

The MEIZU logo is displayed in a blue, sans-serif font. It is positioned in the top right corner of the slide, to the left of the 'IT大咖说' logo. The background features several colorful, rounded rectangular shapes in shades of green, orange, and red, scattered across the top and bottom edges.The 'IT大咖说' logo is located in the top right corner. It consists of the text 'IT大咖说' in a bold, blue font, with 'IT' in a smaller size. Below it, the text '知识共享平台' is written in a smaller, grey font. The logo is set against a white background with a thin grey border.

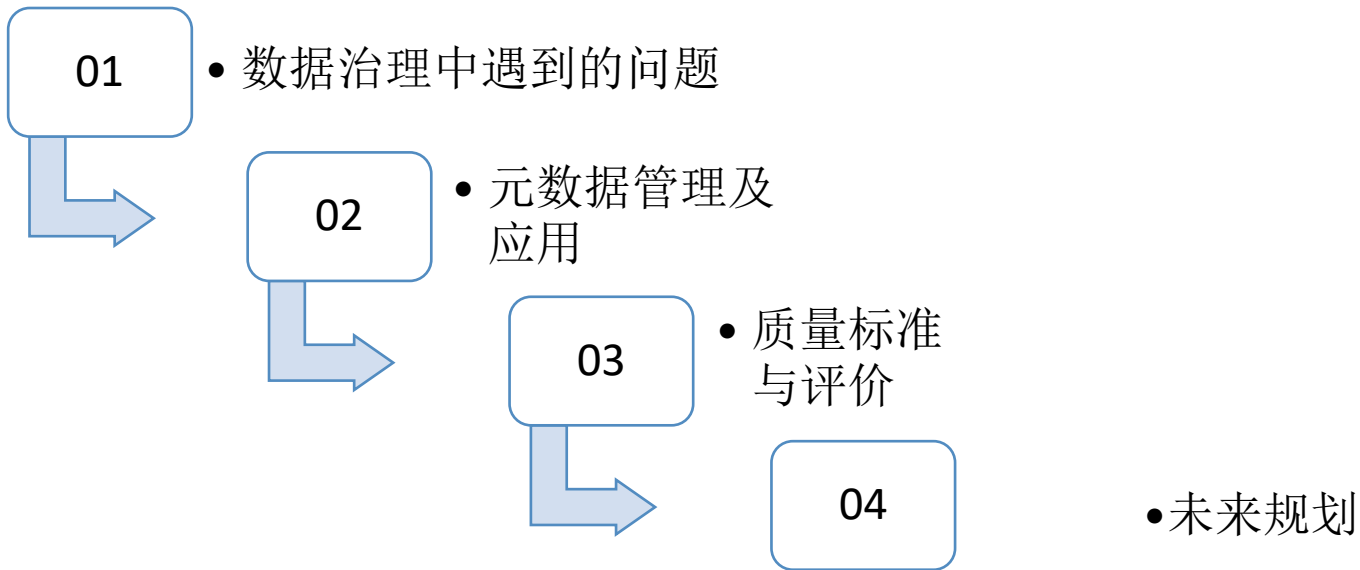
魅族数据治理探索之路

魏战松

The MEIZU logo is located in the bottom right corner. It is a blue, sans-serif font, positioned above the text '魅族技术学院出品'. The background features several colorful, rounded rectangular shapes in shades of green, orange, and red, scattered across the bottom and left edges.

魅族技术学院出品

大纲



大纲

01



- 数据治理中遇到的问题

问题背景

MEIZU

IT大咖说
知识共享平台

业务线	80+	
日新增数据量	> 70TB	每年4~6x增长
技术栈	Hadoop + Spark + Hbase+Hive+TEZ……	
计算平台	调度平台、自助分析平台、机器学习平台、权限……	

面对的问题

这个表有没有
其它地方用到?



ETL人员

今天的xx指标
怎么不正常?
数据怀疑?



