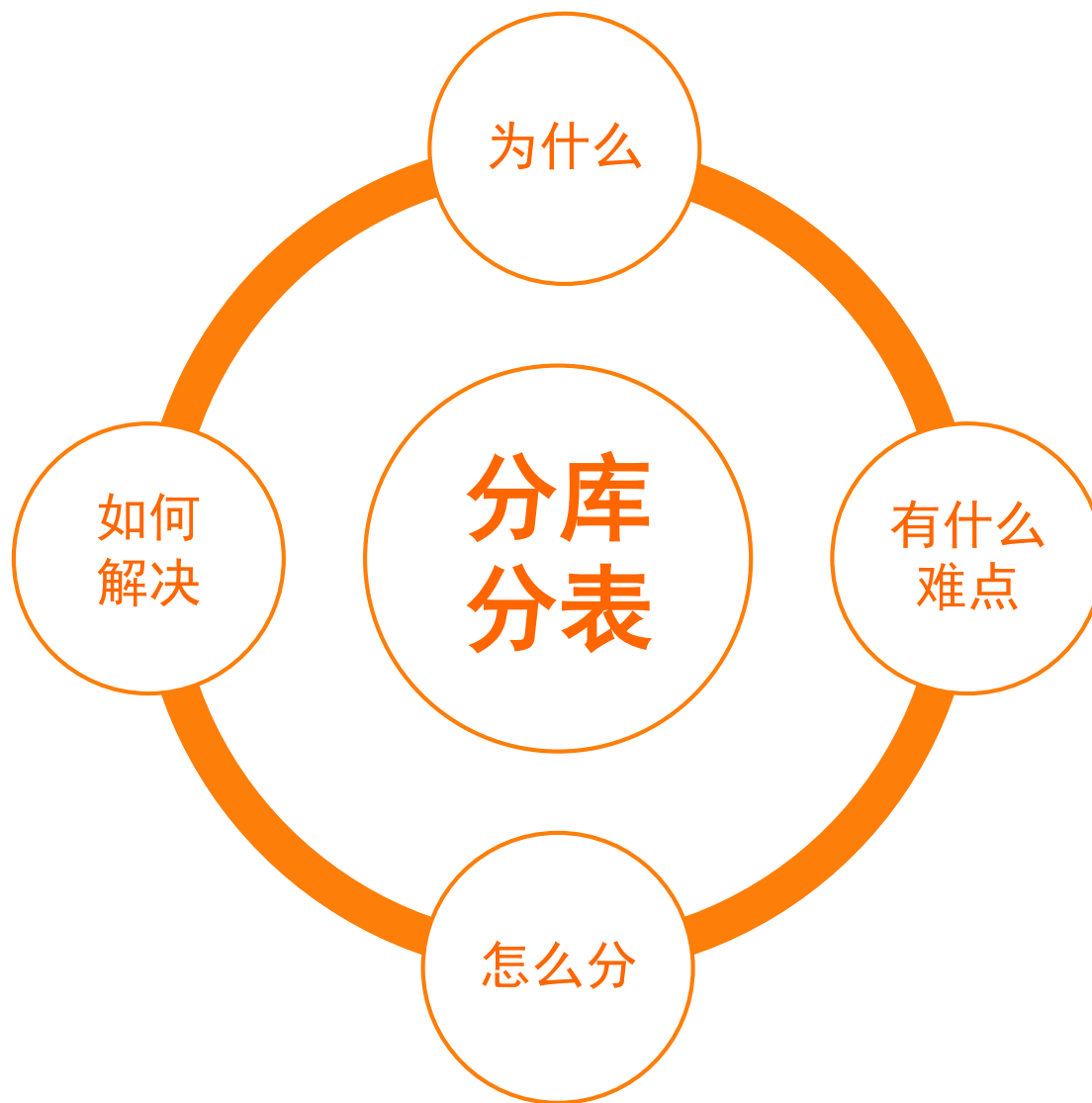




# 基于客户端的分库分表实践

# 目录

- 分库分表理论
- CDS实践
- 如何接入CDS



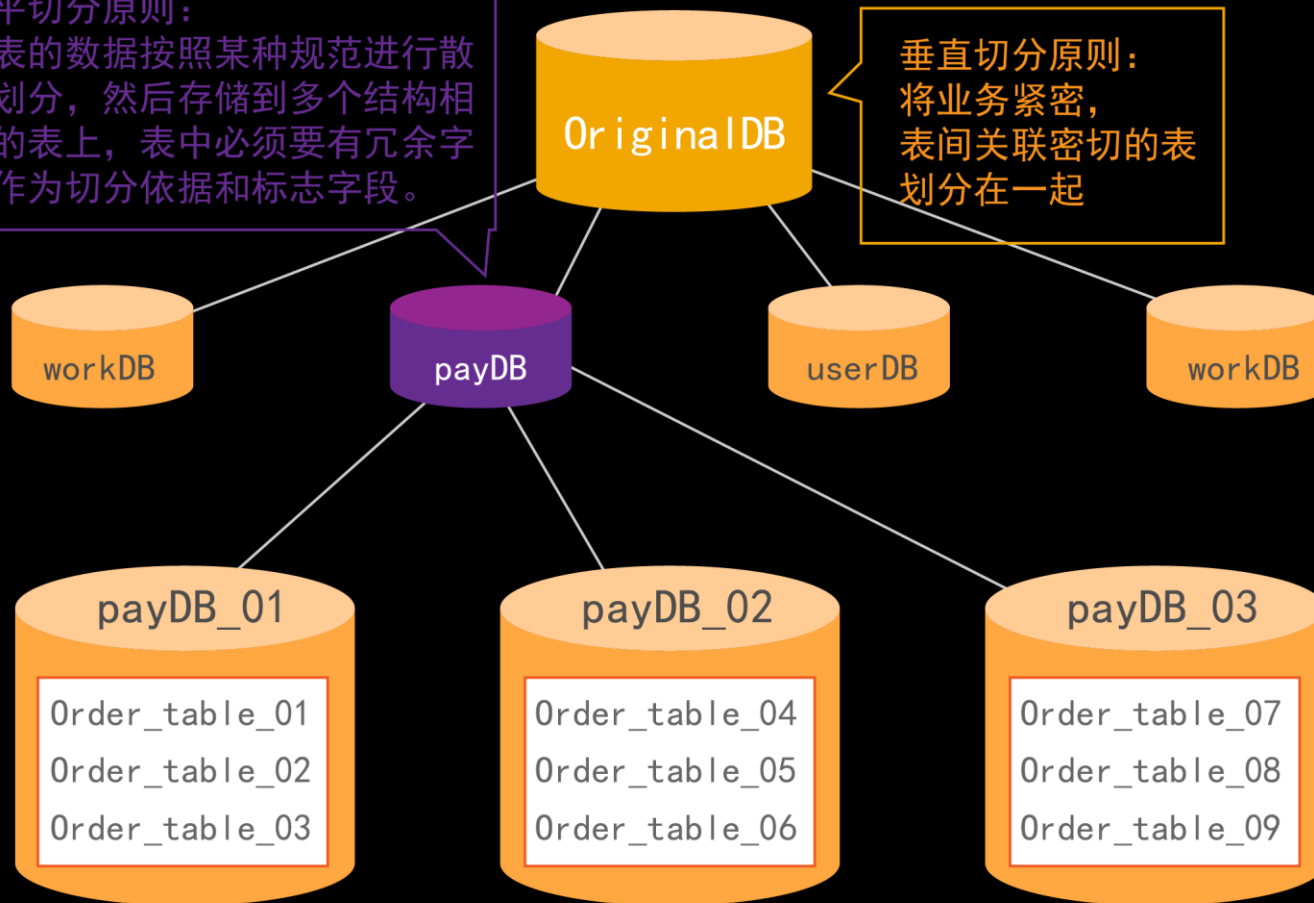
- 为什么要分库分表?
  - ① 业务数据量剧增
  - ② 访问并发量剧增
  - ③ 业务系统可用性需求



