

网易严选大数据架构演进

网易严选 左琴

左琴

2013年毕业于浙江大学，先后任职阿里/华云数据/网易等公司，主要从事大数据存储/计算引擎开发和优化工作

目前就职于网易严选数据及风控部，同时担任网易集团数据委员会委员，是网易严选数据和算法工程团队的负责人，负责严选数据技术体系的建设(大数据平台/中台/数据产品和算法工程)

社区的积极分享者，先后在SACC，QConf研习社等做过多次主题分享，举办过Alluxio、Pulsar杭州的Meetup.



1

数据驱动技术体系

数据分析-> 数据决策

2

数据中台

DataLake

AutoWarehouse

3

数据平台

智能任务调度

Cloud Native

辅助决策 统计 分析 诊断

VIPAPP
移动数据工作台

轩辕
数据大屏

神相
流量分析产品

谛听
舆情洞察

伏羲
营销数据分析

大麦
商品数据运营

哈勃
评论数据产品

统计分析型 数据产品

用户分析

商品销售

渠道销售

逆向分析

供应商管理

流量分析

营销分析

仓储物流

品控分析

财务分析

严选报表分析体系（有数）

DIS
科学实验平台

阶段1 2017~2020: 重点打造分析体系

智能决策 策略 模型

用户运营

用户增长

风控

人

选品

榜单

货

推荐

搜索

栏目

广告

场

辅助决策 统计 分析 诊断

VIPAPP

移动数据工作台

轩辕

数据大屏

DIS

科学实验平台

大麦

商品数据运营

伏羲

营销数据运营

谛听

舆情洞察

数据服务化 破除数据门槛 开放数据+能力

河洛

供应链协同决策

业务财务一体化平台

严选商品中心

刑天

市场投放与归因系统

阶段2 2020~: 打造数据驱动工程体系

全域
用户运营

用户增长

风控

标签 特征 策略 模型

选品

榜单

推荐

赛马

货

场

河洛
供应链协同决策

业务财务一体化平台

开放数据+开放能力

刑天
市场投放与归因系统

严选商品中心

时效: T+1, T+10m, 实时

质量: 准确率, 异常分析

定义: 统一的定义

数据

访问: 一致的访问服务

存储: 对外不直接暴露

计算: 服务化的计算能力
Spark, flink etc.

