



贝壳一站式大数据开发平台实践

人工智能技术中心@贝壳找房

日期：2020年9月

个人简介

仰宗强

- 贝壳找房 人工智能技术中心-大数据平台
- 趣店 数据应用
- 百度 流量平台
- 当当 数据平台
- 西安科技大学 学士

目录

01 背景介绍

02 探索历程

03 整体介绍

04 总结与展望

背景介绍-贝壳数据业务

- 将贝壳数据 分为物的数据、人的数据、行为数据三大块来进行研究。



物的数据

- 楼盘字典：2亿+套房屋信息



行为数据

- 线上行为
- 多样的线下行为



人的数据

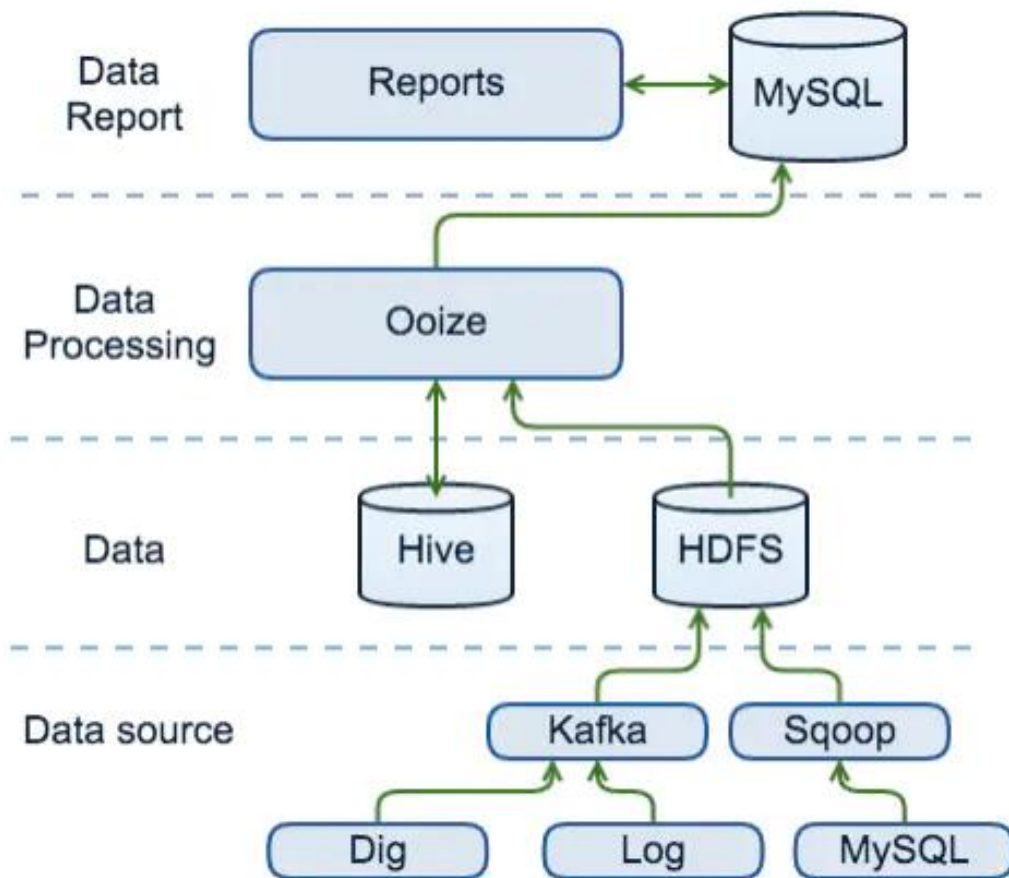
- 买家 -> 上亿级别
- 业主
- 经纪人 -> 40万+
- 品牌

