

基础数据建设实践

美丽说 高玉石

DTCC

2015中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2015

大数据技术探索和价值发现



互联网+

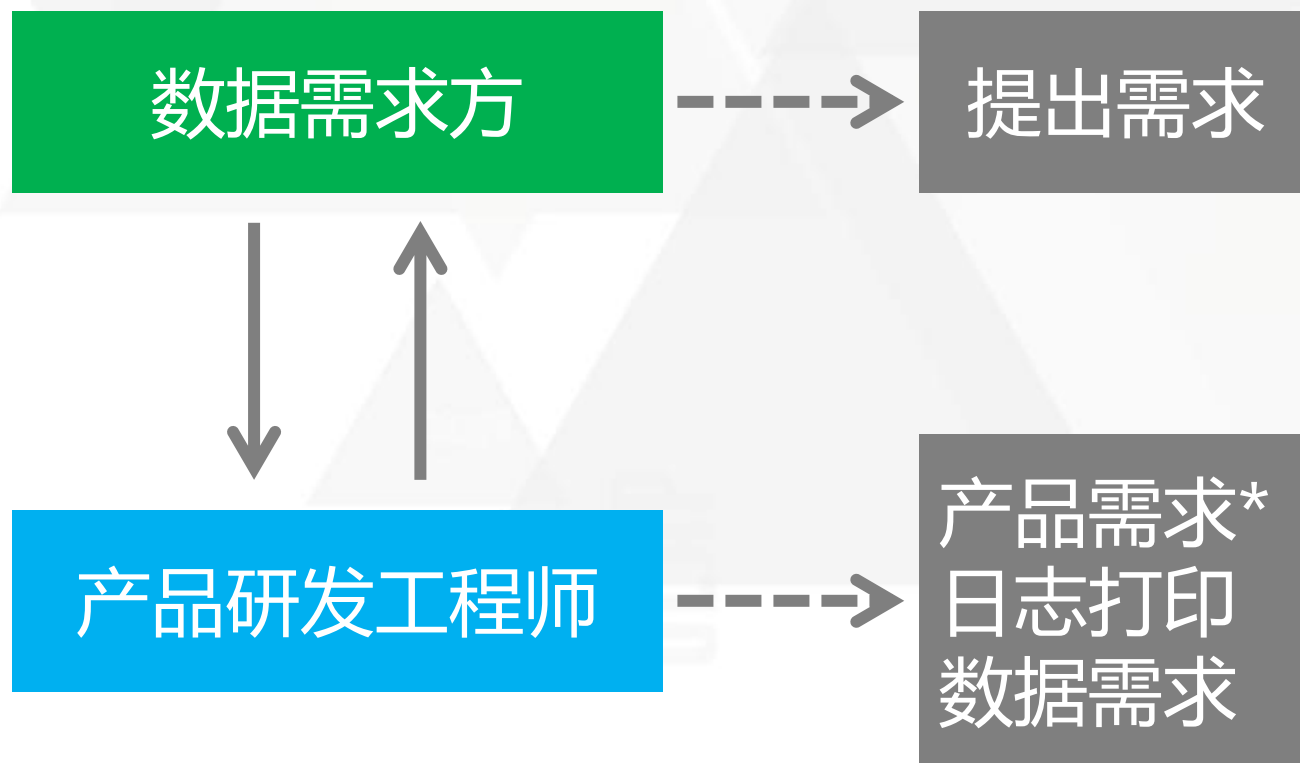