稳定: 满足在线业务SLA要求

能力

标签: 3个不同的服务

特征: 无统一访问服务

存储: 10+ 存储引擎

推荐: 3个推荐服务

临时方案和烟囱开发积累的债务

债务

无ACID能力
无法应对在线业务要求

数仓模式过于重度
无法应对标签与特征的大量需求

开放表, 没SLA的服务
无法应对业务系统对数据体系的要求

演进

面向分析 向 数据驱动 的工程演进

目标

业务目标

翻译

工程目标

分析 | 结合

工程现状

打造数据驱动技术体系

1 数据驱动技术体系

数据驱动工程体系建设思路

业务角度切入，技术视角划分



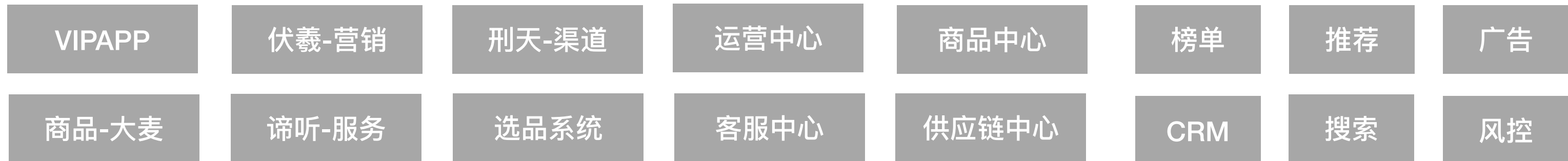