- 如何选择拆分策略?
  - ① 垂直拆分
  - ② 水平拆分
  - ③ 两者结合

水平切分原则：  
将表的数据按照某种规范进行散列划分，然后存储到多个结构相同的表上，表中必须要有冗余字段作为切分依据和标志字段。

垂直切分原则：  
将业务紧密，  
表间关联密切的表  
划分在一起



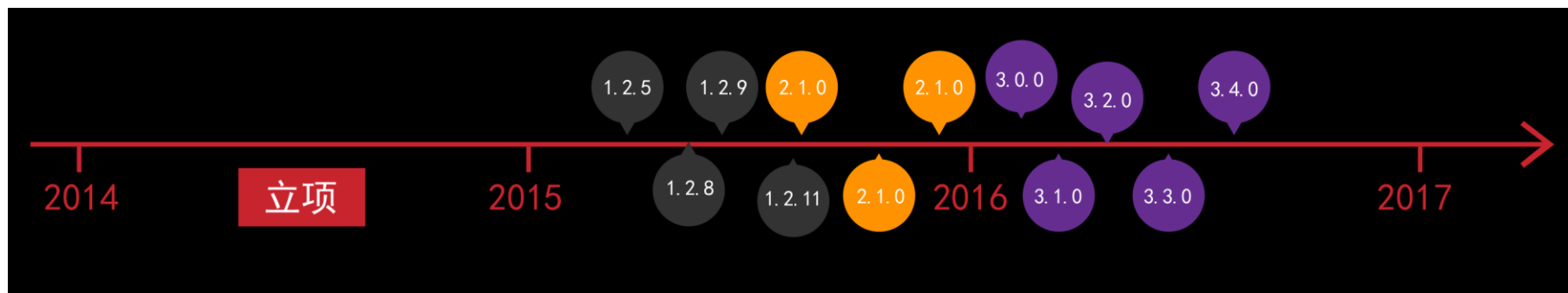
- 分库分表带来了哪些难题？

- ① 分布式事务
- ② 复杂SQL查询
- ③ 跨库跨表Join
- ④ 数据运维

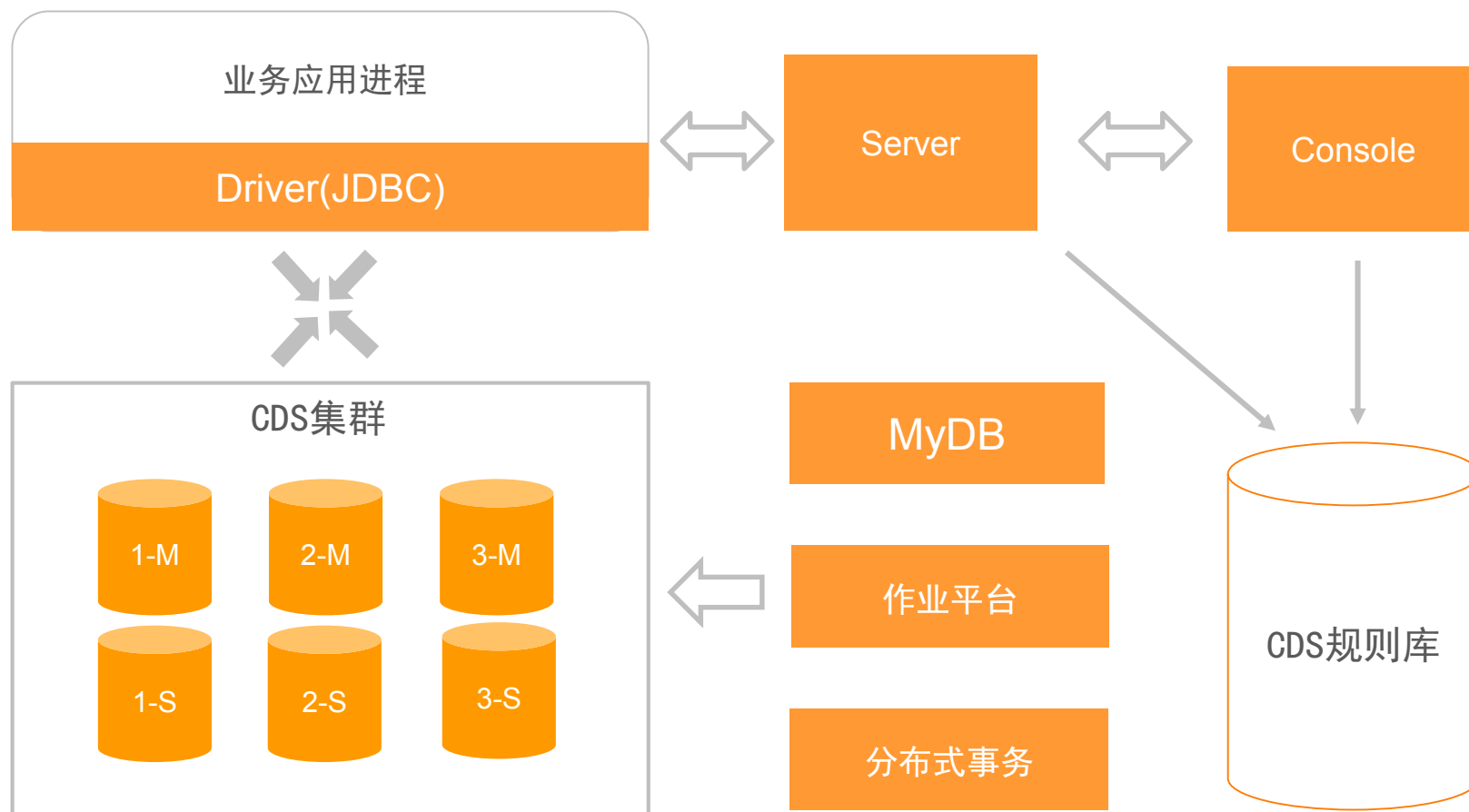
# Completed Database Sharding

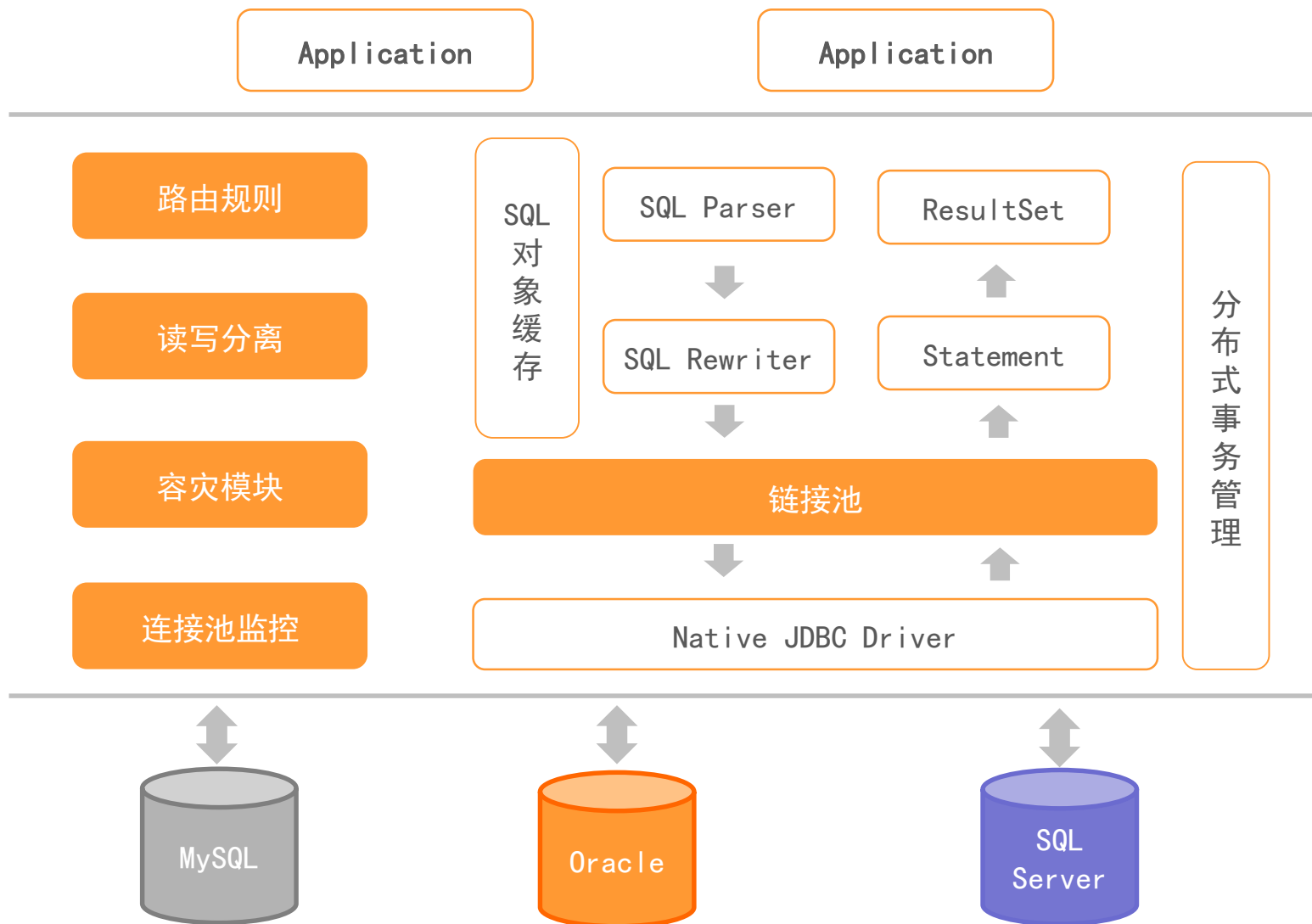


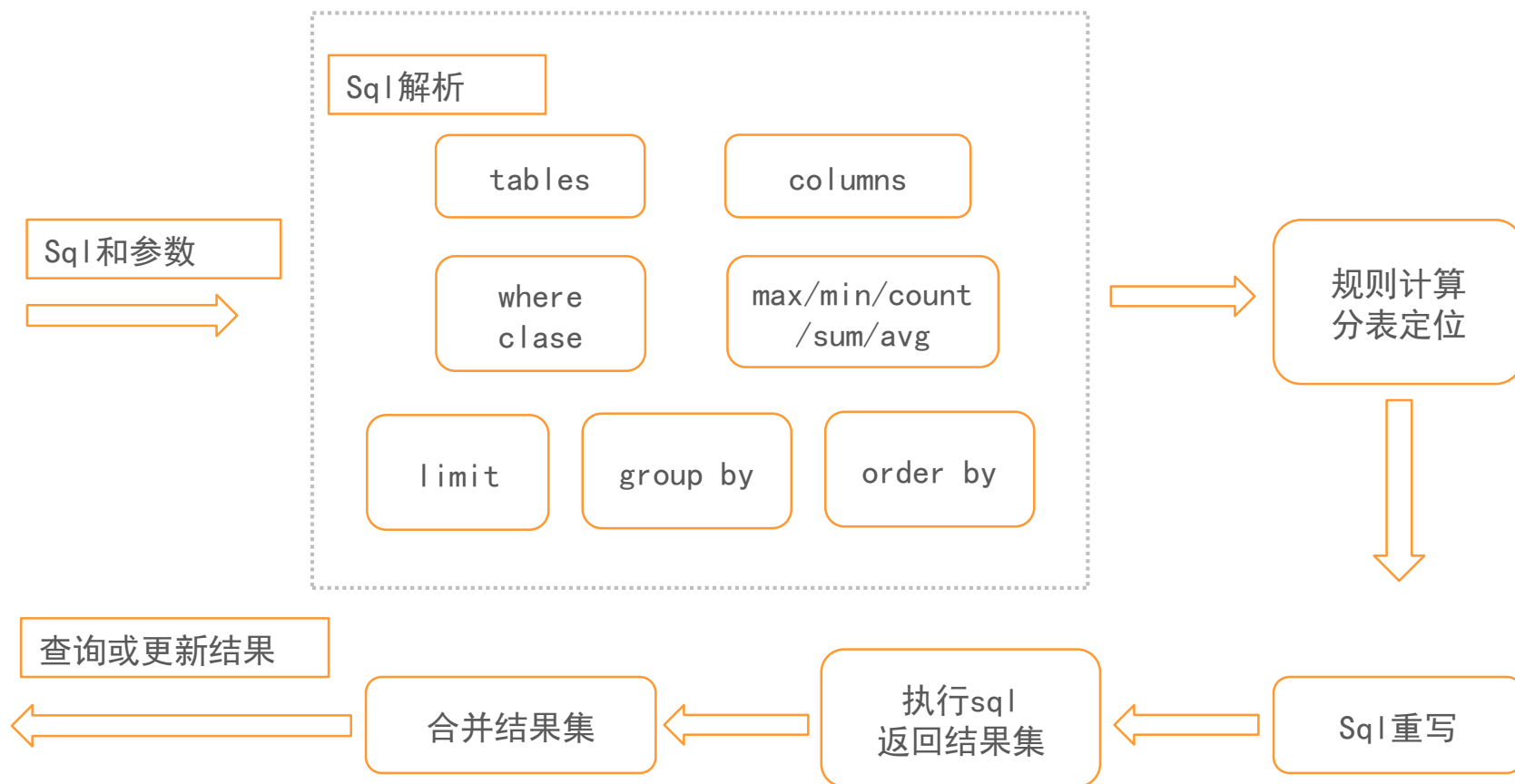




- ❶ 分库分表，读写分离，Failover。
- ❷ 运维监控，自定义路由，级联路由，重读双写，集群配置在线推送，弱XA事务支持，作业平台，离线数据集中平台。
- ❸ 优化sql解析模块，丰富sql语法，groovy路由算法，数据管道服务，分布式补偿事务服务。







- 群组类型
  - ① 全局组
  - ② 工作组
  - ③ 空备组
  - ④ 外挂组
  - ⑤ 重读组
- 表类型
  - ① 切分表
  - ② 全局表
  - ③ 孤立表
- 切分键类型
  - ① 级联切分键
  - ② 复合切分键
  - ③ 查询切分键