- 为了满足各类数据获取/计算等场景，业内涌现出了很多解决方案

- 降本（降低数据处理成本）
- 增效（合理提高数据使用/计算效率）
- 规范化（提供统一的编程范式）



探索历程-初期的架构



优点

1. 开源组件，方便扩展和运维
2. 业界成熟的数据仓库方案，分层模型设计
3. 有利于技术人员培养

缺点

1. 需求迭代跟不上节奏
2. 数据仓库工程师累成狗。
3. 集群作业运维困难
4. 数据泄露



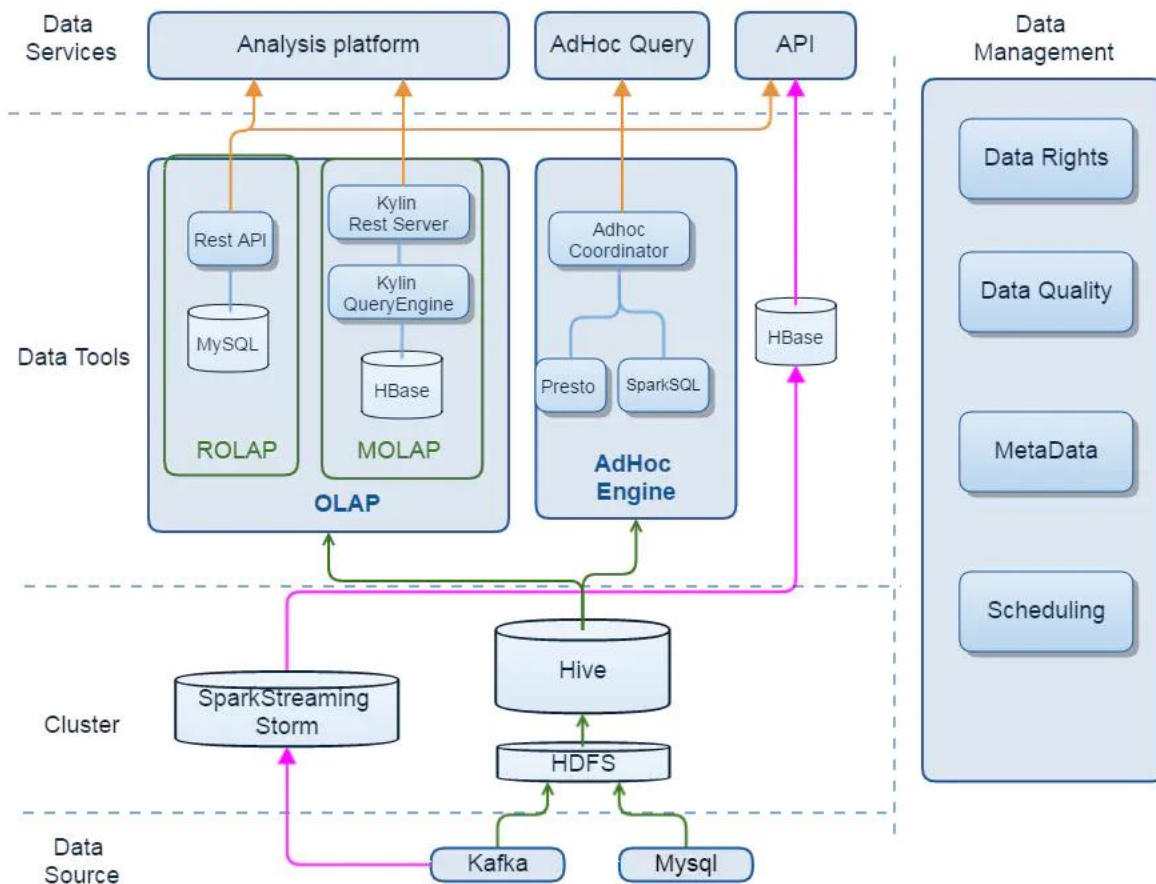
探索历程-平台化架构

优点

1. 解决了数据仓库的瓶颈问题。
2. 业务数据产品可以直接使用数据服务平台提高效率。
3. 自助运维，快速定位

缺点

1. 数仓表/任务量剧增，造成集群资源相当紧张。
