

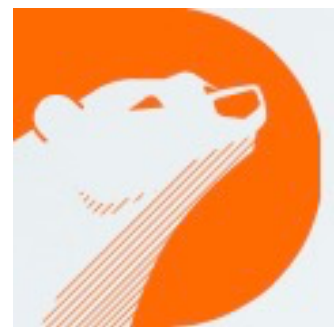
# PolarDB In-Memory Column Index HTAP解决方案

蔡畅  
阿里云数据库内核研发

# 自我介绍

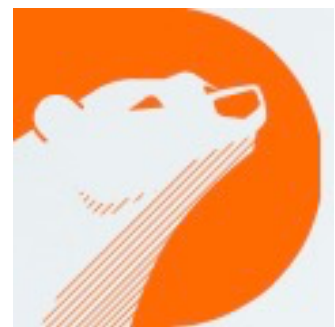


- 达梦数据库公司
- 淘宝数据库团队  
AliSQL分支, RocksDB引擎, X-Engine引擎
- 阿里云数据库团队  
PolarDB, X-Engine引擎, 列存引擎

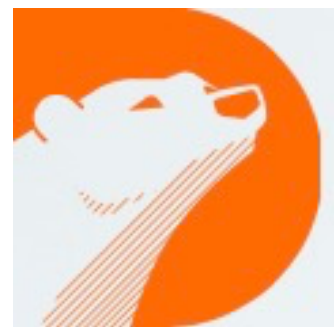


# 目录

- MySQL 在OLAP场景的挑战
- PolarDB In-Memory Column Index 整体架构
- PolarDB IMCI Query Engine
- PolarDB IMCI Hybrid Row/Column Store
- PolarDB IMCI 客户案例
- 总结与展望



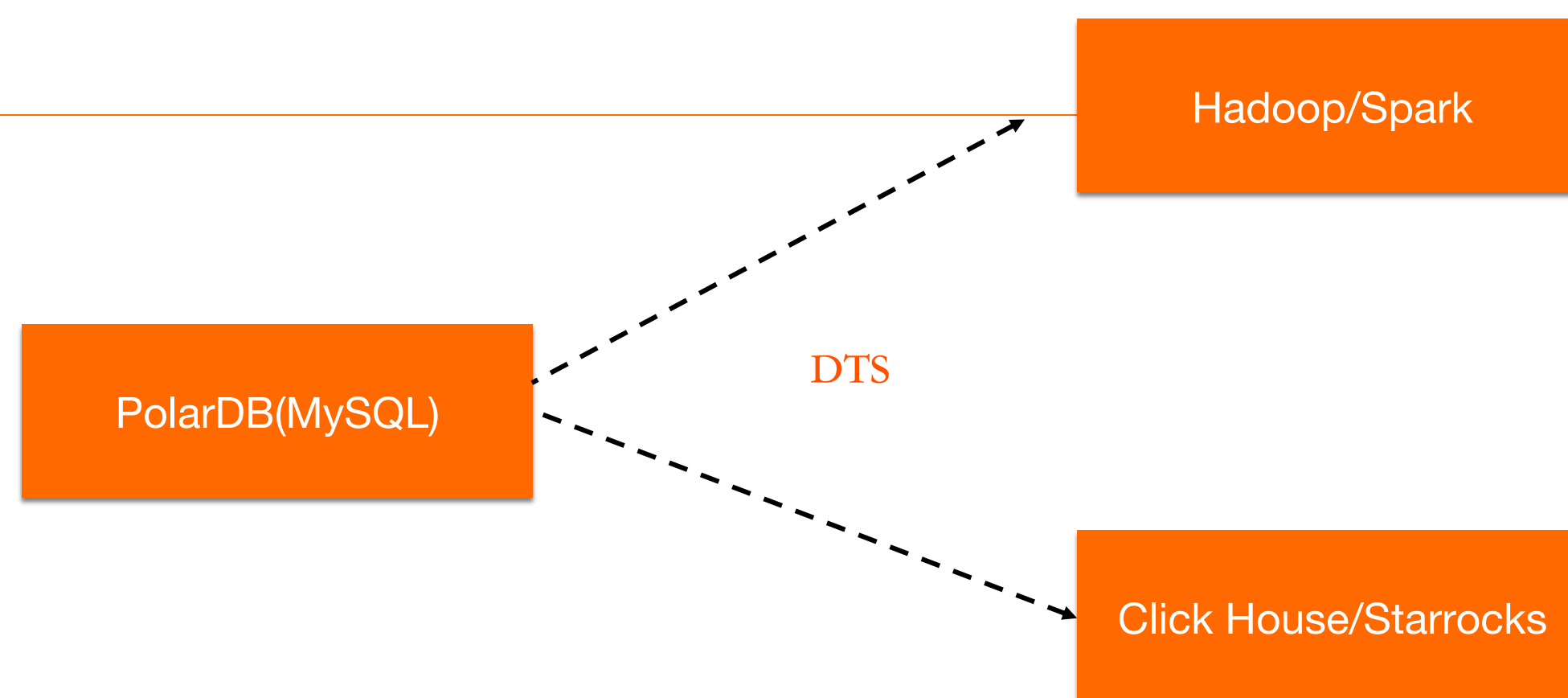
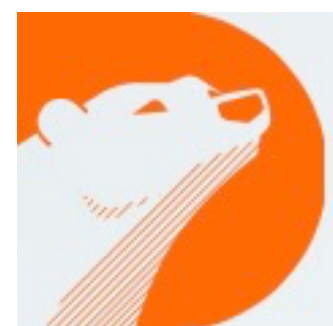
- **MySQL 在OLAP场景的挑战**
- PolarDB In-Memory Column Index 整体架构
- PolarDB IMCI Query Engine
- PolarDB IMCI Hybrid Row/Column Store
- PolarDB IMCI 客户案例
- 总结与展望





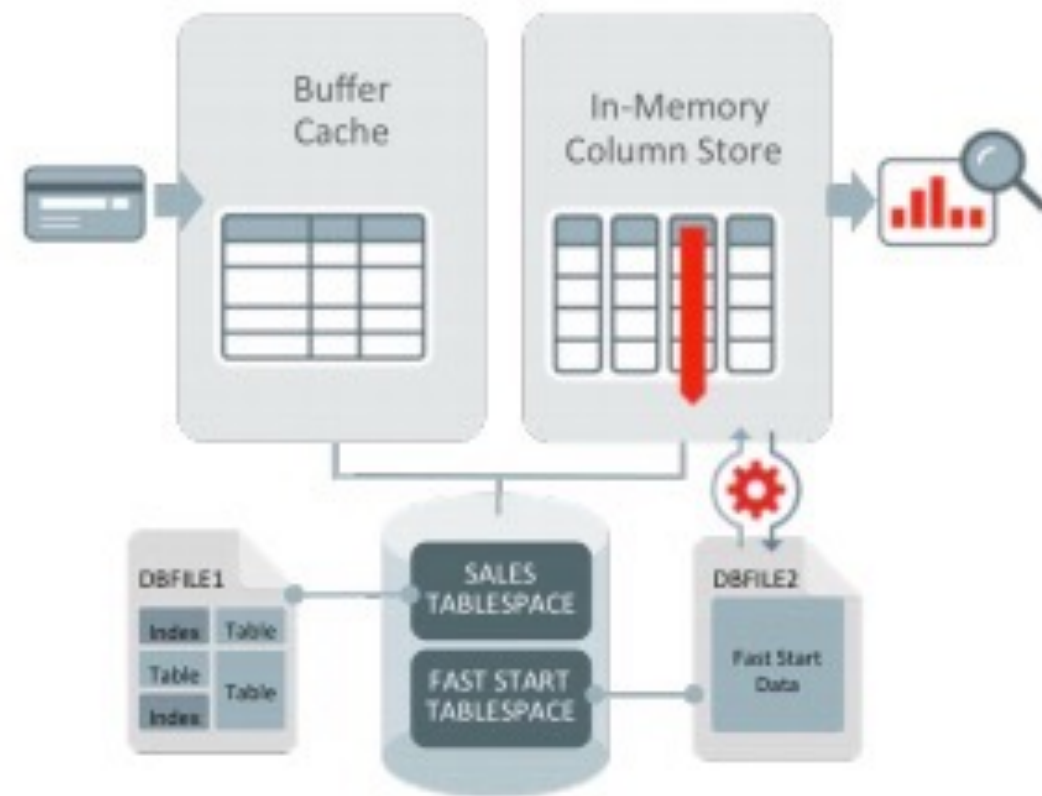
# PolarDB MySQL 的客户诉求: HTAP

- PolarDB MySQL 被大量在线业务所采用。
- PolarDB MySQL 主要面向 TP 场景优化，分析型查询由于执行器及存储引擎限制，性能不佳。
- PolarDB 数据导出至外部专用 OLAP 系统支持大数据的复杂查询，成本高，架构复杂，运行维护麻烦。
- 更易用的一体化数据库，一套 DB 同时满足 OLTP 及 OLAP 需求。

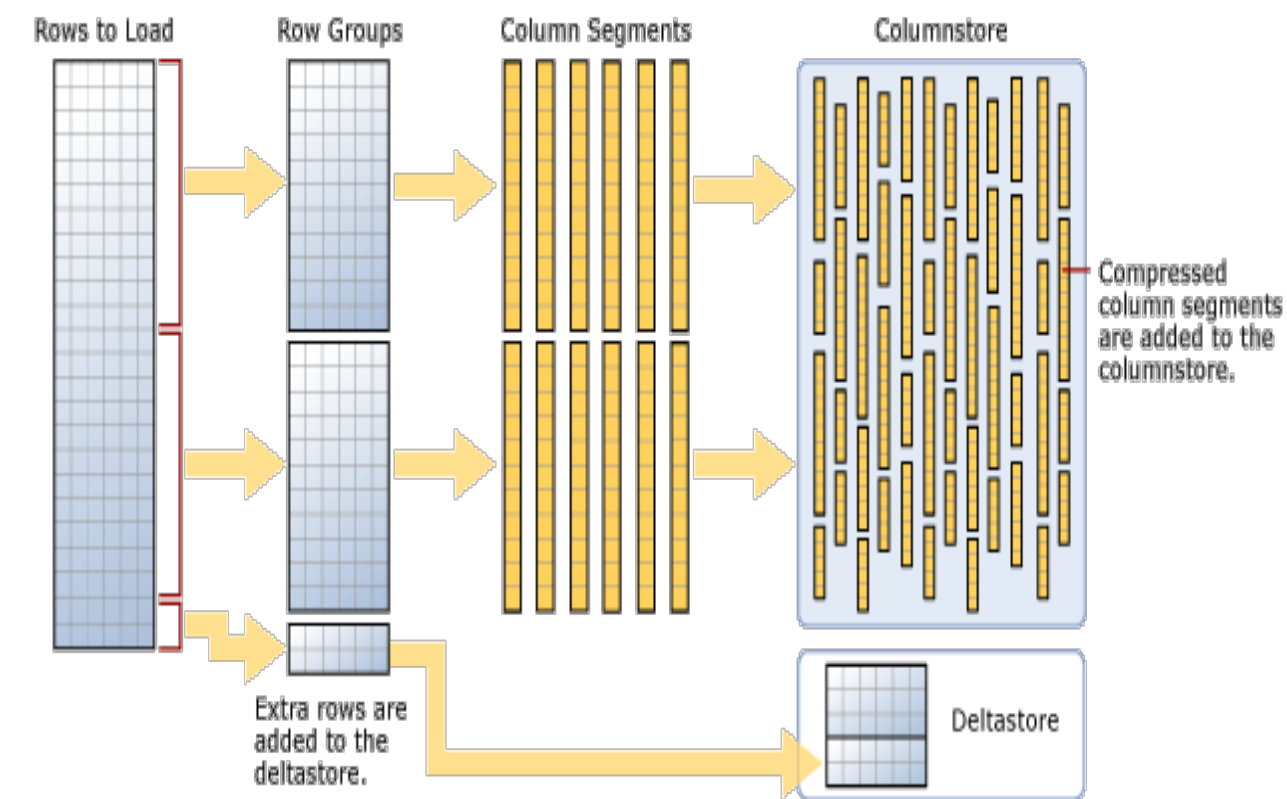
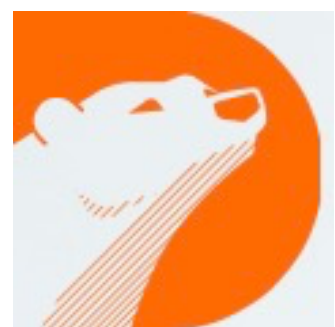


# 传统流行数据库产品

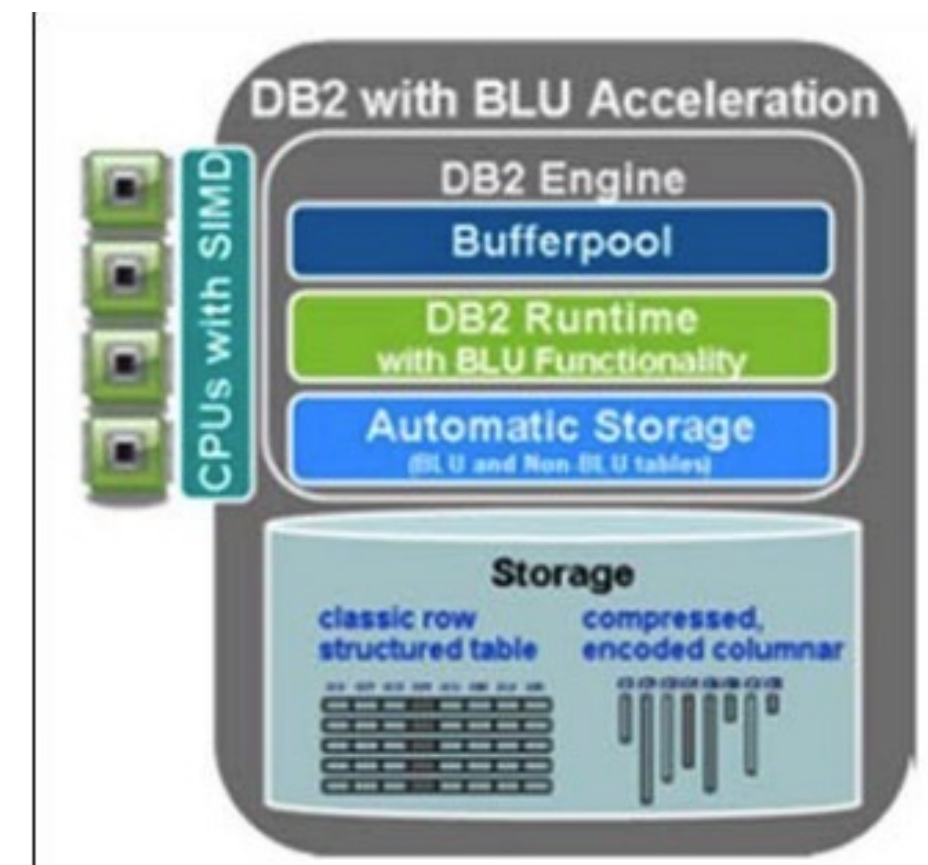
- OLTP与OLAP需求在一套系统里面被满足
- 云原生数据库追赶
  - 容量&弹性
  - AP+TP技术融合



