



网易数据中台建设实践

网易大数据平台负责人 郭忆



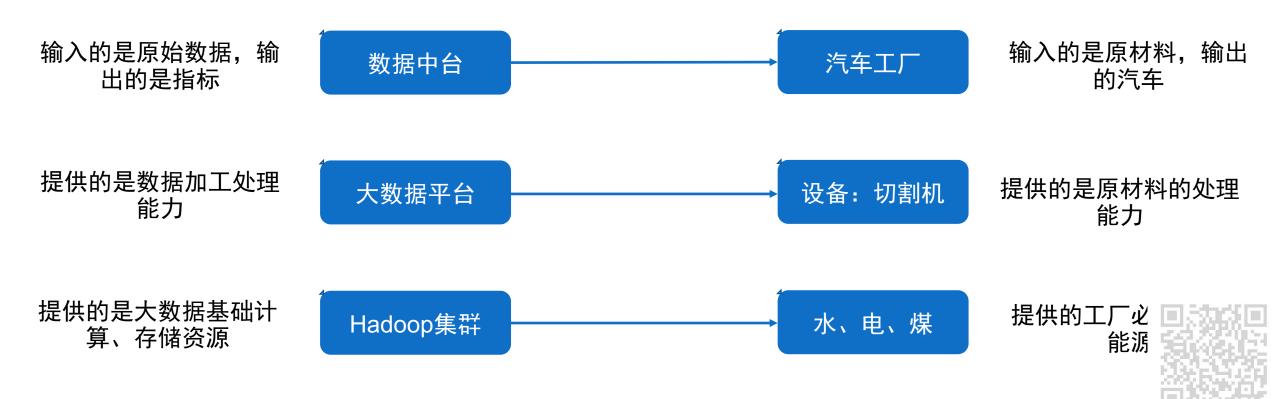
Agenda

- 1 什么是数据中台?
- 2 元数据中心:数据中台的基石
- 3 数据治理:效率、质量、成本
- 4 数据服务:数据中台的门户
- 5 数据中台治理效果

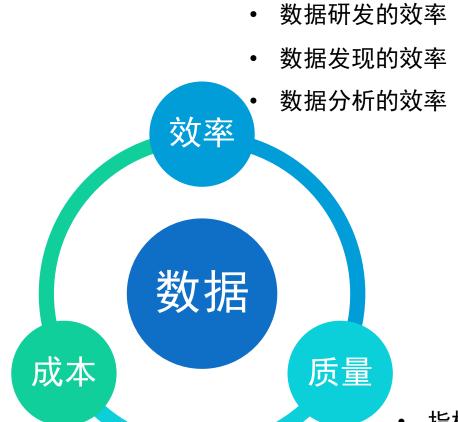


什么是数据中台

• 如果我们把数据中台比作一个汽车工厂



数据中台需要解决什么问题?



- 计算、存储资源成本
- 研发人力成本

- 指标一致性
- 数仓设计质量
- 数据质量



网易数据产品体系:以电商为例

市场运营 供应链 管理层 商品运营 业务场景 用户运营 供应链决 用户行为 商品运营 推广渠道 策协同系 高层看板 管理系统 分析系统 系统 统 数据产品 商品與情 用户精准 活动实时 Vipapp 投放系统 系统 直播

\$2517f25f25l48638200

网易在做数据中台前面临的挑战













全新IT技术AJI交流平台

数据中台支撑产品:网易猛犸

行业数据产品 行业业务系统 网易有数 网易大屏 自助分析 数据服务 数据服务 数据集成 数仓设计 数据开发 数据治理 运维安全 数据传输 指标系统 离线开发 成本治理 任务运维 数仓设计 智能报警 日志采集 实时开发 质量治理 权限中心 数据填报 数据测试 数据地图 埋点管理