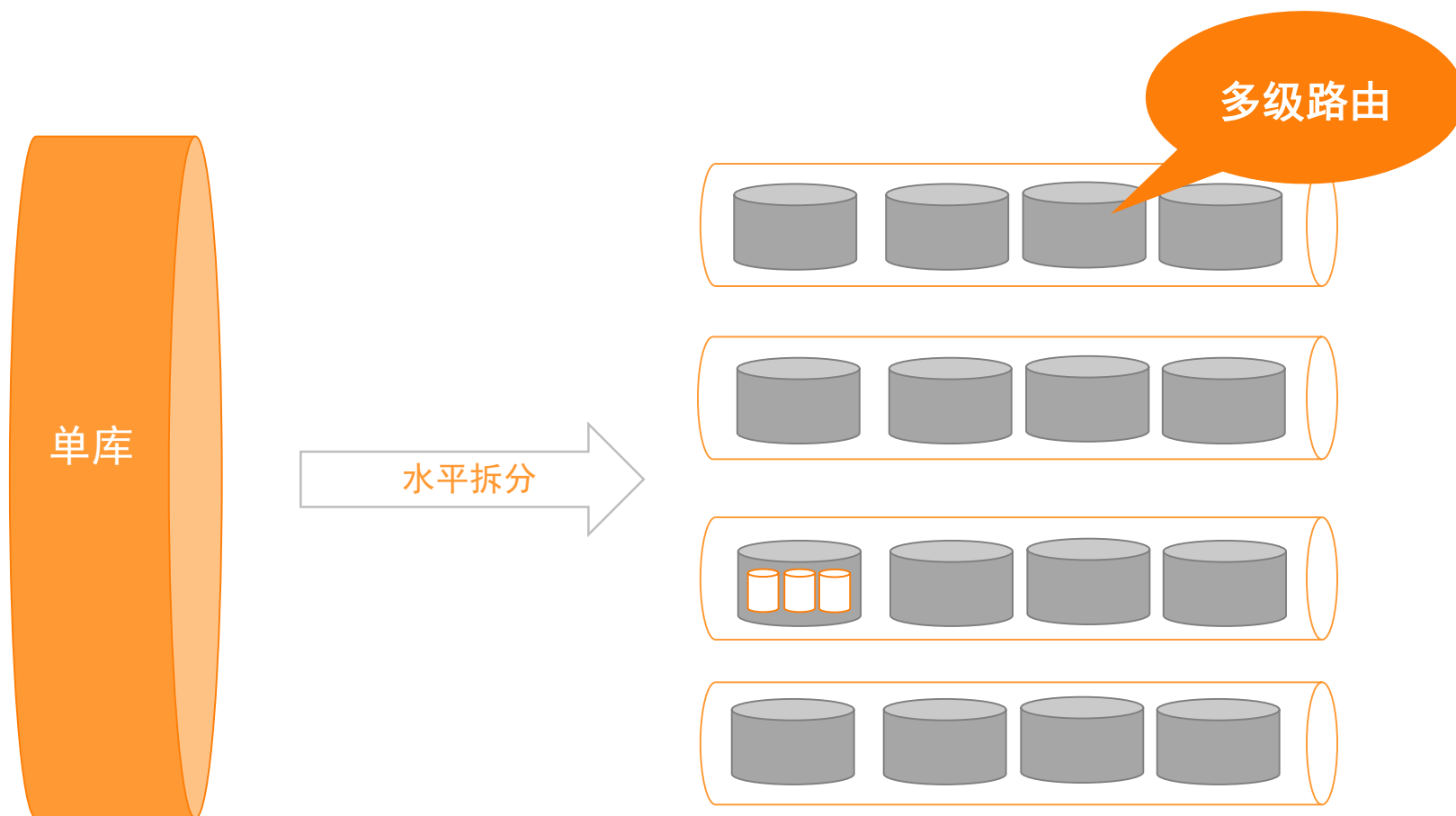


## CDS集群

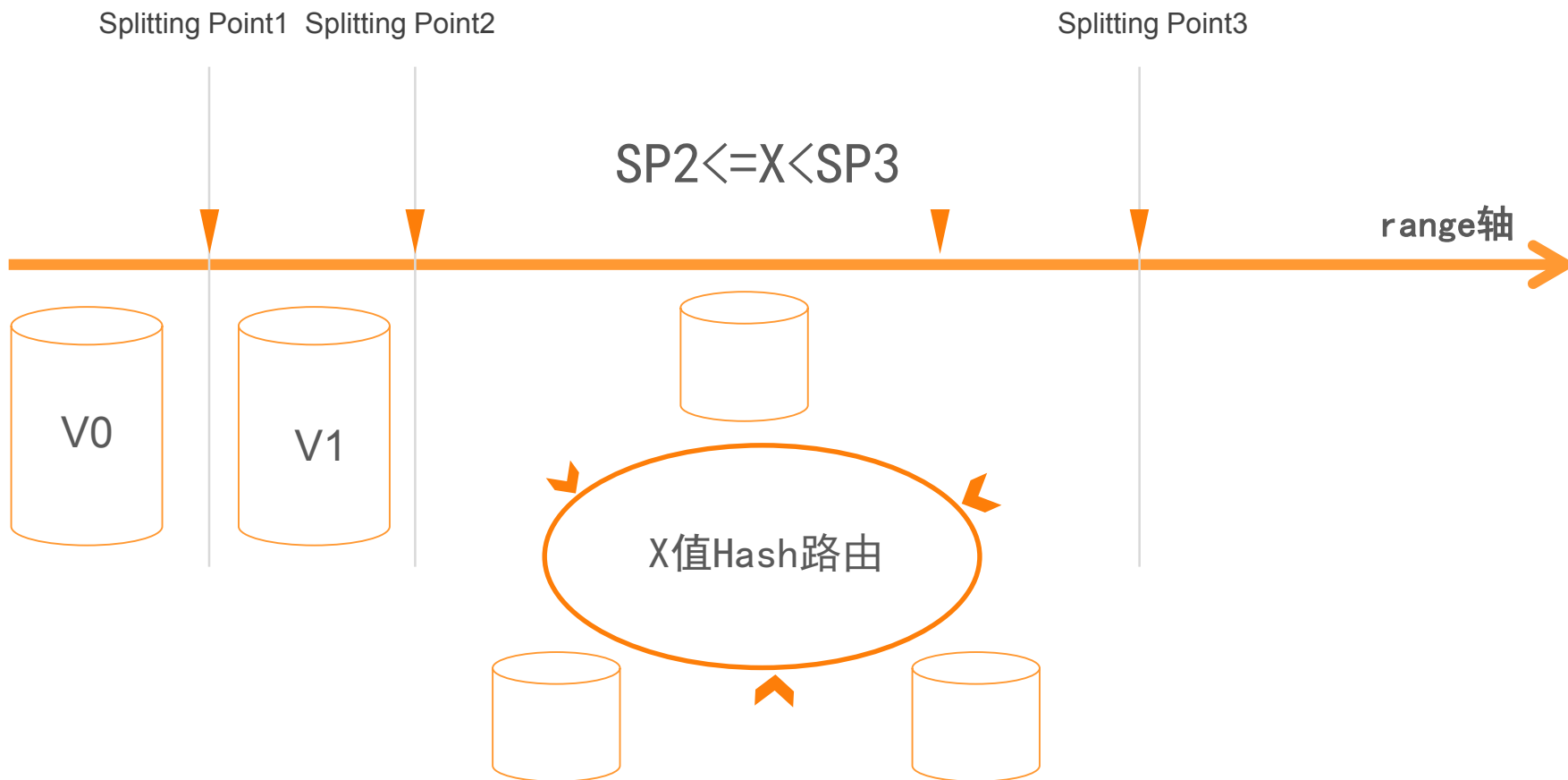
**全局表：**每个工作组都存在相同的数据备份（比如. 城市、类别）

**切分表：**插入数据时根据切分键插入到不同分表（比如. 业务流水表）

**孤立表：**可以存放在任意工作组包含数据表（比如. 配置表）



- Range路由策略 + Hash路由策略





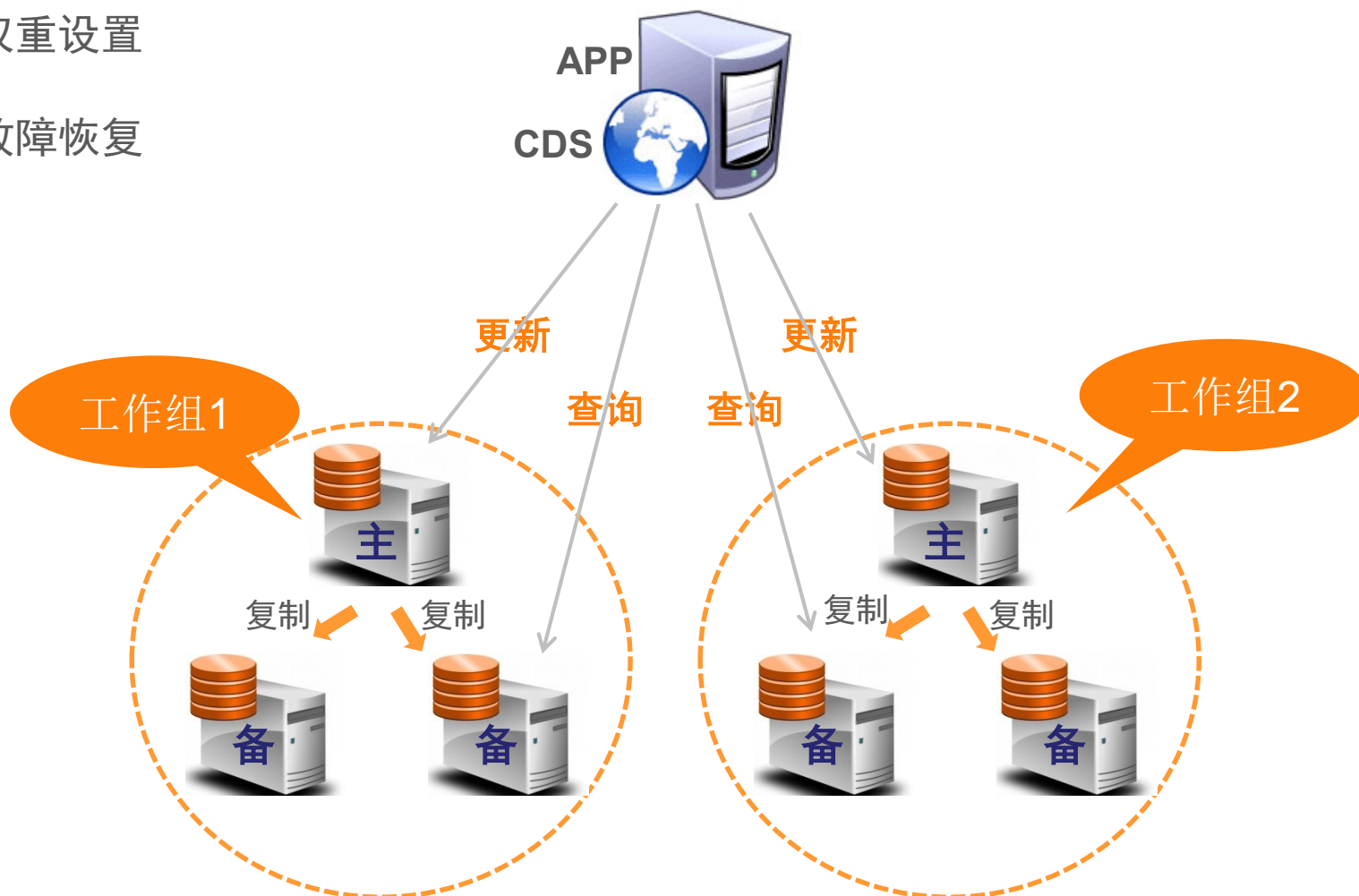
- 分库分表

- ① 分区规则支持：Hash, Range, List, 自定义（groovy）
- ② 支持oracle、mysql、sql server数据库
- ③ 支持复合切分键
- ④ 支持级联切分键
- ⑤ 支持DDL变更
- ⑥ 支持Sum/Max/Min/Count/Avg聚合函数
- ⑦ 不支持：复杂子查询、更新切分键字段