大数据



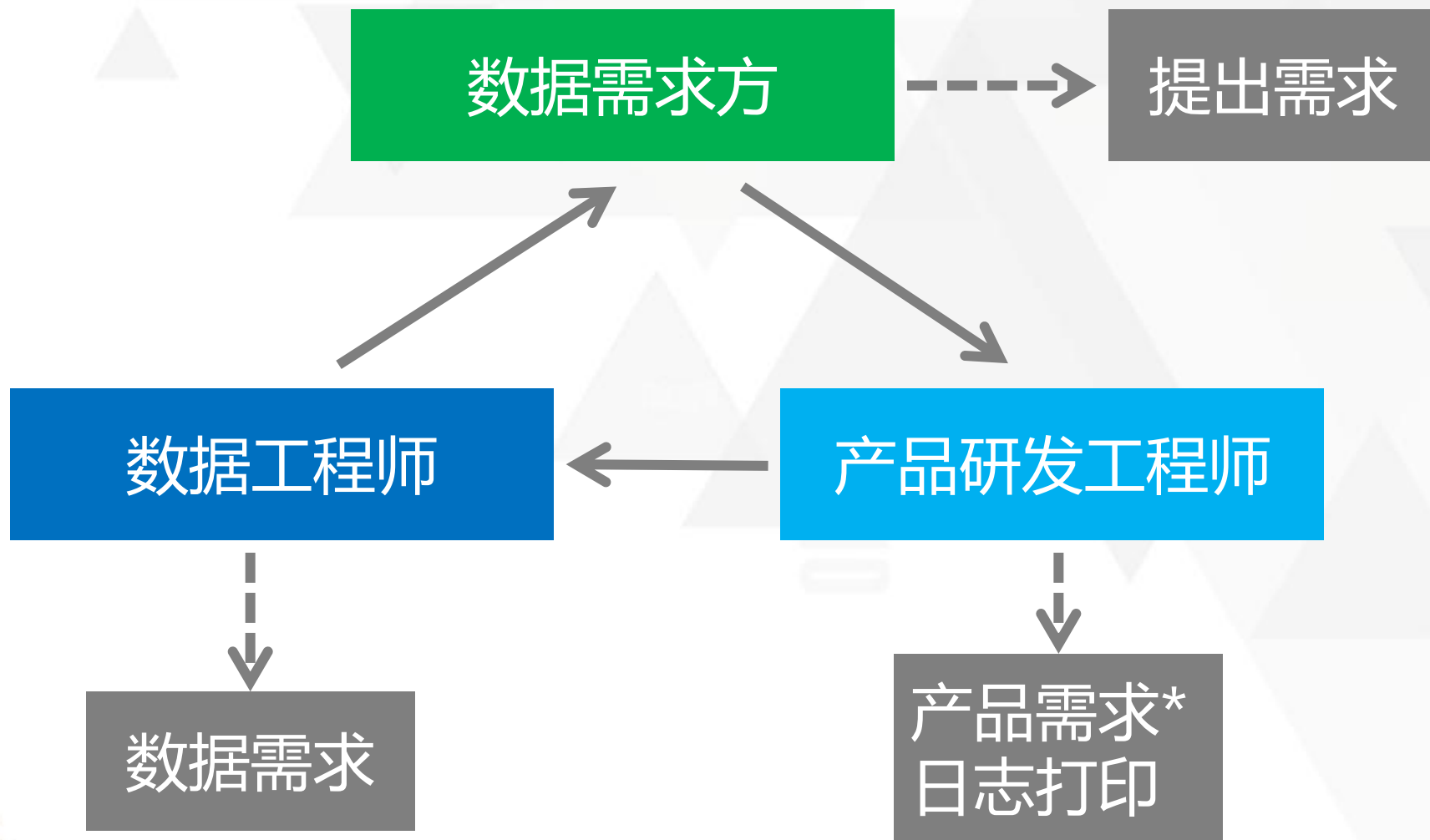
基础数据建设演进



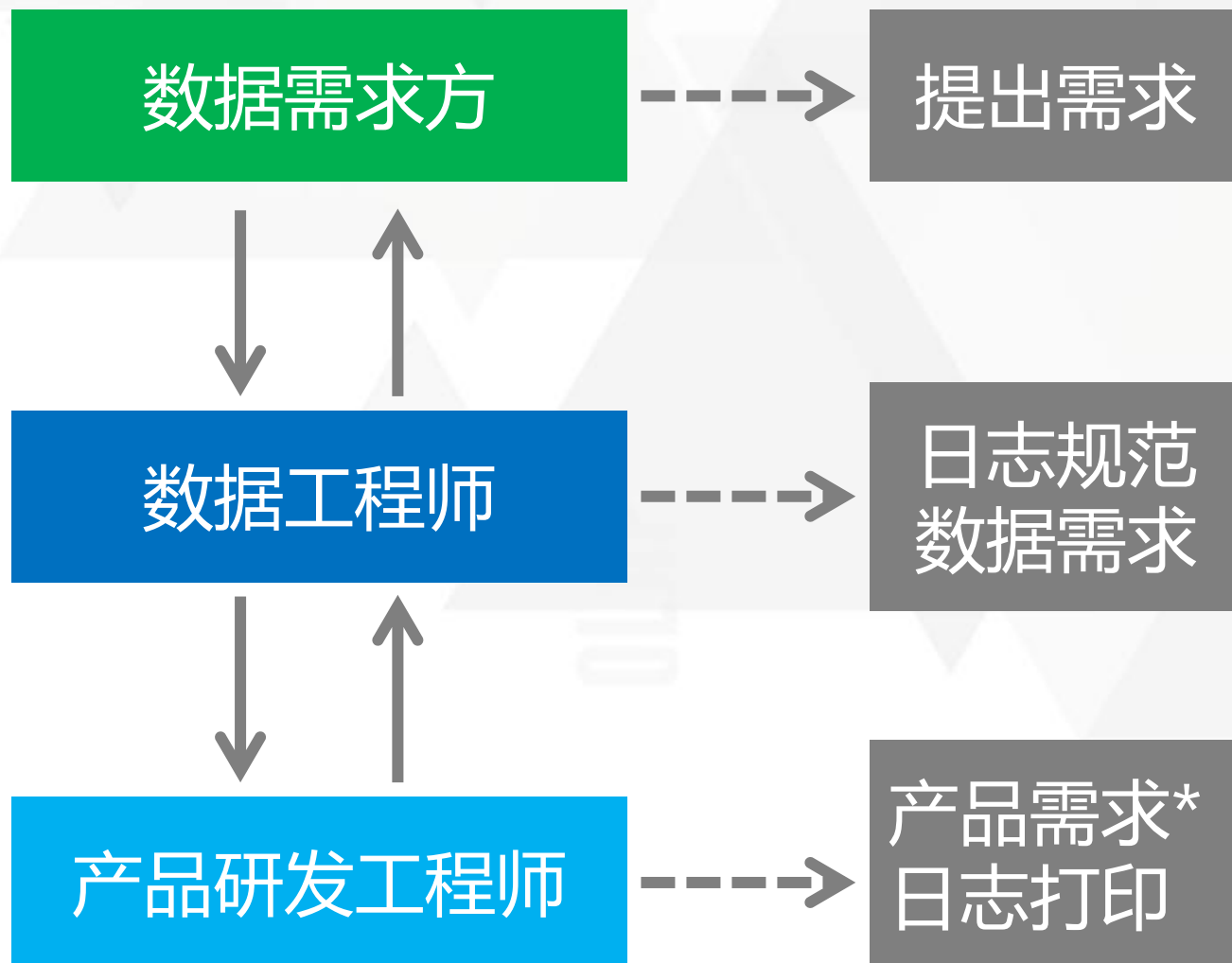
原始社会



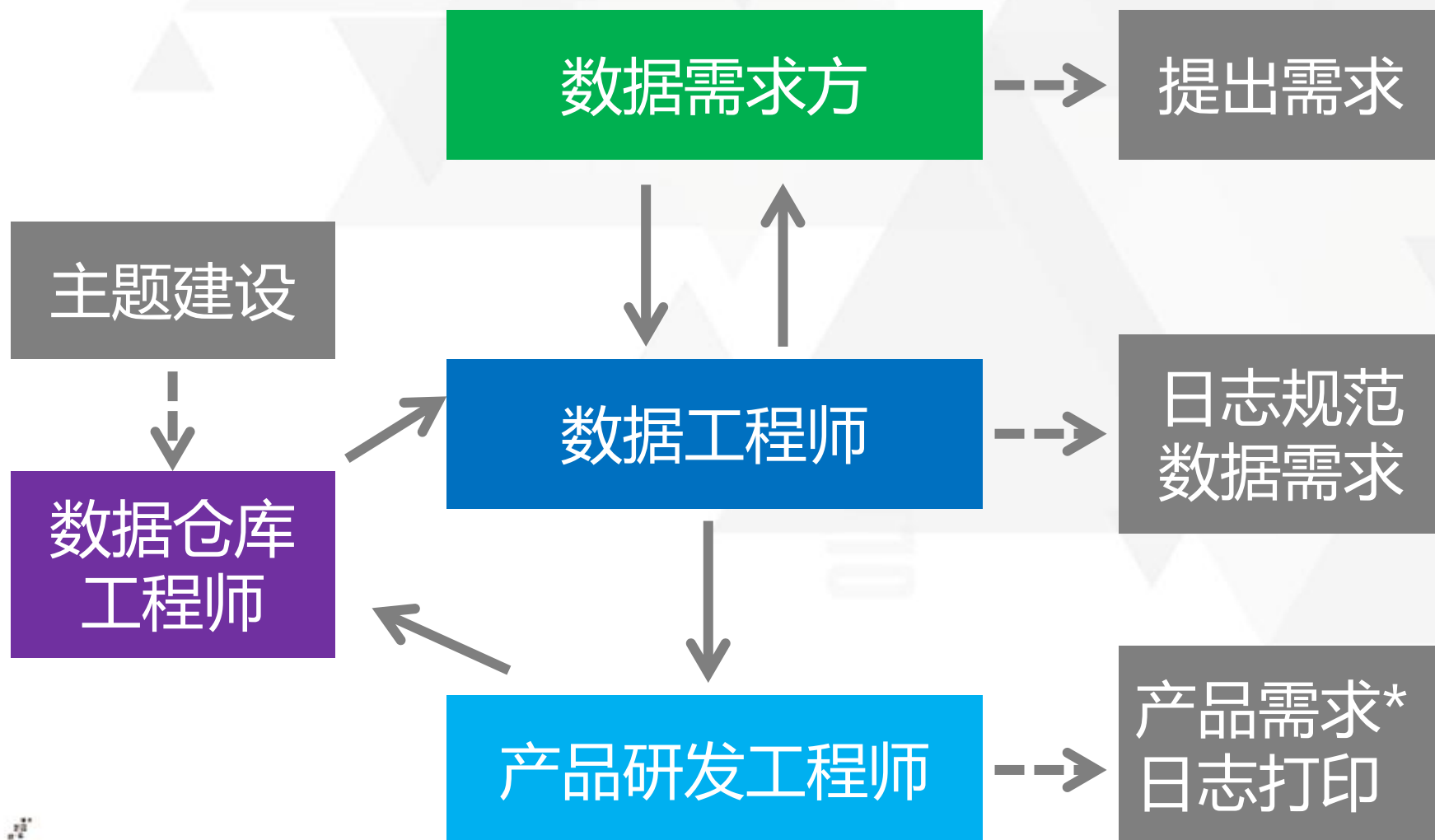