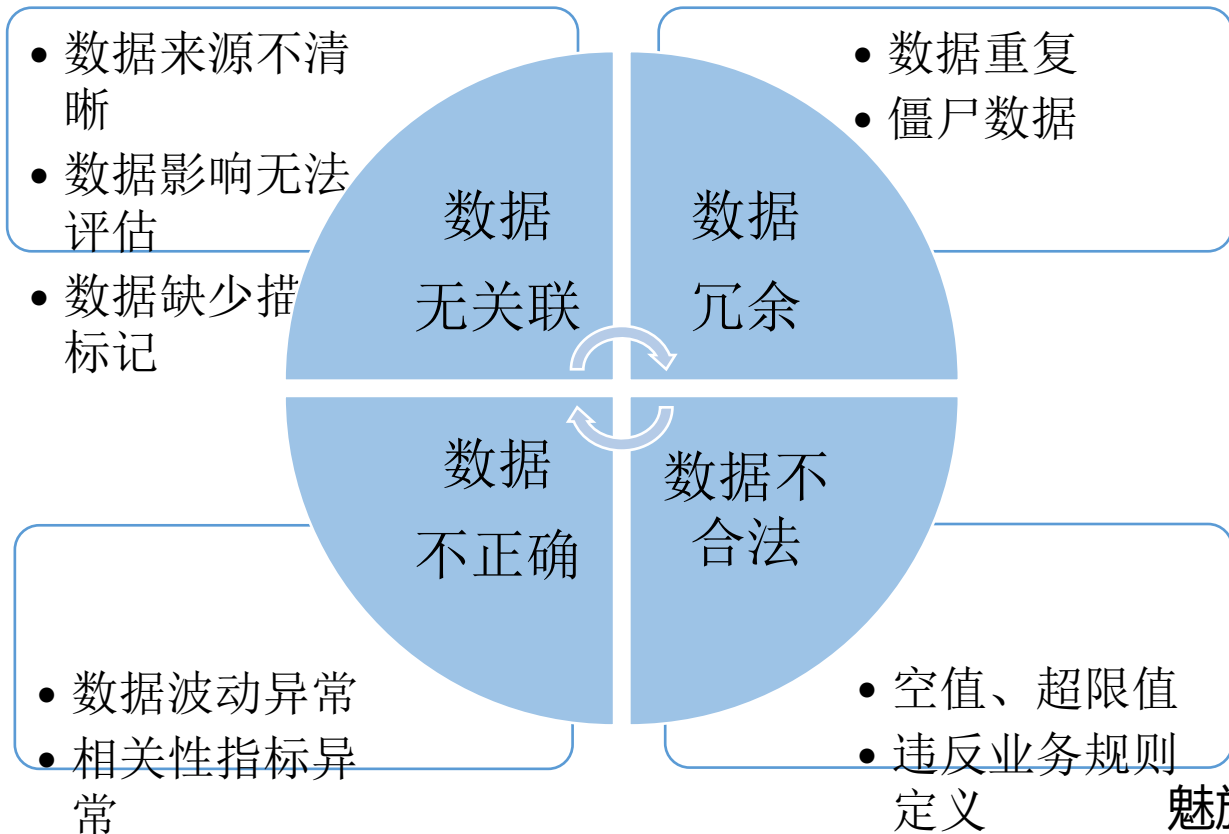
业务分析人员

存储不够了,
哪些数据可以
清理?