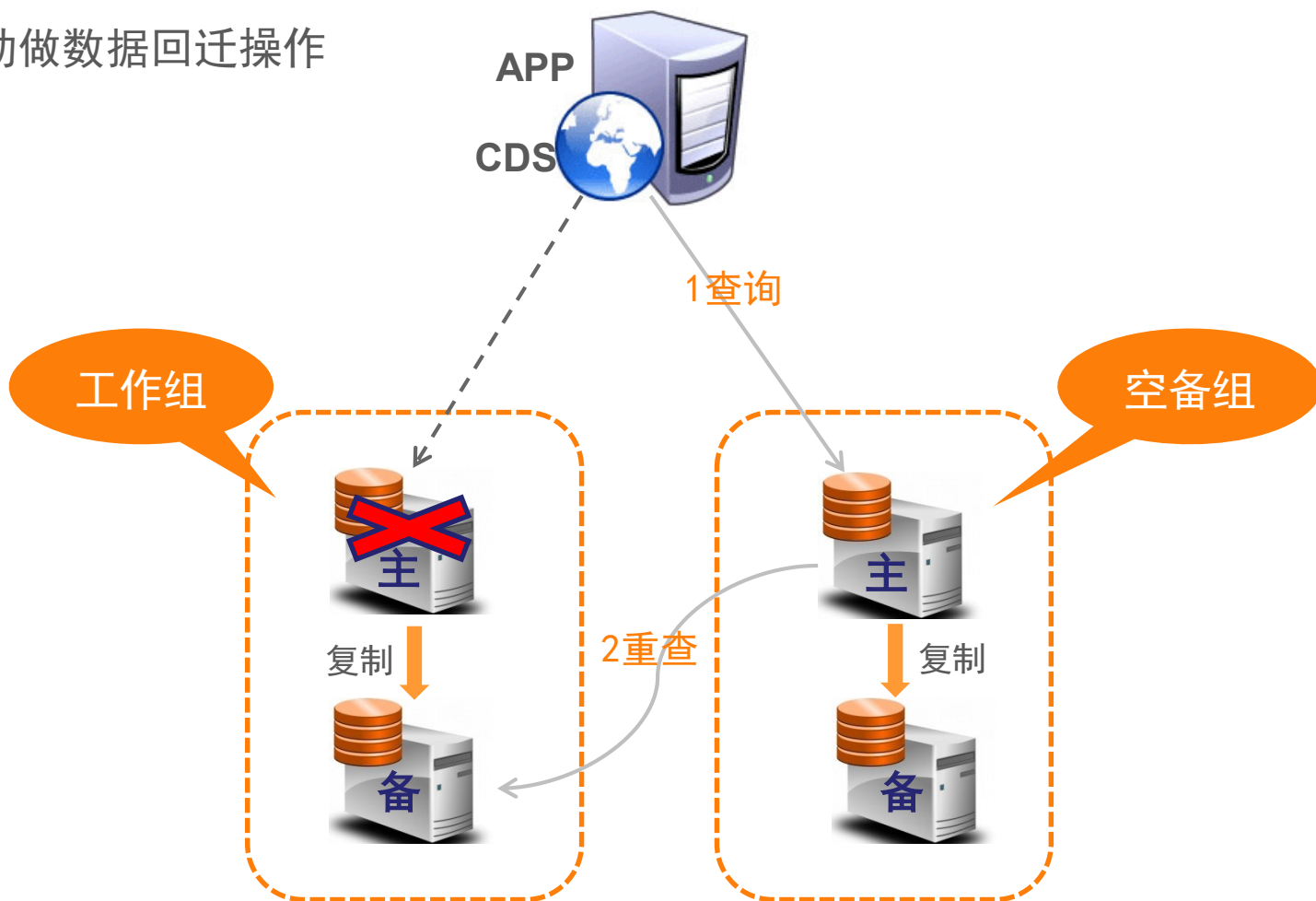
- 非分库分表

- ① 支持存储过程
- ② 支持大部分DML、DDL

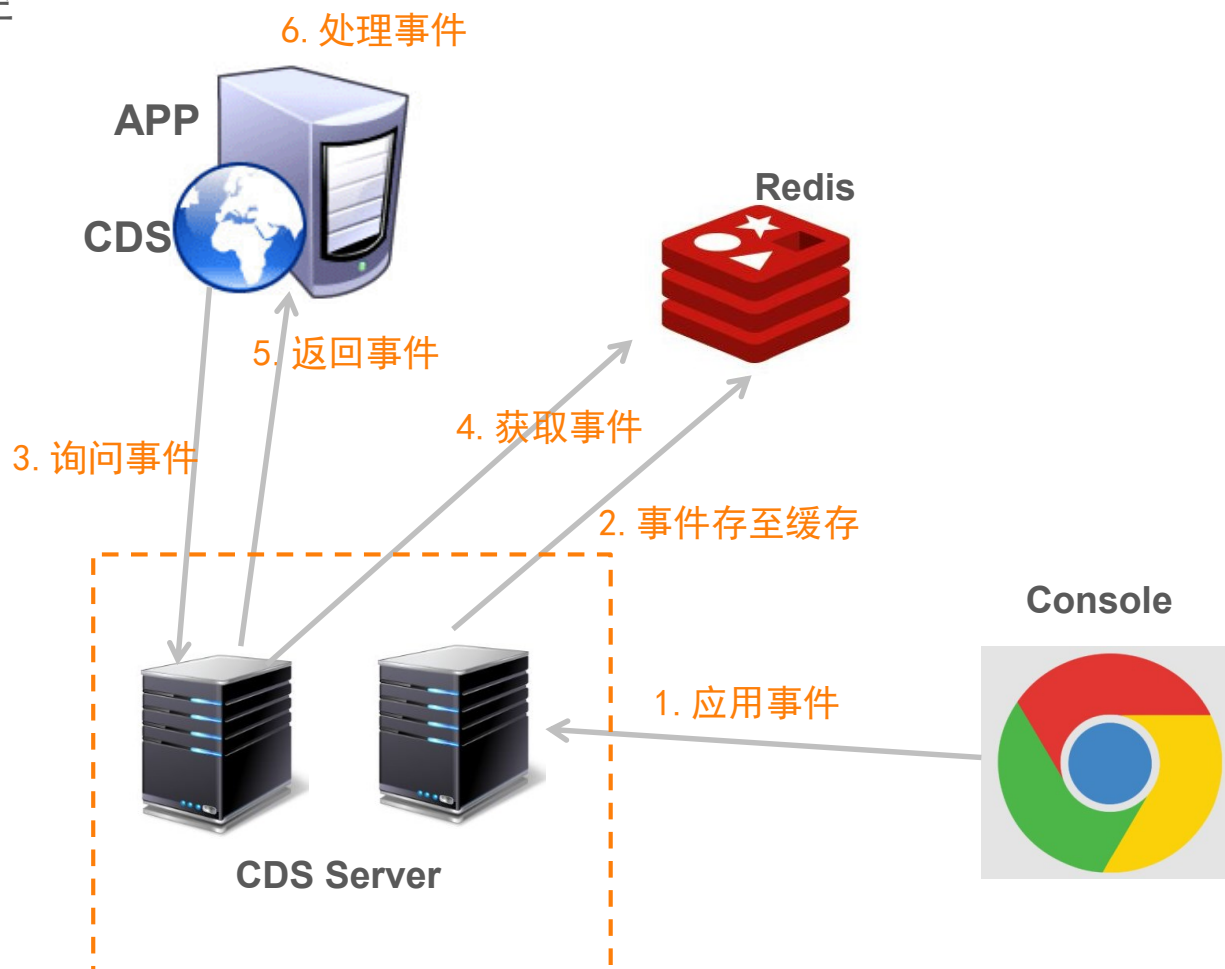
- ① 负载均衡
- ② 权重设置
- ③ 故障恢复



- ① 适用于流水型业务表，记录变更与查询不依赖之前的信息
- ② 主库故障由生产DBA手动触发推送事件，将新请求指向failover库，主库恢复后需要手动做数据回迁操作



- ① 连接池监控
- ② 应用进程信息监控
- ③ 集群配置变更



- 全局Sequence
  - ① 全局唯一
  - ② 不保证全局有序
  - ③ 单点问题
- 分表Sequence
  - ① 解决单点问题
  - ② 与分表一一对应
  - ③ 唯一性由业务保证

- MyDB
- 环境检查工具
- 作业平台
- 数据管道服务
- 分布式补偿事务
- 数据集中平台
- ... ..

## CDS数据库源在线查询，订正，统计，导出，支持工单和复杂用户权限控制

MyDB (1.4.0) 首页 数据字典 统计 关于 欢迎,钱青松! 退出

数据源 trade\_sharding(CDS:oracle)

trade\_sharding(CDS:oracle)

表信息

- CLUSTER\_AUTOINC\_TABLE
- CLUSTER\_TRANSACTION\_REDO
- CLUSTER\_TRANSACTION\_REDO
- TRADE\_AGREEMENT\_CONFIG
- TRADE\_ATTACHED
- TRADE\_BASE
- TRADE\_FUND
- TRADE\_GOODS
- TRADE\_LOGISTICS
- TRADE\_ORDER
- TRADE\_ORDER\_MAPPING
- TRADE\_PAY\_PI\_DETAIL
- TRADE\_PAY\_RELATION

执行(F8) 收藏 格式化 查看执行计划

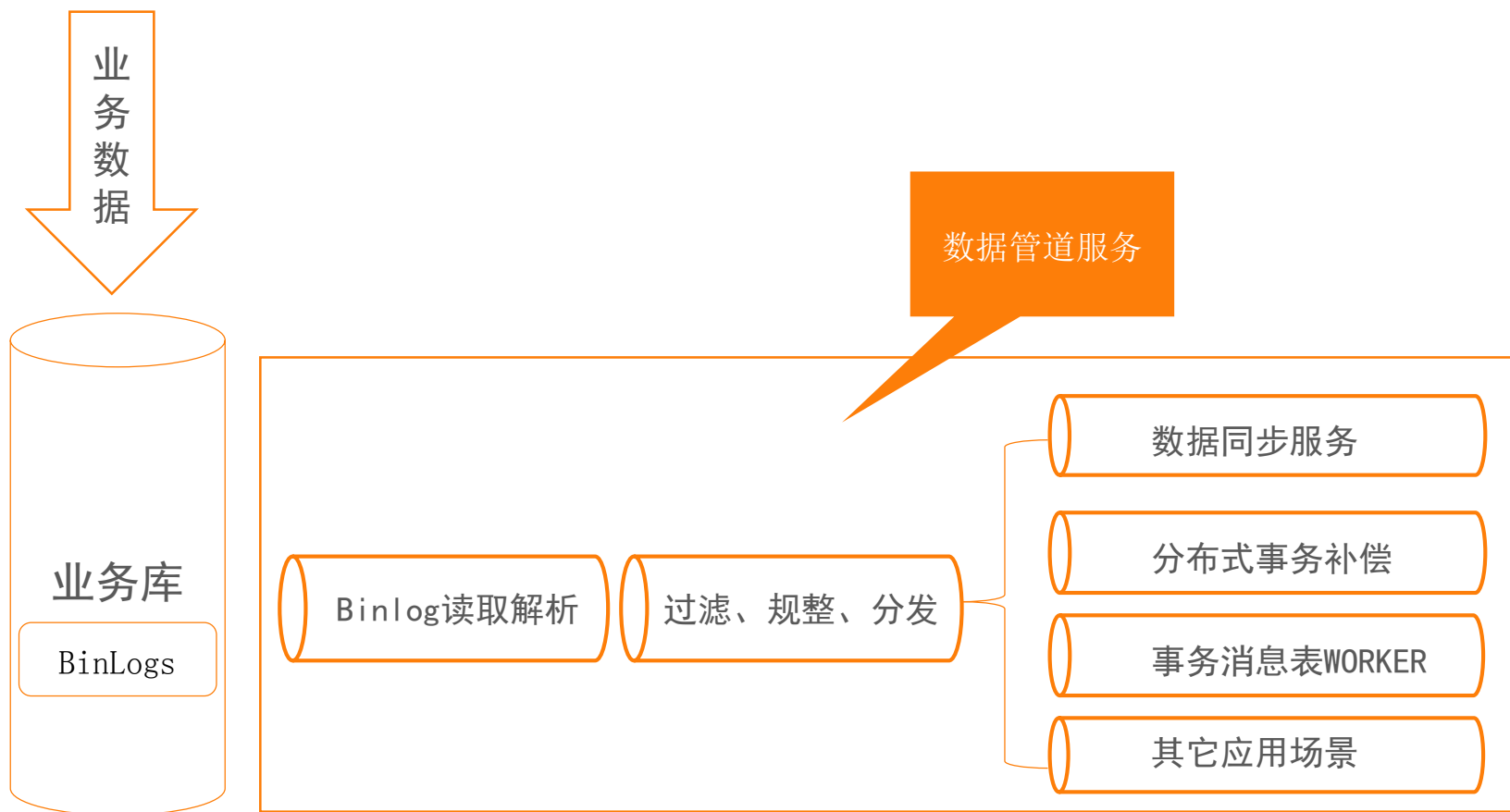
1

Ready

查询 历史 收藏 快捷

- 海量数据迁移
- 增量数据同步
- 自动建表
- Failover数据回迁
- 批量数据导出
- 集群数据库对象比对
- 索引表初始化
- .....