奴隶社会



封建社会



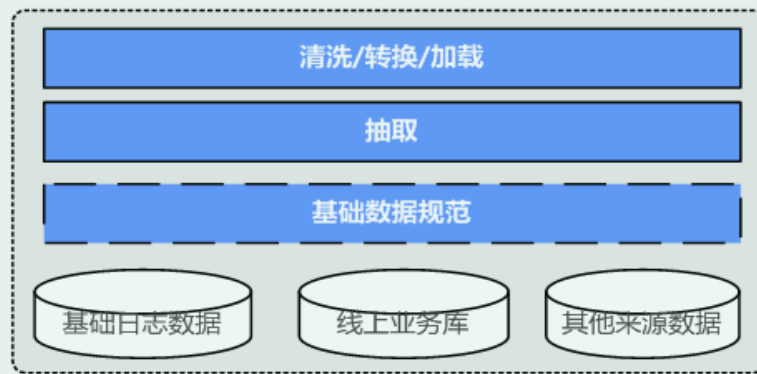
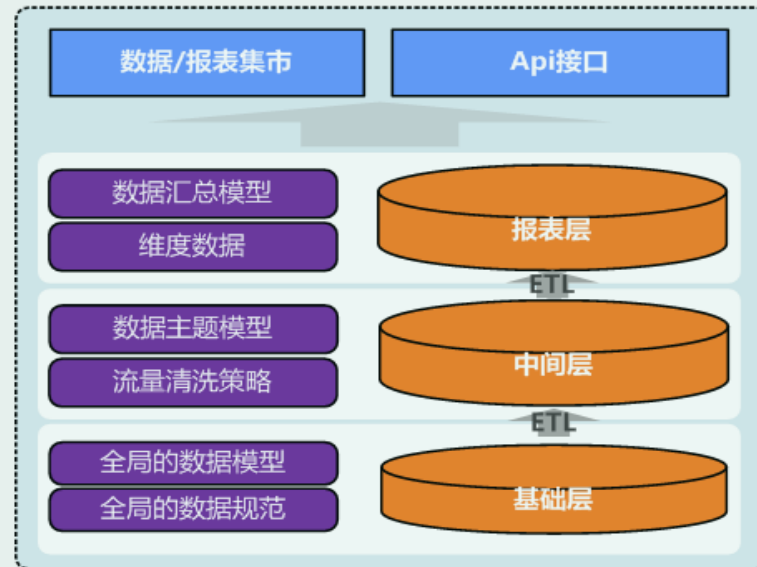
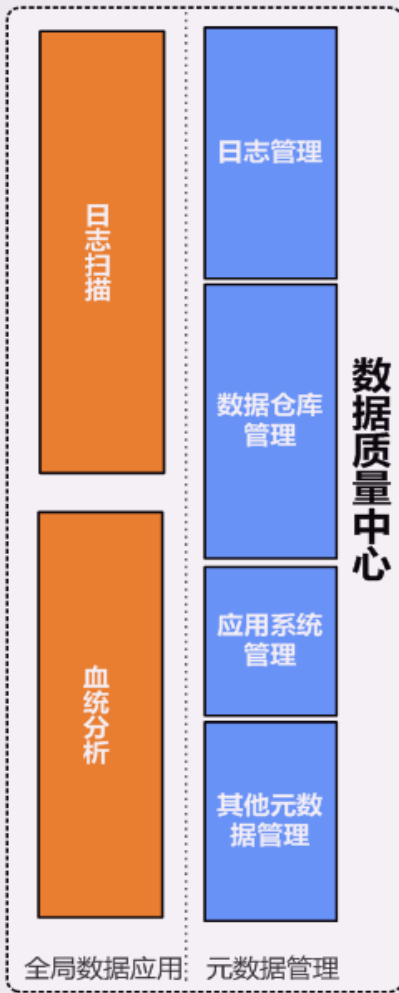
资本主义社会



