

中国数据智能管理峰会

万亿数据库核心存储引擎实现与应用

演讲人: 母延年

演讲人简介-母延年

录信数软



职业生涯的前十年 任职于 新浪、阿里、腾讯

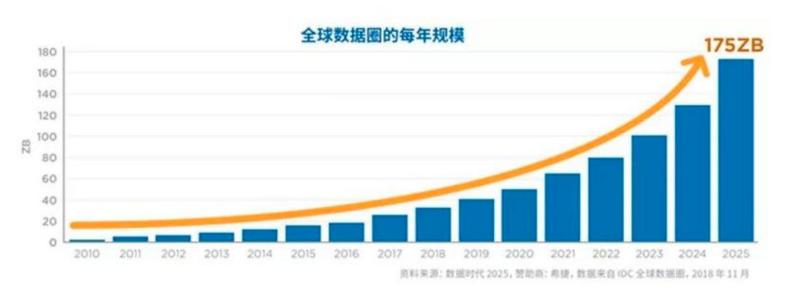
- ▶ 精通hadoop、lucene、storm等多种主流软件源码
- ▶阿里负责海狗 (higo)、黄金策等多个项目的研发,开源2个项目, Mdrill、Jstorm
- ▶ 腾讯负责腾讯数据平台hermes项目开发,单日数据规模支撑3600亿+

最近五年 南京录信创始人&CTO

录信是lucene的中文发音,作为粉丝,以此命名公司

- ▶ 有200多套在生产系统使用超1年
- ▶ 其中1个系统数据规模过10万亿 (1000+节点)
- ▶ 其中10+个系统数据规模过万亿 (300+节点)
- ▶ 其中100+个以上系统数据规模过千亿
- ▶ 获12项技术发明专利

时代背景——数据爆炸式增长



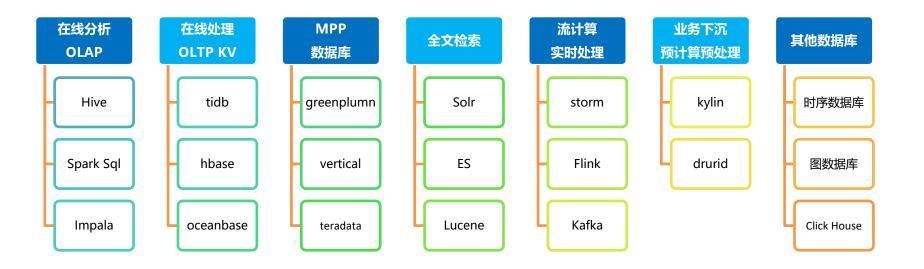
据IDC发布《数据时代2025》的报告显示,全球每年产生的数据将从2018年的33ZB增长到175ZB,2025年全球每天产生的数据量将达到491EB。

行业痛点

- 一、大数据里的产品种类很多,但每个种类内都很单一。
 - 1.绝大部分系统采用单一的"暴力扫描",性能低下。
 - 2.少量系统有索引但功能受限只能KV,或全文检索,做不了复杂的统计。
 - 3.还有一部分系统采用预计算处理,不灵活也不能查看原始明细数据。