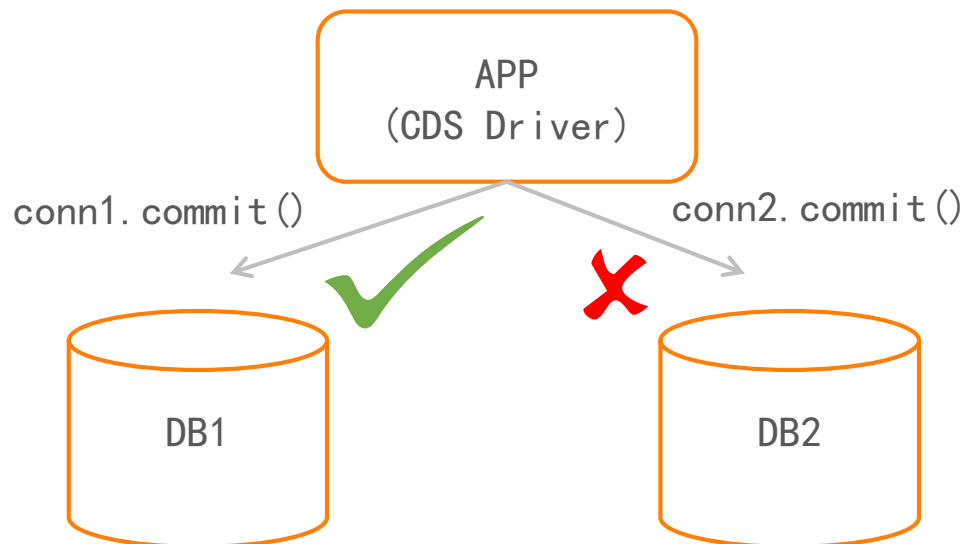
- 分布式系统事务处理的三种方式：

- ① 分布式事务
- ② 基于Best Efforts 1PC模式的事务
- ③ 事务补偿机制

- 小结：

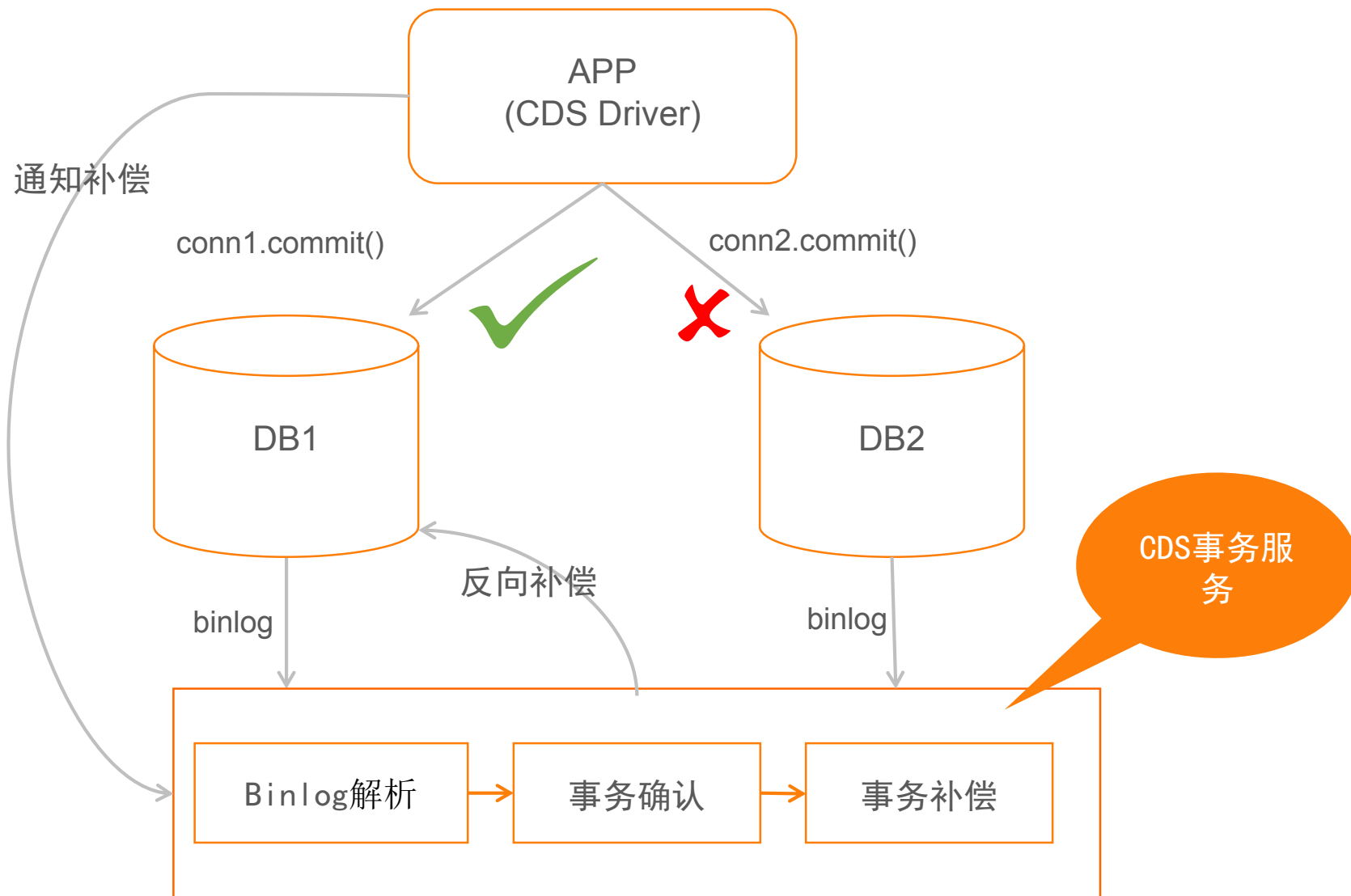
- ① 分布式事务，最严格的事务实现，性能是个大问题；
- ② Best Efforts 1PC模式，性能与事务可靠性的平衡，支持系统水平伸缩，  
大多数情况下是最合适的选择；
- ③ 事务补偿机制，最终一致，牺牲实时一致性，获得最大的性能回报。

- 基于BASE模型，通过解析Mysql的binlog来实现事务补偿，对用户完全透明，不一致时间窗口控制在秒级；
- 最终状态为事务失败状态，业务能感知事务失败异常；



- 问题：

- ① 补偿哪个binlog事务？
- ② 怎么建立conn1.commit()和conn2.commit()的联系



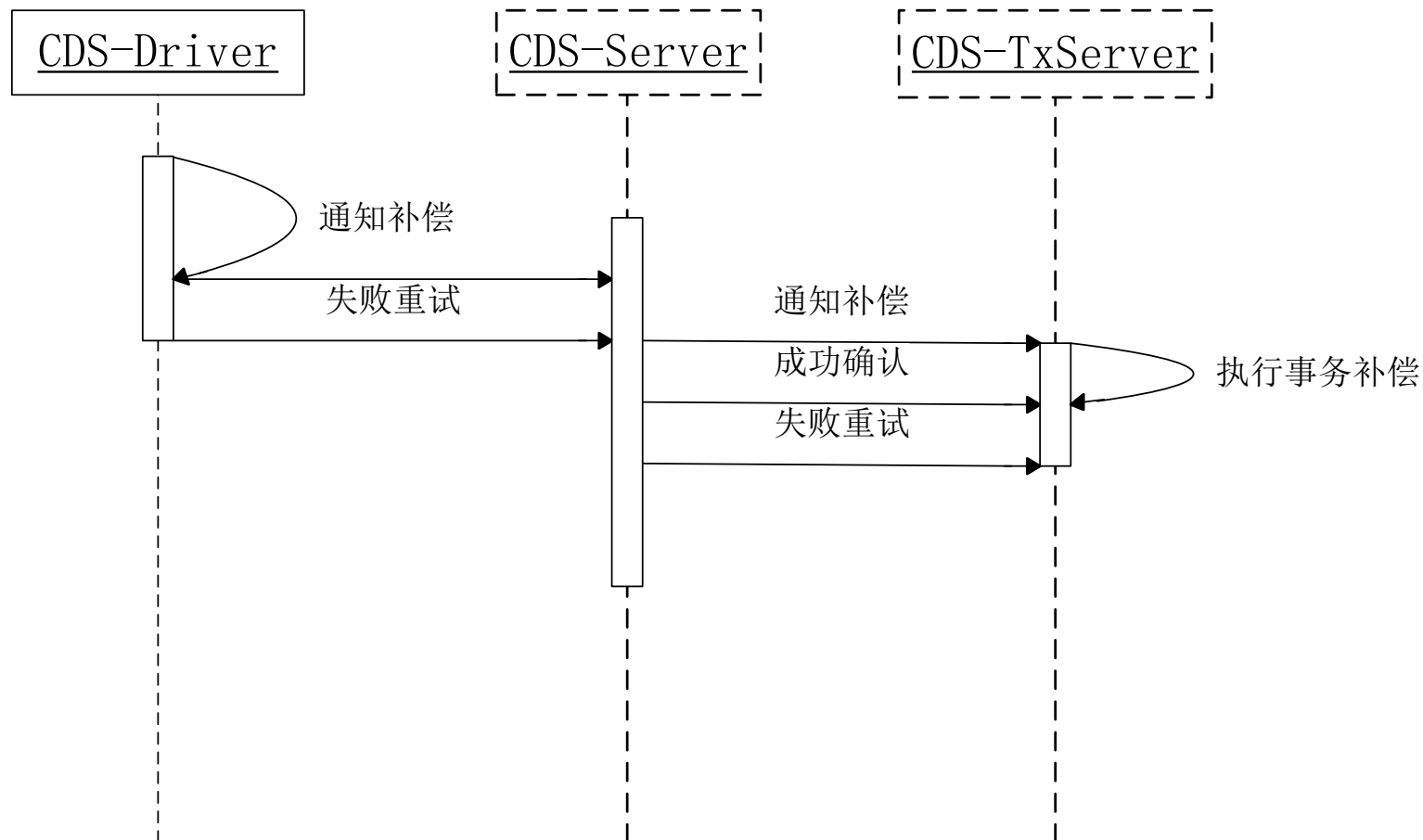
```
try{  
    connection1.insert、 update、 delete...  
    connection2.insert、 update、 delete...  
    ... ..  
    try{  
        insert cluster_xa_1 ✓  
        insert cluster_xa_2 ✗  
        connection1.commit();  
        connection2.commit();  
    }  
    catch(Exception e){  
        notifyTxServer(); //通知事务服务  
        .....  
    }  
}catch(Exception e){  
    connection1.rollback();  
    connection2.rollback();  
}
```

CDS Driver向每个分库  
cluster\_xa表插入一条相同  
记录，业务服务完全透明

- Cluster\_XA表

- ① xa\_uuid: 每个外围事务分配一个uuid, 唯一标识分布式事务
- ② xa\_dbunits: 外围事务涉及到的数据库单元
- ③ cur\_dbunit: 当前分库标识;
- ④ timestamp: 时间戳, 用于计算超时时间

| xa_uuid | xa_dbunits | cur_dbunit | timestamp |
|---------|------------|------------|-----------|
| xxx     | db1, db2   | db1        | xxx       |
| ... ..  | ... ..     | ... ..     | ... ..    |