看技术本质
通用能力下沉

屏蔽底层细节
提供积木搭建

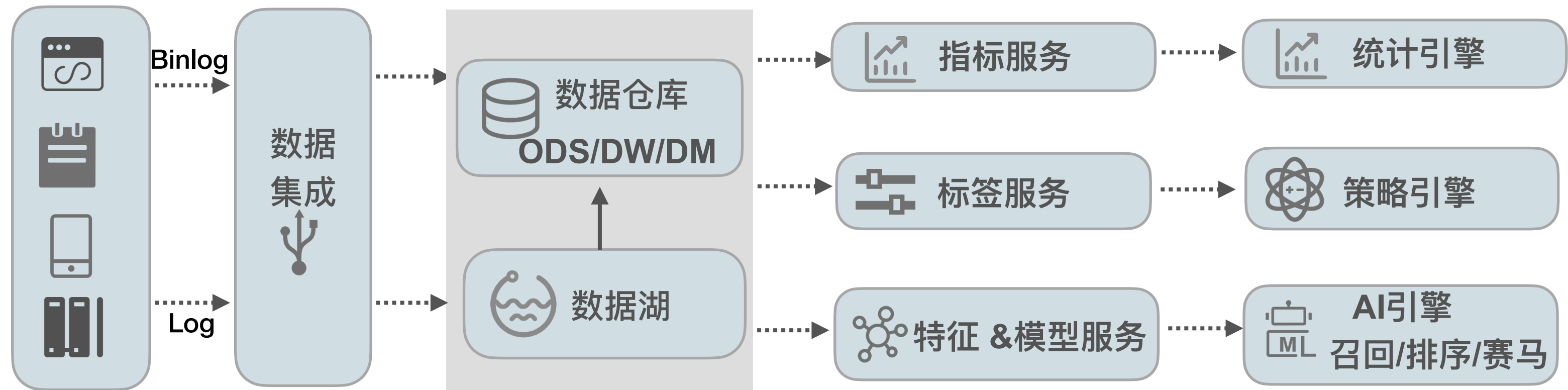
数据驱动技术体系

数据驱动技术体系全景图

前台



中台



平台



科学实验平台

流量控制

效果分析

数据治理

指标治理

标签治理

表治理

任务治理

基础服务

统一元数据

数据血缘

监控服务

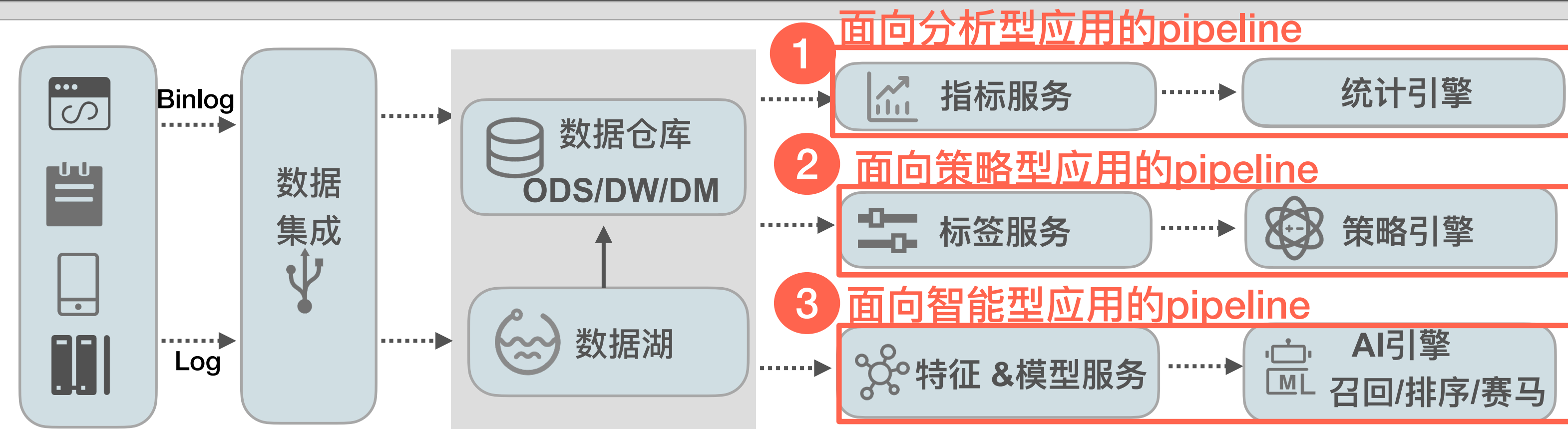
数据驱动技术体系

数据驱动技术体系全景图

前台



中台



平台



科学实验平台

流量控制

效果分析

数据治理

指标治理

标签治理

表治理