二、为了应对这些不完善,需要混合使用多种系统

- 1.数据存储的份数太多,浪费存储资源。
- 2.多个系统之间数据互通很难。
- 3.每个系统接口都不一样,学习与维护成本很高。



用户的期望

能不能将他们整合起来

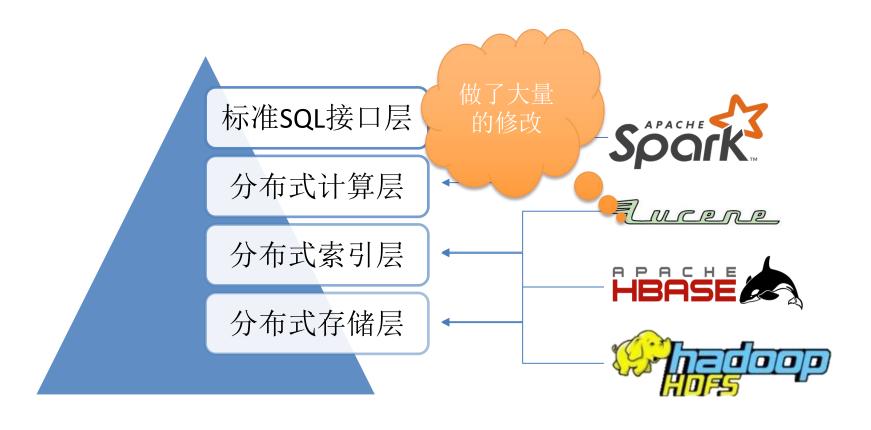
一份数据、一套系统、一种接口



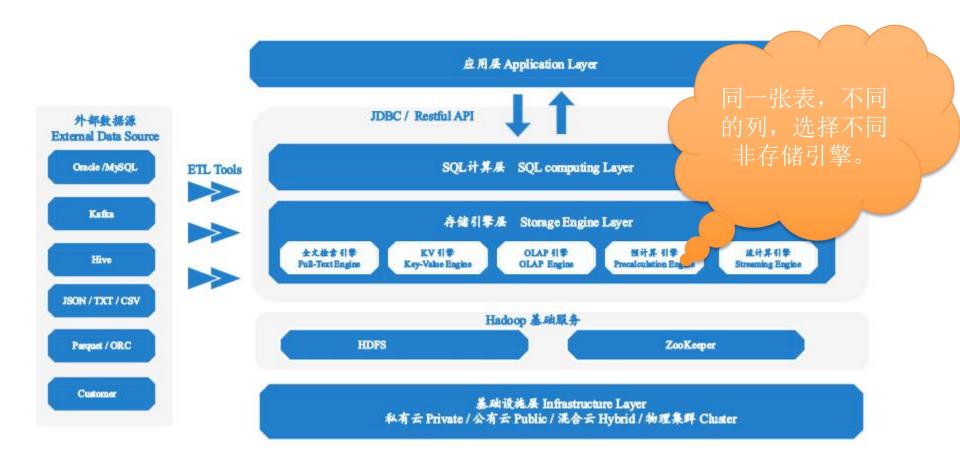


TERADATA.

解决思路



系统架构



分布式索引的存储->HDFS

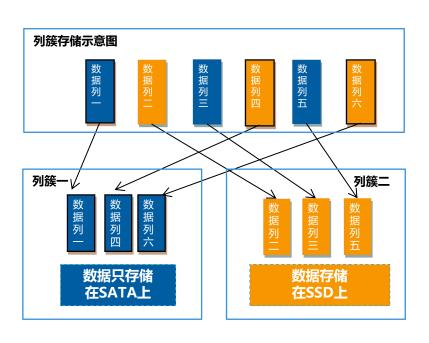
1: 相对于本地文件系统的优点! 2: 面临的主要问题,如何解决? 3: 单库跨多个联邦存储。