运维小哥


要治理什么样的数据




数据治理方案

元数据管理

- 指标定义
- 模型管理
- 血缘追踪
- 影响分析

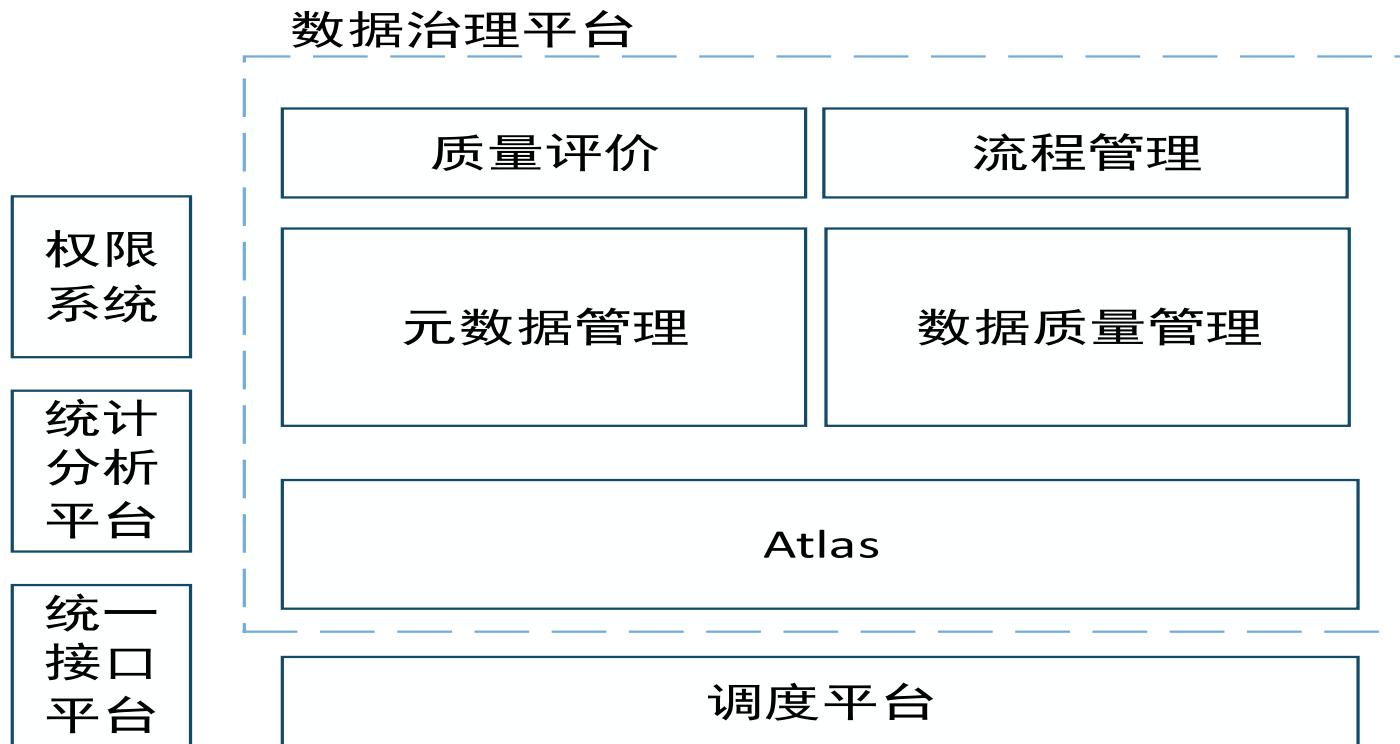
- 
1. 数据无关联
 2. 数据冗余

- 
3. 数据不合法
 4. 数据不正确

- 数据标准
- 规则校验
- 数据监控

数据质量管理

数据治理架构图



大纲

MEIZU

IT大咖说
知识共享平台

02



- 元数据管理及应用

魅族技术学院出品

MEIZU

什么是元数据?

元数据种类

• 技术元数据

基础平台

- HDFS
- Hive
- Hbase
- Spark

计算平台

- 离线任务(CETUS)
- 流平台
- 质量管理平台
- 权限系统(SCT)
- 可视化编码
- 机器学习平台

数据产品

- 统计分析平台
- 埋点系统
- 用户画像&洞察平台
- 统一接口平台

线下数据

- 模型元数据
- 指标定义
- ...

如何管理元数据?

MEIZU

IT大咖说
知识共享平台

自研?