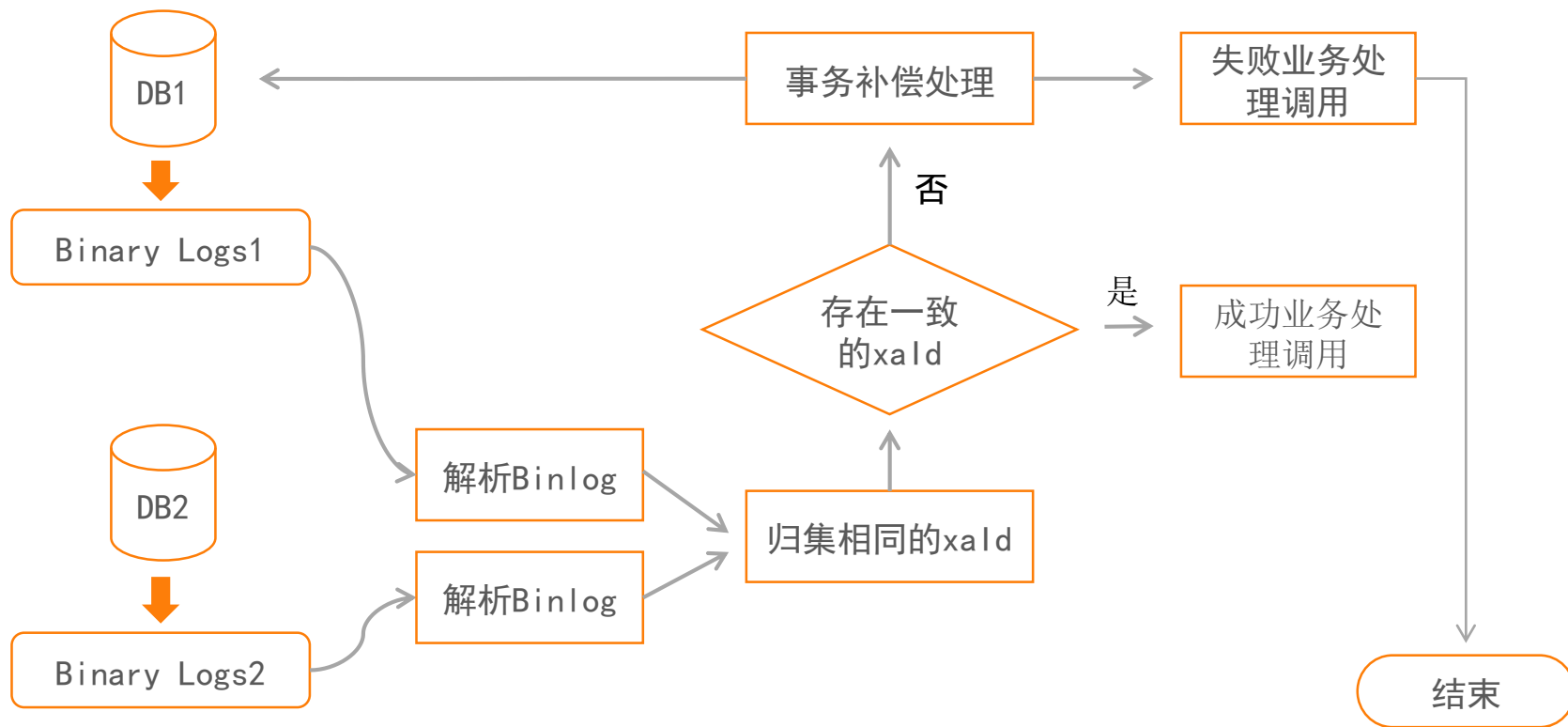




- 解决了DB2宕机事务补偿；
- 解决了到DB2网络中断，无法提交的事务补偿；

?

进程自己宕机了，  
而无法向TxServer发出事务补偿通知...



说明:

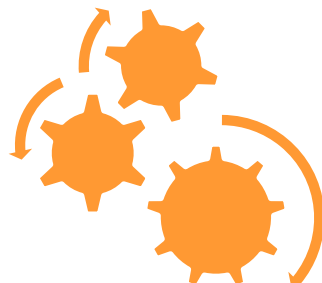
- DB1提交成功了事务, 会存在一条事务id, 并通过cluster\_xa表获悉同事务还包含DB2;
- DB2不存在相同的事务id, Binlog解析的当前时间大于DB1的xaId + 一次IO最大等待时间时, 则判断DB2没有提交成功!