Oracle In-Memory Database



SQL Server Column store Index



IBM DB2 BLU

# PolarDB MySQL HTAP技术挑战

## ➤ 性能

- 处理复杂SQL能力

## ➤ 数据实时性

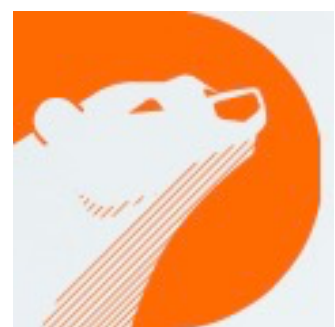
- TP更新负载下，AP请求实时性保证

## ➤ AP&TP隔离性

- TP吞吐和响应时间不受影响
- 混合负载自动分流

## ➤ 弹性扩缩容

## ➤ 100%的MySQL兼容性



# PolarDB MySQL: 并行执行加速

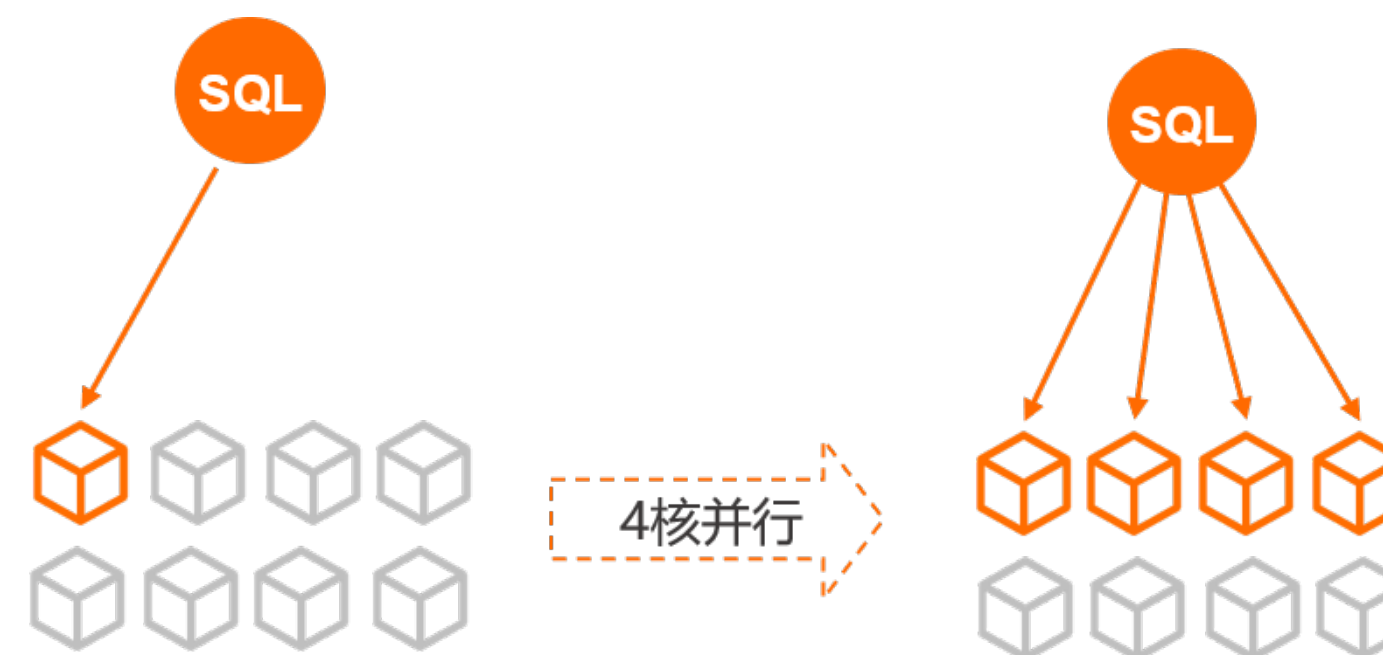
## MySQL-分析场景缺陷

- 行存
- 单核串行执行

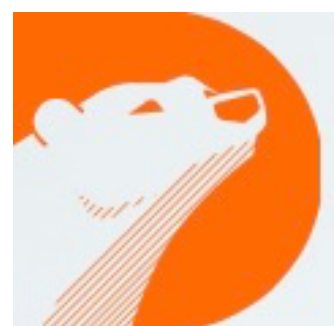
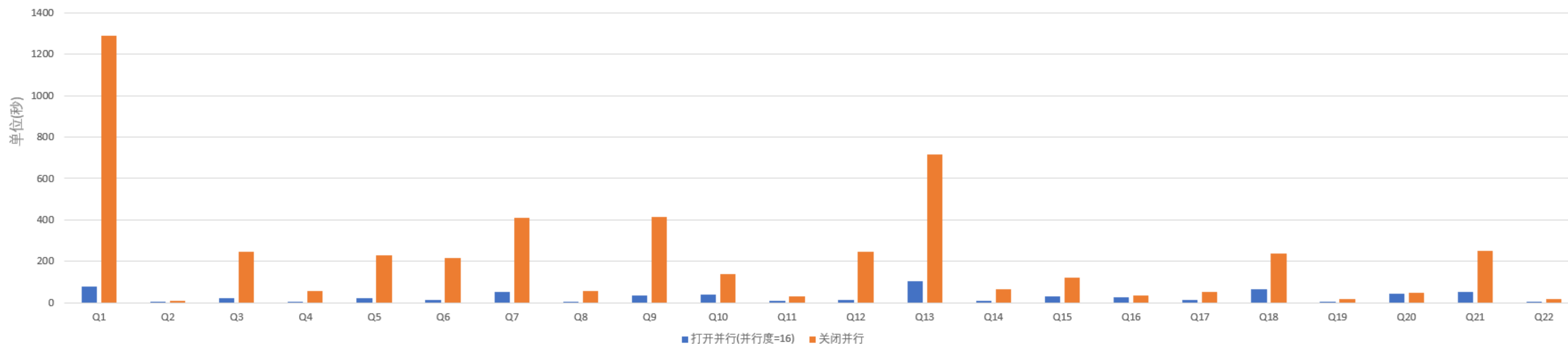
## PolarDB的改进

- 并行查询，利用多核CPU并行处理

## 需要配套的列存PK专用OLAP系统



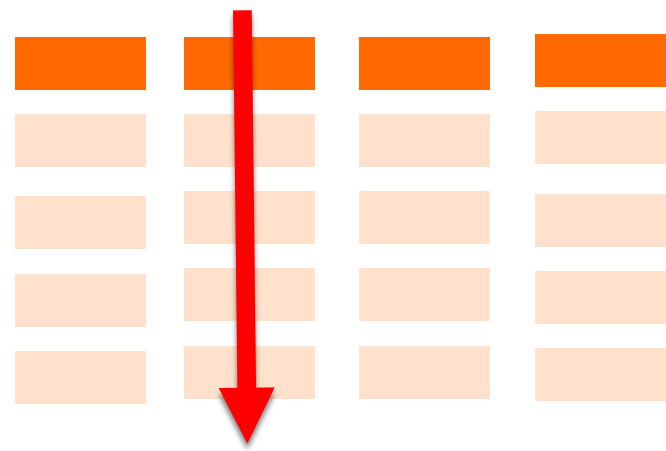
32核实例开关并行查询运行时间(越低越好)





# PolarDB MySQL : Why PolarDB Need Column Store

## 记录按列拆分存储



- 只访问必要数据

## 压缩存储



- 压缩数据及减少IO量
- 压缩数据计算

## 列存稀疏索引



- 大块数据过滤

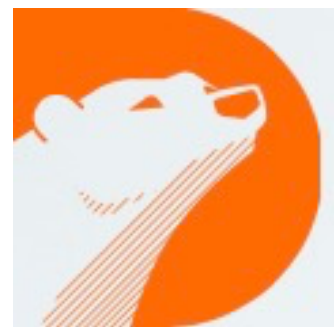
## SIMD向量执行



- 单指令多数据

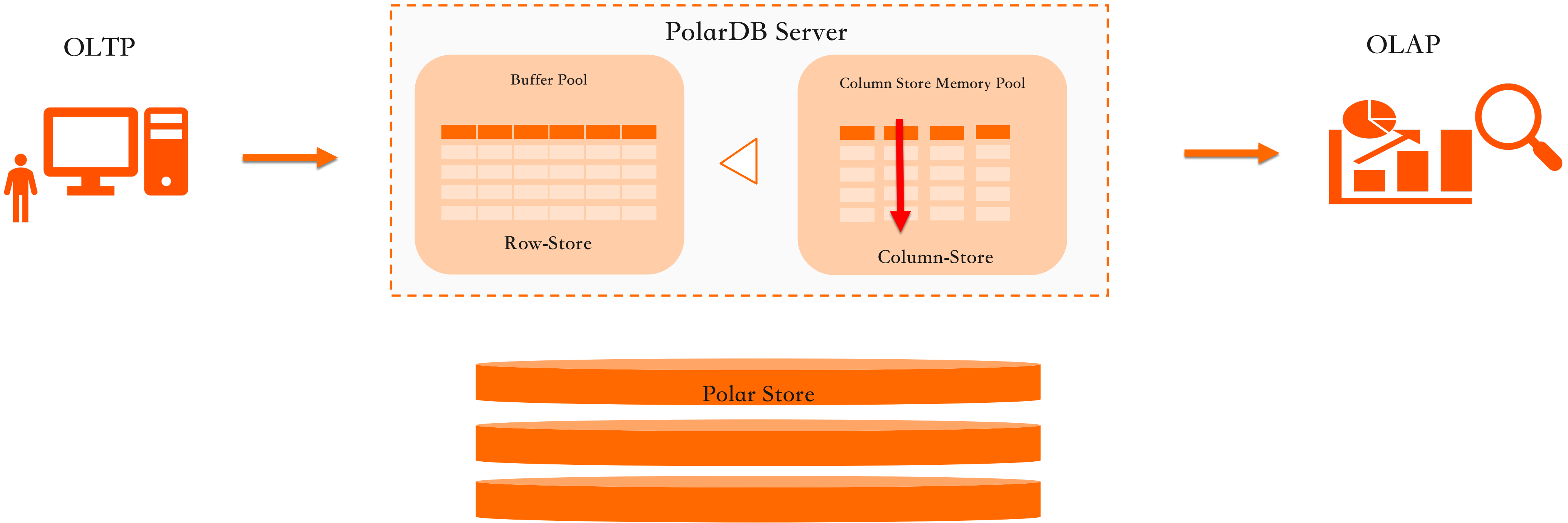
IO加速

CPU加速



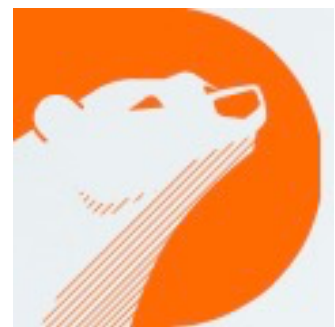
# PolarDB IMCI : 一体化HTAP数据库

- PolarDB IMCI 在一张表上同时支持行存和列存存储
- OLTP负载访问基于Page存储的行格式，OLAP分析型负载访问列存储格式
- PolarDB内部实时维持行列数据一致性，对用户透明

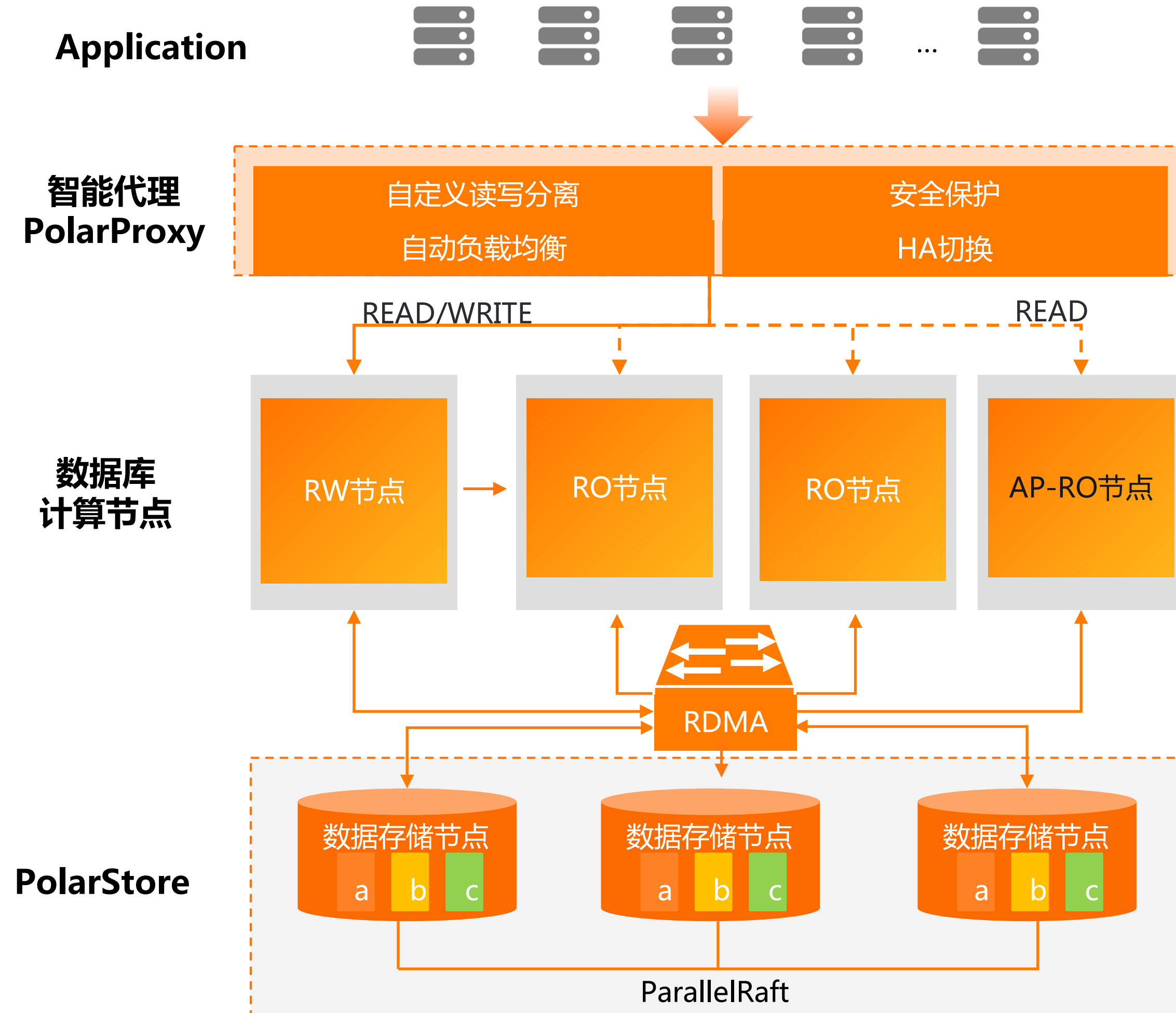


# 目录

- MySQL 在OLAP场景的挑战
- **PolarDB In-Memory Column Index 整体架构**
- PolarDB IMCI Query Engine
- PolarDB IMCI Hybrid Row/Column Store
- PolarDB IMCI 客户案例