产品特色

• "组件式"产品架构,业务可以根据发展阶段选择性搭配

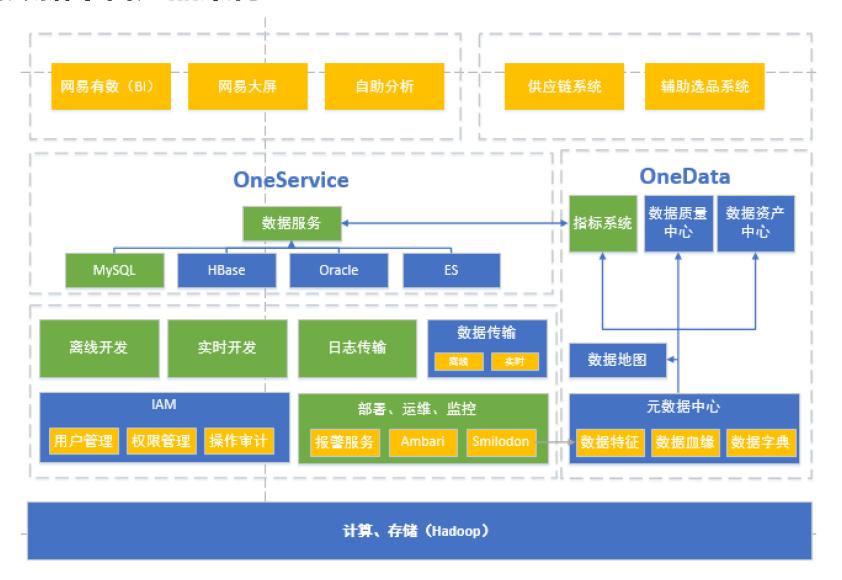
 "开放式"产品架构、聚焦核心通用产品、同时开放基础能力、 允许业务集成新的产品

