

# 海量数据下携程数据库架构的探索和实践

携程旅行网 - 孙天浩

▶ 数据库技术专家

▶ SQLServer / MySQL / Python / Golang

个人简介

数据库混布，无任何Sharding方案

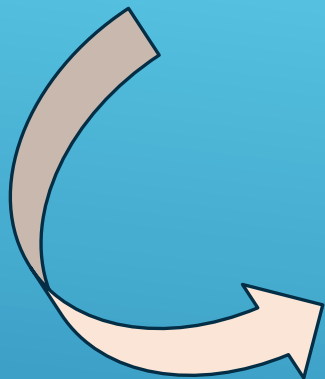
程序日志存放在数据库

用复制分发做读写分离架构，跨BU，网状

BI ETL取数，全量，5点前可以完成

早期现状

2010年，业务快速发展.....



单机存储容量限制

数据库本身处理能力有限

QPS过高数据库性能下降

DB和表过大影响正常运维，如备份，DDL

要求架构支撑10X业务发展！

业务快速增长遇到的问题

## 数据库以及应用的解耦

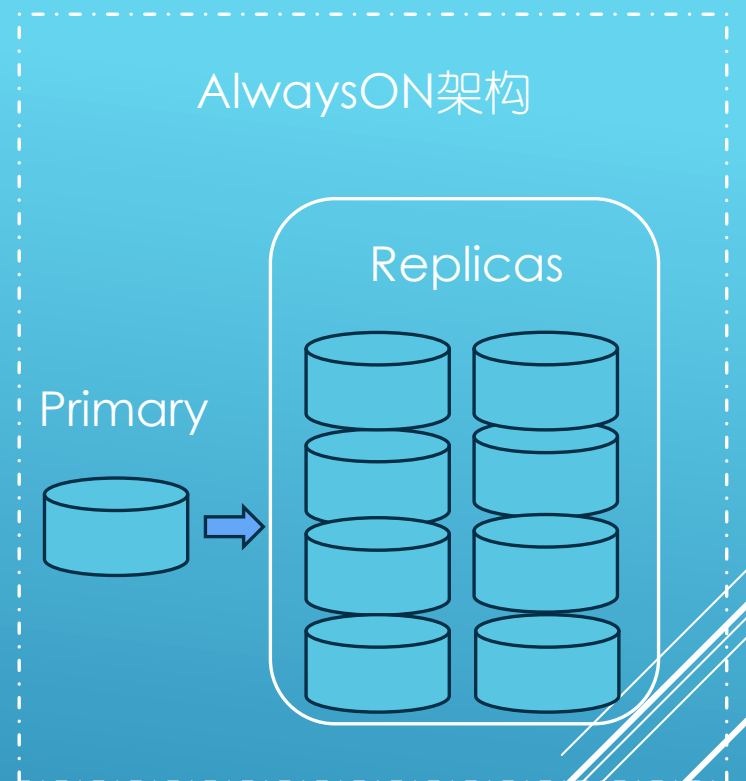
- 复杂业务拆分，降低系统耦合度
- 复制分发消除跨BU链路，改为调用接口
- 立项处理慢SQL，消除存储过程以及视图
- 数据库消除跨库Join，为拆分和数据库sharding做准备