# PolarDB IMCI：云原生一体化HTAP架构



## 高可靠高可用

- 基于分布式共享存储，数据3副本，多可用区部署
- 海量存储，支持100TB
- 秒级快照，可按时间点快速恢复数据
- 主节点故障后快速切换

## 高弹性高性能

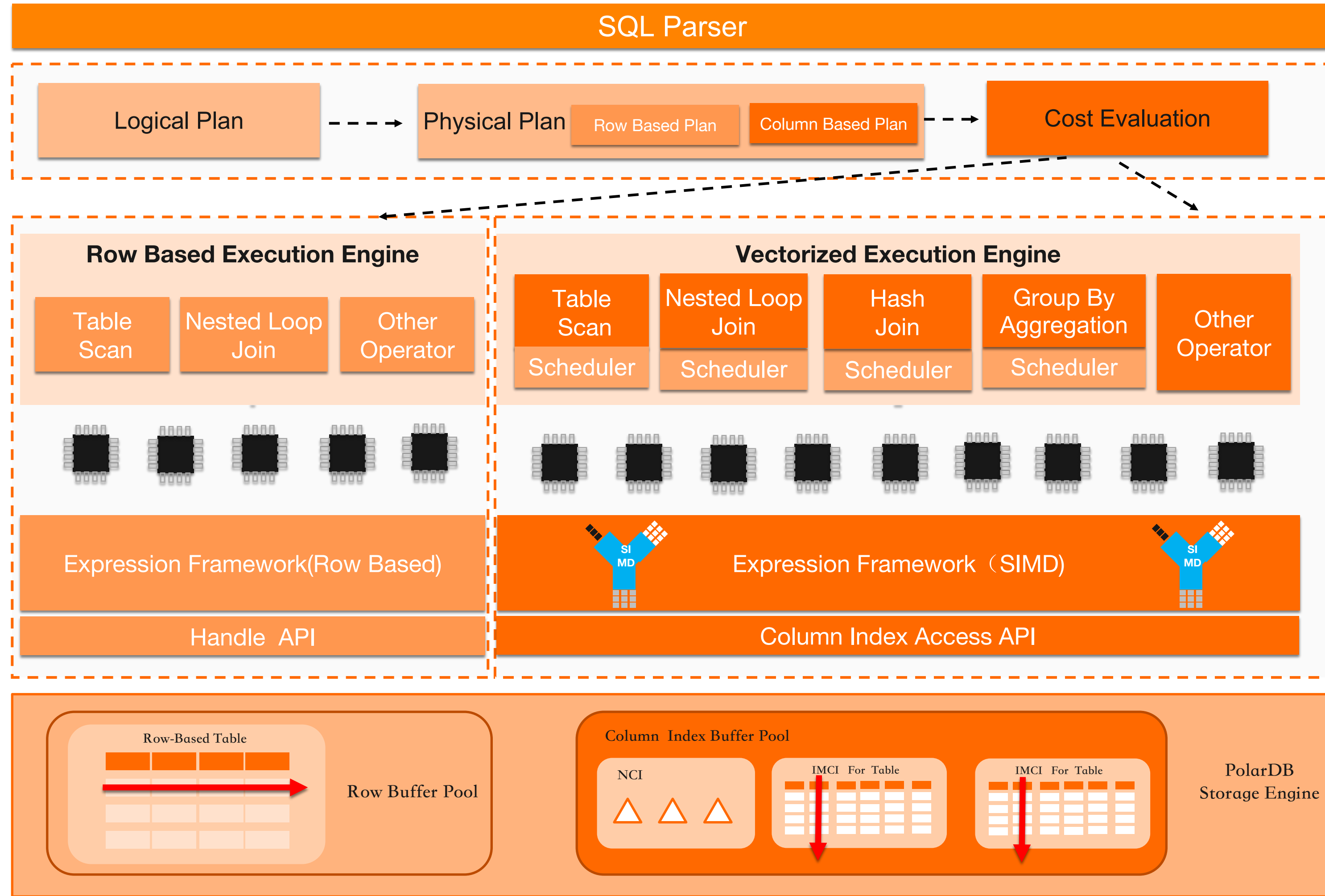
- 计算节点可扩展到**16个计算节点，1000核**
- 计算节点弹性扩展(ScaleUp, ScaleOut)
- 并行查询，充分利用多节点多核CPU
- **AP-RO节点**，并行和向量化技术，复杂查询提供**百倍的加速比**

## 易用性好

- 多个计算节点自动读写分离
- 多个计算节点**AP&TP负载自动分流**
- AP-RO节点单独部署，**AP&TP完全隔离**



# PolarDB In-Memory Column Index 整体架构



## Parser/Optimizer

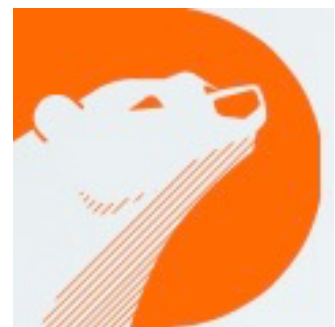
- 100% MySQL兼容
- 串行/PQ/列存混合调度
- Cost-based Optimizer

## Execution Engine

- 向量化处理架构
- 并行执行算子
- SIMD表达式加速

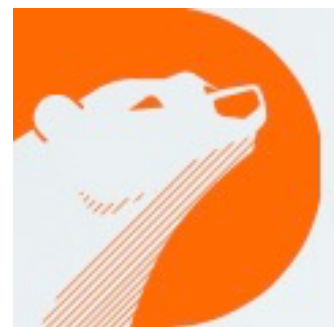
## Storage Engine

- 行列混合存储
- 事务级别的行列一致性
- 数据压缩 & 粗糙集索引
- 列存物理复制及一写多读



# 目录

- MySQL 在OLAP场景的挑战
- PolarDB In-Memory Column Index 整体架构
- **PolarDB IMCI Query Engine**
- PolarDB IMCI Hybrid Row/Column Store
- PolarDB IMCI 客户案例
- 总结与展望



# IMCI Query Engine: Hybrid Execution

## 100% MySQL兼容性如何实现

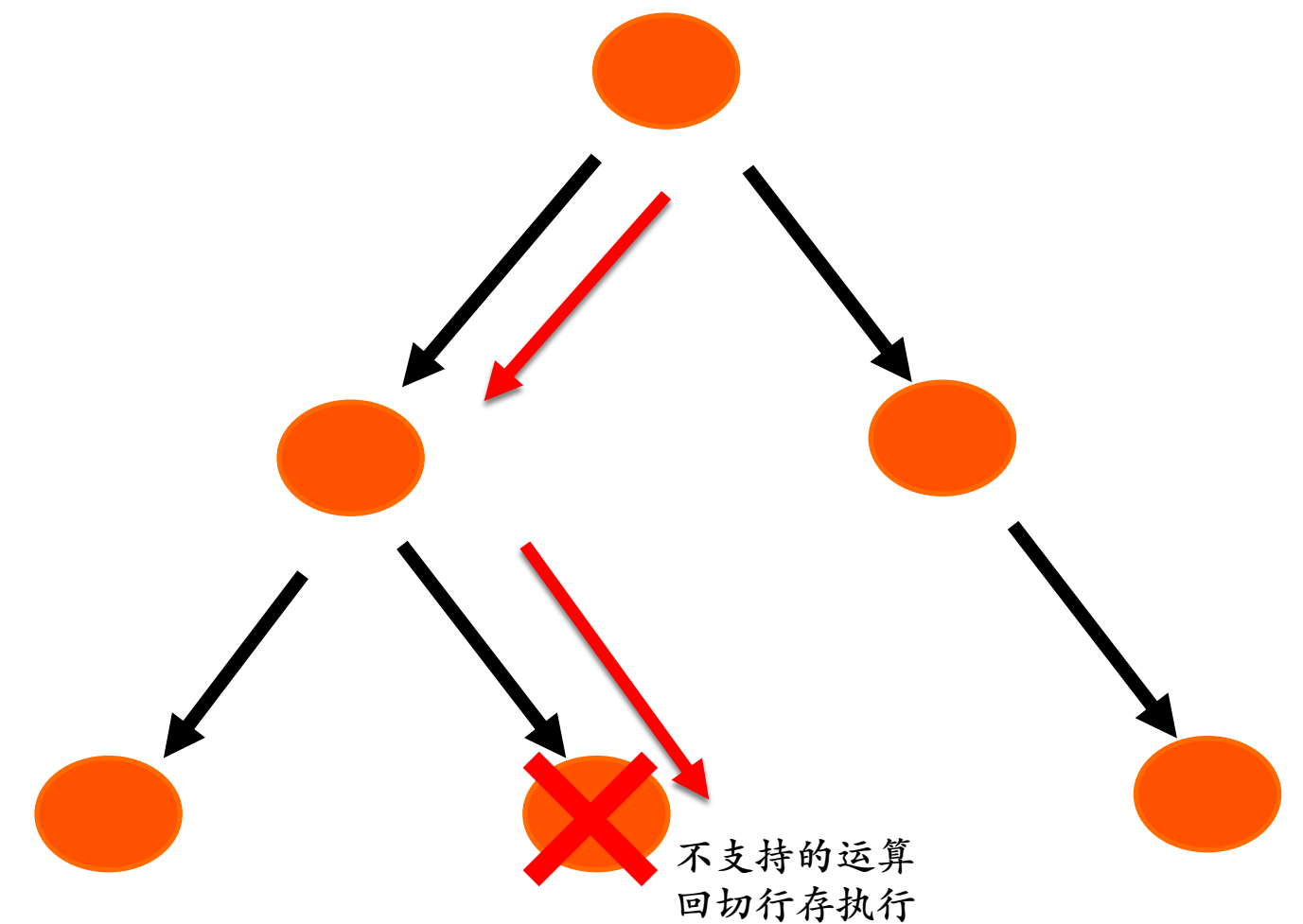
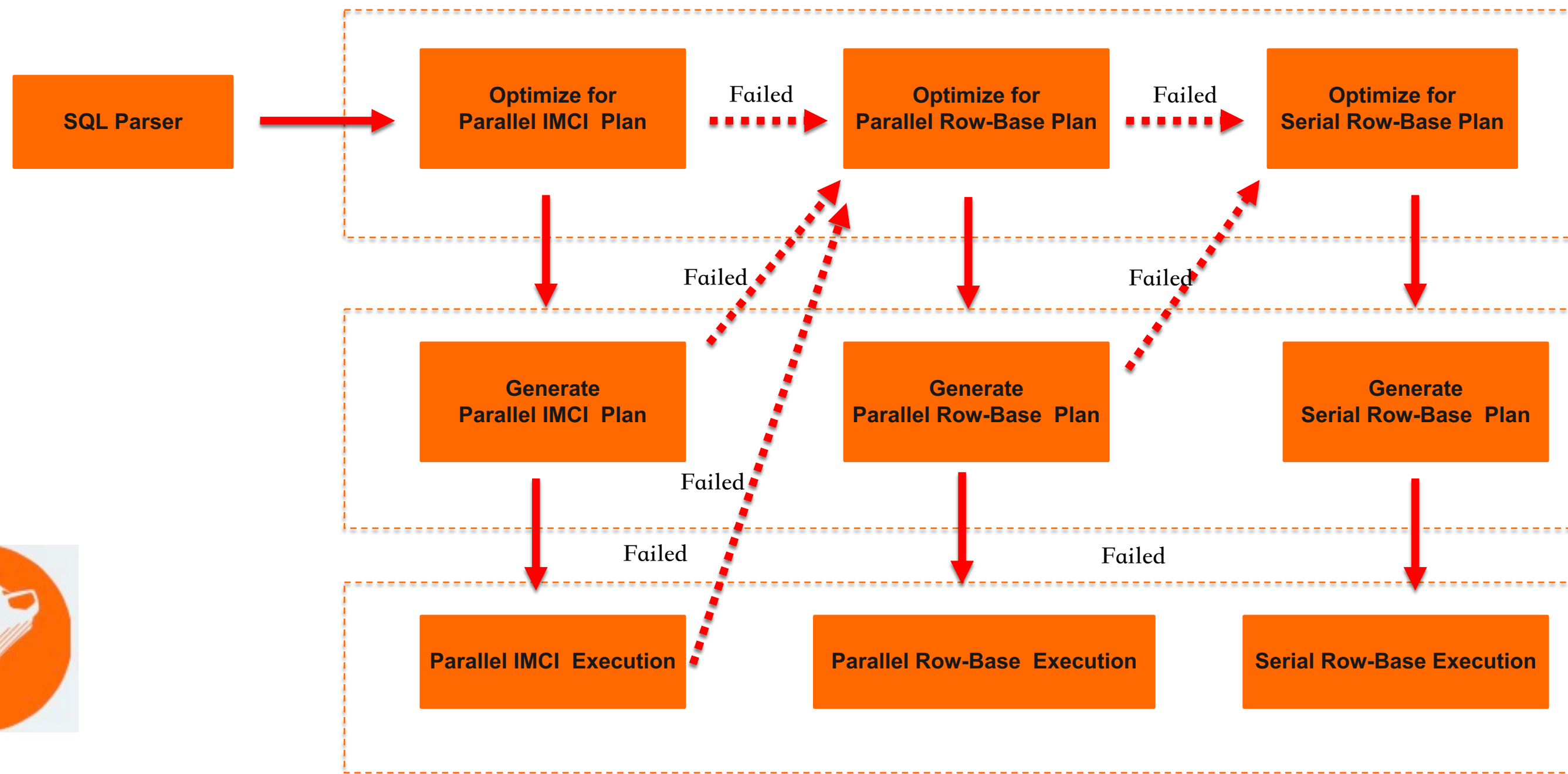
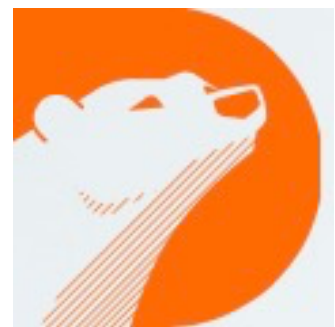
- 行存串行 100%兼容
- 行存并行
- 列存并行