共产主义社会

数据仓库 +





数据仓库建设



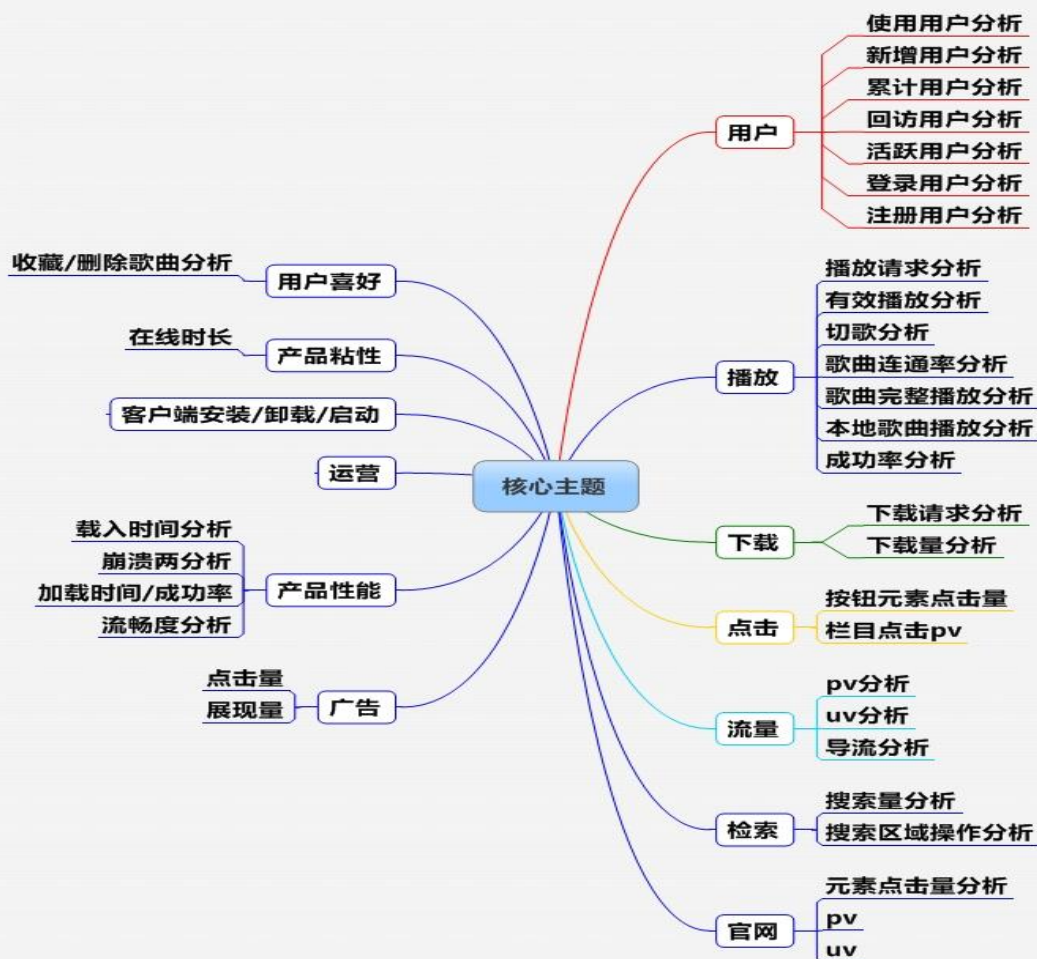


主题划分

一致性

粒度VS灵活

维度模型

数据仓库
建模

主题划分

高度汇总

指标唯一

数据一致

需求驱动

数据分析平台

数据挖掘

线上服务

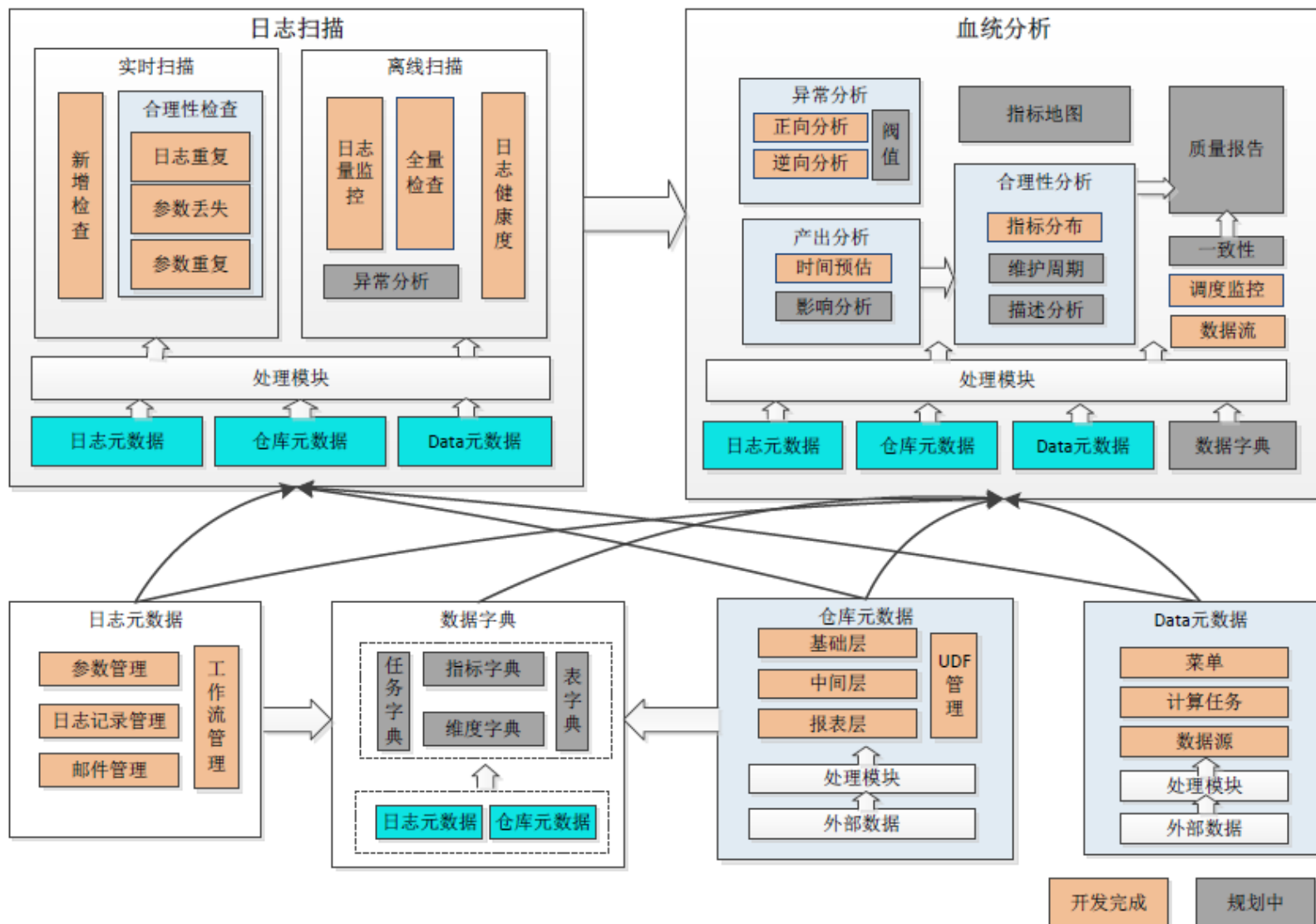
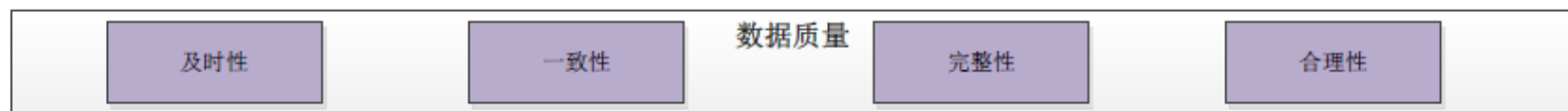
...

报表层



数据质量建设





元数据管理

日志/业务库元数据

日志模块 日志传输、参数结构

业务库表结构 表、属性状态

ETL

任务 层级 依赖 任务执行情况

数据仓库

表结构 属性信息 依赖关系

元数据

前端展现

报表 位置 层级 指标



元数据管理

日志管理

产品线

日志参

日志记

我的曲

日志扫

据仓库

RM元数

统分析

常分析

志扫描

本管理

表自动

人中心

表基本信息 | 表字段信息 | 表维度信息

名称

music_mds_songid_bhv_day

所属层级

中间层

主题

行为主题

终端类型

不分端

抽取周期

天

来源表

music_mds_bhv_play,music_dim_quku_song_si,music_dim_quku_resource_sf,music_dim_quku_song_si

对应程序

处理逻辑

从music_mds_bhv_play统计出每天各端的全量歌曲播放量，以songid作为主键和歌曲全量信息表全关联。歌曲全量信息表是指以songid作为主键将music_dim_quku_song_si,music_dim_quku_resource_sf,music_dim_quku_song_si关联到一起取出每首歌曲歌名，音频，歌词等信息。

功能描述

songid四端全量播放表

修改



日志扫描

类型

离线日志扫描

实时日志扫描

技术

实时计算

规则引擎

分析内容

日志错、漏、重复、异常等

日志使用情况

日志冗余



日志扫描

扫描结果报告

一. 扫描结果 (下载日志)

扫描完成时间: 2014-05-05 10:19:01

扫描成功,结果正常: 成功扫描到该字段, 详情请关注《三》《四》《五》!