## 读写分离由复制分发改造为AlwaysON

- SQL Server AlwaysON，同步快延迟小，维护简单
- 理想的读写分离架构方案，支持跨IDC

## BI ETL取数调整为增量

- 通过添加时间戳
- 取数源从快照调整为可读副本



虽然有些效果，但是随着数据量持续增长，系统瓶颈还是存在☹

## 短期方案

数据库垂直拆分

## 业务梳理

- 分析现有系统架构以及主要业务流程

## 划分类别

- 分析各数据表进行分类

## 容量预估

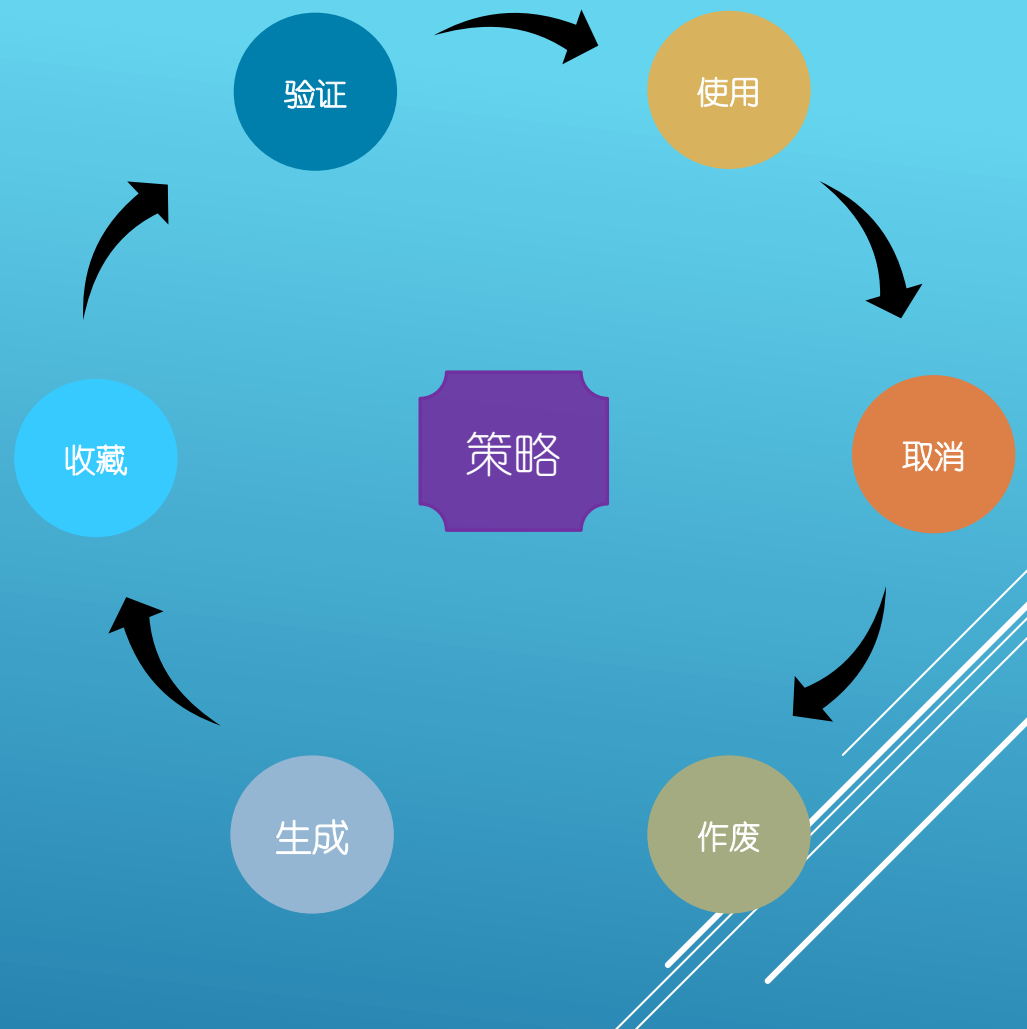
- 分析数据库表结构以及数据量
- 分析增长趋势

## 系统设计

- 数据库设计
- 数据异构设计
- 系统架构设计

垂直切分

水平切分



# 消费券架构演变案例

## Original Table

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
1	TAEKO	OHNUKI	BLUE
2	O.V.	WRIGHT	PURPLE
3	SELDA	JOHN	RED
4	JIM	JAMES	GREEN

## Vertical Partitions

### VP1

CUSTOMER ID	FIRST NAME	LAST NAME
1	TAEKO	OHNUKI
2	O.V.	WRIGHT
3	SELDA	JOHN
4	JIM	JAMES

### VP2

CUSTOMER ID	FAVORITE COLOR
1	BLUE
2	PURPLE
3	RED
4	GREEN

## Horizontal Partitions

### HP1

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
1	TAEKO	OHNUKI	BLUE
2	O.V.	WRIGHT	PURPLE

### HP2

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
3	SELDA	JOHN	RED
4	JIM	JAMES	GREEN