• "轻型易用"平台,通过"增强分析"降低用户使用的门槛

• 完美的支撑数据中台建设,减少重复建设,提高数据共享能力



网易数据中台产品架构



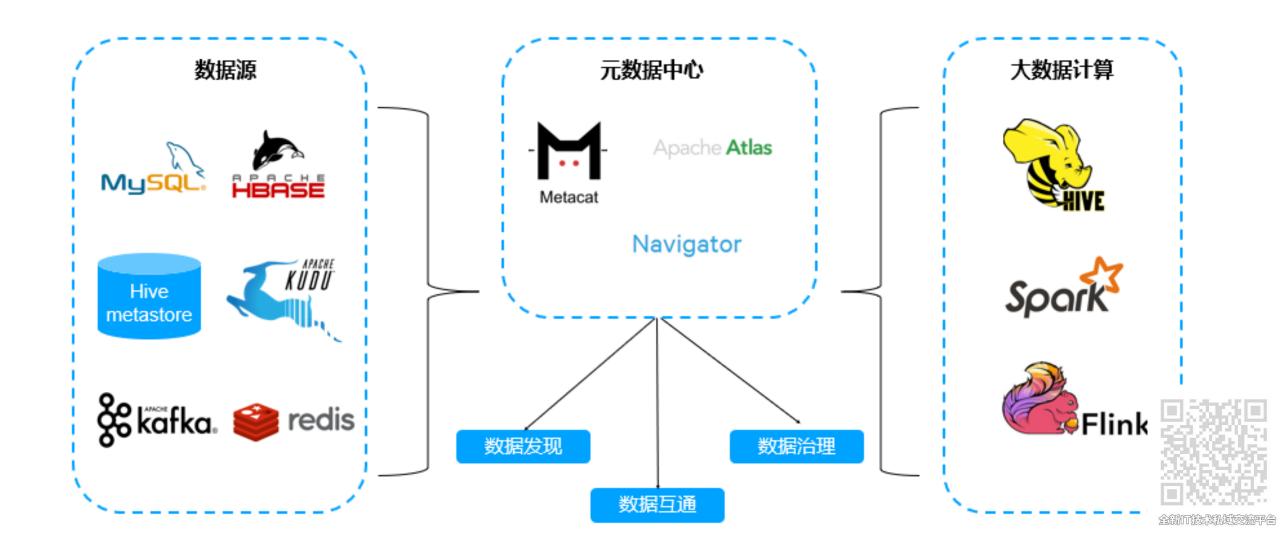


Agenda

- 1 什么是数据中台?
- 2 元数据中心:数据中台的基石
- 3 数据治理:效率、质量、成本
- 4 数据服务:数据中台的门户
- 5 数据中台治理效果



元数据中心:数据中台基石



元数据中心

数据标签

- 通过丰富的不同类型的标签,完善数据特征体系
- 指标标识、数仓的主题域、分层信息,是否是数仓维护的推荐表都以标签形式存在

多租户, 多业务线

能够支持电商(考拉、严选)、互娱(音乐、 游戏)、传媒、教育



多种数据源支持

能够覆盖网易所有的数据源,甚至包括 Kafka, Redis, Hbase等Schema less KV系 统

与大数据系统集成

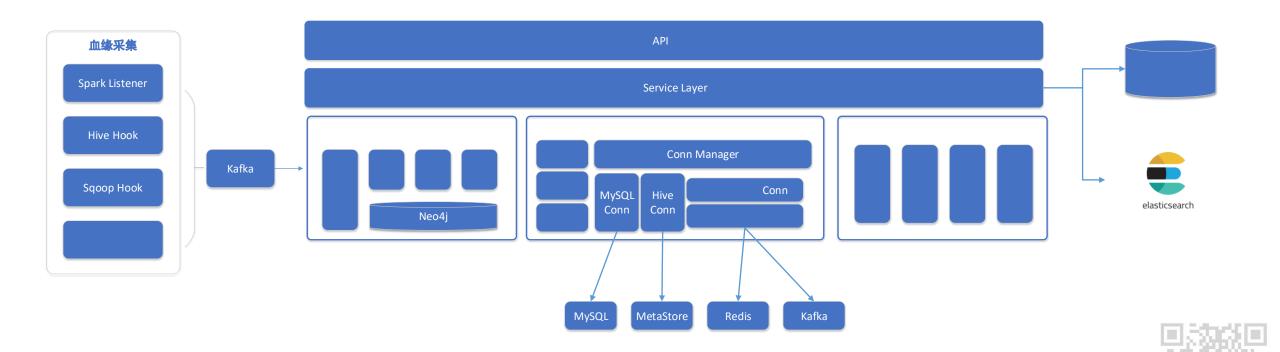
- · 与Ranger结合,允许通过自定义标签的方式对数据进行动态授权
- 数据传输、自助分析与元数据中心集成
- 基于元数据中心、构建数据质量中心、数据资产 管理中心、数据地图

数据血缘

- · 静态血缘 & 动态血缘
- · 血缘支持时间戳,可以按照时间戳读取,过期