二. 扫描基本情况

主扫描字段: < album

扫描规则: 无明确

扫描开始时间: 2014-0

三. 初级扫描结果详情

成功扫描到主字段(

详情日志如下:

第1条:

61.172.247.198 -- 05/May/2014:01:01:01
hello&singer_id=1&singer_id=1
_version=&linkCode=220000
ga_201404251715&userna
xcode=26cb151a204b141f8
ertype=lowflow&prestatus=
__m=mboxCtrl.playSong&
NBID=6D0DEC6BB793873

四. 高级扫描结果详情

1. 暂未发现有重复字段的现象发生!

2. 暂未发现有重复记录的现象发生!

五. 基础扫描字段结果分析

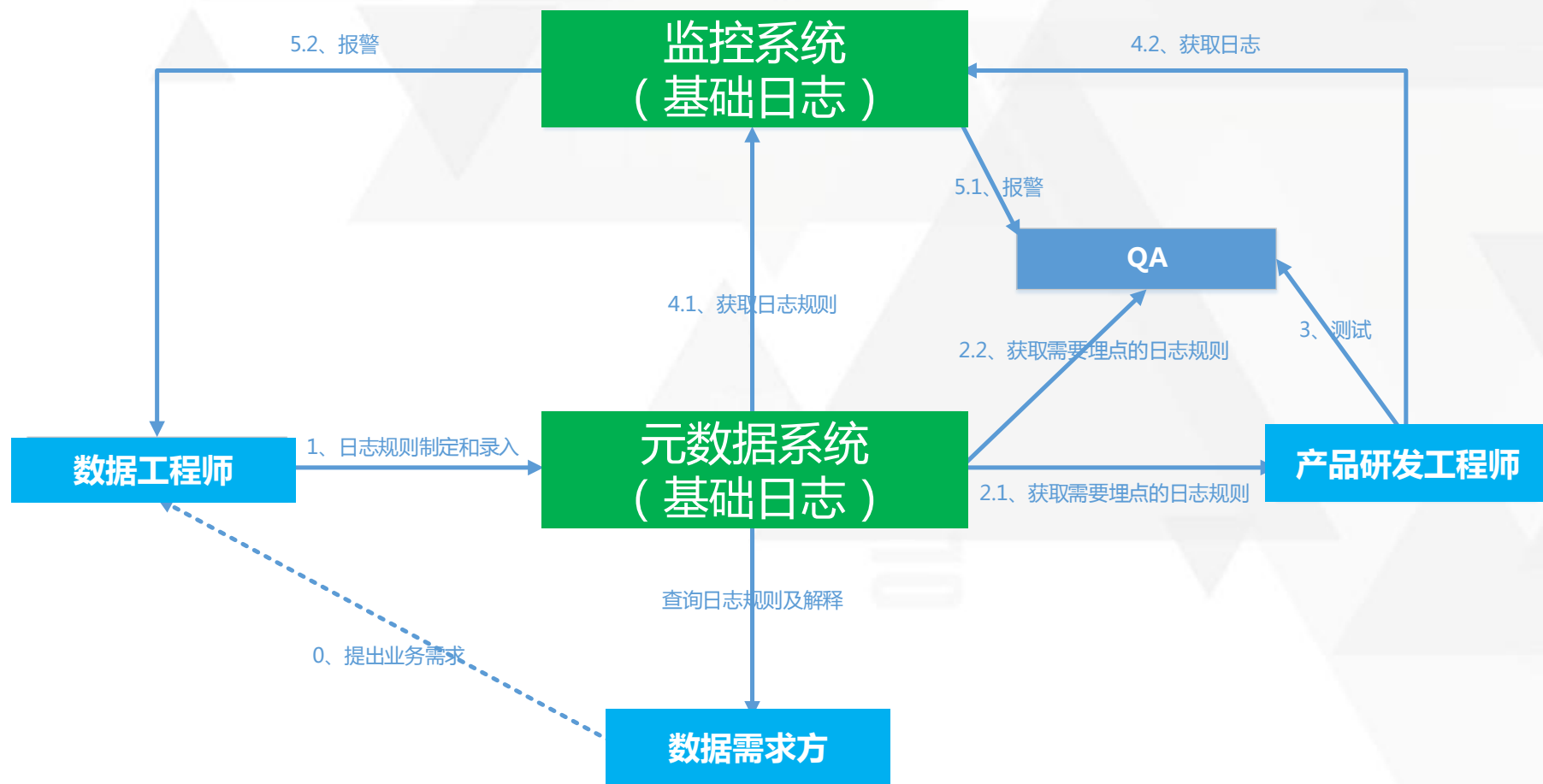
总条数为: 501条,基础字段的详细扫描结果分析如下表:

参数	规则	空值率	扫描结果分析 (占比=存在数/总数)	操作
pid	=	0	正常(value=304)的条数: 501条——占比: 100%	查看详细日志
flow	in	0	未扫描到该字段的条数: 98条——占比: 19.56%	查看详细日志
			正常(value=ad_voice_1)的条数: 403条——占比: 80.44%	查看详细日志
ref	in	0	正常(value=tingbox)的条数: 403条——占比: 80.44%	查看详细日志
			正常(value=radio)的条数: 85条——占比: 16.97%	查看详细日志
			正常(value=aldsinglebox)的条数: 12条——占比: 2.4%	查看详细日志
			正常(value=pc_widgethao123)的条数: 1条——占比: 0.2%	查看详细日志