全文检索->ES场景

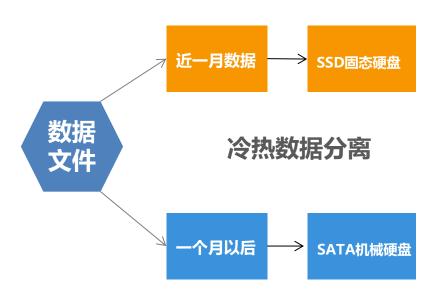


列簇+异构

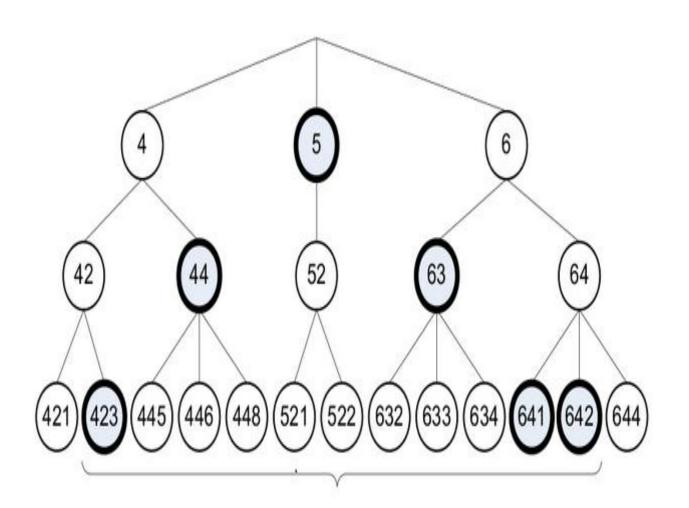
列簇存储



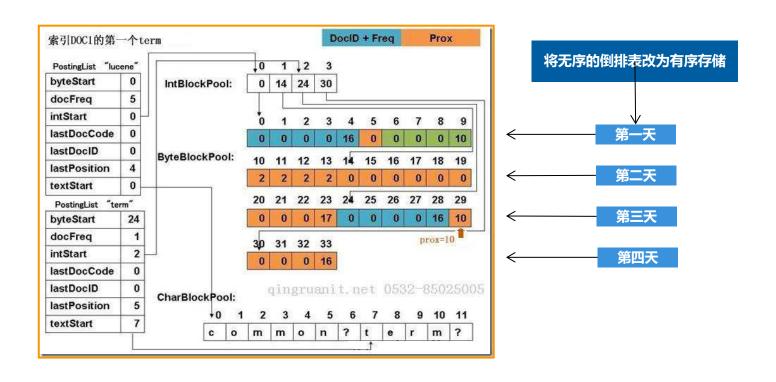
异构存储



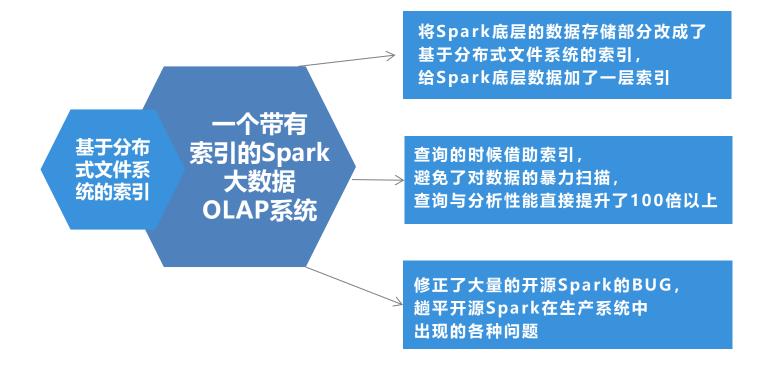
多层次索引-前缀与排序



倒排表数据按时序存储



计算框架->基于索引的Spark



>>> 与spark融合后,查询与统计分析功能更强大



支持SELECT语法:

group by, order by, case when, sum, max, min, avg, count

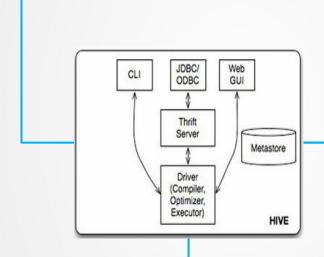
join[left join[right join]]、in、not in、like、not like、with子句、 union/union all、嵌套子查询等常见语法。

支持1000+个条件组合联合匹配查询。

支持丰富的函数操作:

支持数学函数(sin、cos、round、floor等)调用。 支持字符函数(substring、concat等)调用。

支持分析函数(row number、rank、lag、lead等)调用。



支持DDL语法:

create table drop table

支持DML语法:

支持数据插入insert语法 支持数据删除delete语法。

- 支持数据分区清理truncate table语
- 支持数据的大规模导出操作。

支持用户自定义(UDF)函数调用。

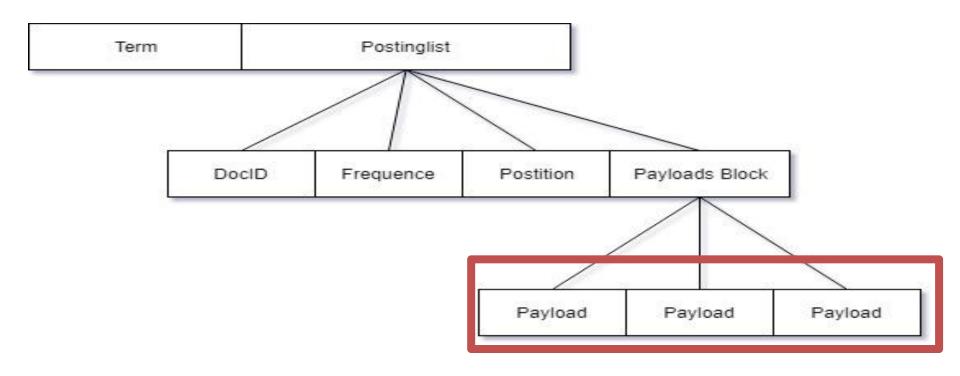
支持丰富的数据类型(string、long、int、double、char、 like、text、geopoint等),可以针对不同的业务场 景进行表数据类型设计:同时也可以支持自定义分 词数据类型。





- ◆ 对手机号进行1*的检索。
- ◆ 对有数据倾斜的列,进行多表关联。
- ◆ Limit 10000000000
- ◆ Partition by 一个不均衡的列
- ◆ 过载控制