SHARDING 拆分方式

# Benefits



- 通过扩展解决容量瓶颈
- 分散负载，提高吞吐
- 加快查询响应时间
- 易于维护，备份恢复时间，表结构变更时间等

# Drawbacks



- 加大架构复杂度
- 分片的扩容/缩容困难
- 全局唯一分片key问题
- 数据记录的维度问题
- 跨库事务问题
- 同组数据跨库问题

SHARDING优势以及存在的问题



## 数据记录的维度问题

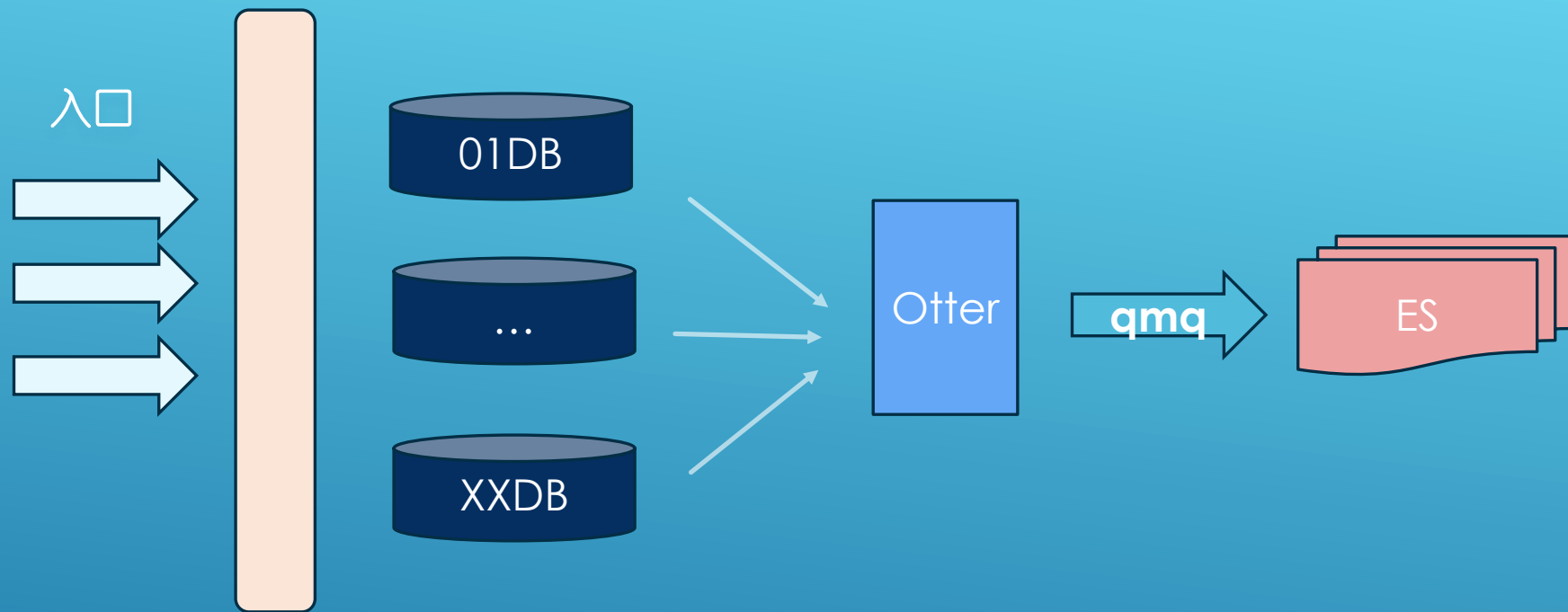
### ► 数据异构

例:客户从携程上订了机票, 酒店, 景点门票等订单后, 需要查询刚刚下的订单?

思路: 各BU产品订单, 会放在以订单号为分片键的shard库内, 直接查的话需要从各BU的订单库内以用户ID查找其下的订单, 这样性能肯定不符合要求。

做法: 携程通过OrderIndex项目实时将各BU的订单信息同步OI库内, 以用户ID作为分片键, 最终用户看到“我的携程”内的数据实际是基础部门以用户ID查询后的数据。

# SHARDING 数据异构



ES:

面向文档型数据库，可以作为全文搜索引擎  
数据可各维度查询

# SHARDING 数据异构

## 跨库事务问题

- ▶ 从架构设计上尽量避免
  - ▶ 同shard片内所有表需以相同维度的分片键；
  - ▶ 或表之间无耦合关系；
- ▶ 特殊情况下通过消息队列来做补偿机制，尽最大可能保证事务最终一致性；