任务治理

基础服务

统一元数据

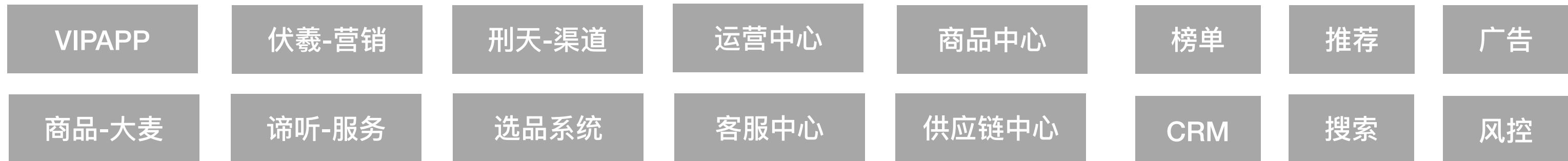
数据血缘

监控服务

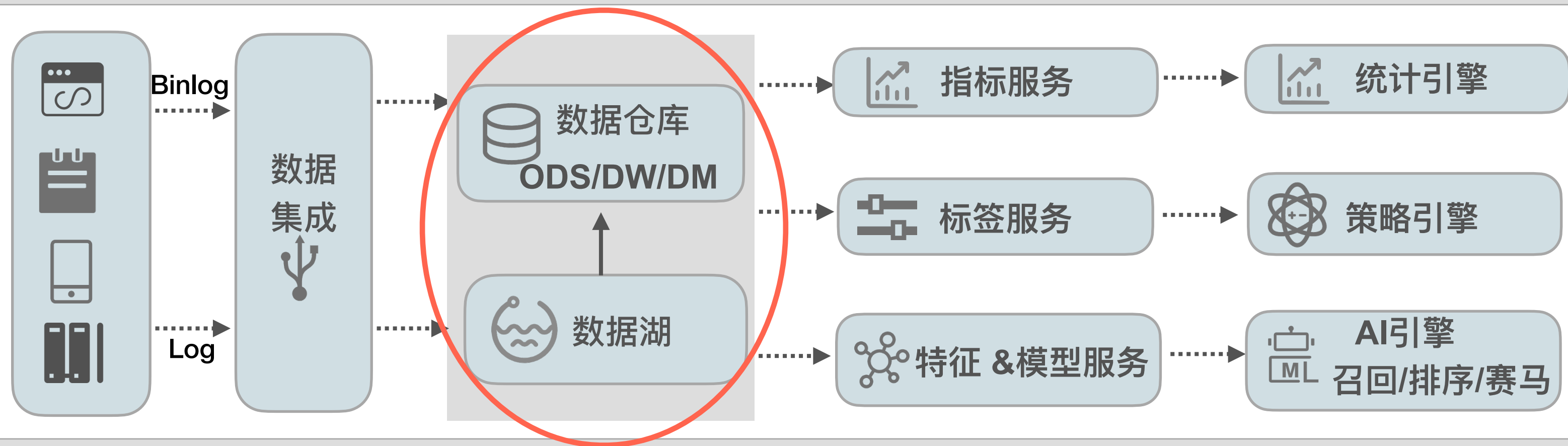
数据驱动技术体系

数据驱动技术体系全景图

前台



中台



平台



科学实验平台

流量控制

效果分析

数据治理

指标治理

标签治理

表治理

任务治理

基础服务

统一元数据

数据血缘

监控服务

1

数据驱动技术体系

数据分析-> 数据决策

2

数据中台

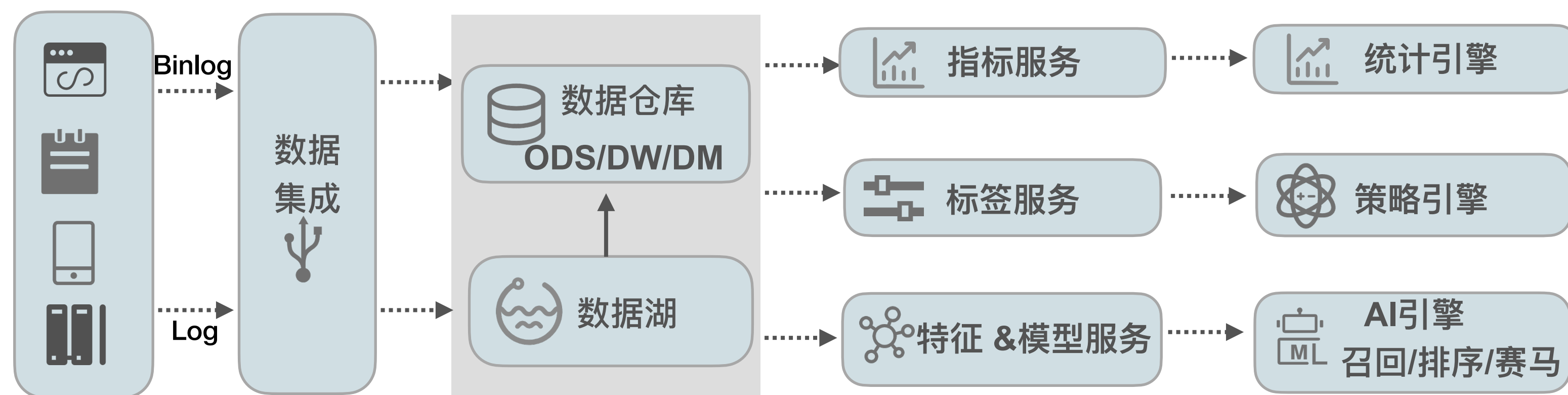
DataLake