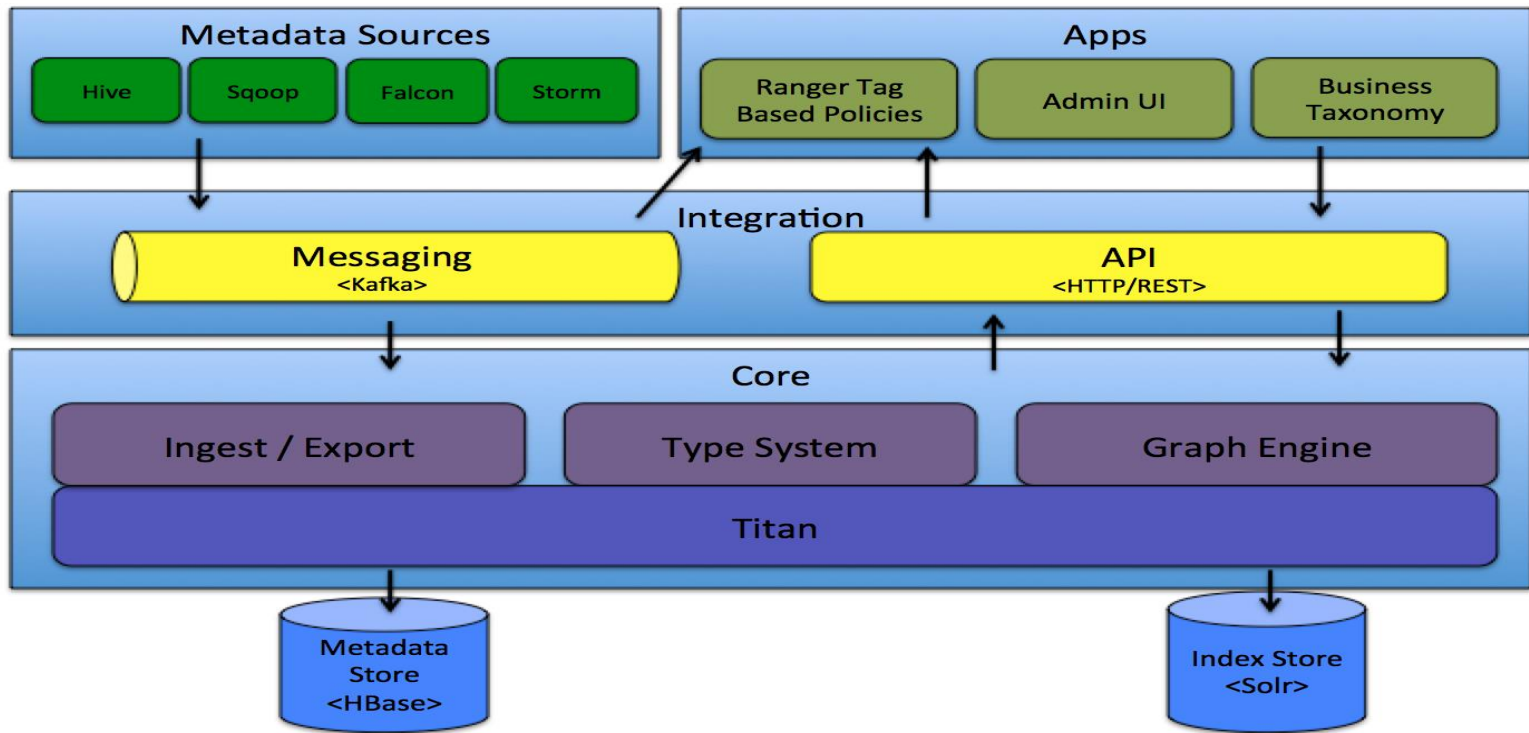
OR

寻求数据治理的开源解决方案?