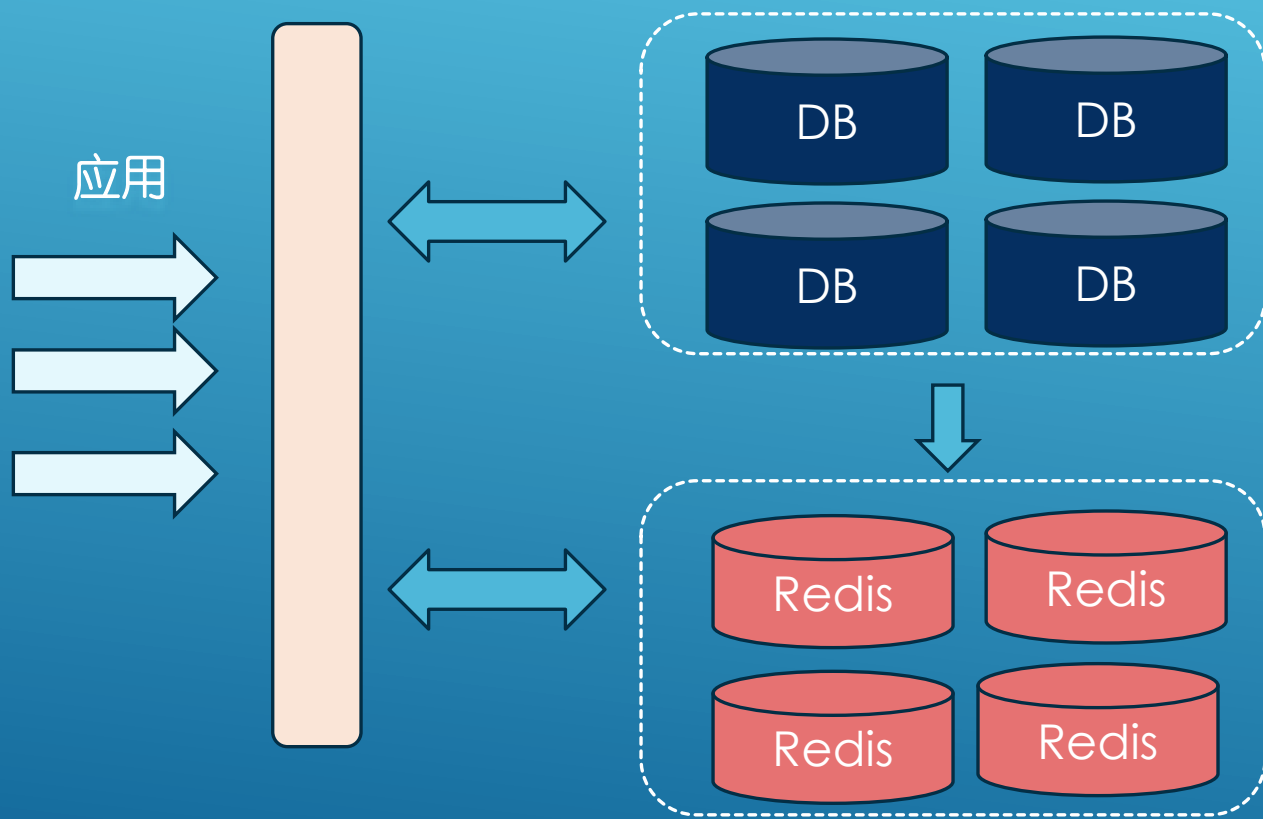
## 同组数据跨库问题

# SHARDING 数据跨库问题

► Memcache

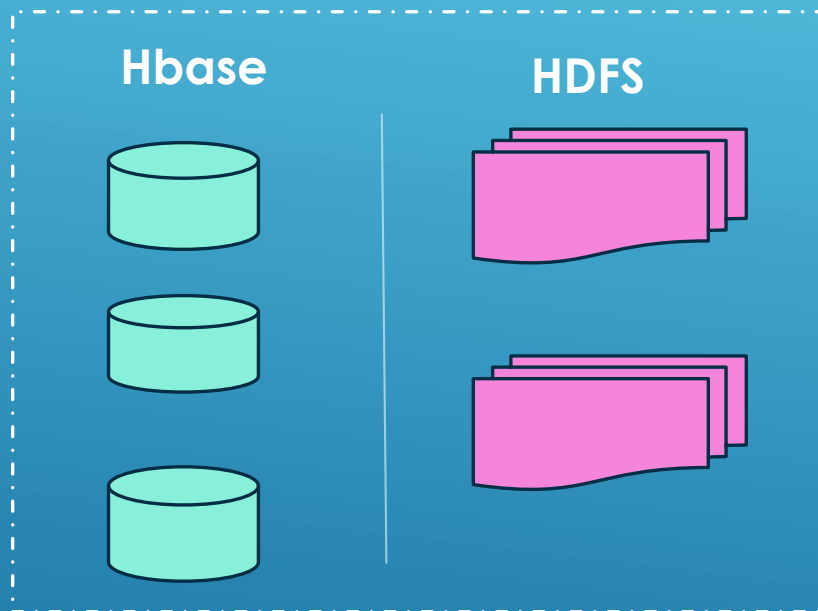