2. 集中的几个人 -> 分散到业务，质量无法把控



探索历程-历程回顾



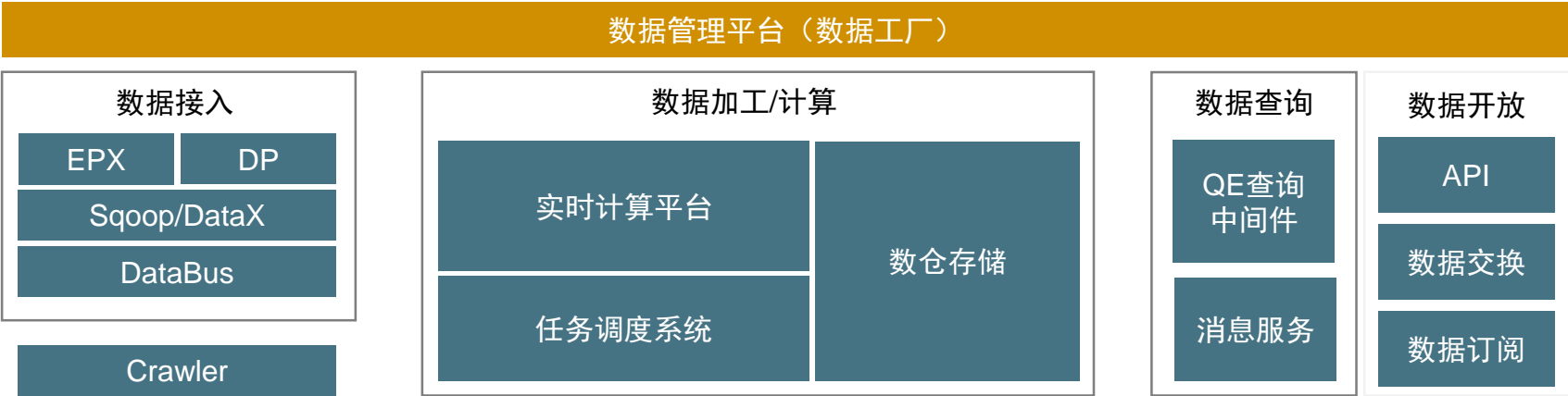
平台架构



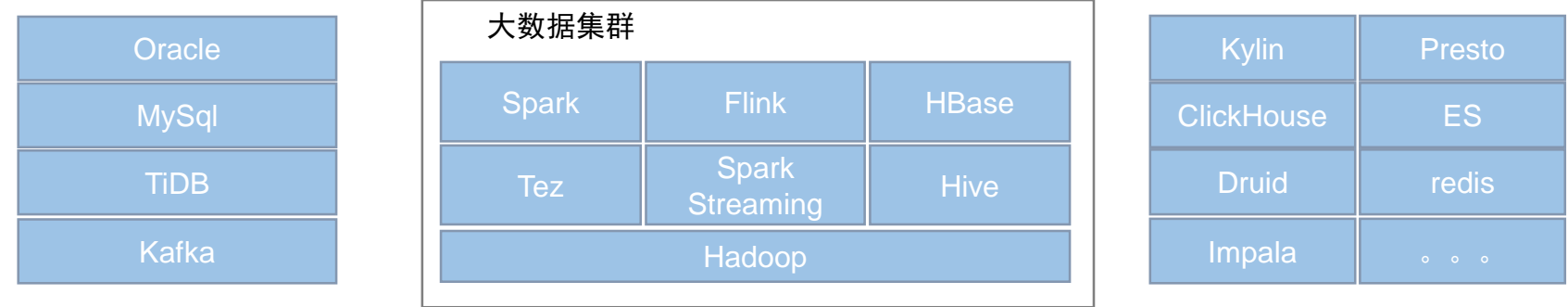
数据视角



平台视角



基础设施



平台介绍-功能划分



平台介绍-1. 数据管理

数据管理：以**元数据为基础**来构建贝壳数据资产的统一管控平台

元数据服务

能力开放

对内应用

对外价值释放

向上对内进行数据资产管理，对外提供价值开放

元数据中心

技术引领

技术

库、表、字段
存储位置 ...

业务

产品线、分类
数据域 ...

管理

负责人、审批人
安全等级 ...

血缘

数据血缘、
字段血缘...