- · 血缘覆盖率以及血缘采集性能

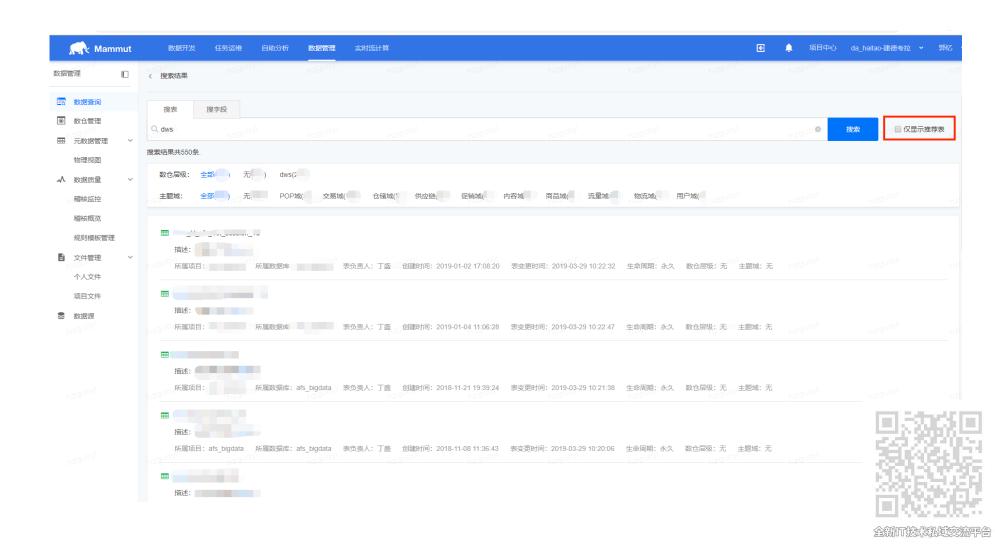


元数据中心



数据地图

解决"有哪些数据可用?","到哪里找数据?"



Agenda

- 1 什么是数据中台?
- 2 元数据中心:数据中台的基石
- 3 数据治理:效率、质量、成本
- 4 数据服务:数据中台的门户
- 5 数据中台治理效果



如何评价一个数仓设计好坏?



- 大量的表没有明确的主题域、业务过程,分层信息,数仓组织混乱
- 超过50%的任务直接引用ODS 层原始数据,30%的表存在跨层引用,DWD建设完善度较低
- DWS 层表复用性差,平均表引用系数低
- 依然有查询ODS 层原始数据的Query, DWS, ADS Query 覆盖率低, 取数效率差
- 表、字段命名规范混乱,数据发现困难



规范化数仓设计

数仓设计度量

EasyDesign

规范化管理

团队协作

- 各层表的分布以及各层被下游表和任务 引用情况, Query 查询覆盖率
- DWD: ODS 被跨层引用的表的数量
 DWD 平均被下游表引用系数
 ODS 被Query 查询情况
- DWS: DWS 平均被下游表引用系数
 DWS Query 覆盖率
- 度量管理
- 维度管理
- 基础字典管理
- 模型设计