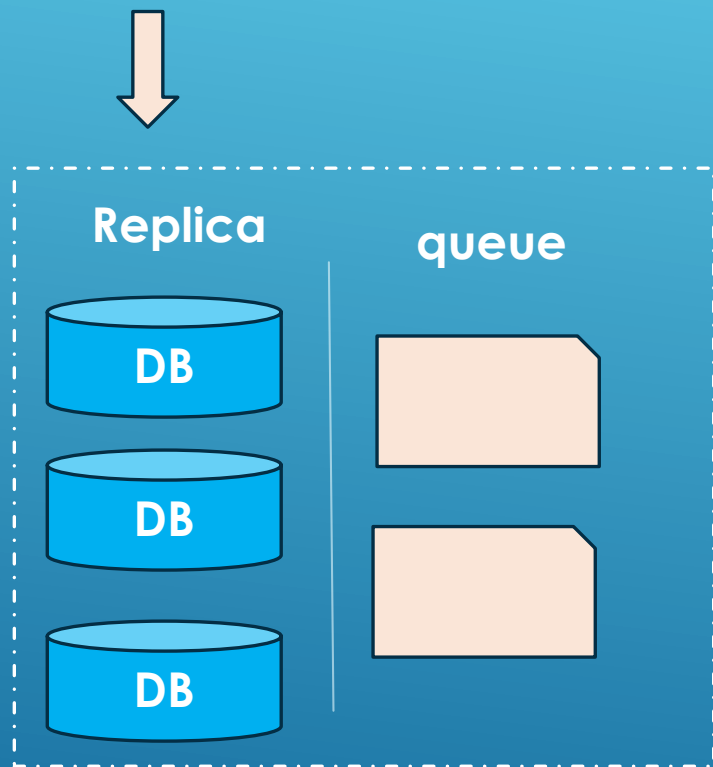
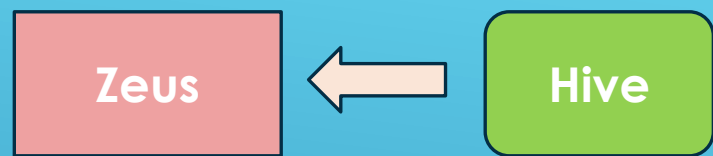
► Redis