资产化管理

贯穿全域

能力化输出

统一元数据模型

向下完成4大方向，全域元数据汇聚

数据源

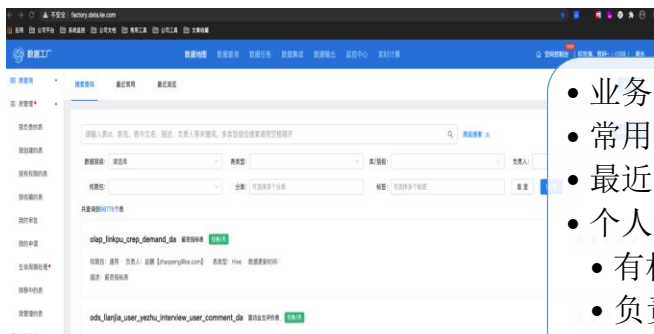
数据获取与存储

新数据

数据分析与应用

产品与服务

平台介绍-1. 数据管理

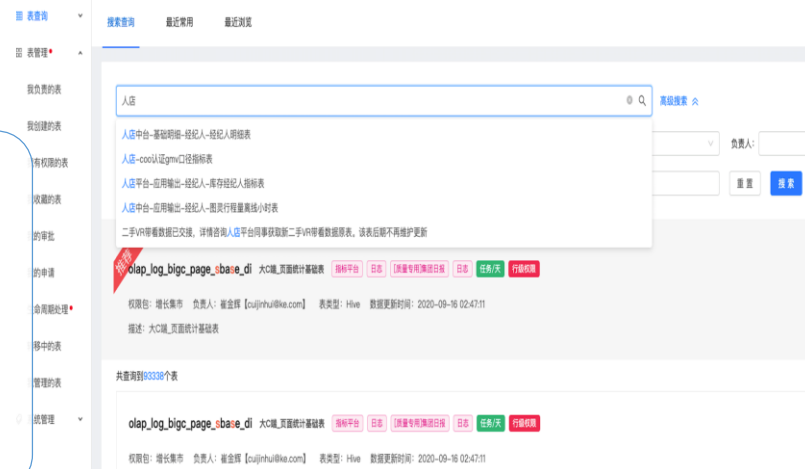


- 业务图谱
- 常用的表
- 最近浏览的表
- 个人管理模块
 - 有权限的表
 - 负责/创建的表
- 申请/审批
- 表交接转移

数据
导航

智能
搜索

- 根据表中文搜索
- 根据表英文名搜索
- 根据描述、标签、负责人搜索



建表类型

创建数据表是一种在数据仓库中进行建表开发的操作，本权限仅对总部部分产品研发和BI同事开放，**请相关同事勿申请。**

Hive

适用于多种离线数据分析，如数据ETL、数据存储管理、大型数据集的查询和分析等，地图上多为hive表。

ClickHouse

适用于OLAP实时分析查询，如实时指标构建、罗盘用户行为分析、用户画像人群圈定等。

Kafka

适用于各种日志的存储，数仓的存储等，对各个数仓层级间进行解耦。

- Hive
- CK
- 物理表
- 视图
- Kafka

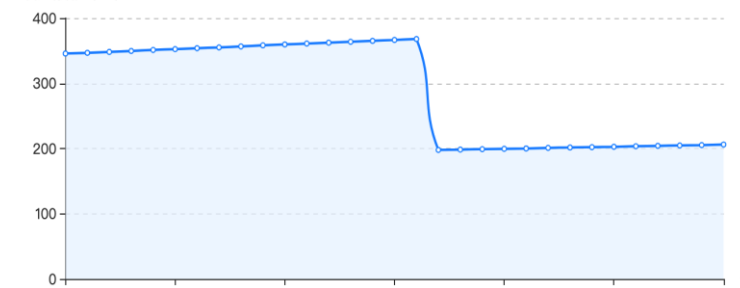
数据
录入

资产
管理

- 数据生命周期管理
- 空间信息展示
- 空间资产管理分析
- 团队协作



已占用存储量(GB)