URL/JDBC

OLTP



- a) 海量数据存取，提高TPS
- b) 读写分离
- c) Failover高可用

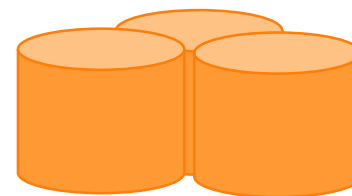


作业平台

- a) 实时的流式数据处理
- b) 作业调度
- c) 执行监控

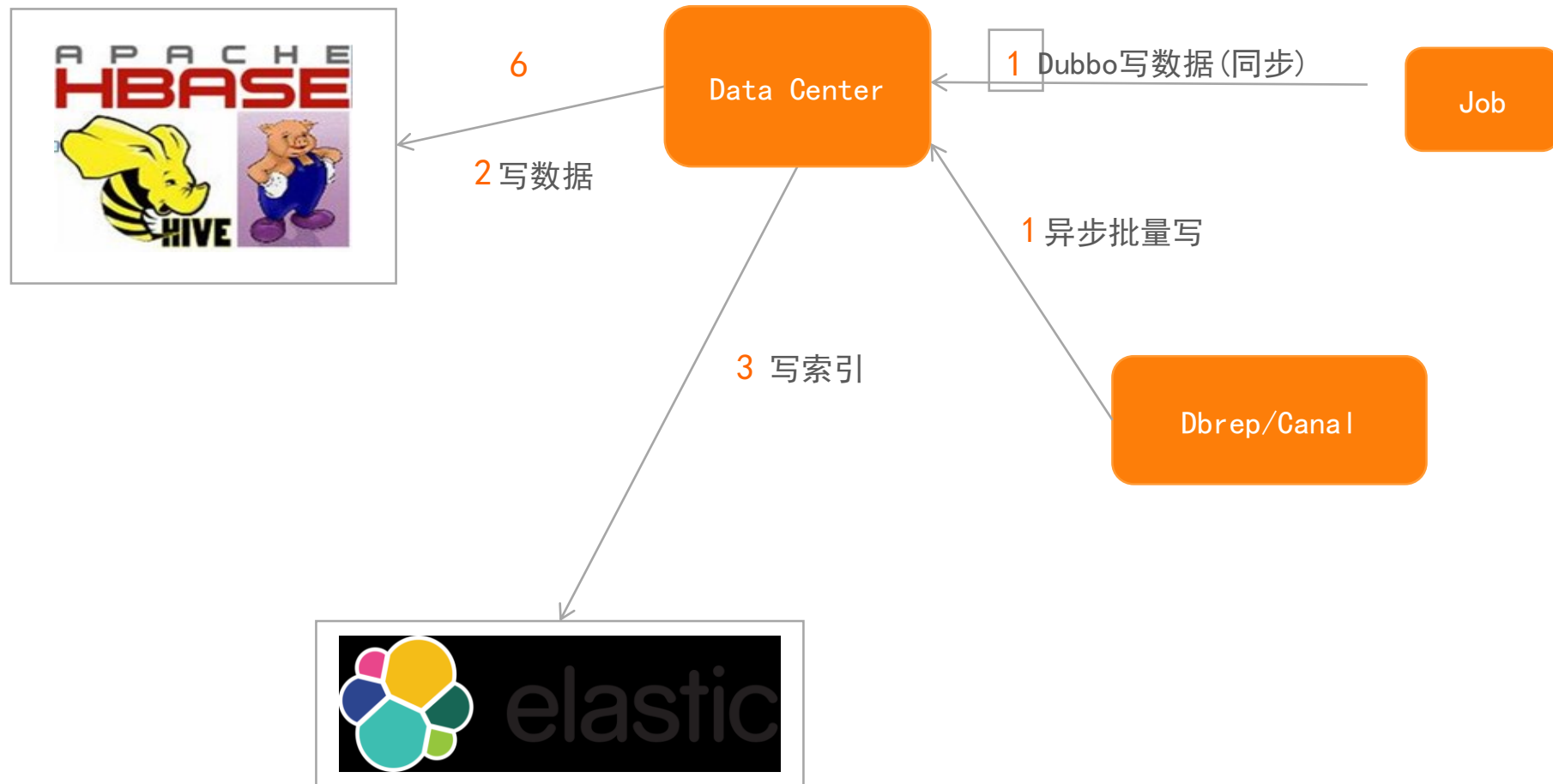
SOA

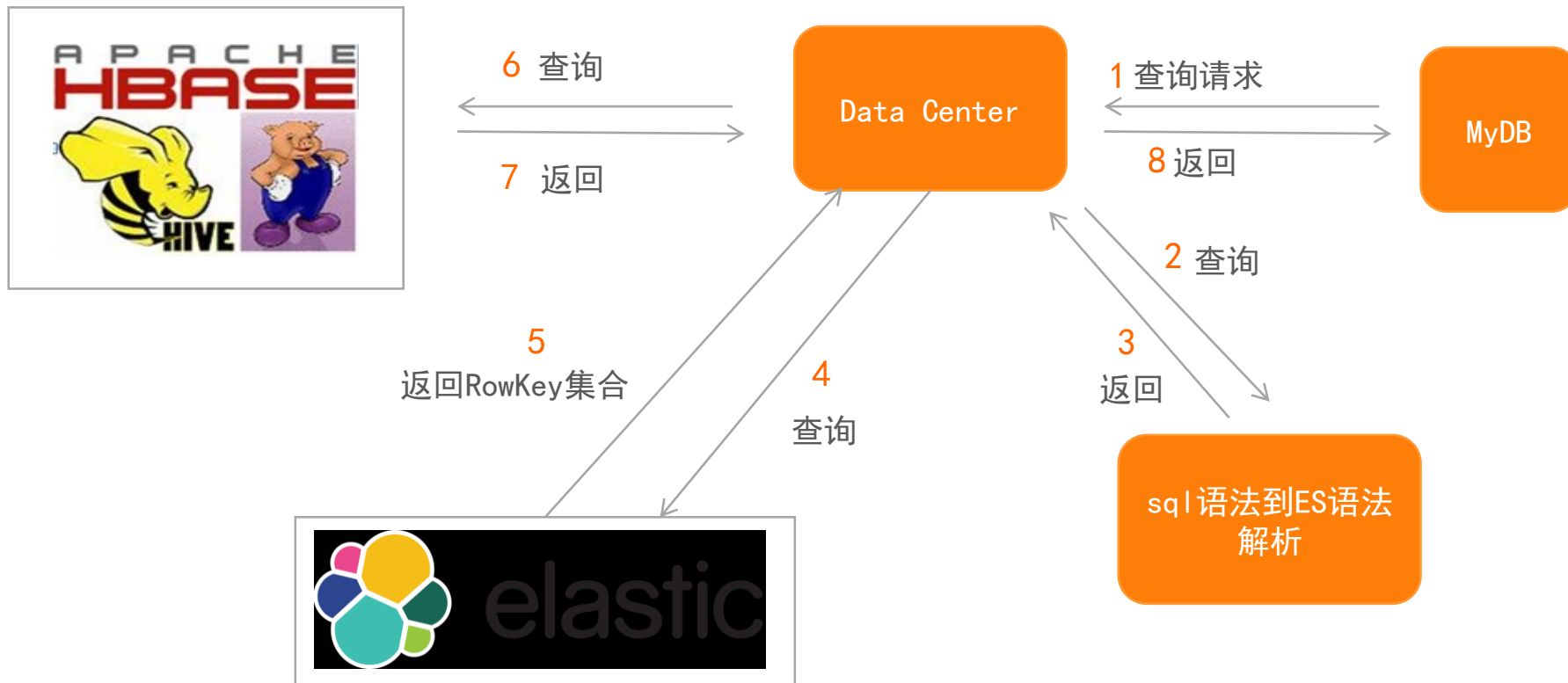
OFF LINE

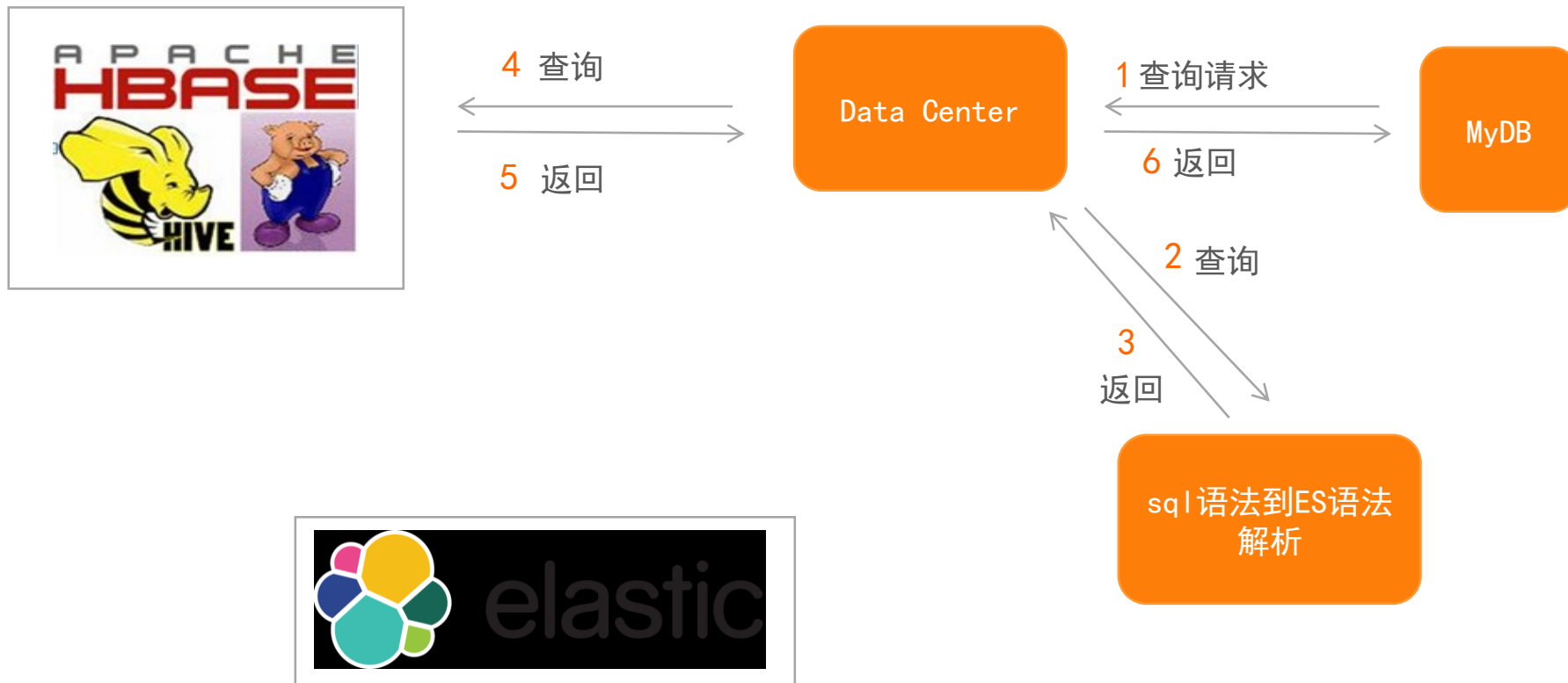