► 公司开发了Credis/Xpipe, 支持分片扩展以及IDC间数据同步



数据缓存

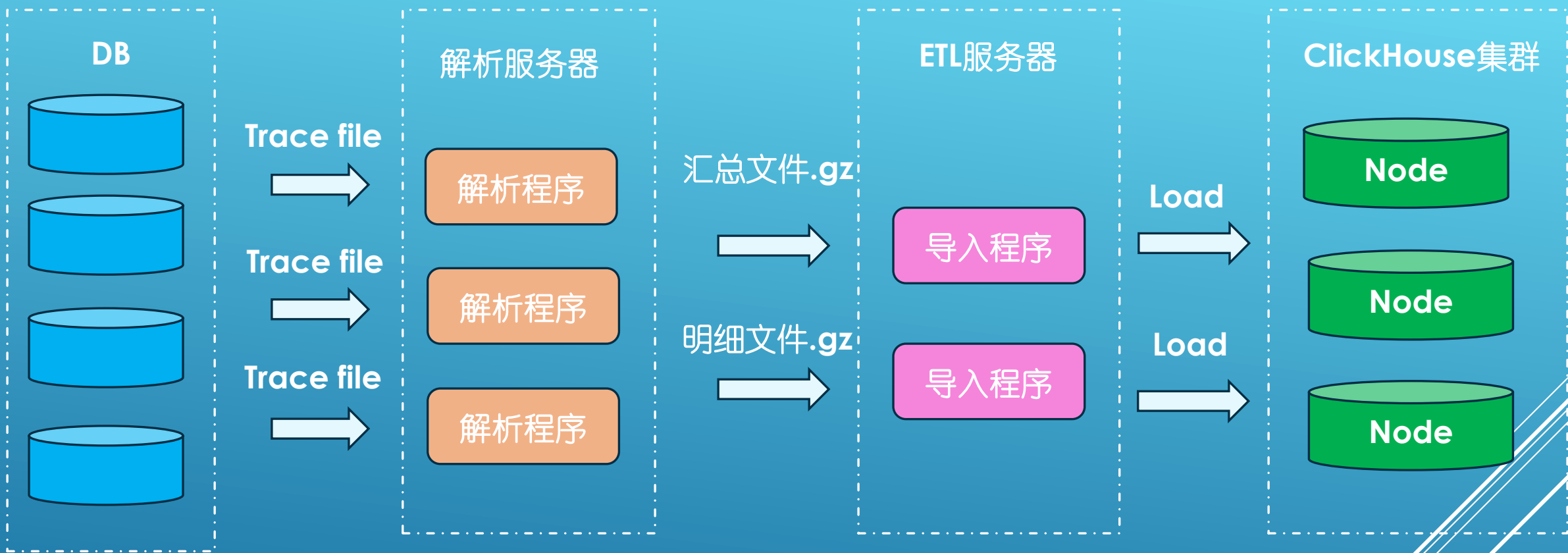
# 分布式数据库解决方案



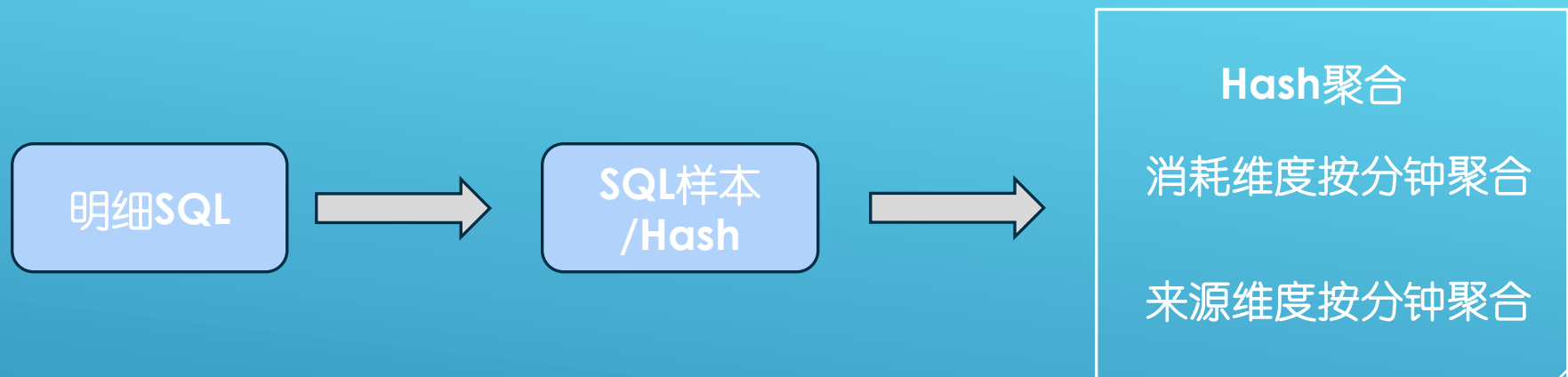
## NoSQL - HBase:

- ▶ 解决了扩展问题
- ▶ 缺少ACID, SQL 支持
- ▶ Region Compaction

数据离线计算



## TRACE数据流式处理



- ▶ 日均大小从**2T** 增长到**40T**
- ▶ 数据存放从**SQLSERVER** 到 **GreenPlum** 再到 **ClickHouse**的演进
- ▶ 整体延迟在**2分钟**以内