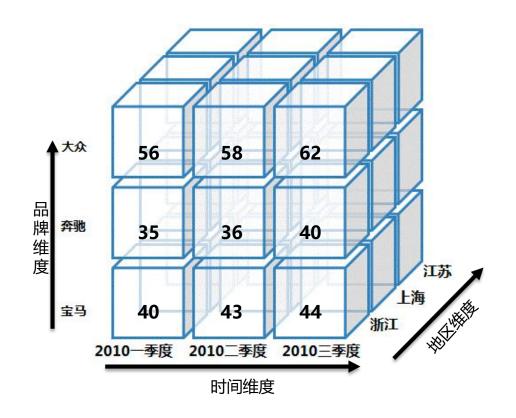
预计算索引->kylin+流计算



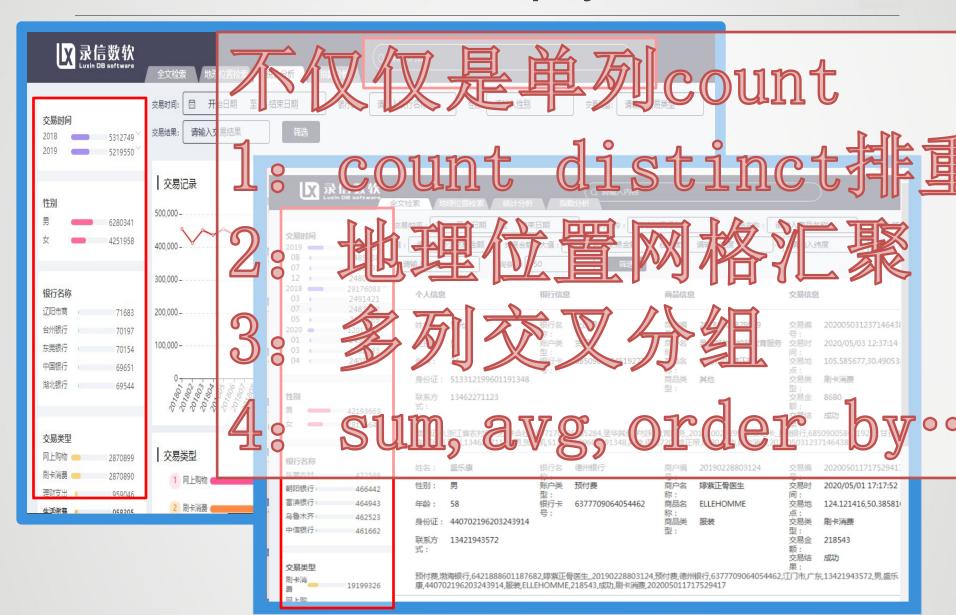
多维统计-> clickhouse场景

1: 95%以上的值为null 值。

2: 碎片化数据多维统计分析。



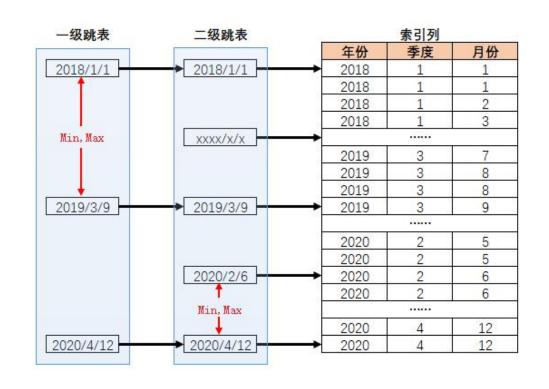
>>> 干人干面统计分析->vertical的projection



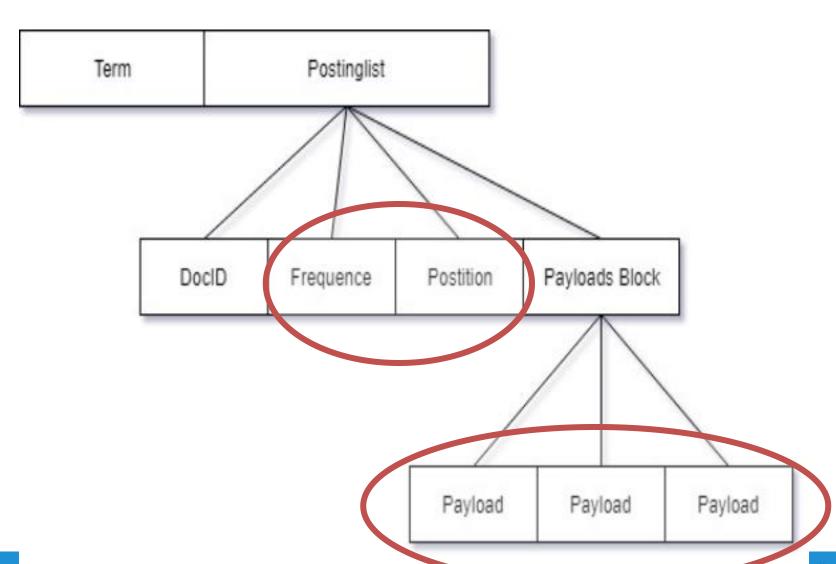
多维统计方案——多列联合索引

多维统计

- 1、每列之间采用列存储。
- 2、干预数据的排序分布,让列存储的 压缩更有效。
- 3、依据查询构造顺序读取。
- 4、多个列之间有层次关系。
- 5、结合分块存储。