AutoWarehouse

3

数据平台

智能任务调度
Cloud Native

数据&算法中台



面向业务场景的, 数据(指标/标签/特征/模型), 工程组件的技术解决方案



科学实验平台

流量控制

效果分析

数据治理

指标治理

标签治理

表治理

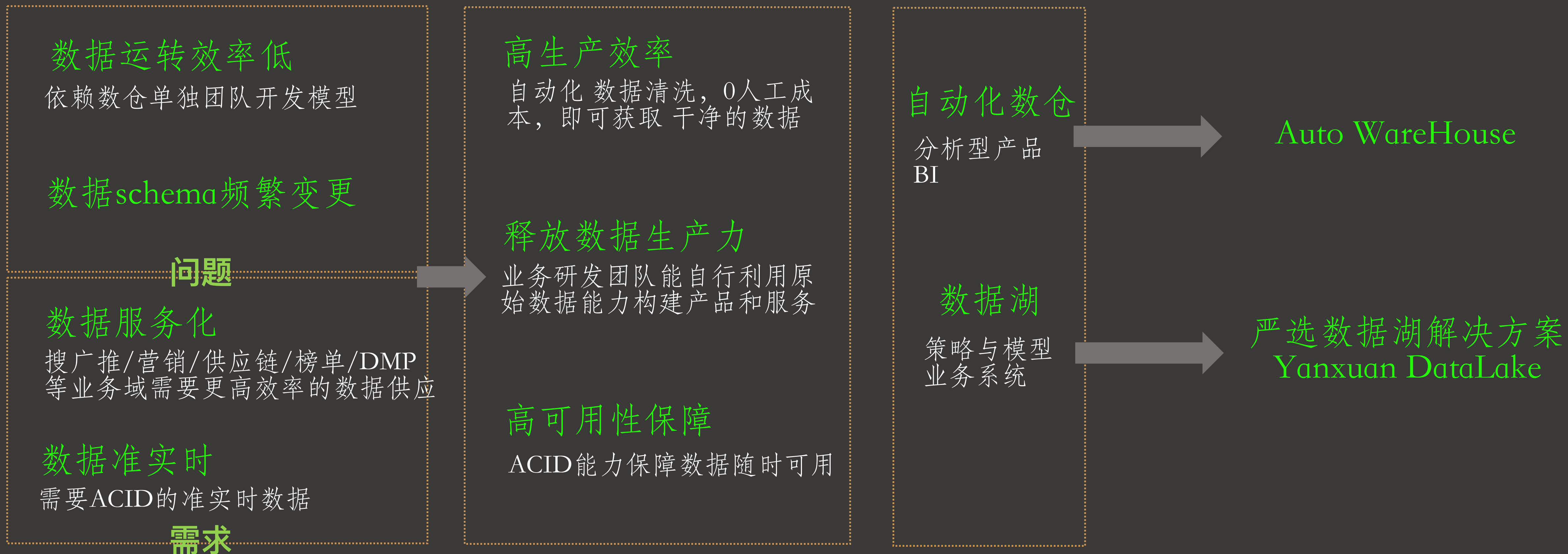
任务治理

基础服务

统一元数据

数据血缘

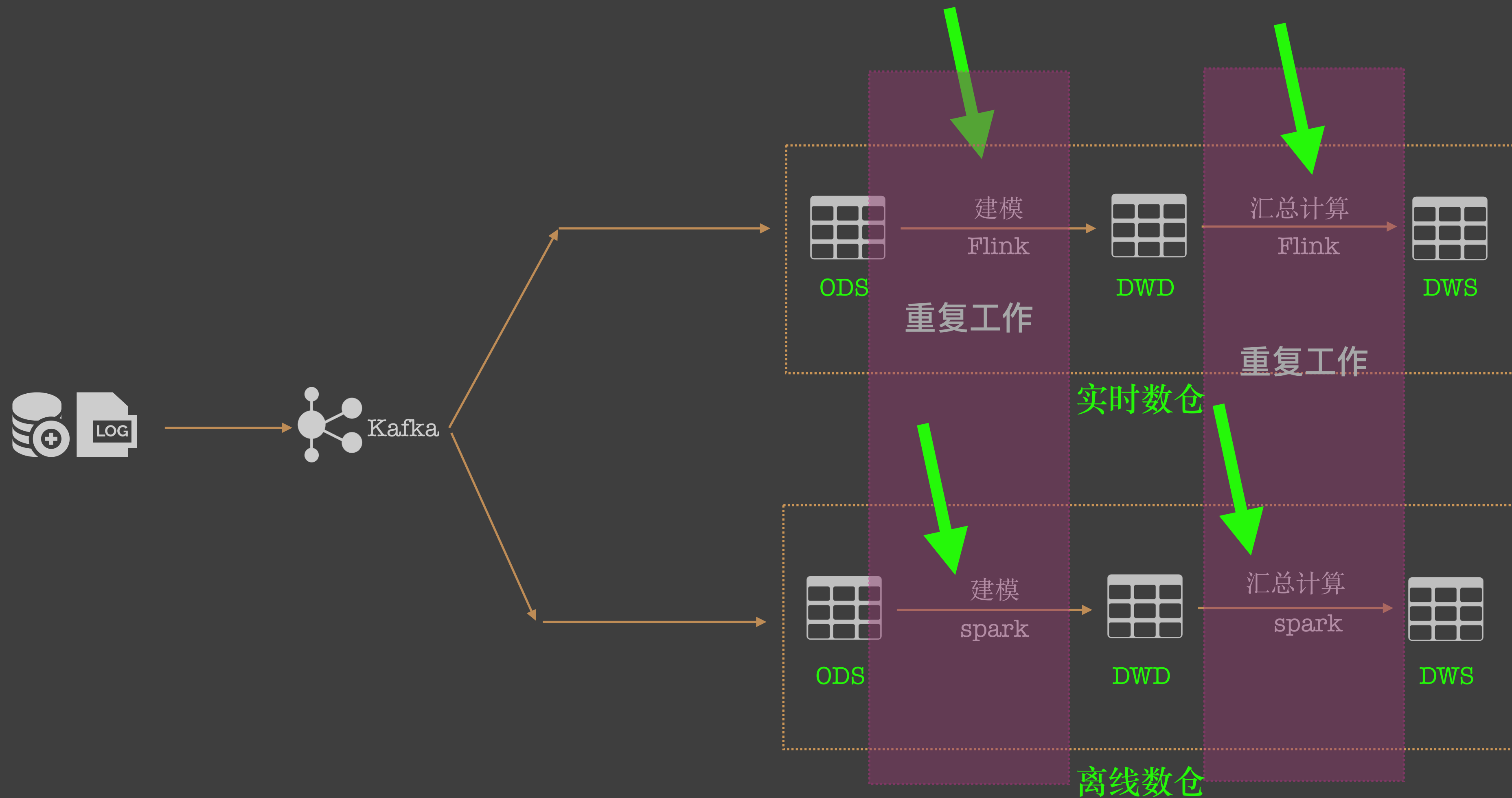
监控服务



问题和需求

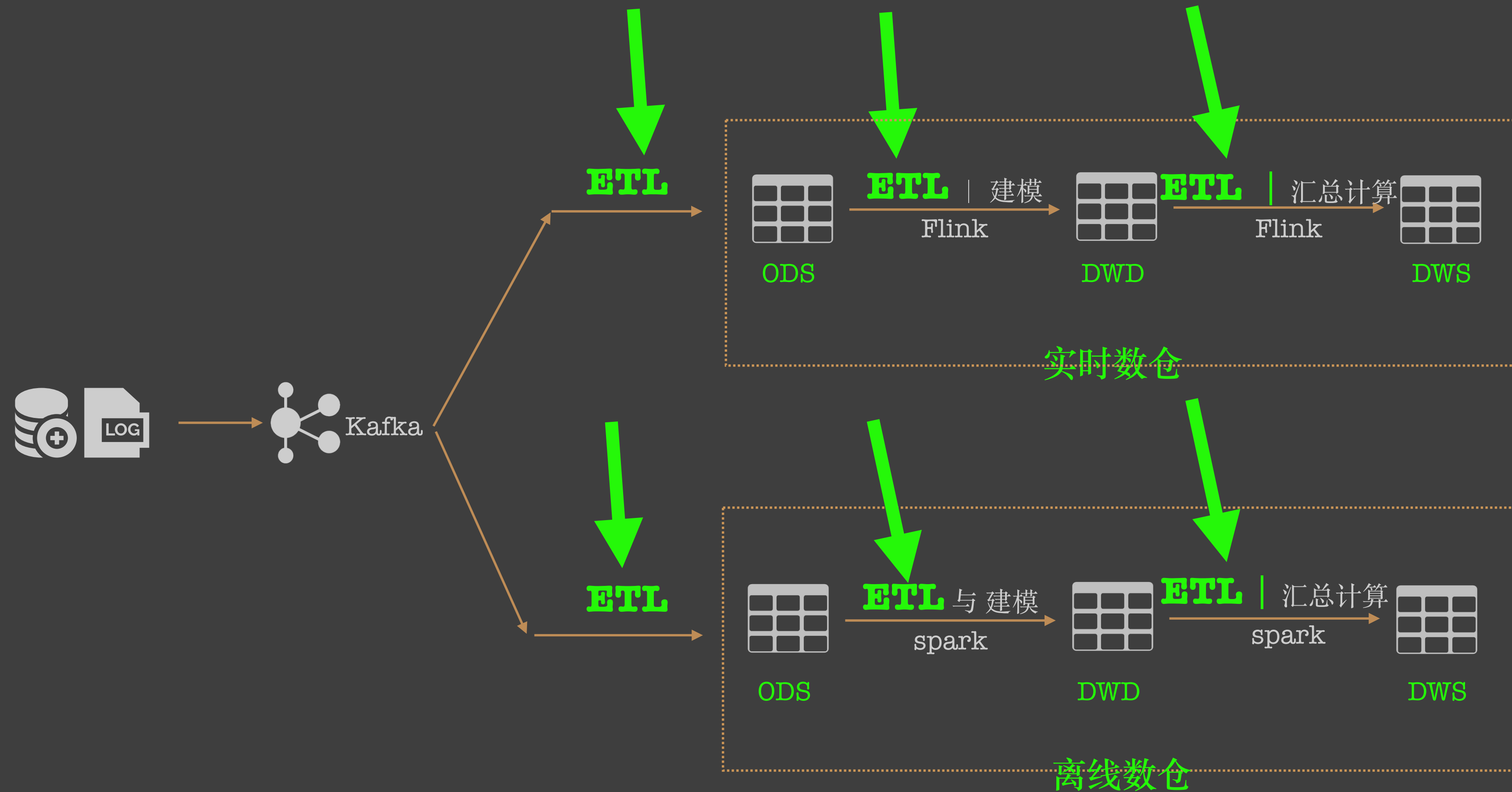
价值

策略



实时和离线数仓:

两条数据流存在着许多重复
但又不完全一致的工作



数仓有很多规范和口径
规范与SQL实现并没有一致

eg. 规范

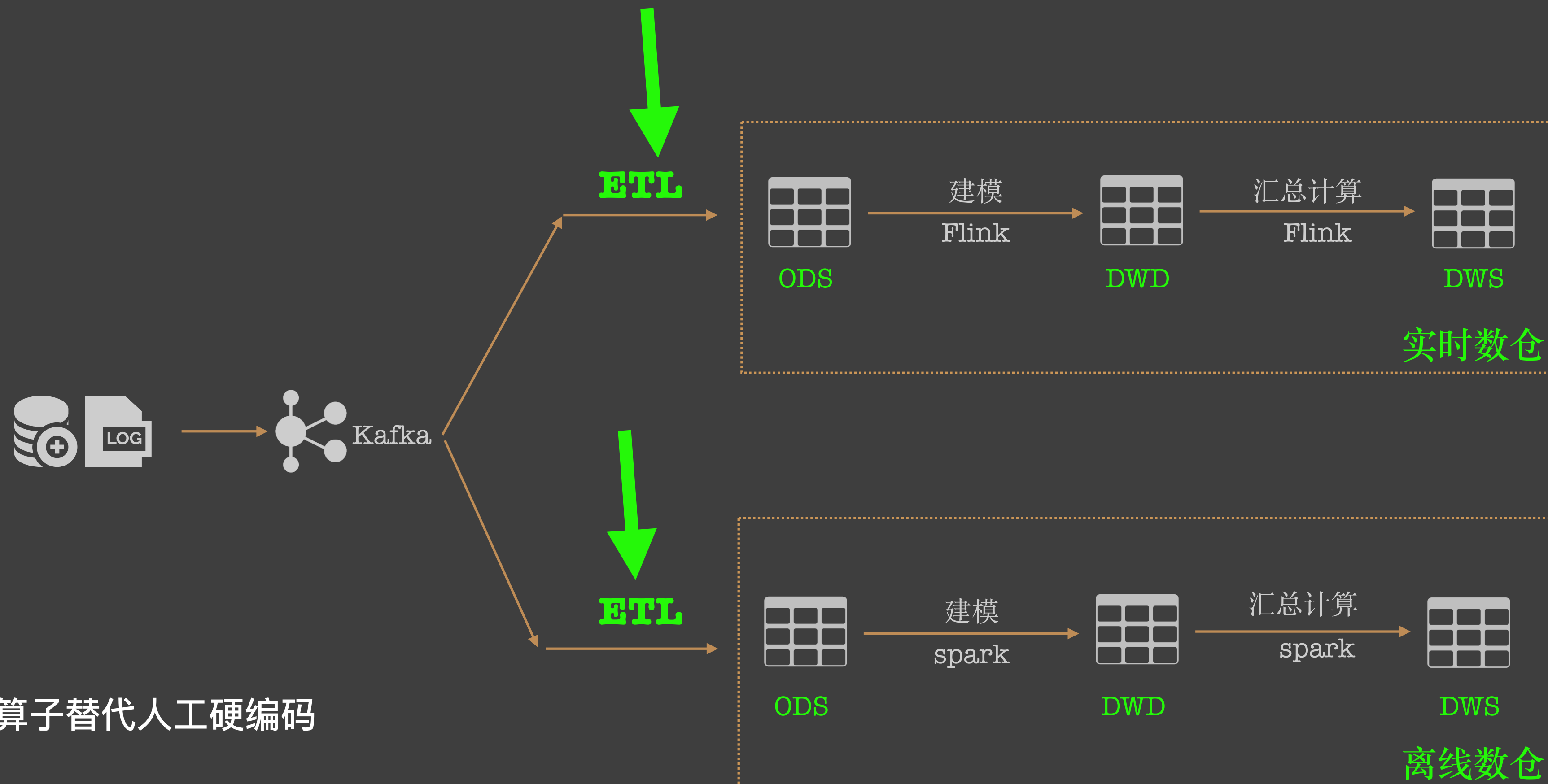
ODS要完成的数据清洗, 在DW和DM都存在

ODS不能有join