## 基于白名单的执行加速

- 遍历执行计划树以探测IMCI兼容性
- 不兼容的Plan回退行存执行

## 基于代价的执行计划选择

- 行存Cost VS PQ COST VS 列存Cost
- 整体代价较低时选择行存



# IMCI Query Engine: Parallel Execution

## 行存执行器

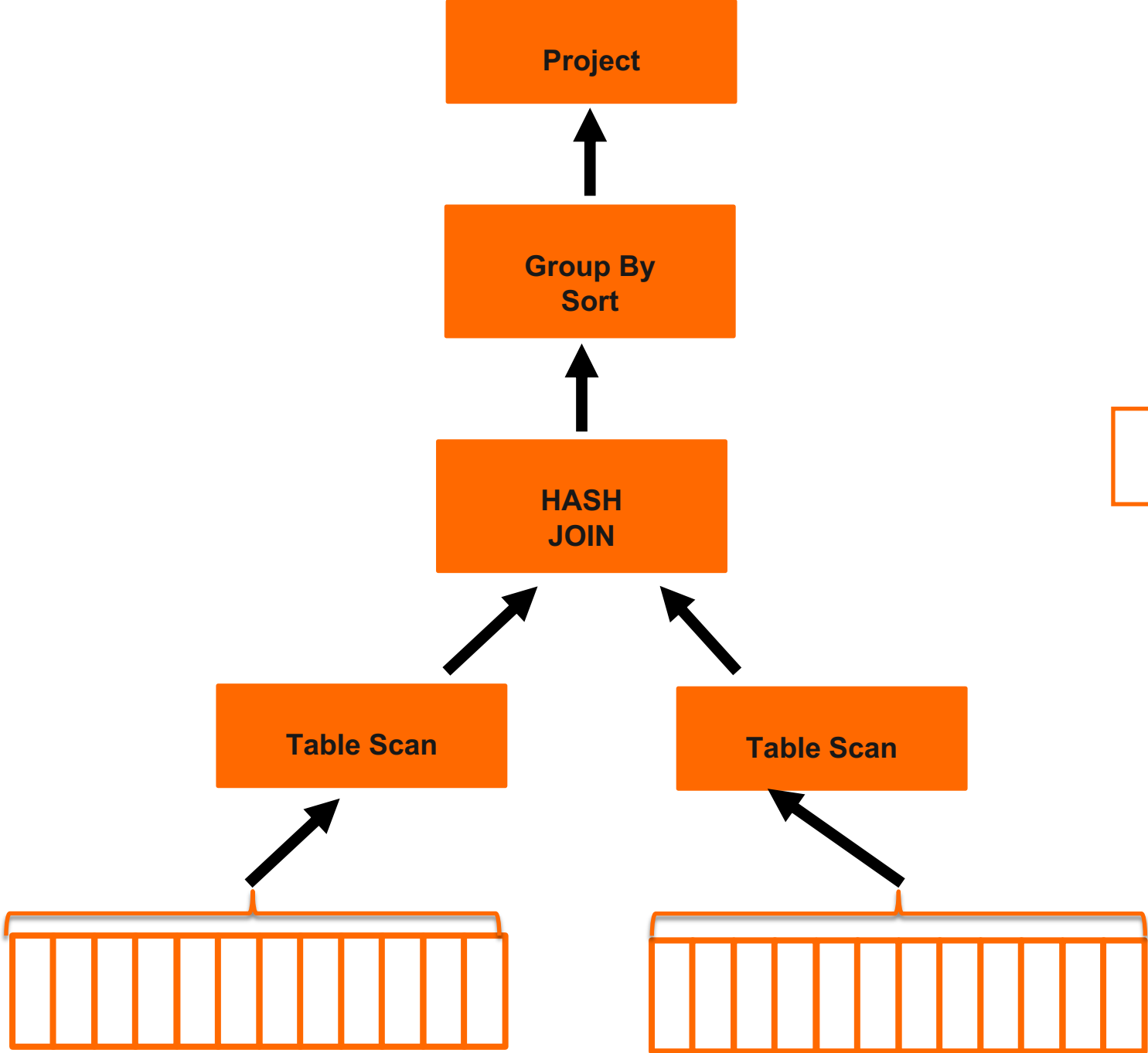
- 火山模型执行模式
- 按行迭代
- 算子串行

## 列存执行器

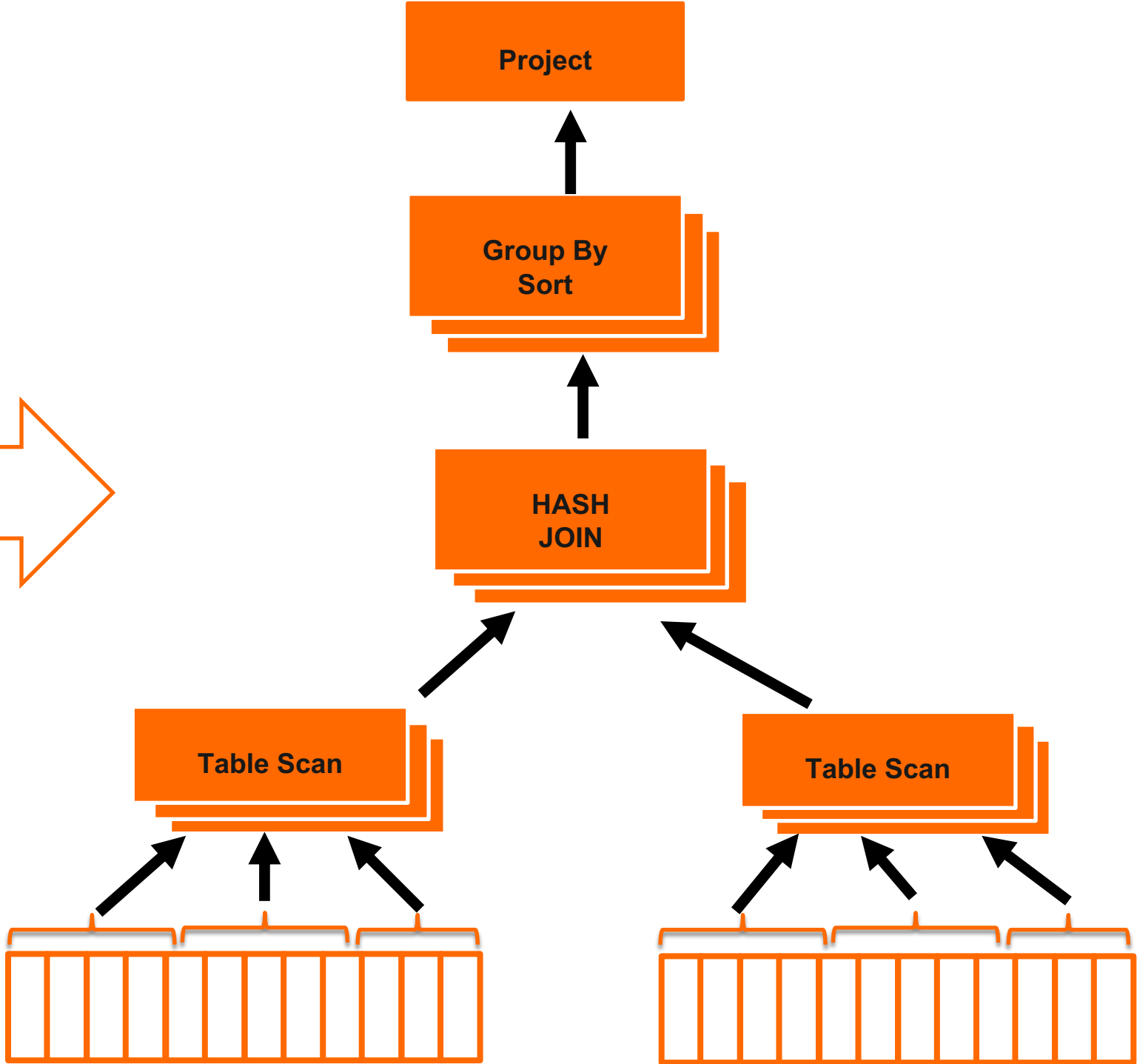
- 火山模型执行模式
- 批处理
- 算子并行



Serial Execution Plan



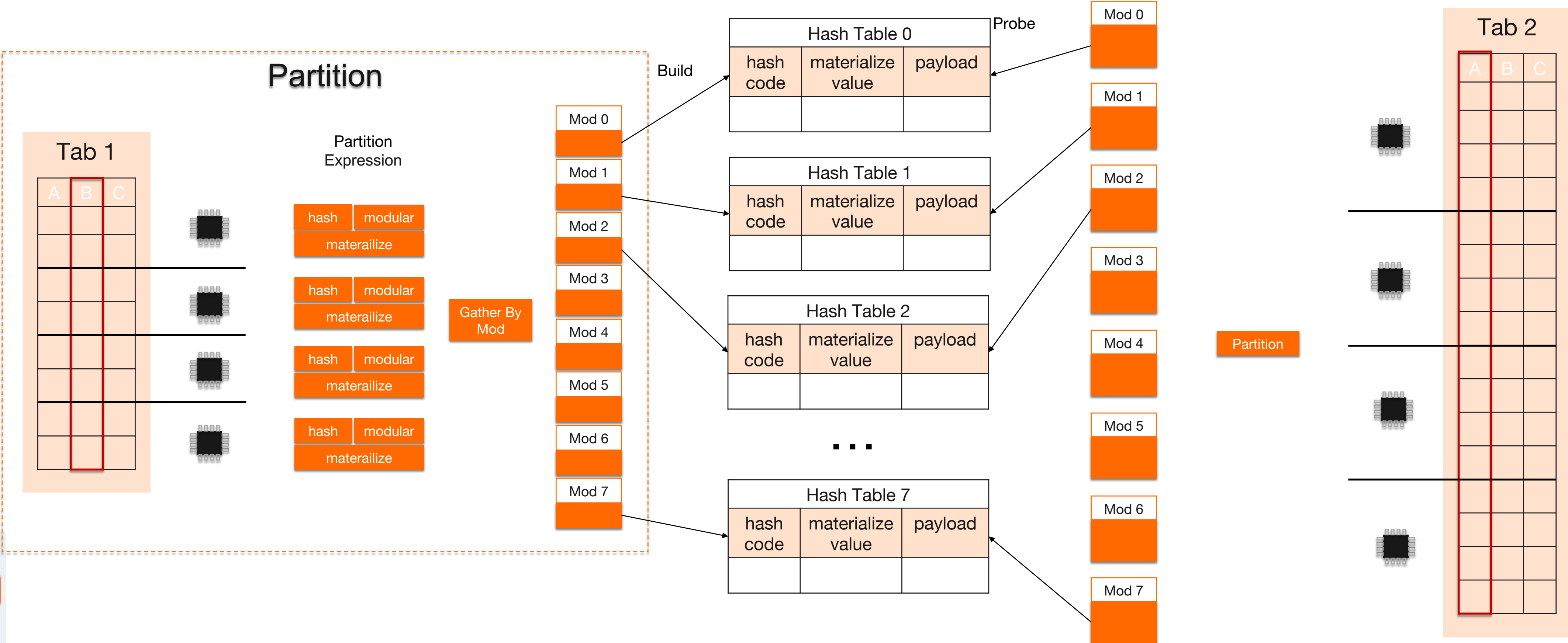
Parallel Execution Plan



```
SELECT customers.cust_first_name, customers.cust_last_name, MAX(QUANTITY_SOLD), AVG(QUANTITY_SOLD)
FROM sales, customers
WHERE sales.cust_id=customers.cust_id
GROUP BY customers.cust_first_name, customers.cust_last_name;
```



