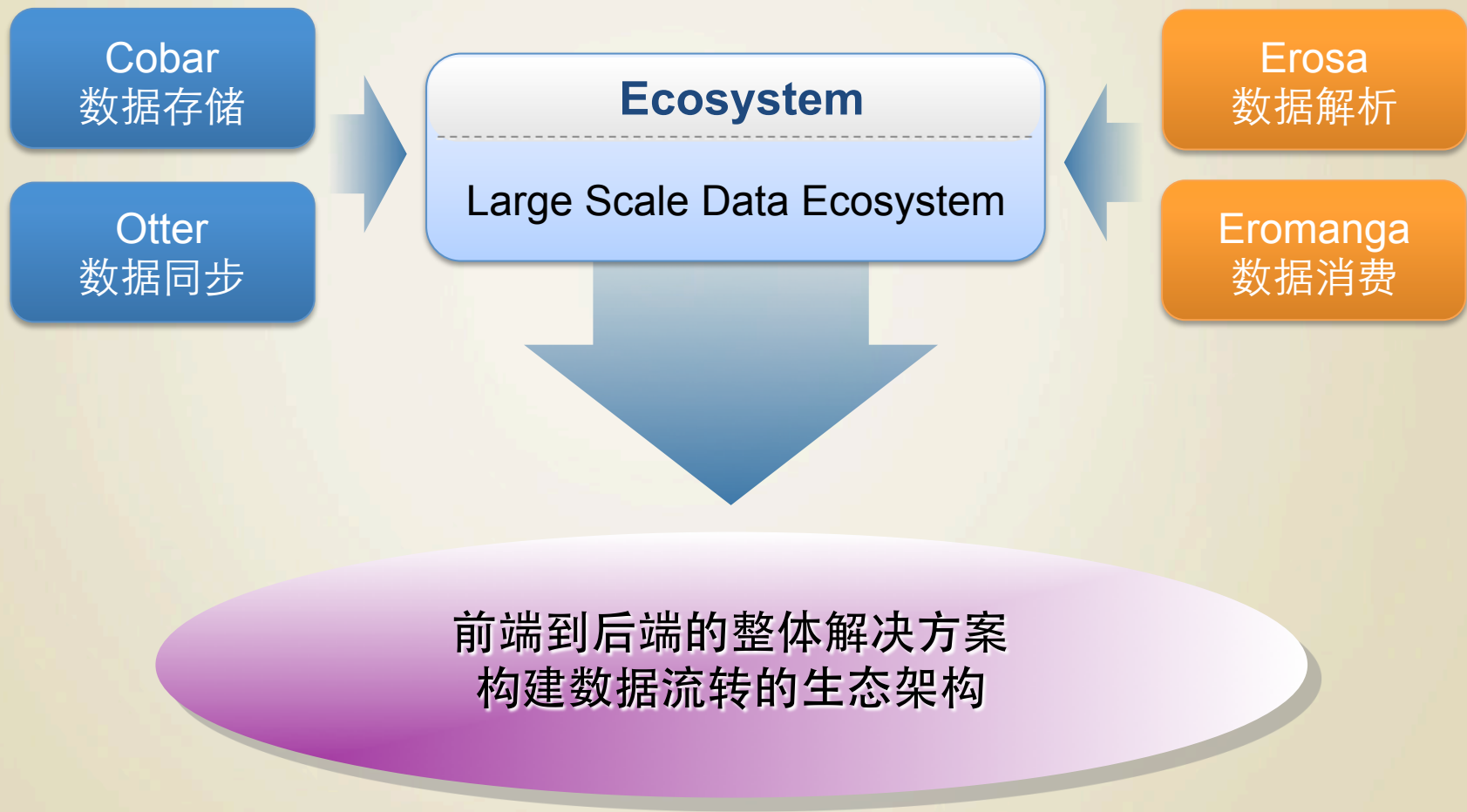


遇到的问题



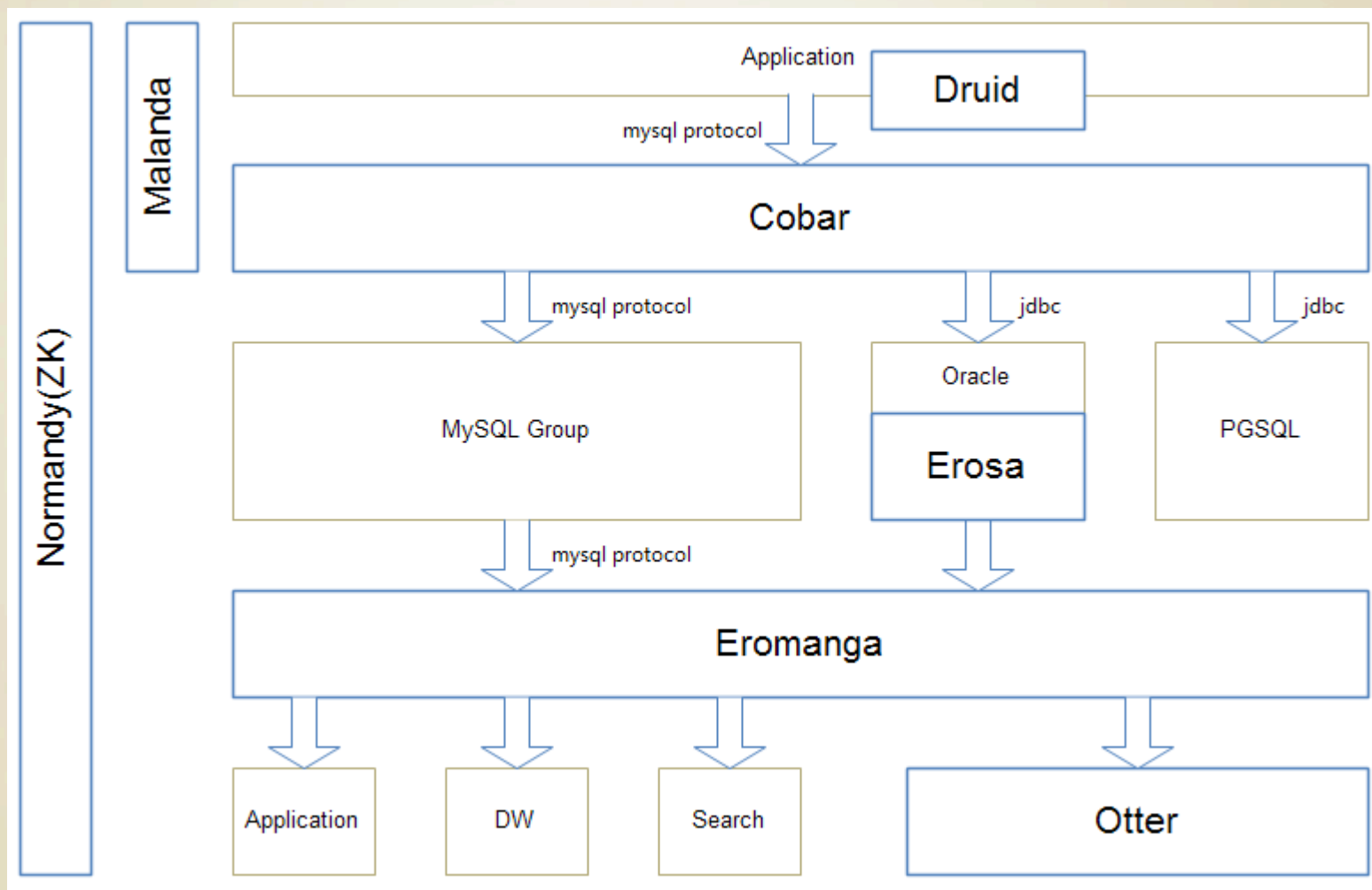


分布式数据库架构生态体系





整体架构





思考与展望

- 产品架构闭环
 - Under constructing
 - OLTP & OLAP
 - NoSQL数据库
- 软硬结合
 - 单机多实例
 - 单机高性能



MySQL优化

单机多实例

- 存储技术飞速发展，IO不再是瓶颈
- MySQL对多核CPU利用率低

单机高性能

- RAID: BBWC/Fastpath
- Fusionio
- Flashcache
- Semi-Sync



Q&A

Thanks!

No best, only the most suitable



北京站 · 2012年4月18~20日
www.qconbeijing.com (11月启动)

QCon杭州站官网和资料
www.qconhangzhou.com

全球企业开发大会

INTERNATIONAL
SOFTWARE DEVELOPMENT
CONFERENCE