DWS 指标的定义和实现不一致

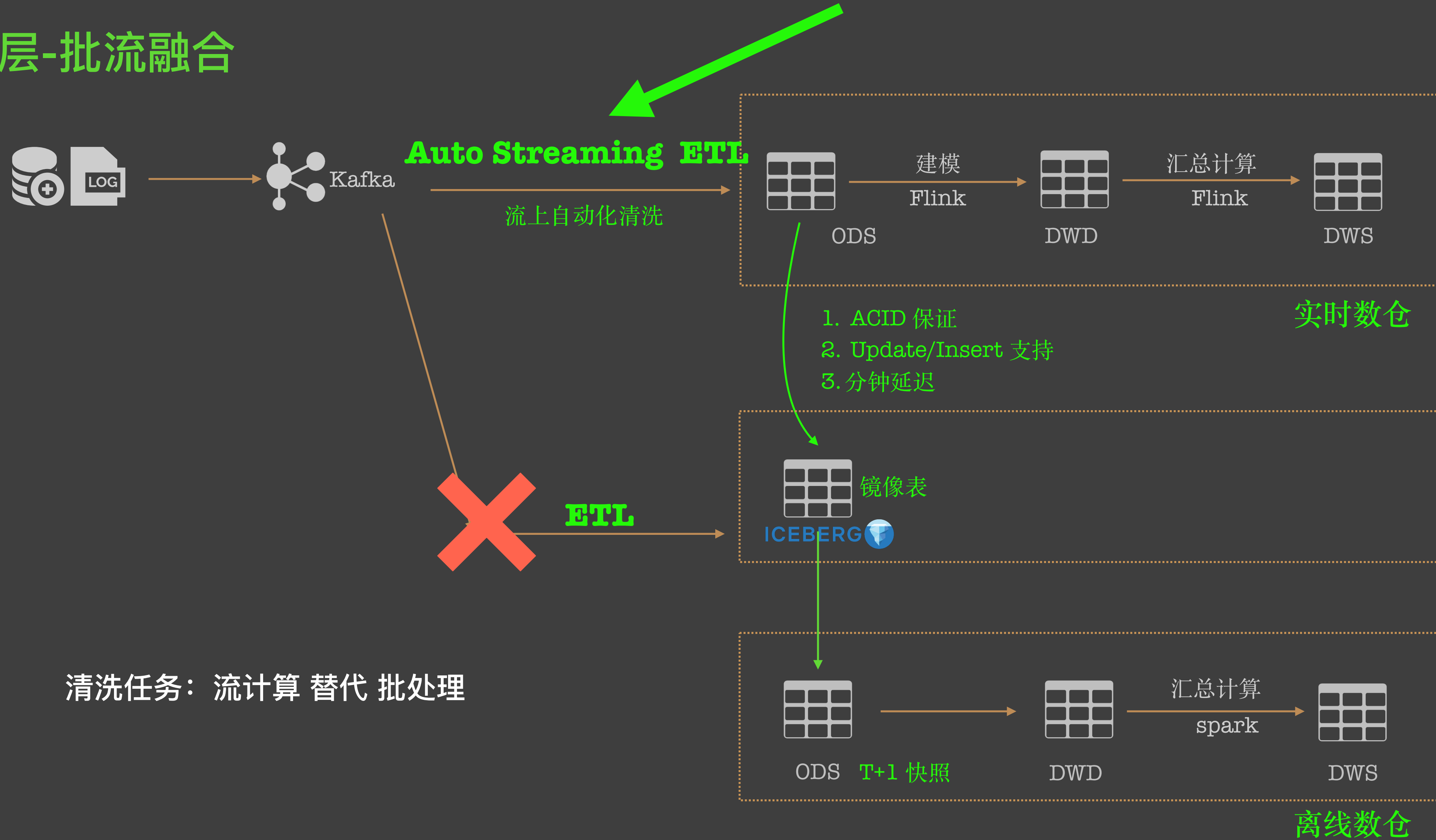
指标的二义性-大量重复计算代码

ODS层-AutoETL

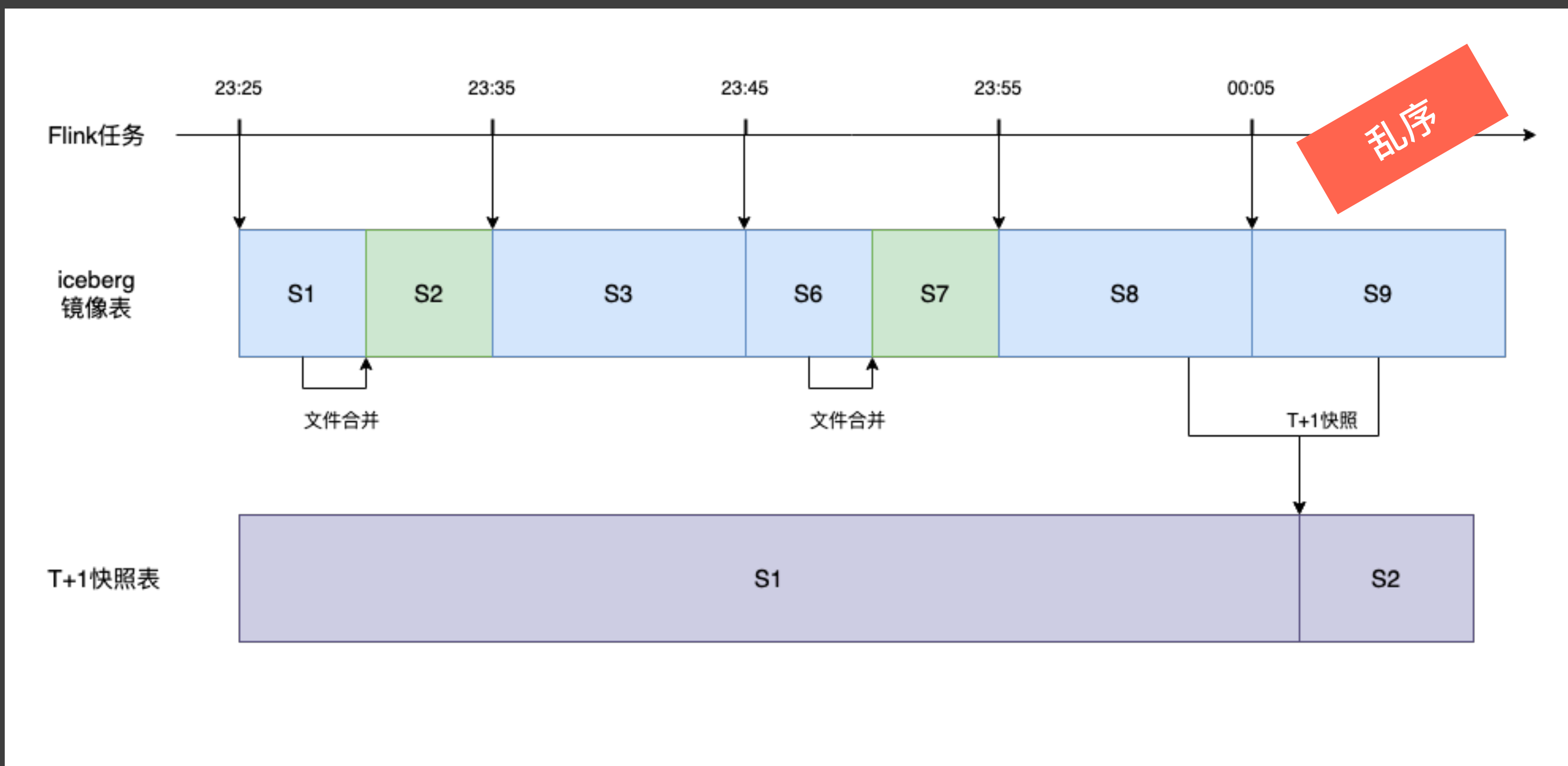