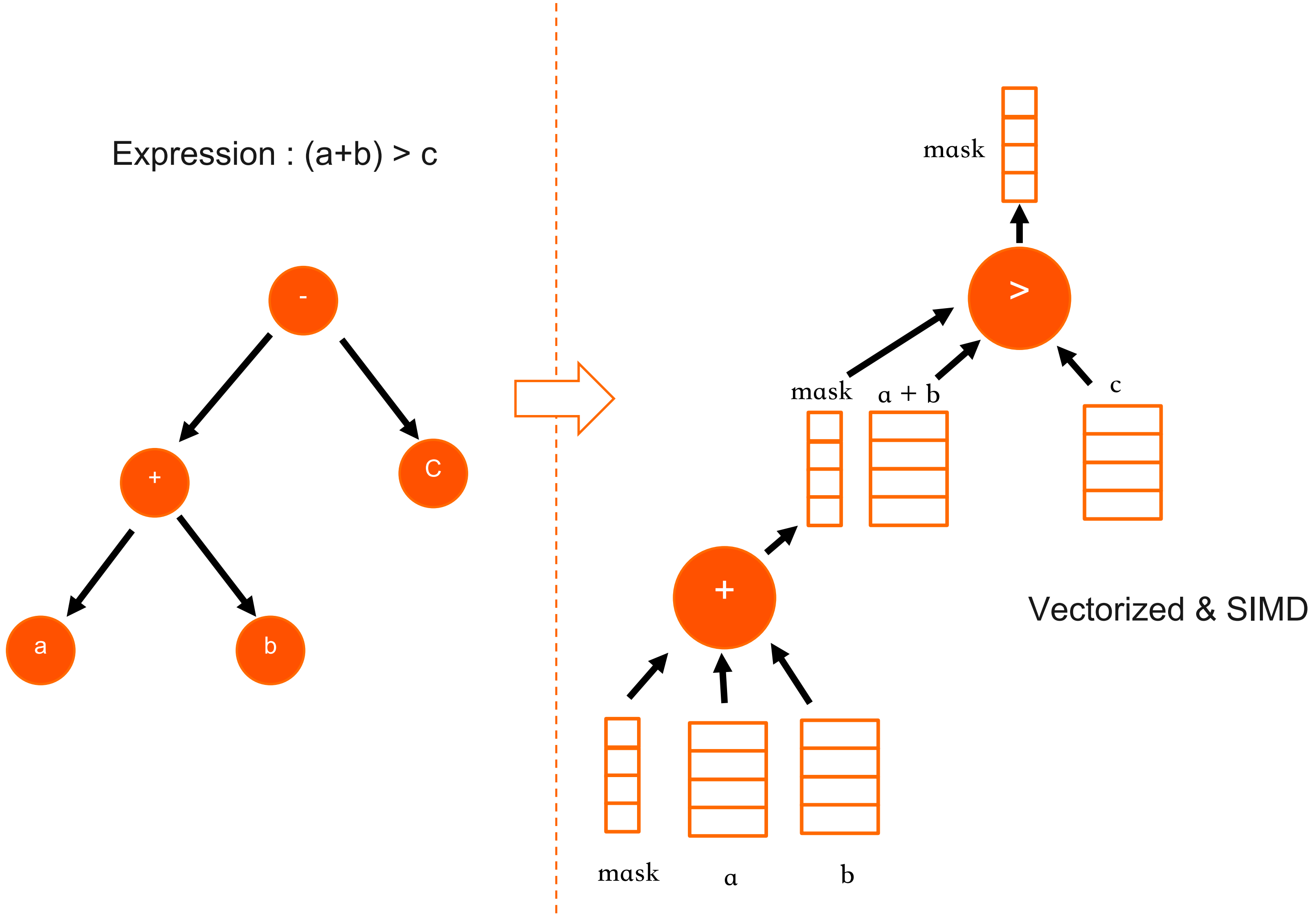
# IMCI Query Engine: Parallel Hash Join



# IMCI Query Engine: Vectorized Expression

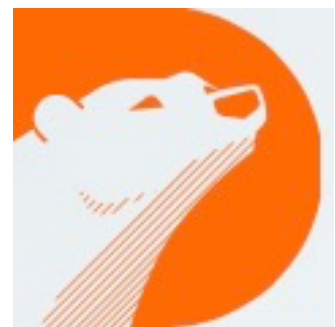
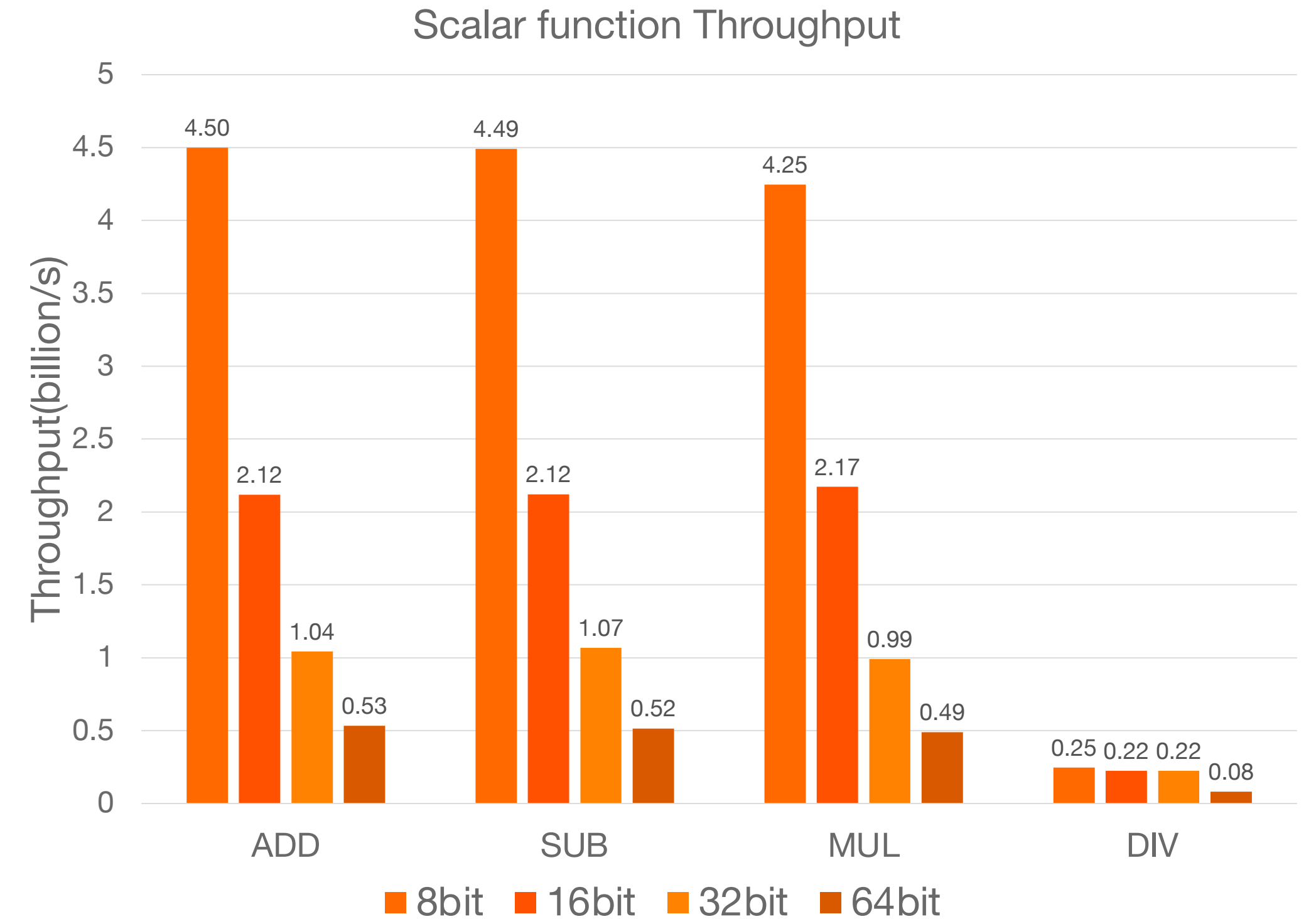
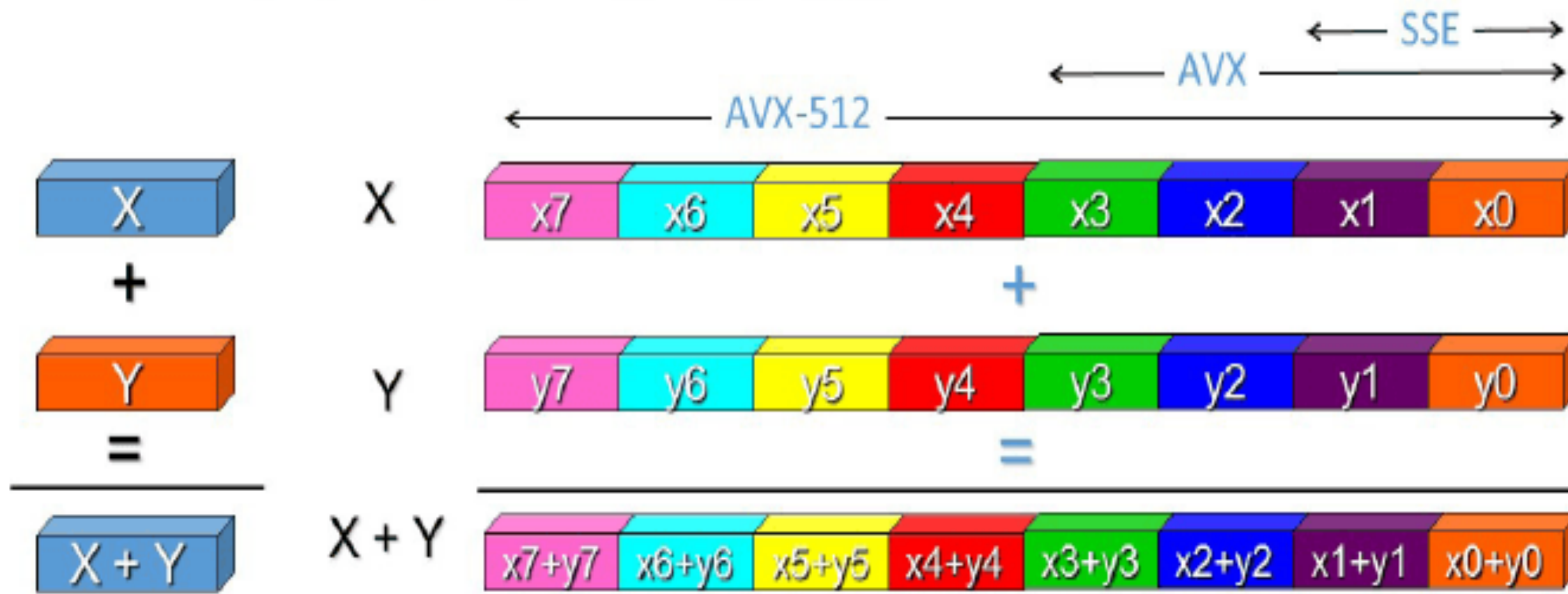
## 向量化表达式

- 批量处理带Mask运算
- AVX 512 SIMD指令加速
- 表达式递归消除提升CPU计算效率



# IMCI Query Engine: SIMD based Expression

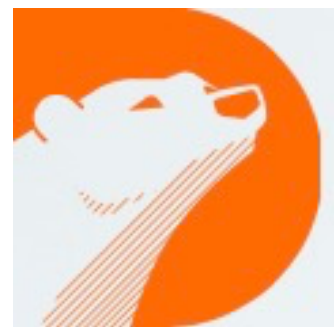
```
double *x, *y, *z;  
for (i=0; i<n; i++) z[i] = x[i] + y[i];
```



使用 SIMD(AVX512) 指令可以获得近10倍的加速效果

# 目录

- MySQL 在OLAP场景的挑战
- PolarDB In-Memory Column Index 整体架构
- PolarDB IMCI Query Engine
- **PolarDB IMCI Hybrid Row/Column Store**
- PolarDB IMCI 客户案例
- 总结与展望





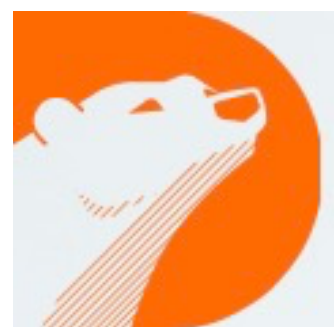
# IMCI Storage : 行列混合存储

## 存储设计，按访问模式聚集数据

- TP场景：点查，查明细，并发实时更新(行存)
- AP场景：大量行扫，少数列聚合，复杂运算
- 行列存储融合

## 行列混存优势

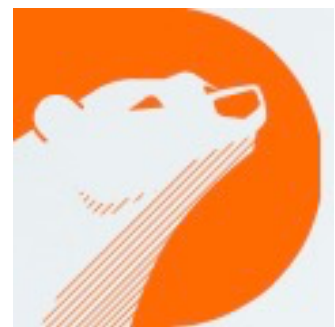
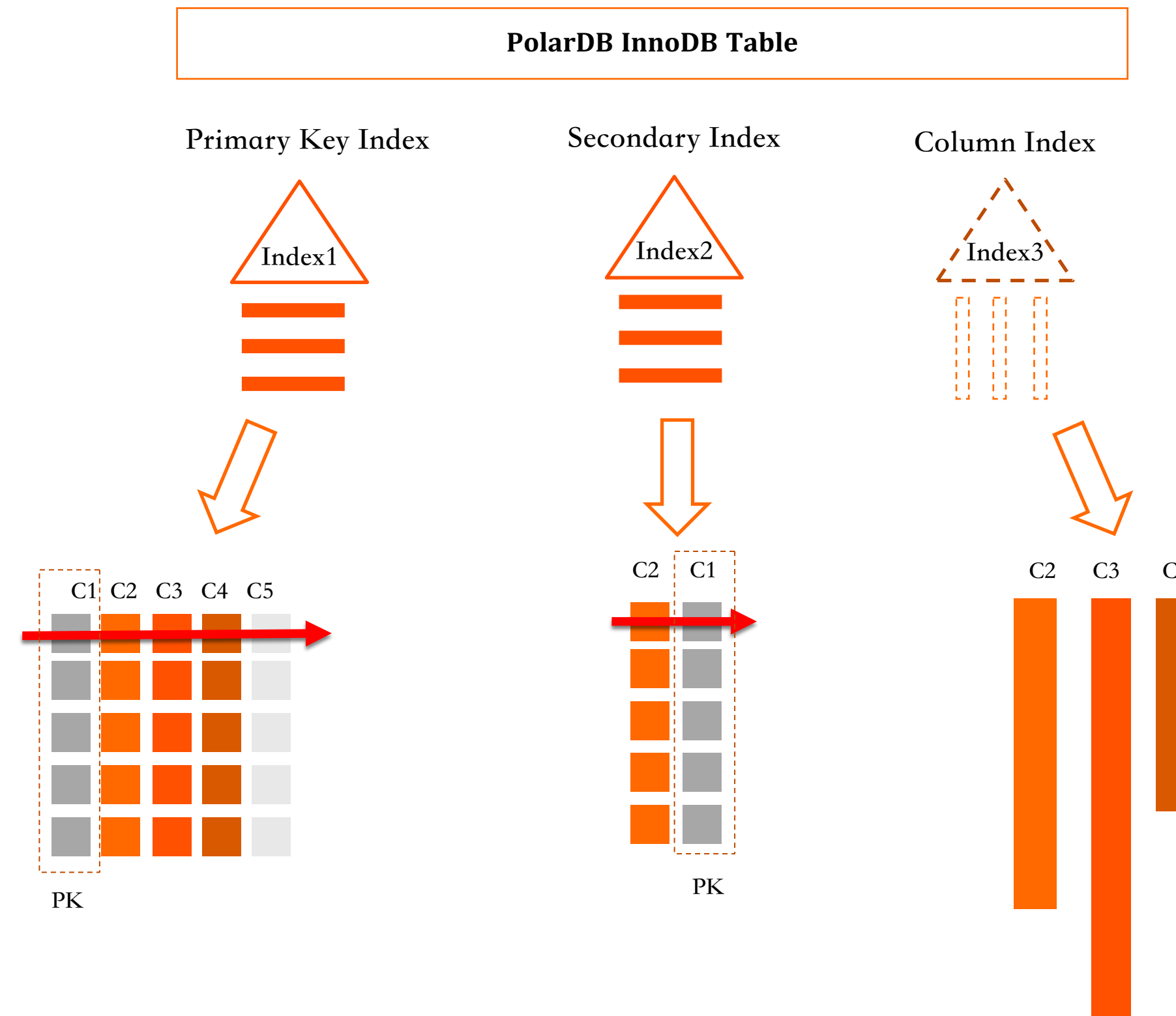
- 数据一致性保证
- 轻量级数据同步，实时性更好
- 行列混合执行
- 对用户透明，管理运维方便