审批流程



数仓升级的目标

覆盖度

- 消灭ADS/DWS 直接 引用ODS 层原始数 据
- 消灭Query直接查询 ODS 层原始数据
- DWS/ADS Query 覆 盖度上升

复用性

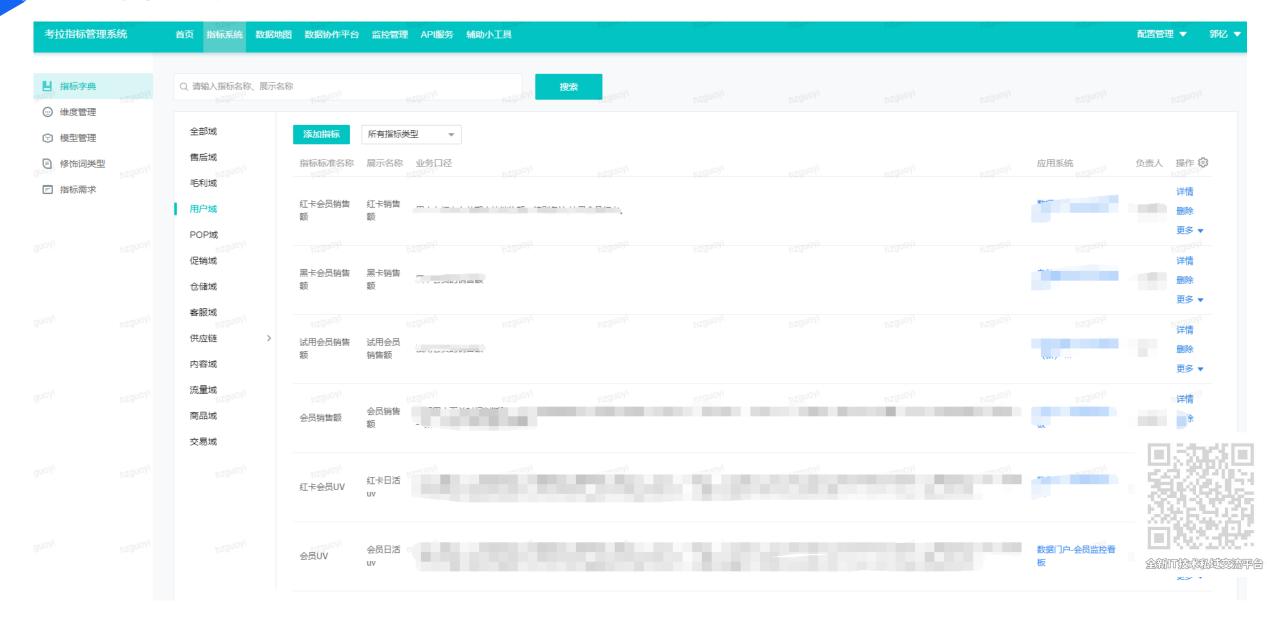
• DWS、DWD 平均每 张表被下游表引用数 量增加 规范性

- 表、字段命名规范统
- · 建表流程审核

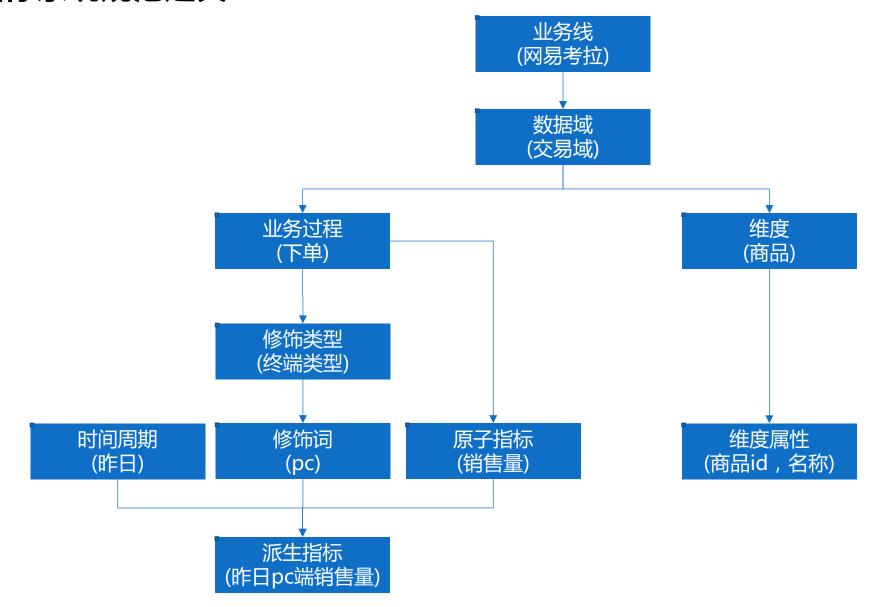
需求响应速度提升、查询速度提升、查询成本降低,数据使用者满意度提升!