平台介绍-2. 数据集成

数据集成最终目标是交付一套有效提升集成效率和覆盖度的工具

数据接入效率&覆盖度

多元集成

支持多种数据类型的集成，满足公司99%以上业务

能力提升

自动化：新建库/表自动接入（DB）

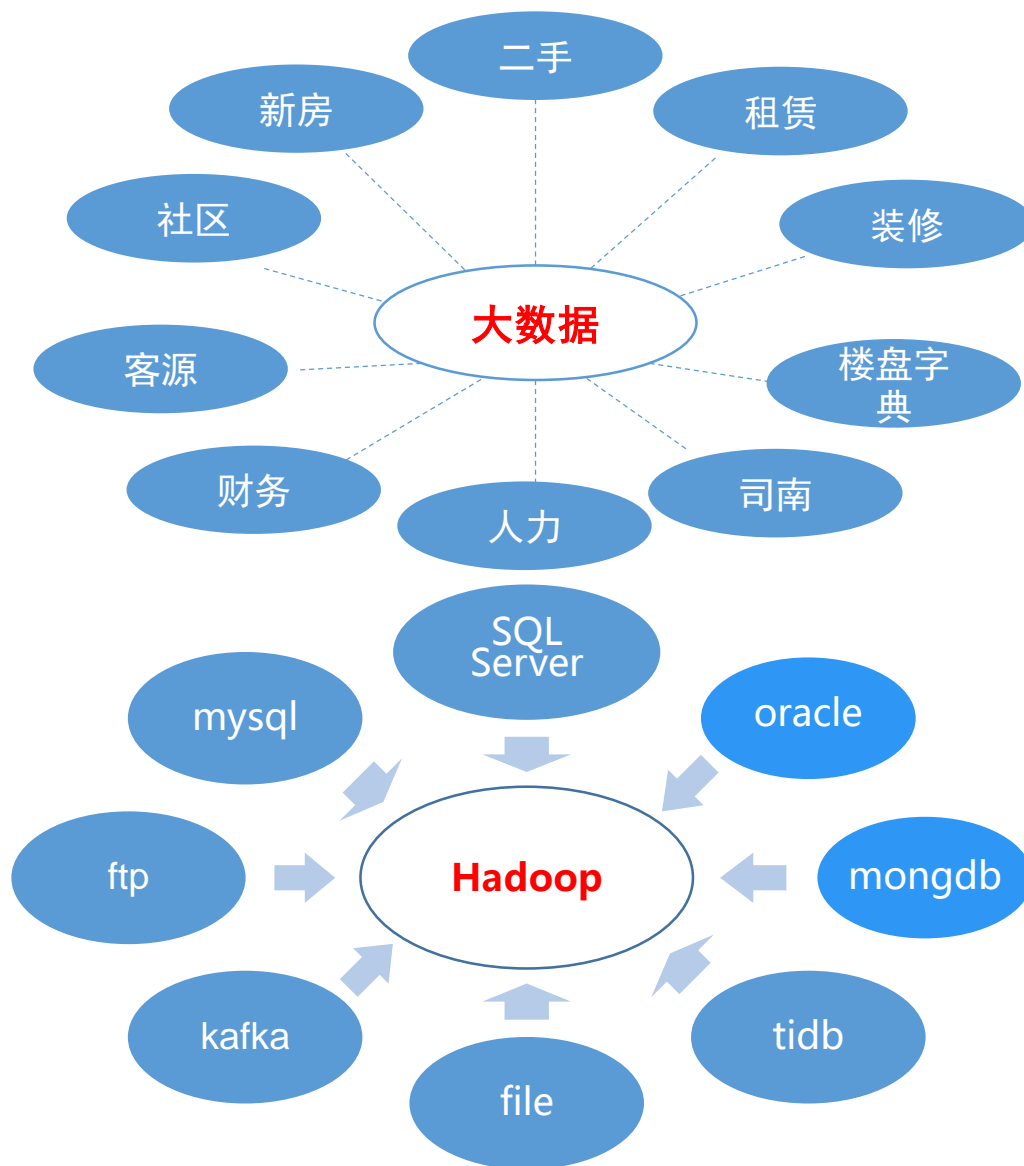
迁移提效：整库/表迁移自助

接入速率：T+1时效

长效保障

变更自动同步&提醒

接入任务监控机制保障



平台介绍-3. 作业调度

创建表

基本信息

字段信息

数据逻辑

任务调度

版本测试

数据上线

1 基本信息

2 字段信息

3 任务调度

基本信息

* 数据层级:

请选择



* 集群名:

hive集群



* 库/层级:



* 表名称:

表名称 (小写英文字母)

* 表中文名称:

表中文名称

* 权限包:



分类:

请选择



* 负责人:

请输入负责人



* 安全等级:

C级



* 审批人:

请输入审批人



* 所属空间:



* 描述:

请输入描述

平台介绍-3. 作业调度

创建表

基本信息

字段信息

数据逻辑

任务调度

版本测试

数据上线



基本信息

2

字段信息

3

任务调度

添加字段方式:

可视化添加

语句添加

序号	* 字段名	参照表	参照字段	* 字段中文名	描述	枚举值		操作
1	pt	请输入参照表 ▼	请输入参照字段 ▼	分区字段	hive分区标识		stri	删除

添加新字段

上一步

仅创建表

创建并配置调度

平台介绍-3. 作业调度



数据源

* insert语句:

1	
---	--

平台介绍-3. 作业调度



基本信息



字段信息

3

任务调度

* 任务周期:

请选择

☐ crontab形式

* 任务有效期:

2020-10-13 19:00:52



2022-10-13 00:00:00



* 触发方式:



时间触发 (按照crontab指定时间运行)



依赖触发 (上游依赖就绪时, 自动触发)

重试次数:

3

任务并发数:

1

* 启动超时:



设置启动超时报警



不设置启动超时报警

* 运行超时

任务运行时长超过 120

分钟后即可推送报警

* 任务队列:



* 需求方:

请选择需求方



依赖数据:

添加依赖

推荐配置

任务失败报警通知设置:

平台介绍-3. 作业调度

创建表

基本信息

字段信息

数据逻辑

任务调度

版本测试

数据上线

基本信息

字段信息

分

序号

版本号

1

215105

【yangzongc

请选择测试内容：

X

☐ SQL执行测试☒ 数据准确性测试

请选择数据验证点：

☒ 表主键、复合主键是否唯一☒ 维度表维度值是否缺失☒ 目标表于源表数据量是否一致源表名：☒ 表字段是否存在空值☒ 枚举值分布是否正常分区：☐ 测试通过后自动上线

取消

开始测试

评论

收藏

已有限

测试报告

操作

开始测试 强制上线 版本记录

平台介绍-3. 作业调度

支撑**批计算任务**的管理、分发和执行，充当数据全链路中的心脏中枢



效果

- ◆ 任务配置简单，在页面上简单配置一个表单即可操作。
- ◆ 提供常用的ETL组件，零编码
- ◆ 依赖关系可视化，一键修复追溯，将排查问题修复数据的时间由一人天缩减到10分钟。
- ◆ 智能调度，错峰运行，保证高优先级任务优先运行。

平台介绍-4. 数据质量

保障数据资产健康，促进数据资产流通

数据准确性痛点

- 数据修改无法测试
- 数据例行无法监测

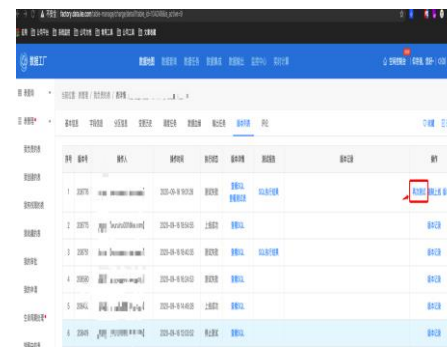
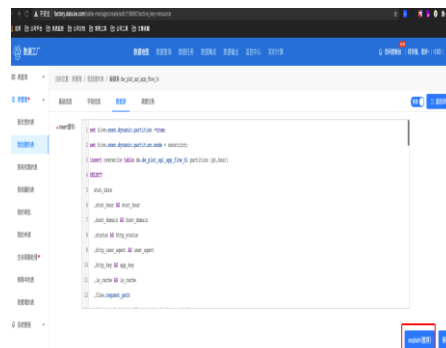
解决方案