# IMCI Storage : 实现为索引的列存

## 列存形态

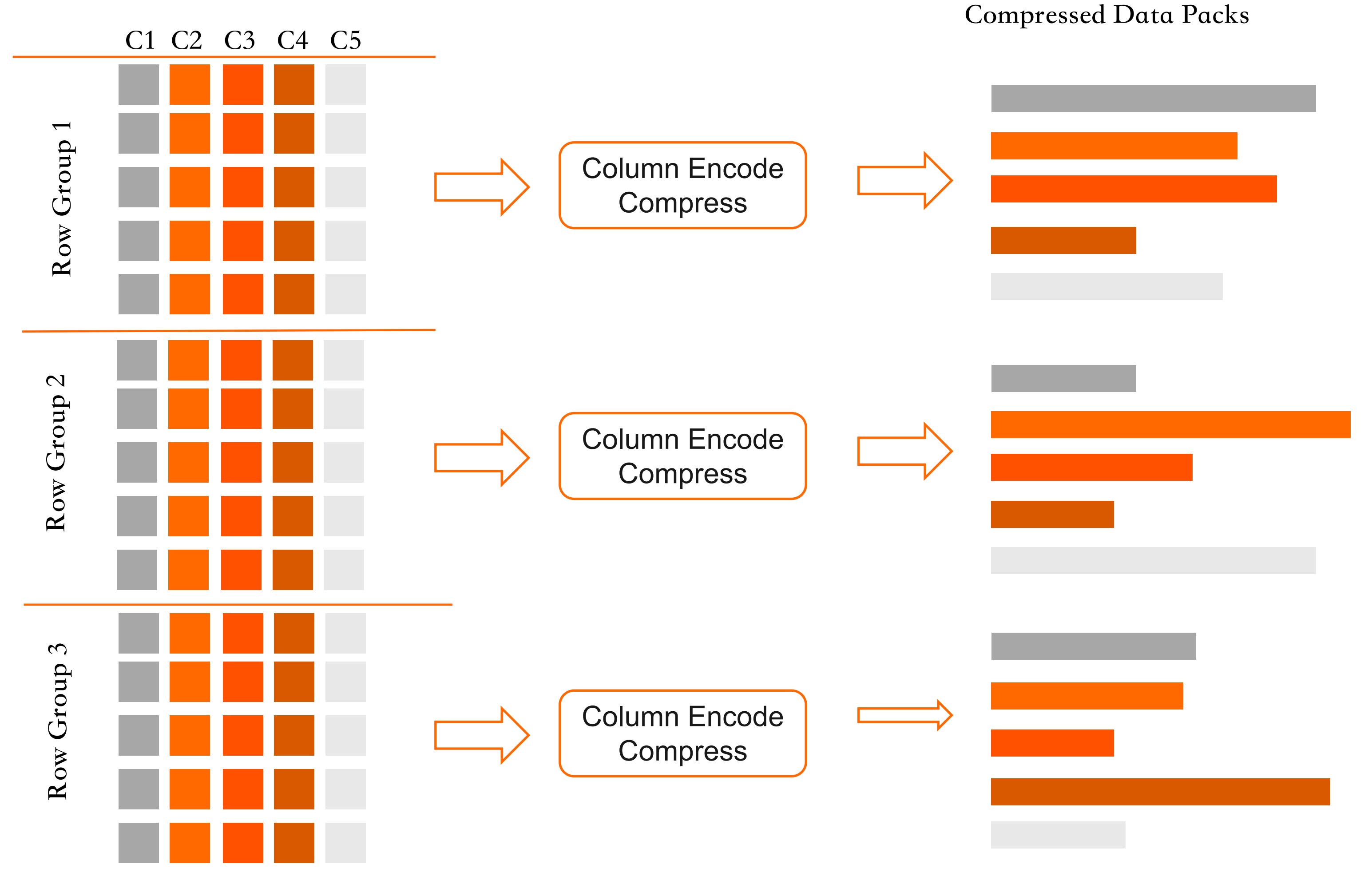
- 列存实现为InnoDB表虚拟列索引(virtual column index)
- 按表维度动态新增/删除列索引, 与二级索引一样
- Column index 覆盖列上的同步更新Apply到列存索引.
- 按需使用行存索引和列存索引对查询加速



# IMCI Storage: 列存数据组织及管理

## IMCI 列存组织

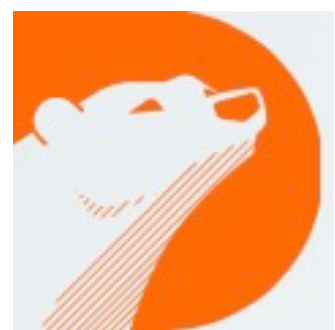
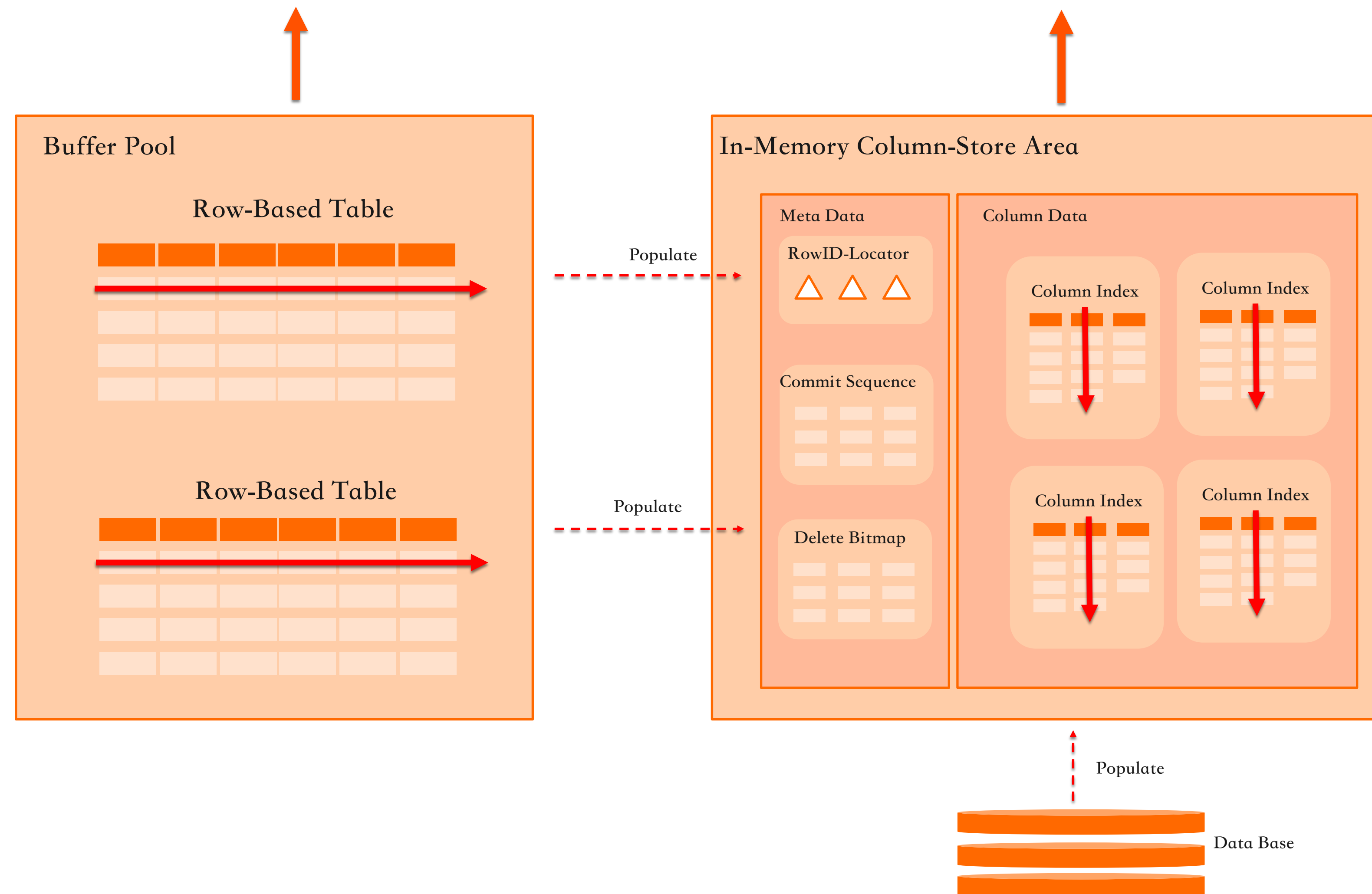
- 堆表，64K行称为一个Row Group
- 逻辑RowID标识唯一一行
- Row Group内拆分按列存储
- 列Data Pack压缩(Encoding, Compress)



# IMCI Storage : 列存管理

## 列存管理

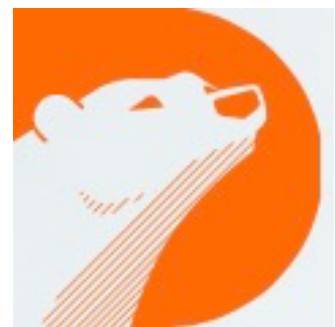
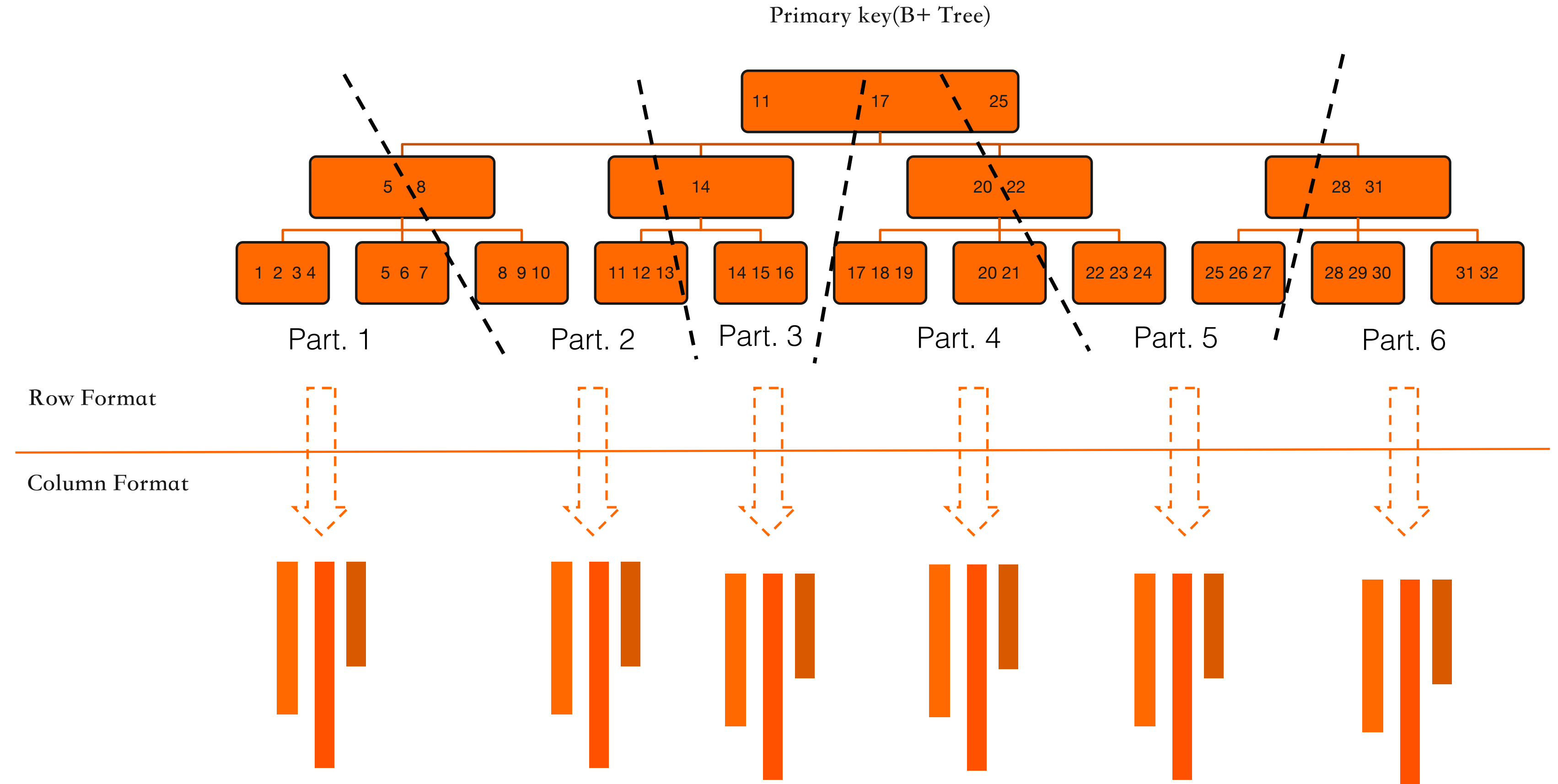
- insert追加写
- delete bitmap , commit-version
- RowID-locator , 快速定位待更新行
- Compaction机制 , 回收空洞
- Checkpoint机制 , fast-recovery



# IMCI Storage: 存量行存表转列存

## 全量构建列索引

- DDL对存量数据新增列索引
- Parallel Scan主表
- OnlineDDL, 不阻塞DML操作



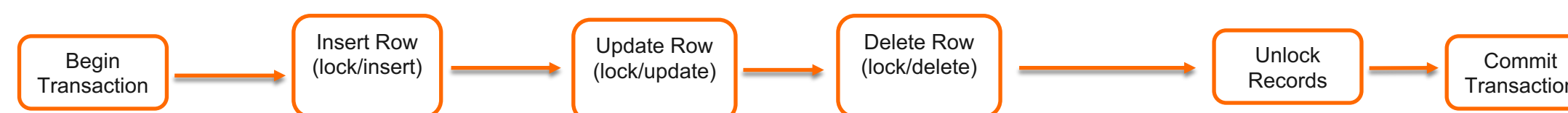


# IMCI Storage: 行列数据的同步更新

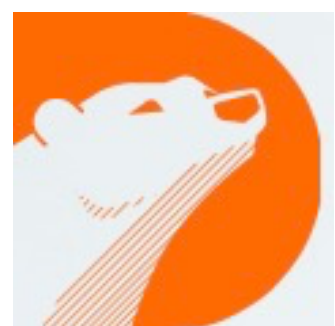
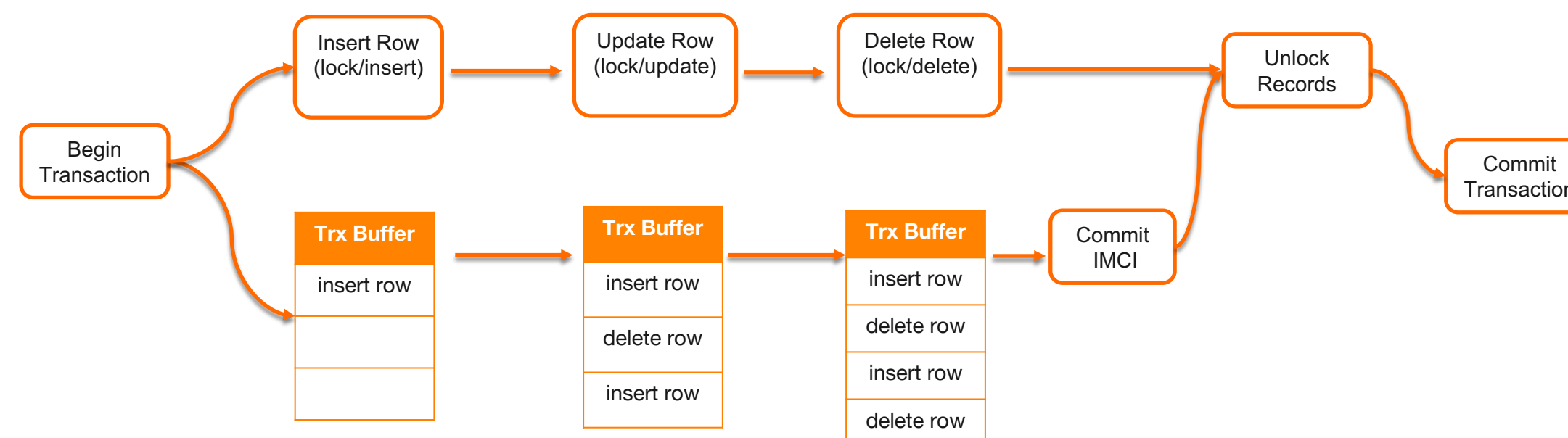
## 列索引增量构建