指标系统



指标系统规范定义

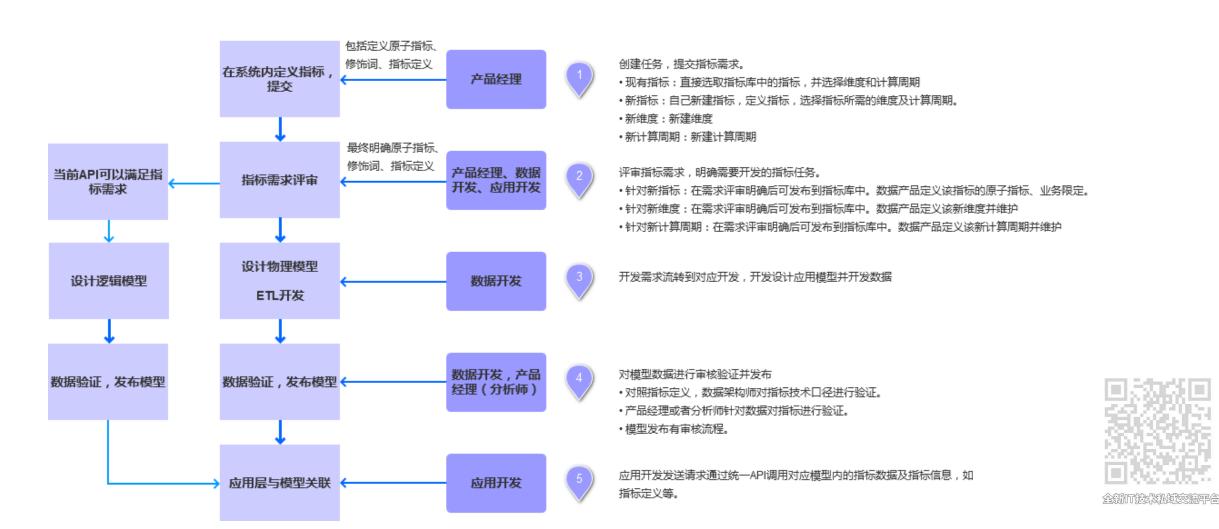




全新T技术规划交流平台

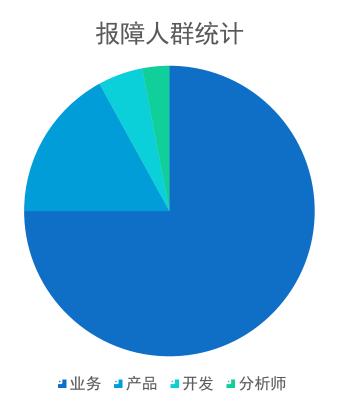
指标开发实施规范

指标开发工作流程

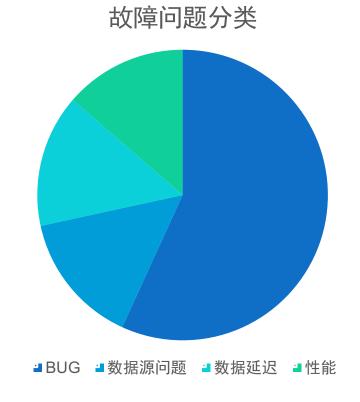


在做数据质量前业务面临的现状

• 超过90%的问题是由业务和产品发现



• 收集的问题中存在研发bug的占比超过50%





数据质量方法论

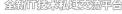


完整性

是指需要的数据已完整记录,可以分为记录数完整性和字段值完整性。

准确性

是指数据在数仓中的值和实际值是否相同,可以分为口径实现结果和数据逻辑合理性。



数据稽核监控规则

完整性

- 表数据量波动监控和绝对值监控
- 主键唯一监控
- 字段为空,为0的监控
- 数据完整性监控,订单24小时,终端覆盖

一致性

- 同一个指标在不同模型不一致监控
- 相关指标趋势监控,比如uv和pv走势一致
- 聚合逻辑一致性监控
- 不同数据源对同一个实体的值一致

准确性

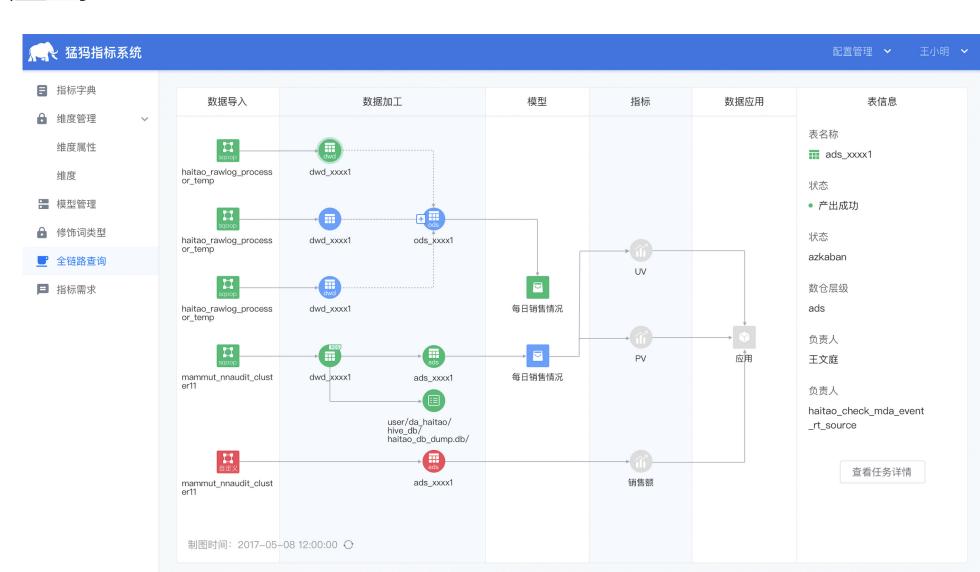
- 数值计算逻辑监控: 商品只能归属一个BU
- · 数据格式的监控,比如IP,URL
- 码表的监控
- 数据异常监控,比如日期还没发生