数据及时性痛点

- 例行任务失败、延迟无法感知
- 任务上游下线无法感知
- 怎么保障任务及时产出

1

设计完善的开发流程，覆盖研发-测试-变更的事前全流程，识别并拦截潜在风险



2

提供完善的任务监控服务，保证任务问题及时预警

- 通过任务失败报警，将失败的任务及时通知用户
- 通过超时报警，保证用户在规定时间内未成功能够及时感知
- 通过表/任务链路报警，保障业务链路能够正常及时产出

3

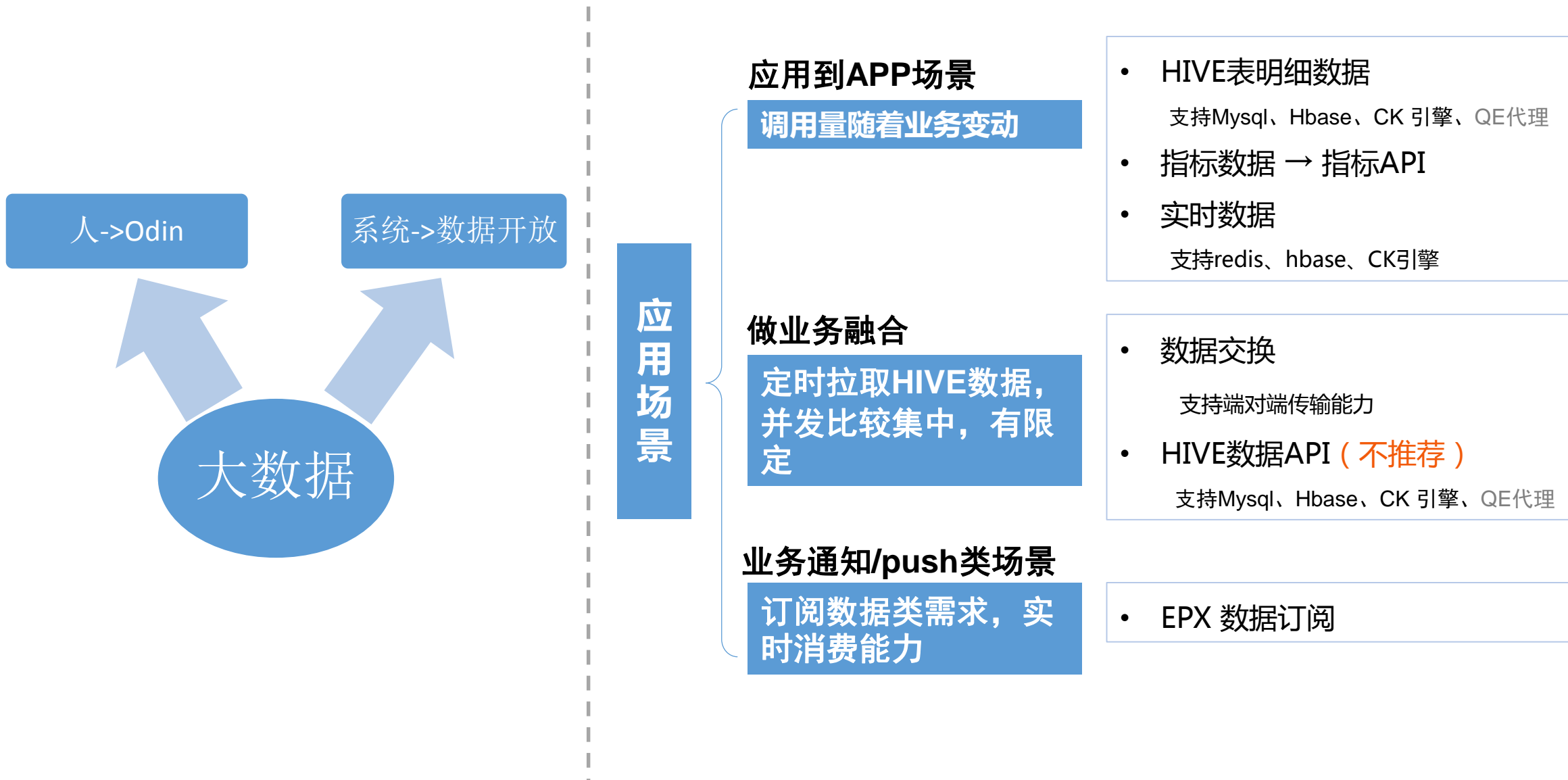
提供多样性的数据质量监控，保障数据资产的准确性



- ☐ 字段变化报警: 字段数量增减报警
- ☐ 采集完整性报警: 采集完整性
采集波动阈值 与业务系统系数浮动 30 %报警
- ☐ 及时性报警: 与该表近 7 次产出时间增加 60 分钟报警
☐ 未在活跃时间内产出进行报警
☐ 是否是核心任务
- ☐ 数据量波动检测: 与该表近 15 次均值增加 30 %报警
与该表近 15 次均值减少 30 %报警
☐ 与上一次比较

平台介绍-5. 数据开放

提供 **HIVE数据API**、**指标API**、**数据订阅**、**数据交换**、**QE代理** 等各种能力和大数据输出的体系化解决方案。



平台介绍-5. 数据开放

创建API

基本信息

请求信息

返回信息

排序分组

Sql预览

生成上线

数据工厂

数据地图

数据查询

数据任务

数据集成

监控中心

实时计算

数据开放

NEW
空间控制台 | 仰宗强, 您好~ | 退出

API项目

API管理

API测试

当前位置：数据开放 > API管理

使用帮助

API名称/地址:



负责人:



API分类: 全部



表名:

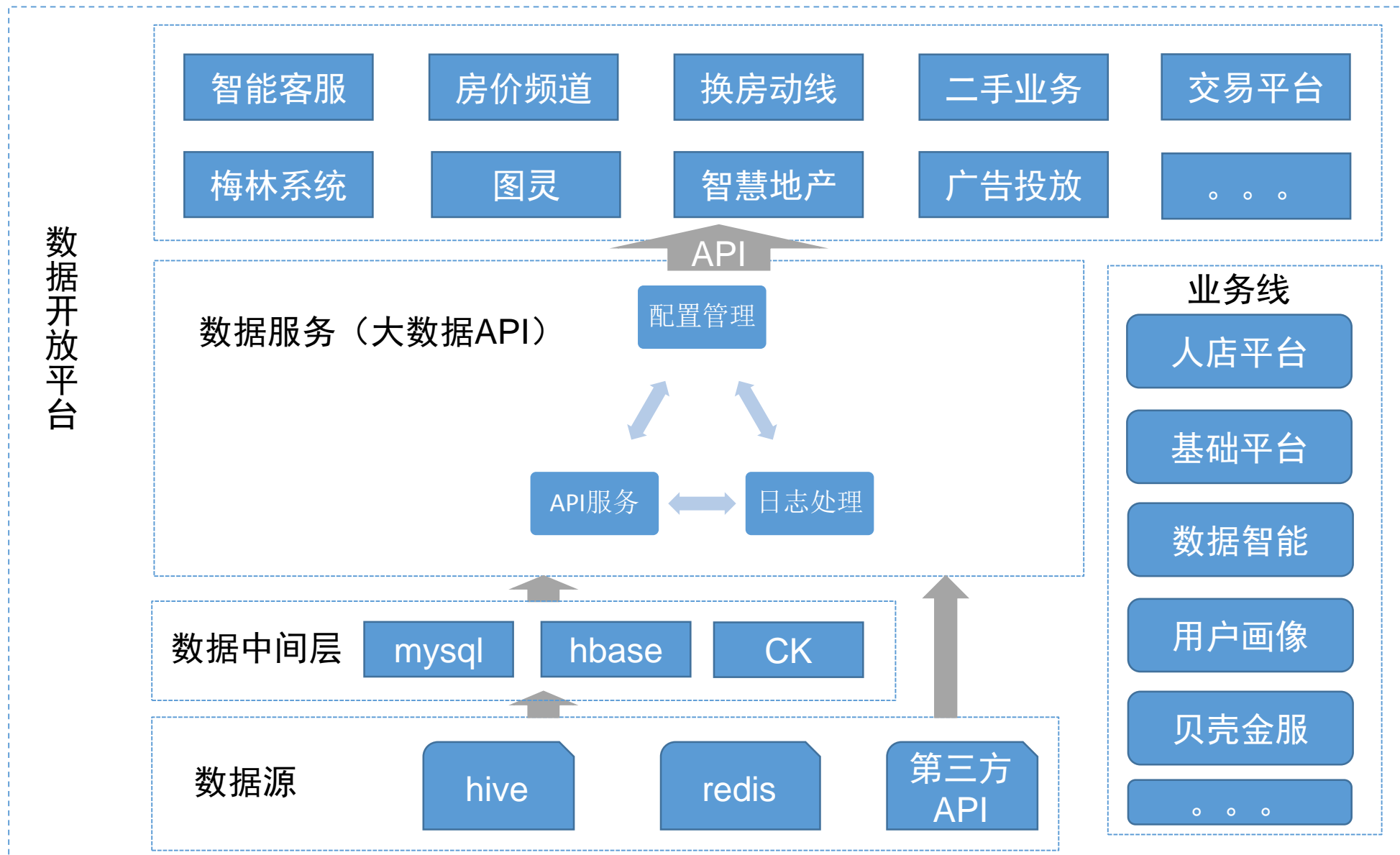
重置

搜索

新建

API_ID	API名称	所属项目	状态	API类型	数据源	负责人	更新时间	操作
6263	test_1233_qa meta/sssette	open专用	● 已上线	公有	mysql	原玉娇 【yuanyujiao001@ke.com】	2020-12-20 16:20	编辑 测试 ...
6262	test1222 meta/ssssssssssssssssssss...	open专用	● 已上线	公有	mysql	原玉娇 【yuanyujiao001@ke.com】	2020-12-20 16:20	编辑 测试 ...
6261	agents_ca_info meta/agents_ca_info	签单小诸葛	● 已上线	公有	hbase	张毅 【zhangyi138@ke.com】	2020-12-16 16:20	编辑 测试 ...
6260	getBizAndDistricts meta/getBizcircleAndDistri...	商业报买项目	● 已上线	公有	mysql	邢光辉 【xingguanghui002@ke.com】	2020-12-11 16:20	编辑 测试 ...
6259	花桥B1认证门店查询 meta/huaqiao/cert/shop	花桥学校项目	● 已上线	公有	mysql	顾贤斌 【guxianbin001@ke.com】	2020-12-16 16:20	编辑 测试 ...
6258	WI分_主行程API_租赁_V5 target/wi/v5/mainStrokeRent	链家支持中心WI分	● 已上线	公有	proxy	牟甜甜 【mutiantian@ke.com】	2020-12-13 16:20	编辑 测试 ...

平台介绍-5. 数据开放



总结与展望

- ◆ 数据资产化管理，全链路数据追踪和分析，提升数据价值
- ◆ 加密、脱敏、敏感监控等多种安全策略，全方位保障数据的存储、访问、传输过程
- ◆ 沉淀一系列的技术能力和组件集合，构建共性能力和通用服务，打造企业级大数据平台



THANK YOU