- 更新操作同时更新行存和列存缓存
- 与行存事务共享相同的并发控制机制
- 轻量级事务提交操作
- 大事务的处理

Transaction Without IMCI



Transaction With IMCI



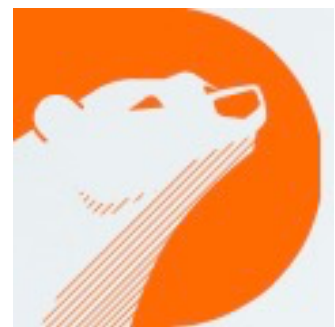
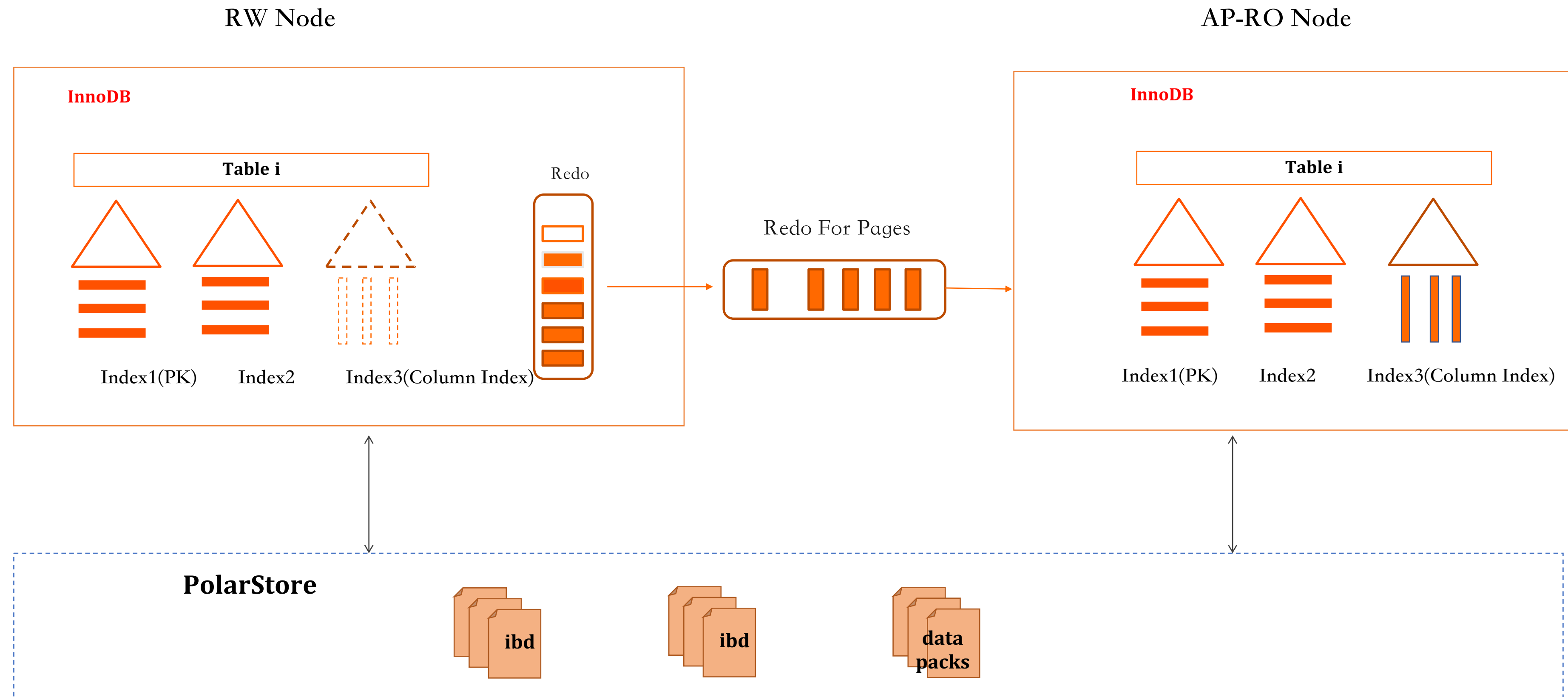
# IMCI Storage: 基于Redo的行列混合物理复制

## 数据复制技术

- binlog逻辑复制
- binlog-in-redo
- reuse-redo物理复制

## 实时同步技术

- 异步复制一致性(Isn位点)
- 多级并行回放



# IMCI Storage : 列存数据压缩

## 数字类型编码

## 字符串编码压缩

## 通用块压缩算法

22, 22, 23, 68, 63, 63

- 原始数据  
Min value : 22

0, 0, 1, 46, 41, 41

- FOR  
存储与区间最小值差值

22, 0, 1, 45, -3, 0

- DELTA  
存储相邻数据的差值

(22, 2), (23, 1), (68, 1), (63, 2)

- RLE  
将重复值用重复数字代替

0, 0, 1, 2, 3, 3

Dict  
0 : 22  
1 : 23  
2 : 68  
3 : 63

- DICT  
字典编码

- Prefix Encoding  
前缀压缩

- RLE  
有序/重复字符串

- Huffman coding  
将重复值用重复数字代替

- FSST coding  
字典编码

- DICT  
字典编码

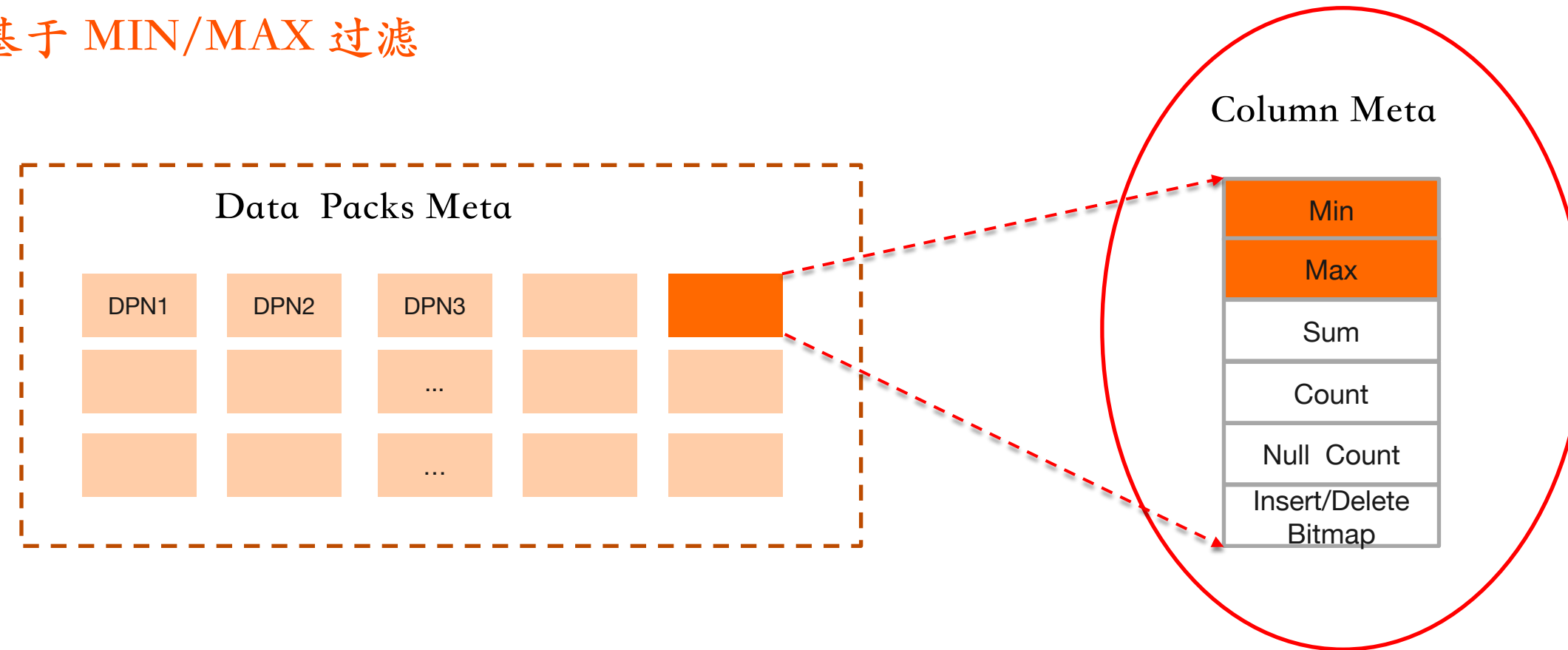
- LZ4

- ZSTD



# IMCI Storage : 列存辅助索引

## 基于 MIN/MAX 过滤



## Column Store 索引加速

### ➤ Min/Max 索引

列存数据块维护MIN/MAX/SUM等统计信息  
基于统计信息实现大块数据的Pruning

### ➤ Histograms

单个pack按数据min/max等分若干分区  
使用bitmap记录对应分区是否有数据命中

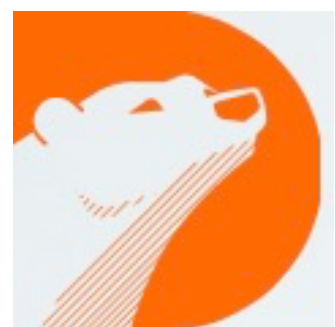
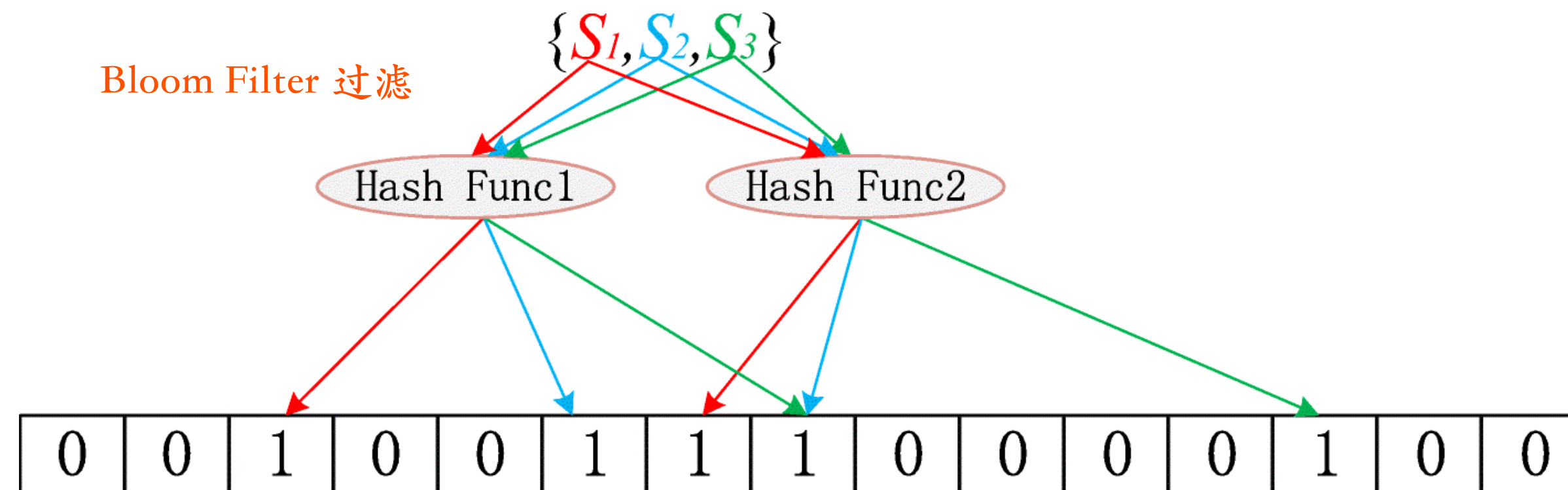
### ➤ Bloom Filter

单个pack中借助多个hash函数构建bitmap  
优化等值/in等精确命中查询

### ➤ Adaptive Radix Tree

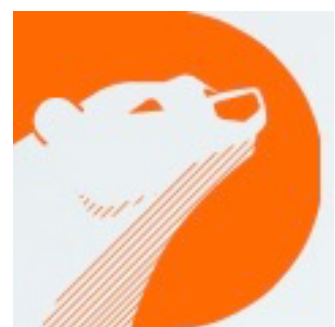
单独列再次基于ART 构建有序索引结构  
优化字符前缀，范围查询

## Bloom Filter 过滤



# PolarDB IMCI Performance: TPC-H 100G PK Click House

- TPC-H对比行存加速超100倍
- 80%的SQL比Click House更快



### TPC-H 100GB MySQL 行存 VS Column Index VS Click House

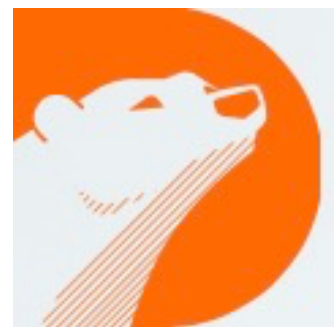


Max DOP = 32    Data All in Memory



# 目录

- MySQL 在OLAP场景的挑战
- PolarDB In-Memory Column Index 整体架构
- PolarDB IMCI Query Engine
- PolarDB IMCI Hybrid Row/Column Store
- **PolarDB IMCI 客户案例**
- 总结与展望