Payloads压缩与按列存储-适合检索后的统计分析



山風粉提賀配管理峰会

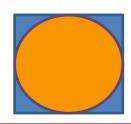
区域检索->数据预分布的变种



临近存储->数据预分布的变种

原生Lucene

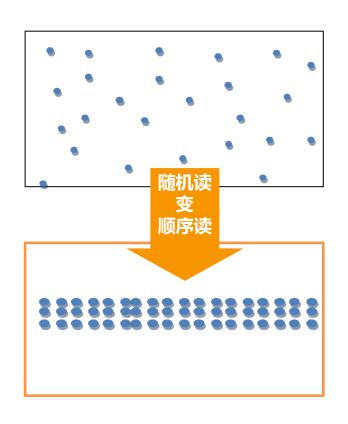
采用GeoHash选择正方形,再根据DocValues进行 二次验证裁剪。



黄色部分 需要剪切验证

优化方案

通过地理位置临近数据**临近存储**的方式构造硬盘上的连续读取,大幅度的减少随机读取的次数, 从而提升查询响应的速度。





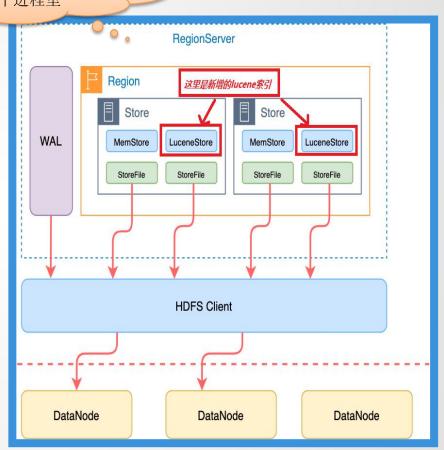
>>> 实现 Hbase 二级索引的方法



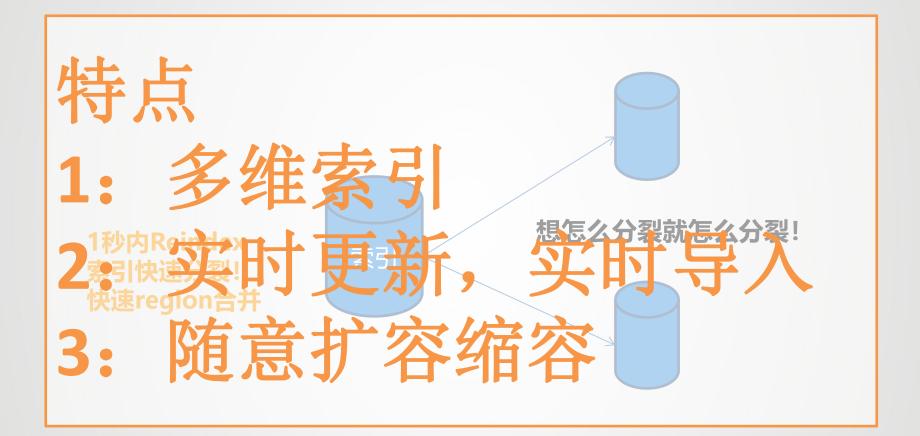
Executor与RegionServer嵌 在同一个进程里

HBase多维度查询

本身只提供基于行键和全 表扫描的查询,而行键索 引单一,需要采用HBase 的二级索引方案来进行多 条件的查询。







>>> 其他要实现的



1: 任务调度, IO调度。

2: 物化视图。

3: 索引加载问题。

应用场景——公安军队





场景描述

公安部门汇集了全网全维度的海量数据,包含互联网数据、社会数据、通讯数据等,通过<mark>实时检索、关联碰撞</mark>,为各警种提供智能研判的关系网络(同行、同 住同飞、同出入境等),大幅增加可用情报线索,提升侦破水平。

应用场景——汽车行业



车联网迅速发展,数亿T-BOX采集的<mark>万亿级别</mark>数据经过分析可让经销商、主机厂、国家监管部门多方受益。尤其随着国六标准出台,国家要对汽车尾气排放进行监管,需要大量采集油耗和排放数据。

录信与中国汽车技术研究中心(中汽研)合作,针对国六标准的推行,对商用车和乘用车的尾气数据、油耗数据、加油站地理位置信息进行监控,并提供整体存储和检索分析解决方案,为一汽解放、江淮汽车等主机厂的创新型应用提供坚实基础

Q&A



南京录信软件技术有限公司



中国数据智能管理峰会

THANK YOU!