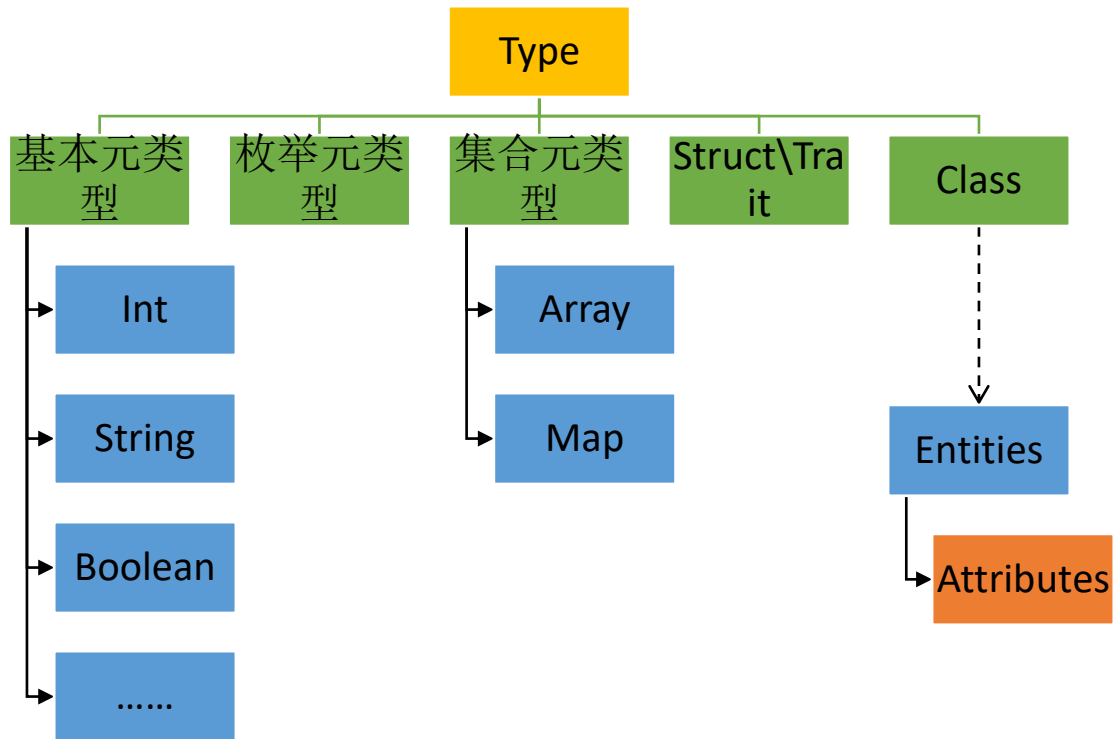
魅族技术学院出品

MEIZU

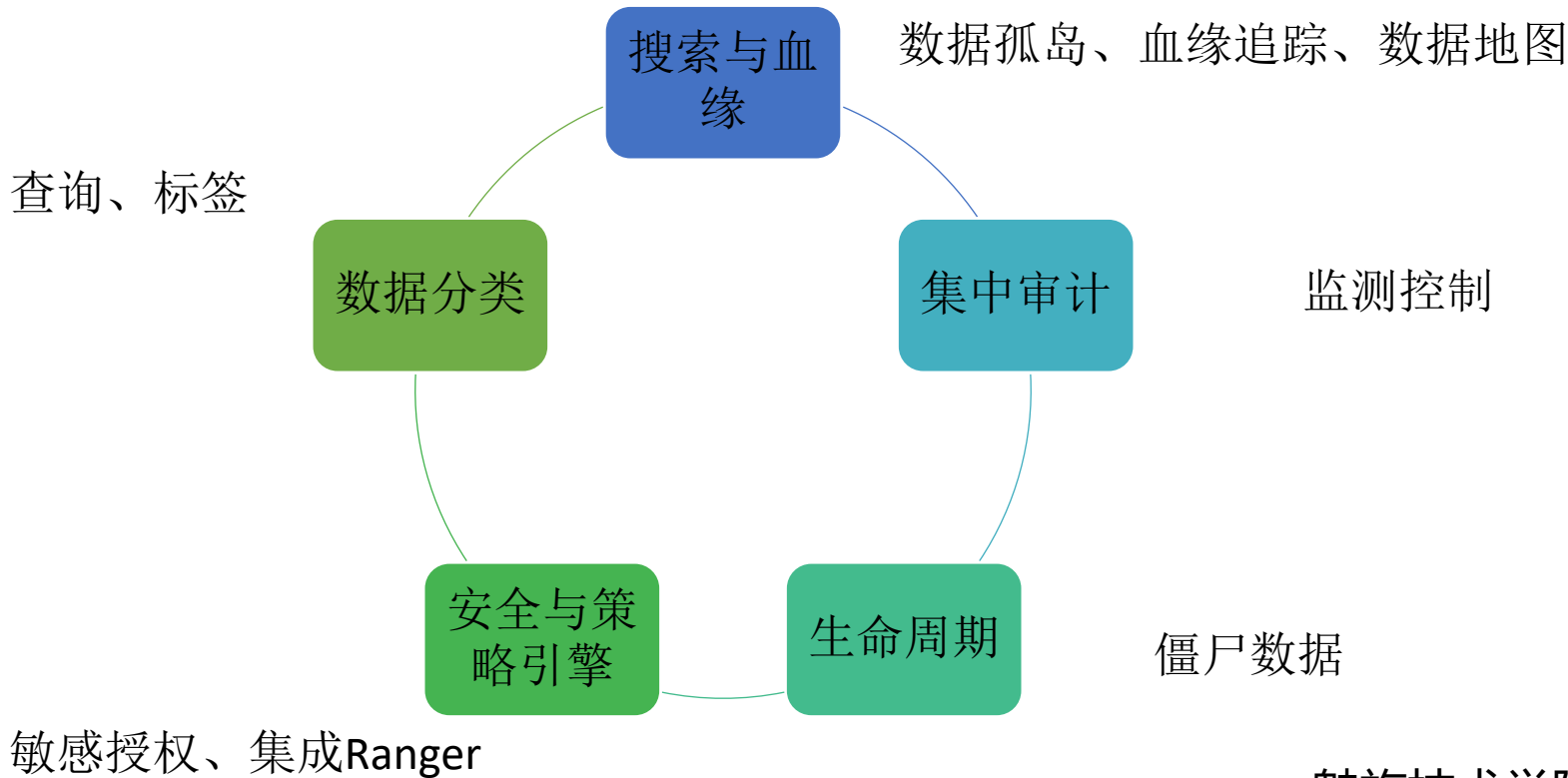
Apache Atlas



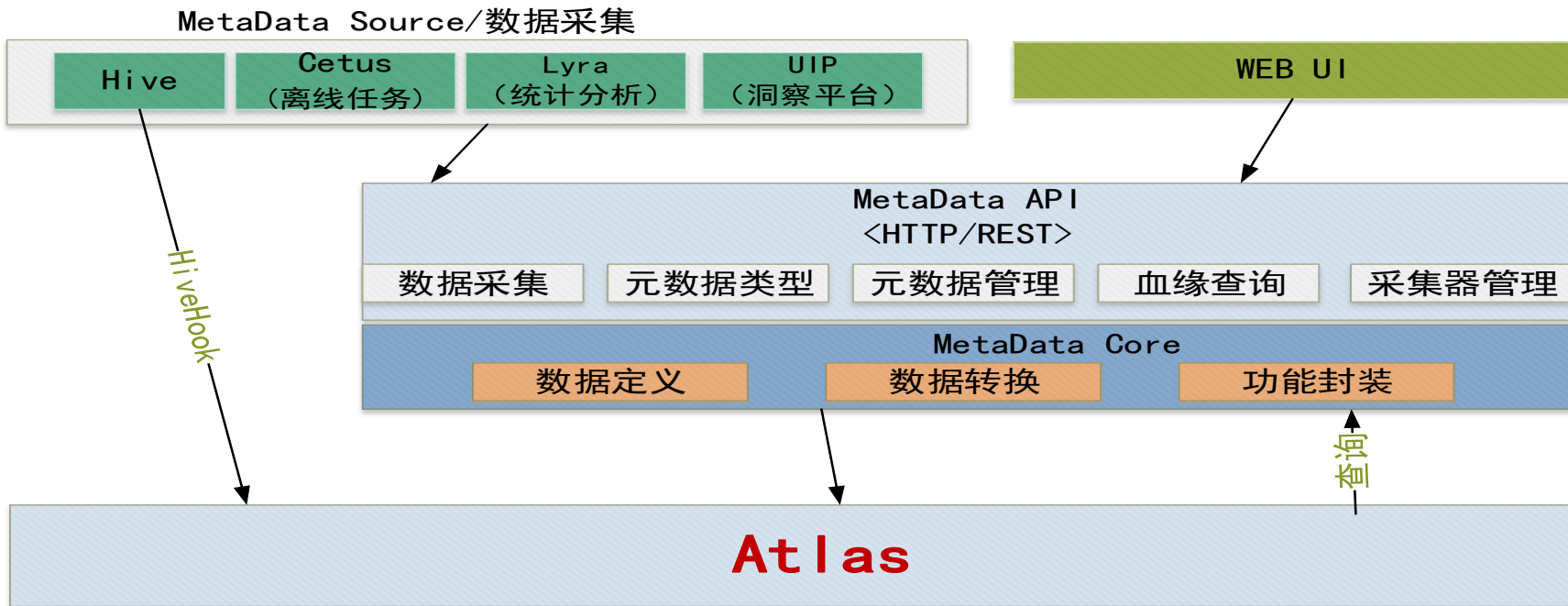
Apache Atlas 类型



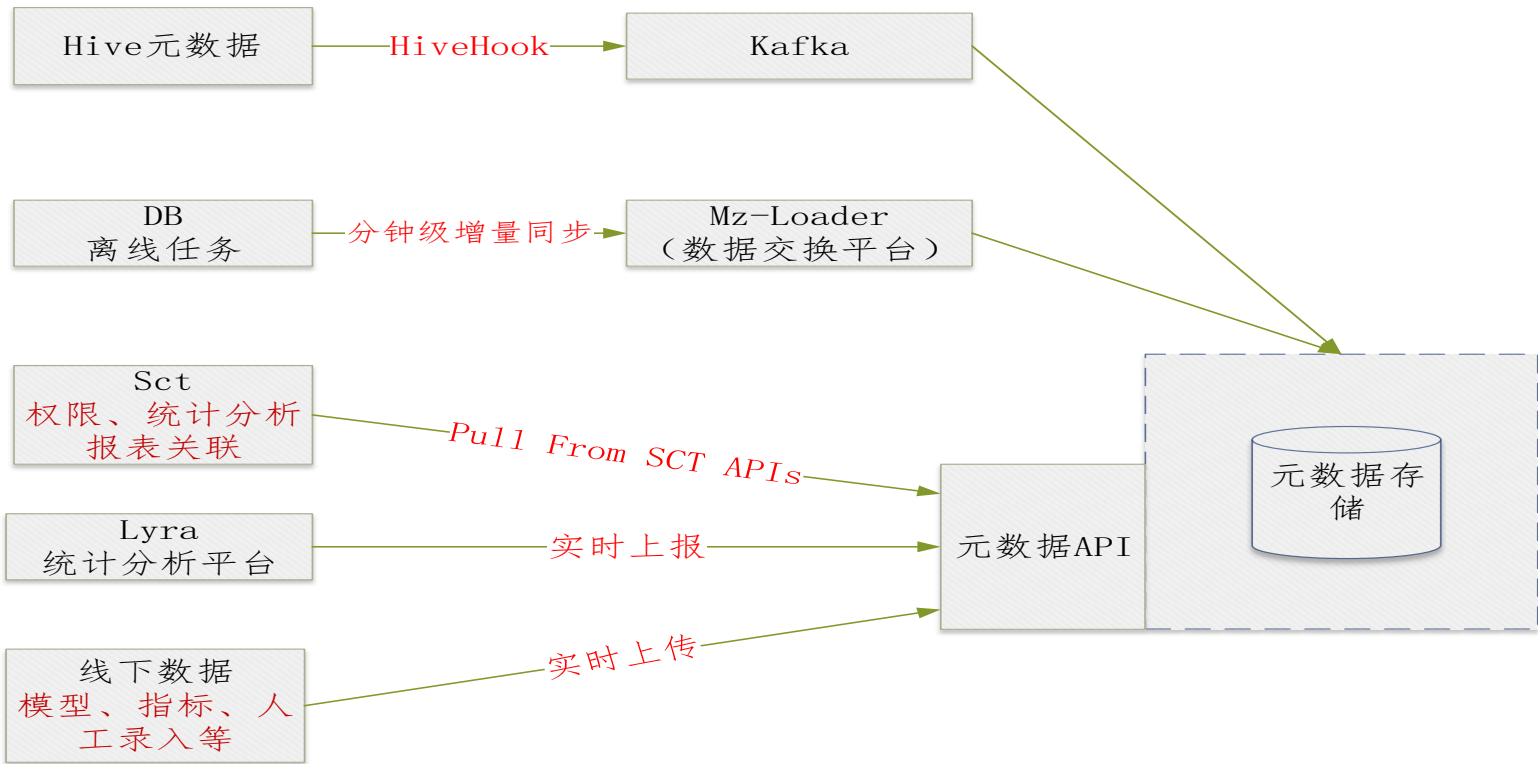
魅族元数据管理需求



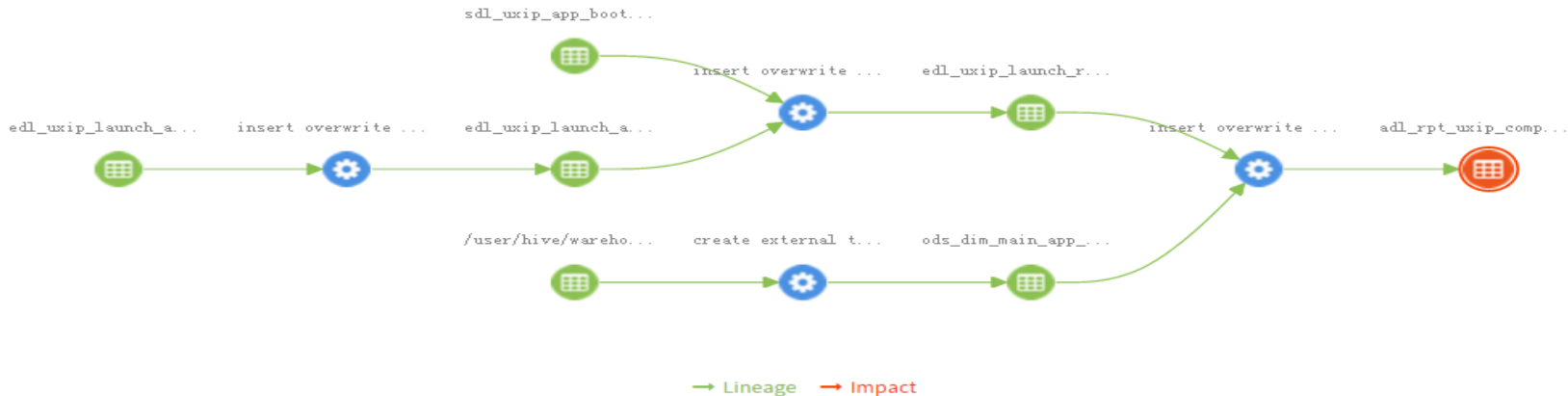
元数据平台架构



元数据采集



元数据应用-血缘分析



hive_process:

```