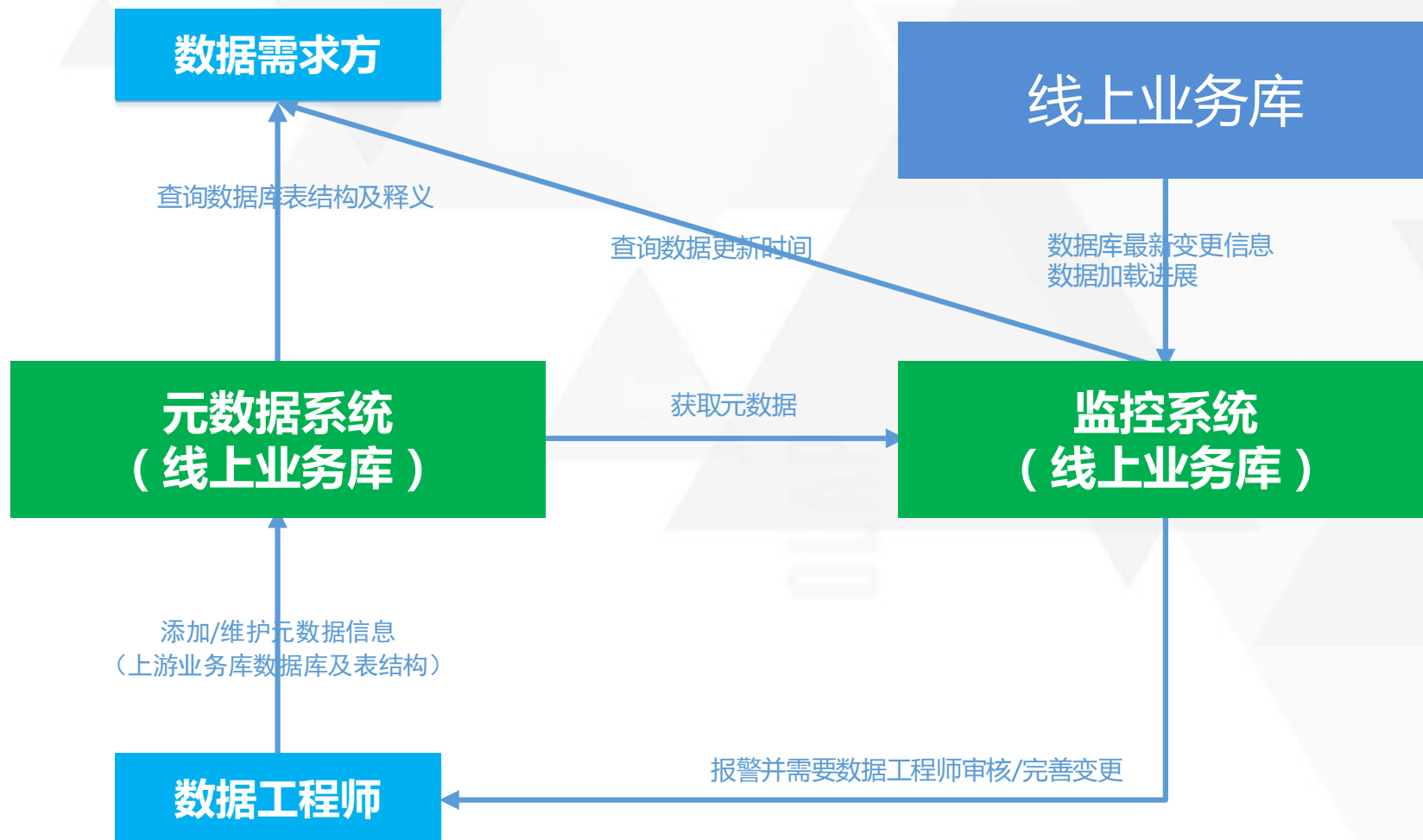
[返回主页](#)