# 跨境电商：易仓科技

## 分析业务

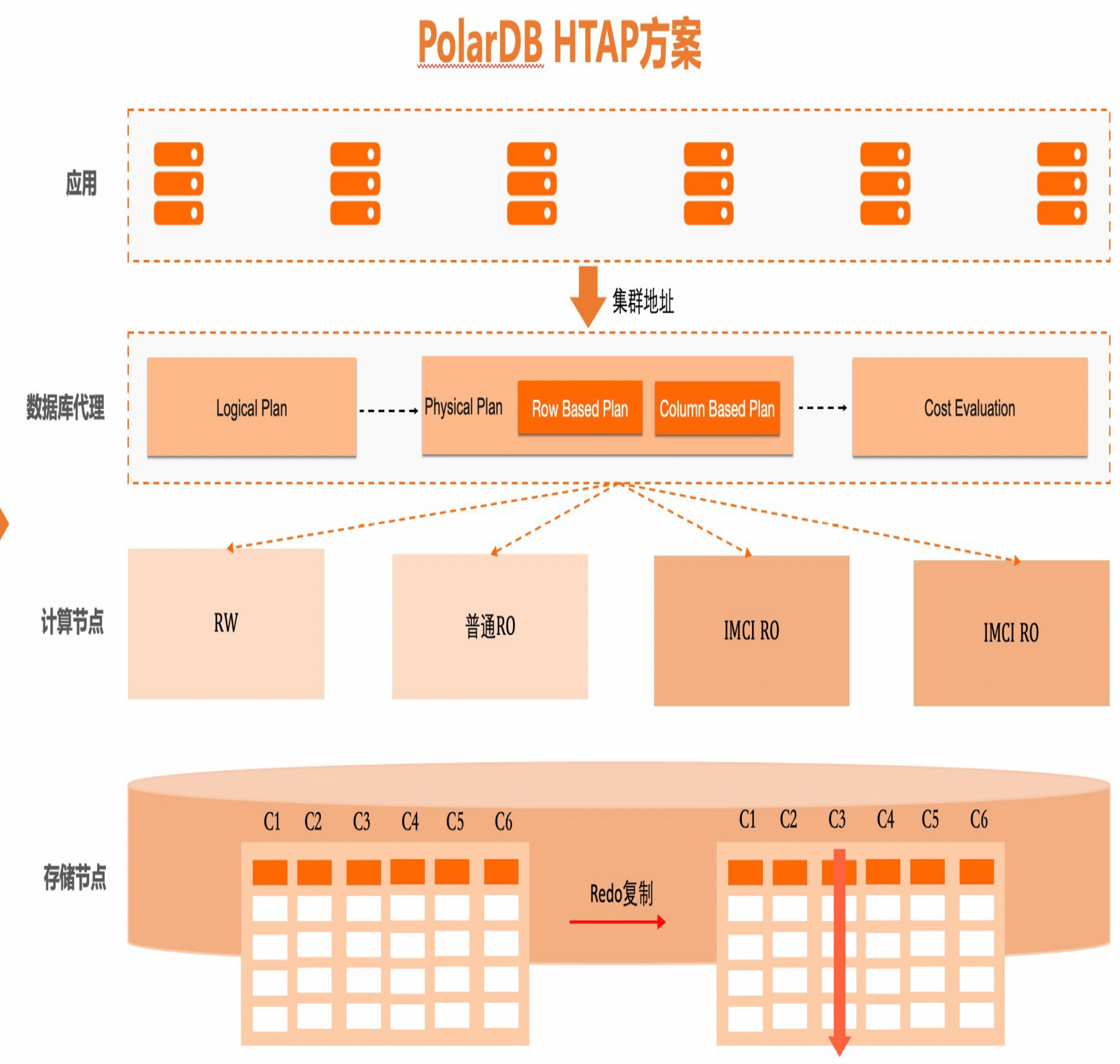
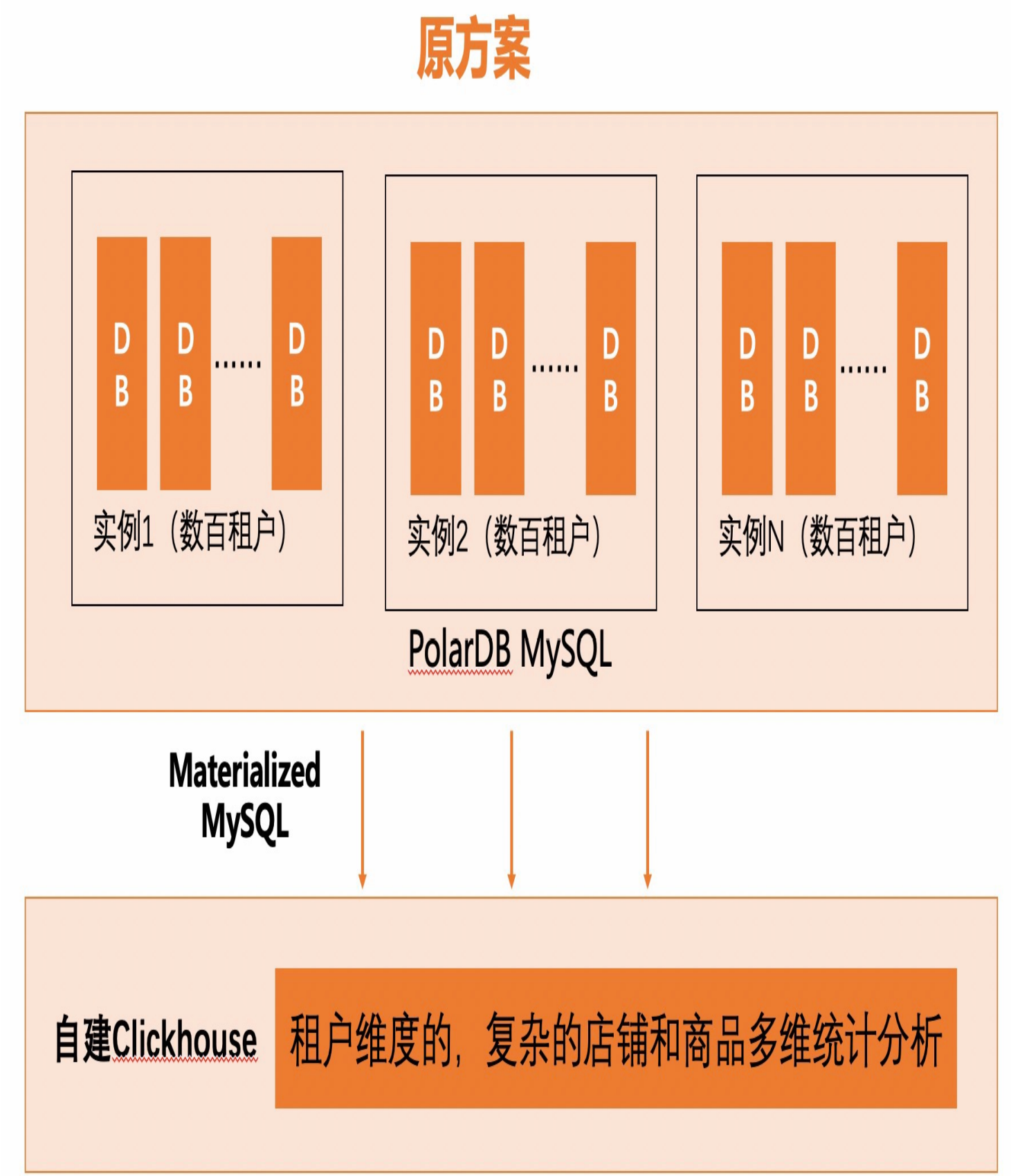
- 店铺仓储管理
- 商品推广分析

## 原有架构问题

- 复制延迟大
- 性能瓶颈
- 管理复杂

## IMCI核心优势

- 简单易用，100%兼容MySQL
- 实时性好
- 成本低



# 数字营销：云想科技

## 分析业务

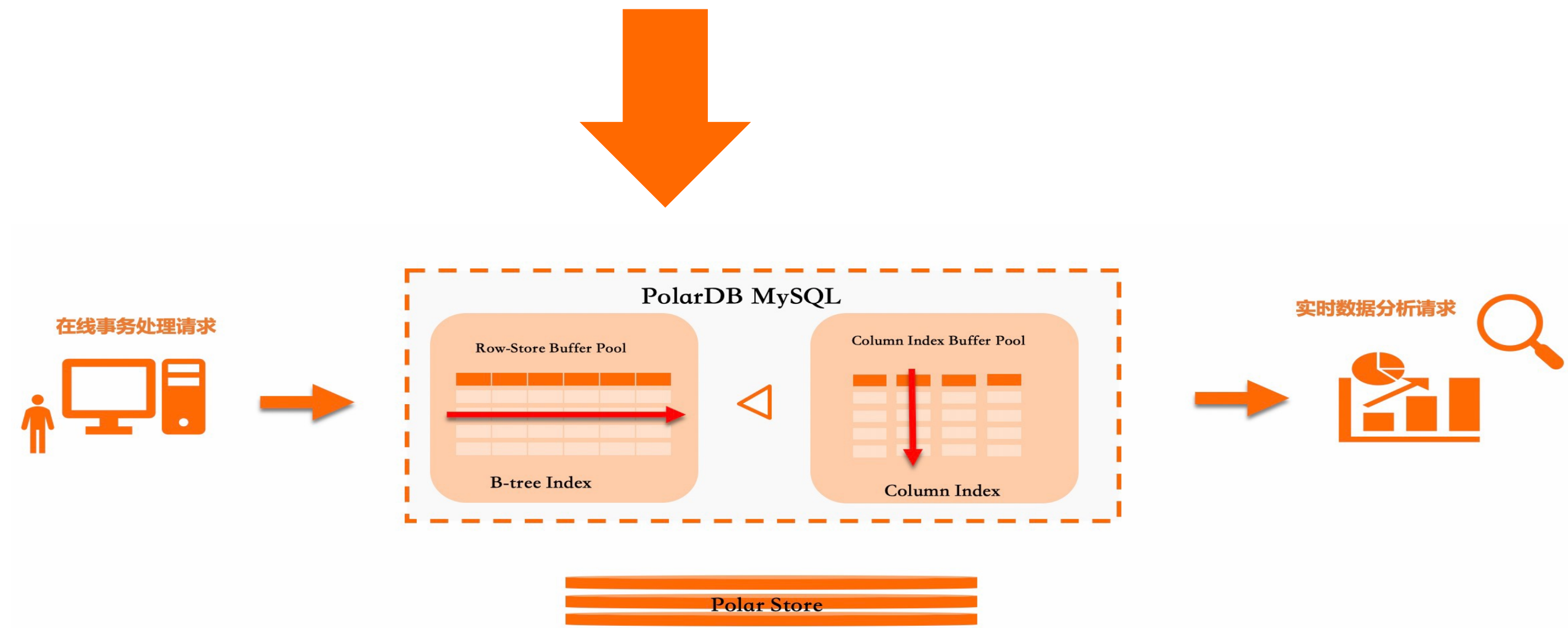
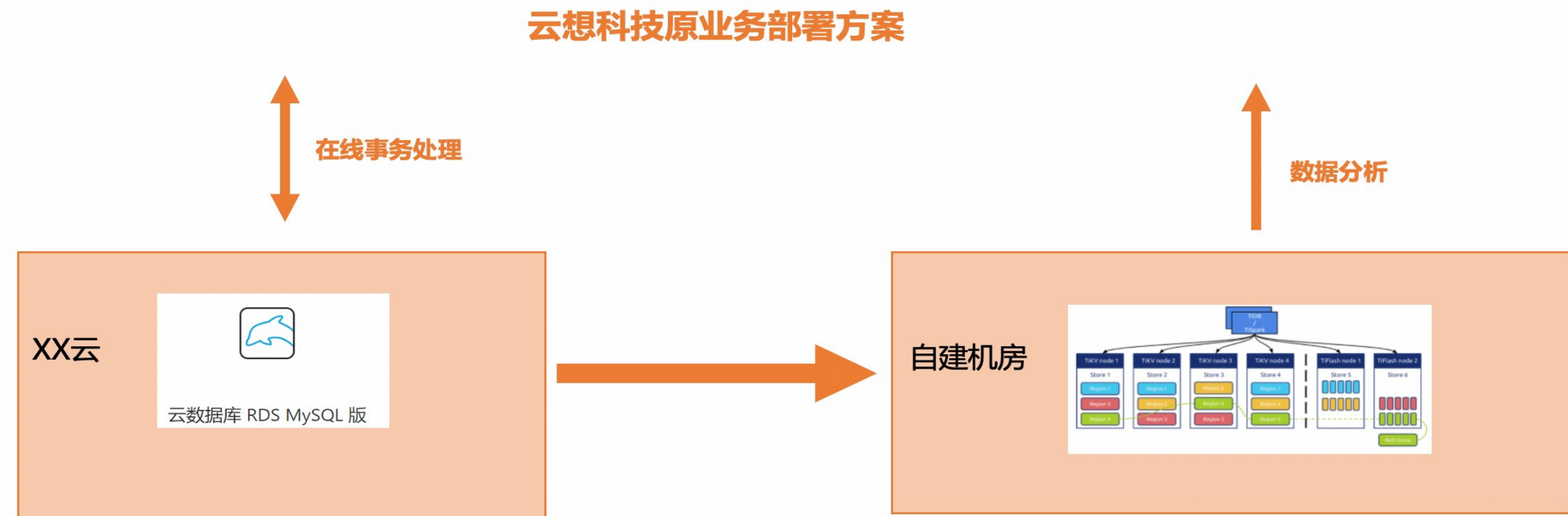
- 短视频精准营销
- 广告智能投放

## 原有架构问题

- 公司自建网络稳定性差
- TiDB 串行DDL变更卡顿
- 运维TiDB人力成本
- 两套系统一致性问题

## IMCI核心优势

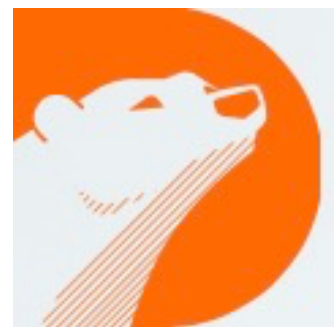
- 一套系统同时解决TP&AP
- Instant-DDL，并行DDL





# 目录

- MySQL 在OLAP场景的挑战
- PolarDB In-Memory Column Index 整体架构
- PolarDB IMCI Query Engine
- PolarDB IMCI Hybrid Row/Column Store
- PolarDB IMCI 客户案例
- **总结与展望**



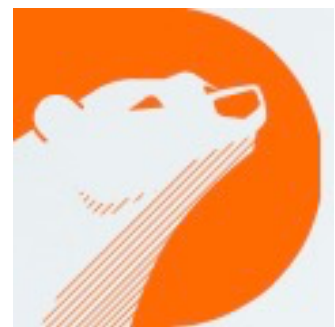
# 总结与展望

## PolarDB(HTAP)不是简单的TP+AP

- AP性能
- 数据实时性
- AP&TP隔离性
- 弹性扩缩容
- MySQL兼容性

## 下一步发展

- MPP 集群，更大容量数据分析
- 行列混合执行
- 列存表，冷热分离，降低分析成本





# Thank You