insert overwrite table edl_uxip_launch_retention_rate partition(stat_date=20170924)
SELECT *
FROM
  (SELECT *
   FROM sdl_uxip_app_boot_stat_d s
   WHERE stat_date = 20170924
        AND event_name = 'boot_app'
   GROUP BY umid,packagename) s
RIGHT JOIN
  (SELECT *
   FROM edl_uxip_launch_app_user t
   WHERE (stat_date >= 20170917
        AND stat_date <= 20170923)
        OR (stat_date IN (20170910, 20170825))) t
ON s.umid=t.umid
   AND s.packagename=t.packagename
GROUP BY t.packagename ,t.stat_date) s
```

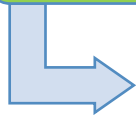
大纲

MEIZU

IT大咖说
知识共享平台

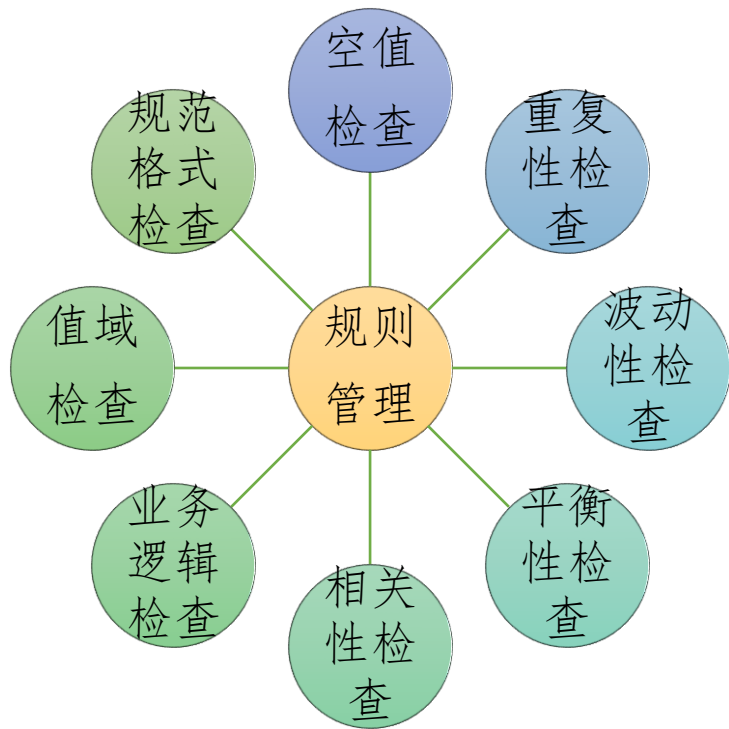
03

- 数据质量管理及评价体系



如何评价数据质量？

度量规则



质量管理系统

MEIZU

IT大咖说
知识共享平台

故障管理平台

告警平台

流程管理平台 (CVN)