# TRACE数据汇总分析

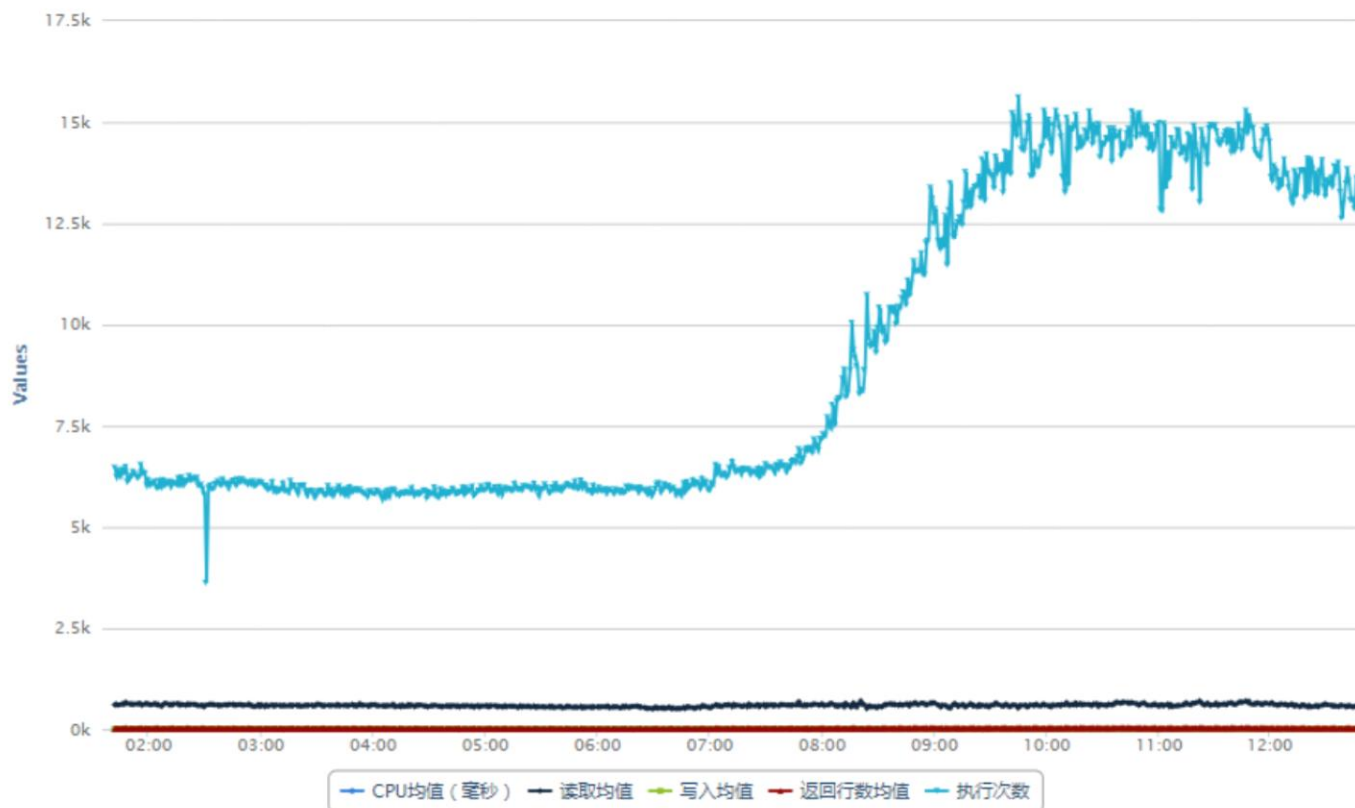


语句开始时间: 2019-03-28 12:41:05 ~ 2019-03-28 12:51:05
 服务器: xxx
 数据库: xxx
 [Q 查询](#)
[0 下载](#)

SQL(Sample)	beginTime	abortcount	execcounts	Duration(ms)	reads	writes	cpu	rowcounts	type	link
su	2019-03-28 12:40		26980	315733 (20.23 %)	30953972 (16.69 2 (0.00 %)		296789 (25.86 %)	1302502	null	
arc	2019-03-28 12:40		819	91241 (5.85 %)	50324122 (27.14 1 (0.00 %)		87389 (7.61 %)	4140	null	
2019-03-28 12:40			2038							
2019-03-28 12:40			9748							
2019-03-28 12:40			5662							
2019-03-28 12:40			7368							
2019-03-28 12:40			126							
2019-03-28 12:40			6678							
2019-03-28 12:40			40							
2019-03-28 12:40			4062							
2019-03-28 12:40			3076							
2019-03-28 12:40			1072							
2019-03-28 12:40			584							
2019-03-28 12:40			22							
2019-03-28 12:40			44							

语句开始时间: 2019-03-28 01:41:05 ~ 2019-03-28 12:51:05
 服务器: xxx
 数据库: xxx
 [Q 查询](#)

性能指标每分钟均值



TRACE聚合展示

## ▶ NewSQL

- ▶ 解决扩展问题
- ▶ 支持ACID
- ▶ 支持SQL
- ▶ Paxos / raft 数据强一致性保证
- ▶ OLTP / OLAT 均有适用场景

谢谢！