时效性

- 任务延迟监控
- 表产出时间的监控
- 源数据延迟监控

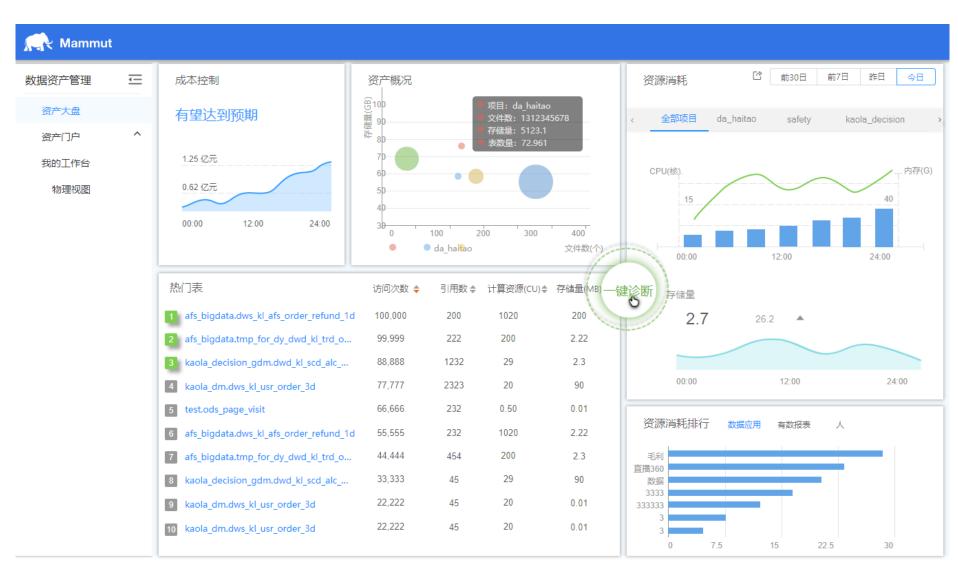
全链路数据质量监控

- 覆盖数据产出的完整 生命周期
- 全链路数据血缘的实时监控
- 快速了解哪些数据产 品的哪些指标异常
- 故障恢复时间的请准 预估

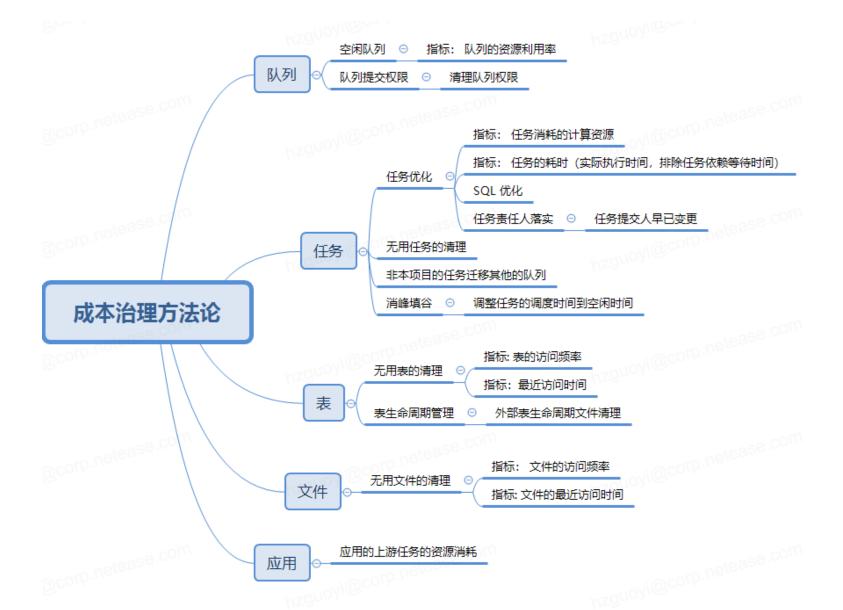


成本治理

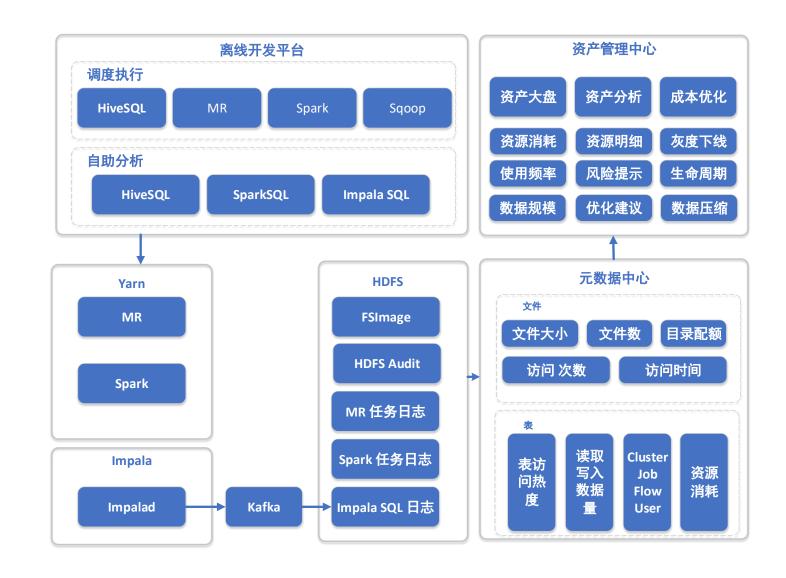
- 每个项目花了多少钱?
- 每张报表花了多少钱?
- 每个数据产品花了多少 钱?
- 每个人花了多少钱?
- 各个业务线预算符合度?
- 根据表的热度,存储空间、加工表消耗的资源, 确认表是否可以优化? 给出优化建议?
- 表的一键下线



