

移动云 (<http://ecloud.10086.cn/>)

数据库产品

rds pg
目前支持rds pg 10、12
后续计划合并到云原生pg
今后发展可能以云原生pg为主

云原生pg
海山数据库 (He3DB)

基于shared-storage构建

受启发于aurora
log is database

He3FS
基于JuiceFS
改造和区别

He3Open
主节点wal写共享存储
从节点接收wal写本地存储
主要价值: 文件句柄到文件的mapping关系

He3Write
在JuiceFS写入逻辑的基础上,
增加解析wal的逻辑

维护wal的映射关系
key —— 数据页相关信息
value —— 对数据页修改的wal描述信息的链表
database的oid
table的oid
forknum
blocknum
其他信息
wal所属文件的inode
wal的长度
wal的offset
其他

数据写入时, 先判断是主节点还是备节点 —— 主节点 —— 解析wal

DataRead
平常查找
增加接口功能
查找数据页相关wal
拼接数据页与wal

He3UnLink、He3Close、He3Truncate、He3Fsync、He3Lseek等其他接口
参考: https://gitee.com/he3db/he3fs/blob/master/docs/zh_cn/desc.md

对标产品
aws aurora
阿里云polardb

锁机制

RegularLock
主要用于事务
table level
row level

LWLock
主要用于共享内存管理
类型
share
exclusive
有等待队列
共享内存对象
shared buffer
wal buffer
clog buffer —— mvcc的特殊性

SpinLock
主要用于其他锁 (如LWLock) 的底层实现
cpu自旋 —— 消耗cpu
依赖主机, 机器cpu指令 —— TAS (Test And Set)
不依赖主机, 依赖操作系统, 信号量 —— PGSemaphoreTryLock