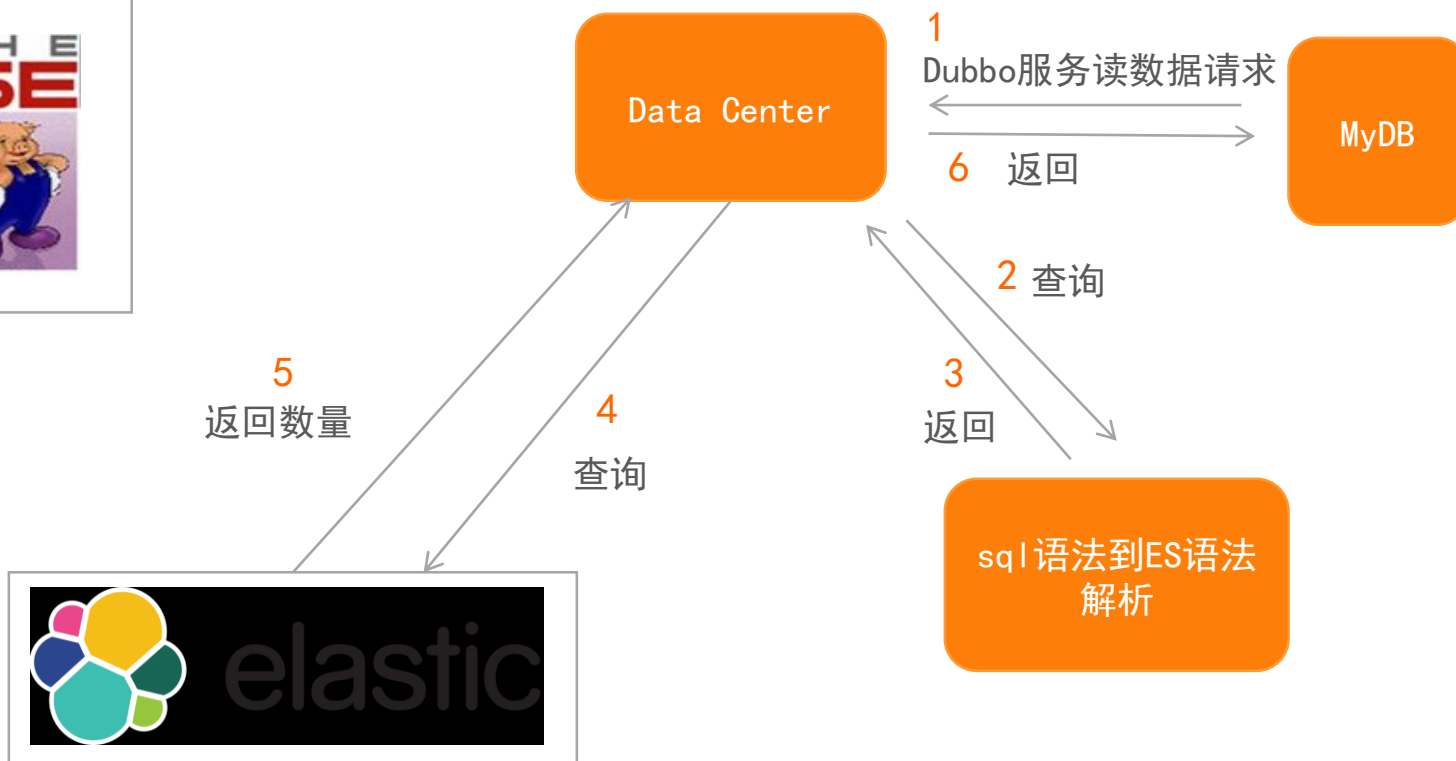
- a) 离线、非实时
- b) 多维度查询
- c) 简单的统计

- 离线数据查询、统计，服务于运维；
- 只关注有限时间内的数据，比如1年内；
- 索引创建，服务多维度查询
- 简单数据统计











## 欢迎使用：**CDS 申请接入流程**

联系方式：

邮箱：wy-cds@jd.com

咚咚：wyxuli quan@jd.com

wywuyue@jd.com

wangzhengzheng@jd.com

**THANK YOU !**