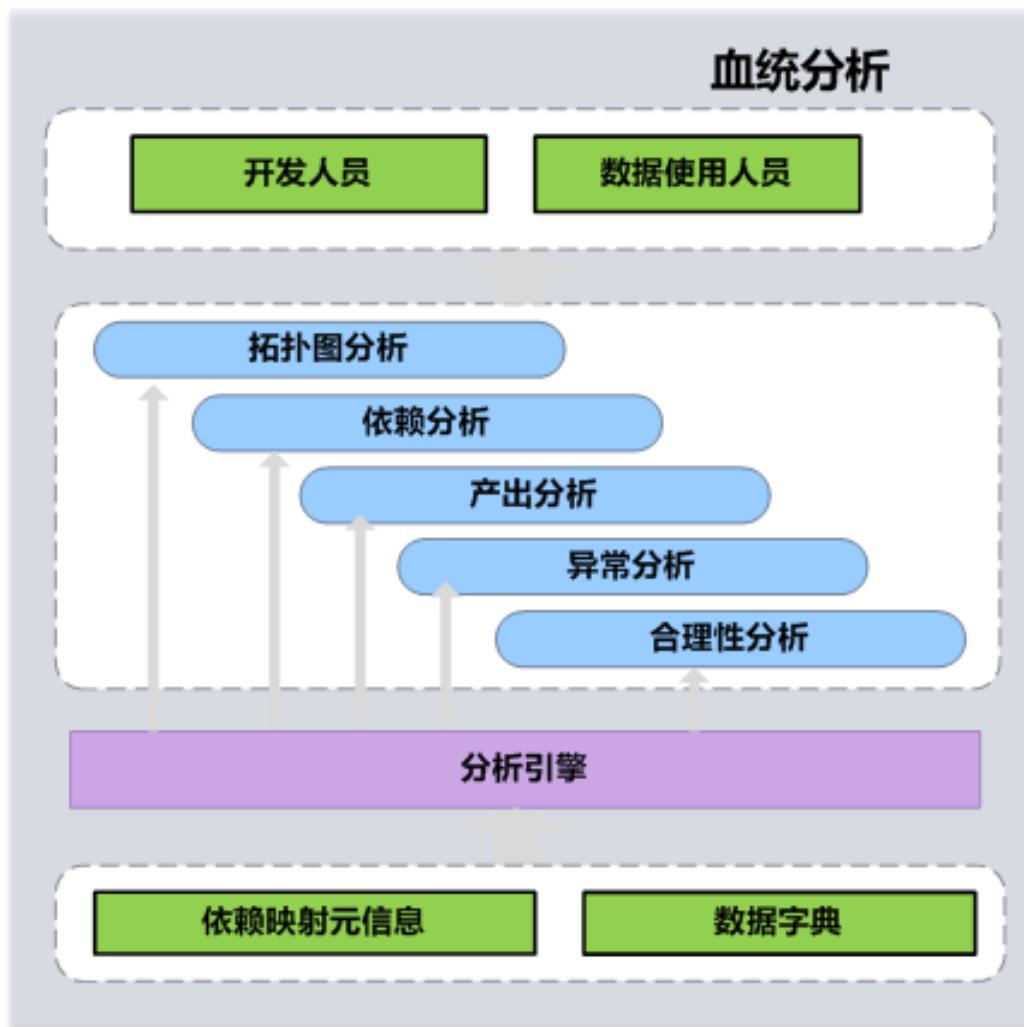
日志规范流程



业务库规范流程



血统分析



拓扑图分析

了解单点对于上下游甚至全局的影响

依赖分析

任务调度/运行异常、依赖异常、重写异常

产出分析

数据产出时间预估等

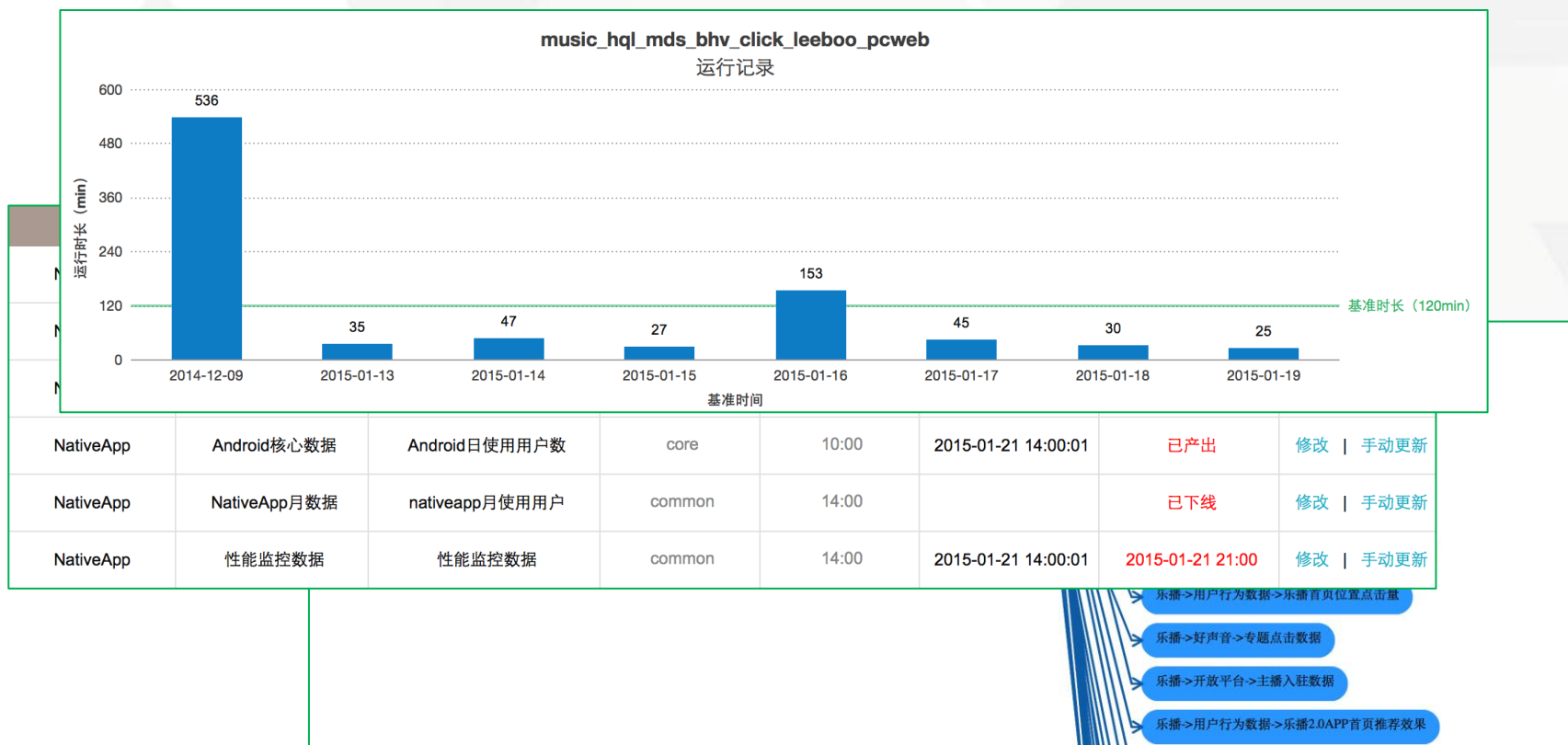
异常分析

合理性分析

...



血统分析



欢迎加入我们！



微信：kioskgao

邮箱：yushigao@meilishuo.com





THANKS