度量规则管理

基础数据度量规则

指标度量规则

质量问题发现

数据核检结果

质量问题提交

质量问题告警

质量问题分析

数据质量分析报告

数据质量对比分析

查询功能

检核规则查询

检核结果查询

规则执行情况查询(执行时间、
时长、资源)

检核规则调度

HQL任务完成后自动执行规则调度

独立任务调度

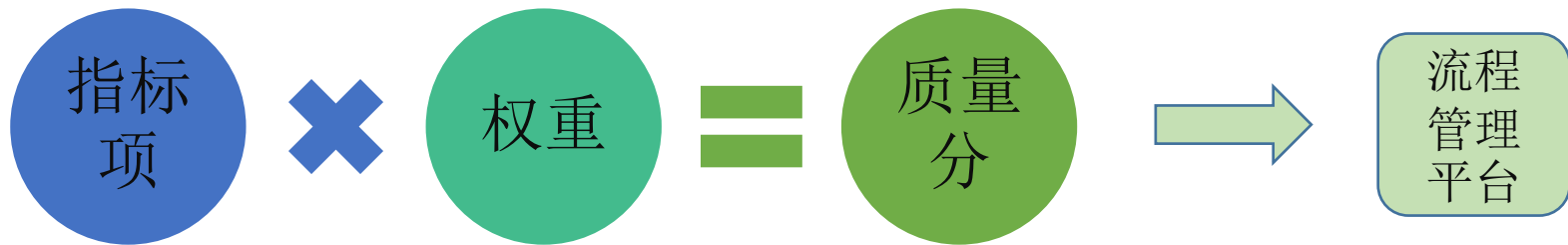
魅族技术学院出品

MEIZU

质量评价体系

MEIZU

IT大咖说
知识共享平台



指标定义

1. 离线任务失败个数
2. 离线任务整体延迟分钟数
3. 九点之前未完成的任务数
4. 流平台告警次数
5. 规则执行失败次数

魅族技术学院出品

MEIZU

大纲

MEIZU

IT大咖说
知识共享平台

04

• 未来规划



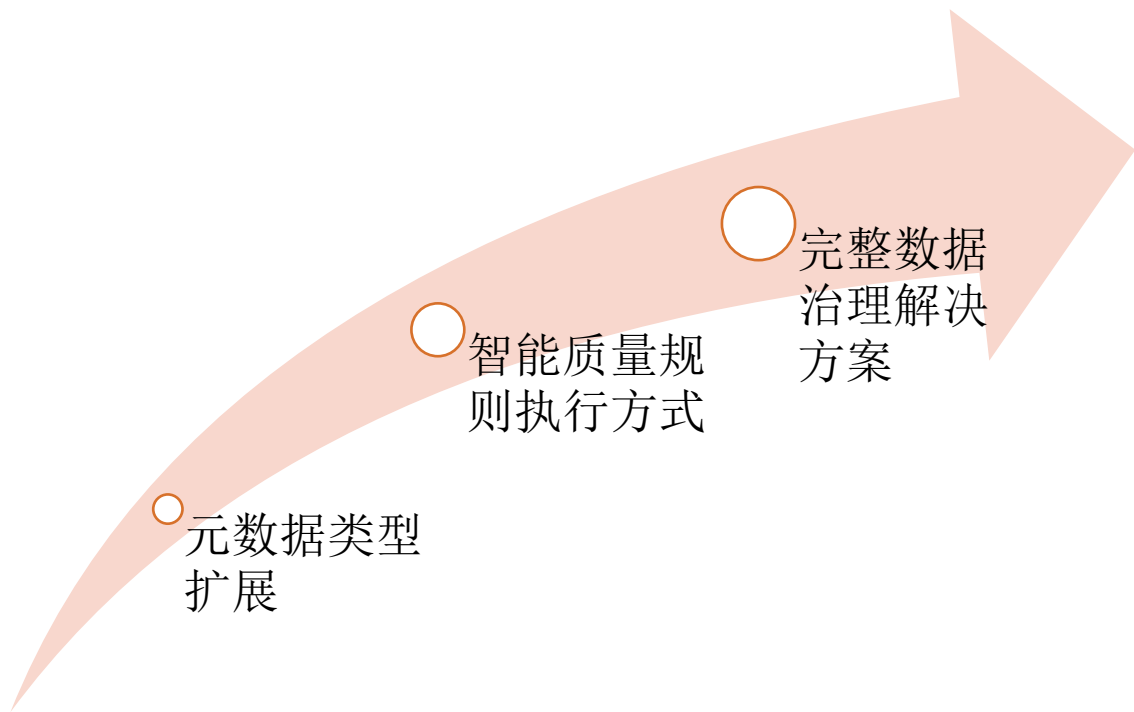
魅族技术学院出品

MEIZU

未来规划

MEIZU

IT大咖说
知识共享平台



The End