成本治理



系统架构



Agenda

- 1 什么是数据中台?
- 2 元数据中心:数据中台的基石
- 3 数据治理:效率、质量、成本
- 4 数据服务:数据中台的门户
- 5 数据中台治理效果

数据服务

指标口径

- 相同指标在多个表中存在
- 指标口径不统一
- 指标复用率低,指标重复加工

效率

- 所有需求一个接口
- SDK 接入效率高

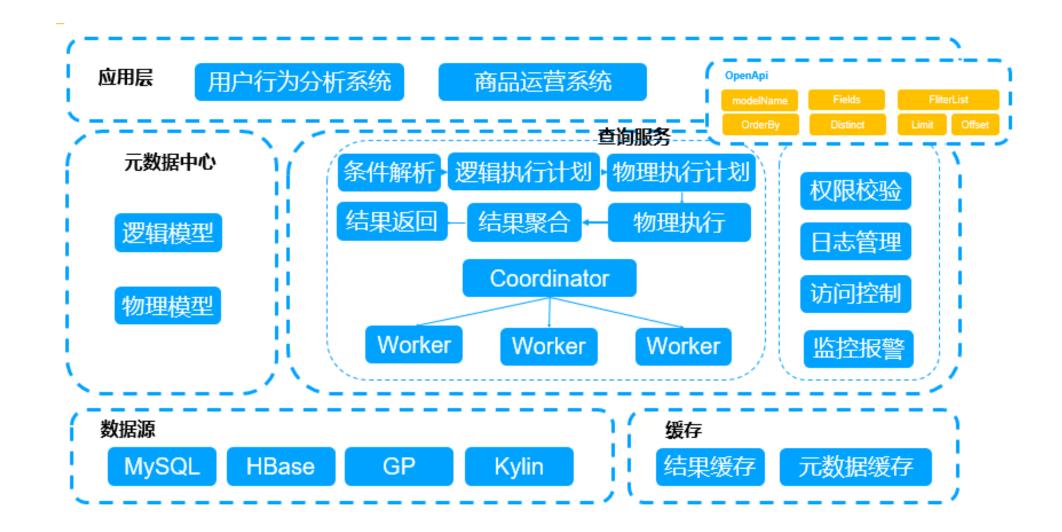
安全隐患

- 底层表直接暴漏
- 所有对数仓的访问权限控制

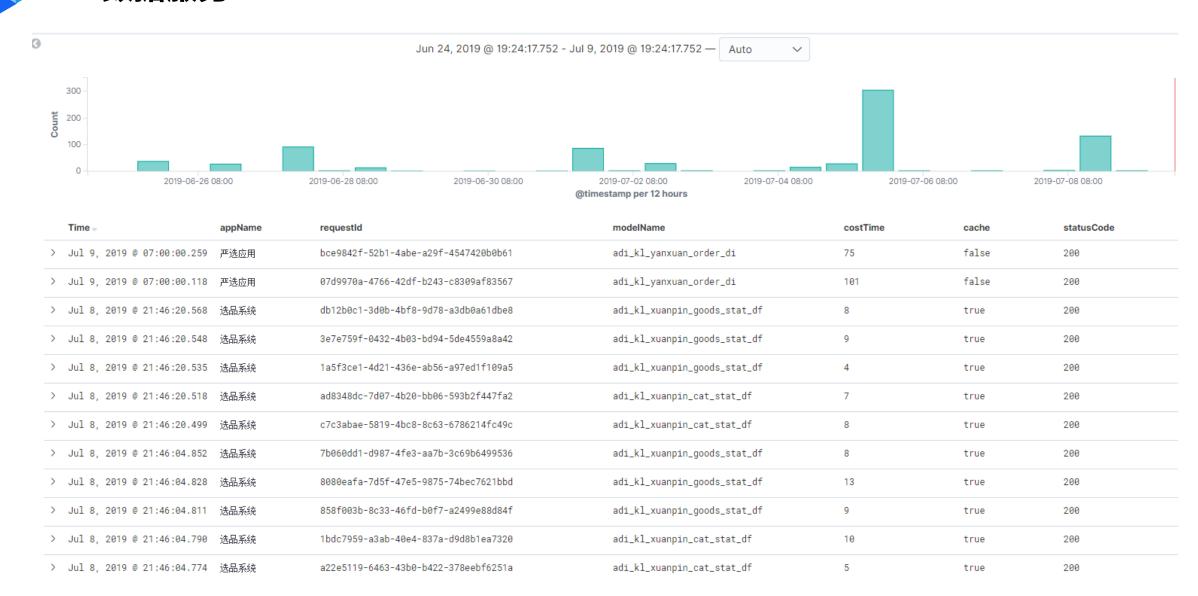
维护性

- 打通了数据应用-指标-数仓表的全链路监控
- 底层表变更,需要应用层修改

数据服务



数据服务



Agenda

- 1 什么是数据中台?
- 2 元数据中心:数据中台的基石
- 3 数据治理:效率、质量、成本
- 4 数据服务:数据中台的门户
- 数据中台治理效果

数据中台的实施价值

100% 数据产品的 指标通过系统化管 理,有明确的业务 口径、计算逻辑、 数据来源,完全消 除指标二义性

基于数据地图,实现100%自助取数,取数效率提升300%!分析师满意度提升!

所有数仓维护的表都有明确的分层、主题域,业务过程,ODS 层引用任务降低,100%消灭ODS 层Query,DW层引用任务增长,DWS/ADSQuery 99.9%覆盖,表的引用系数提升

全链路数据跟踪,回答"数据准不准?""哪些数据故障?""什么时候恢复?",加速数据故障的排查定位,助力99.8%SLA达成

通过全链路资产分析,消灭低价值的资产,为业务节省20%的资源成本

总结