构建清洗算子替代人工硬编码

ODS层-批流融合



ODS层- 批流融合



如何解决乱序问题

Flink-AutoETL任务

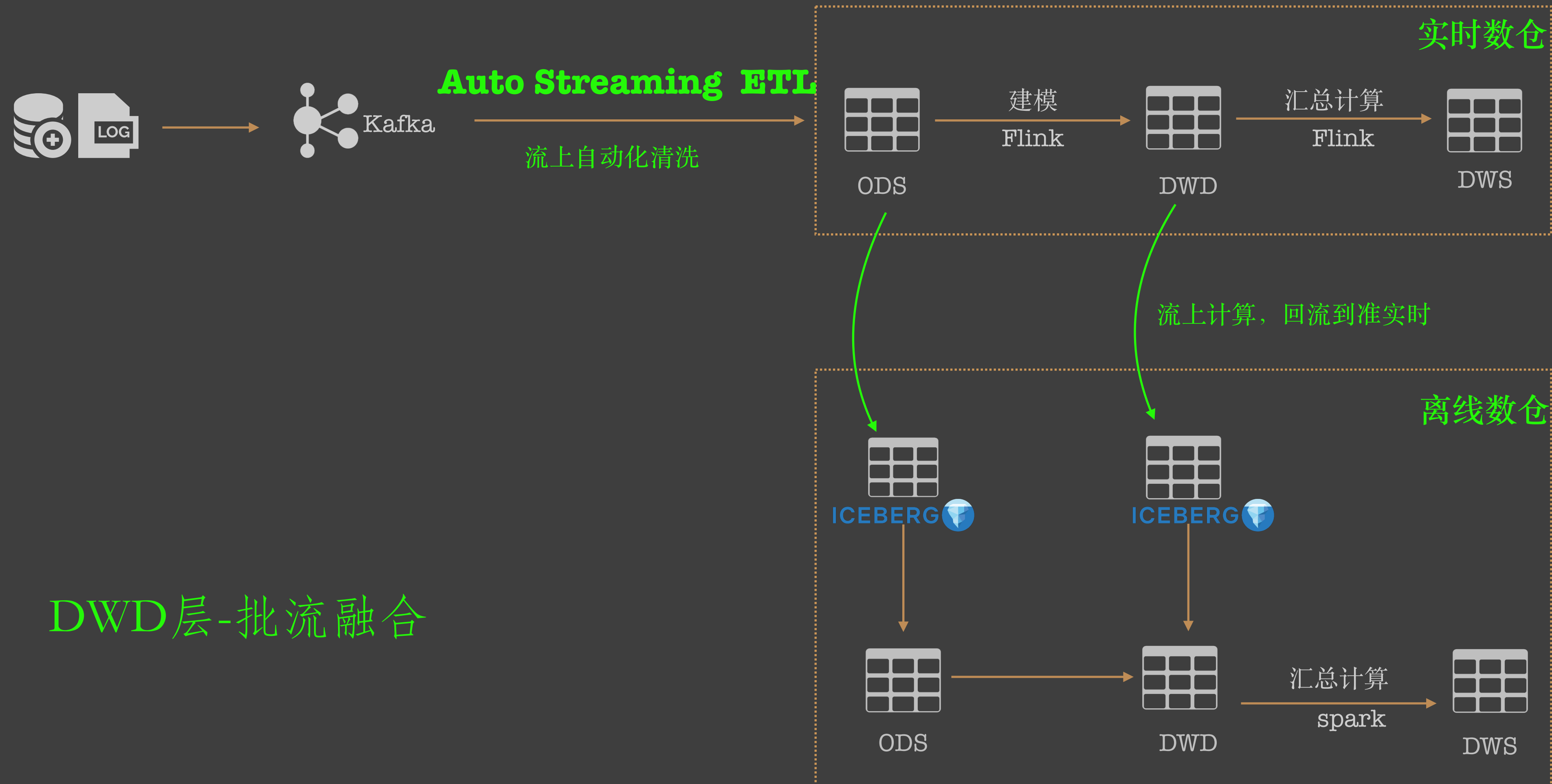
自动化清洗逻辑
保证数据顺序写入Iceberg

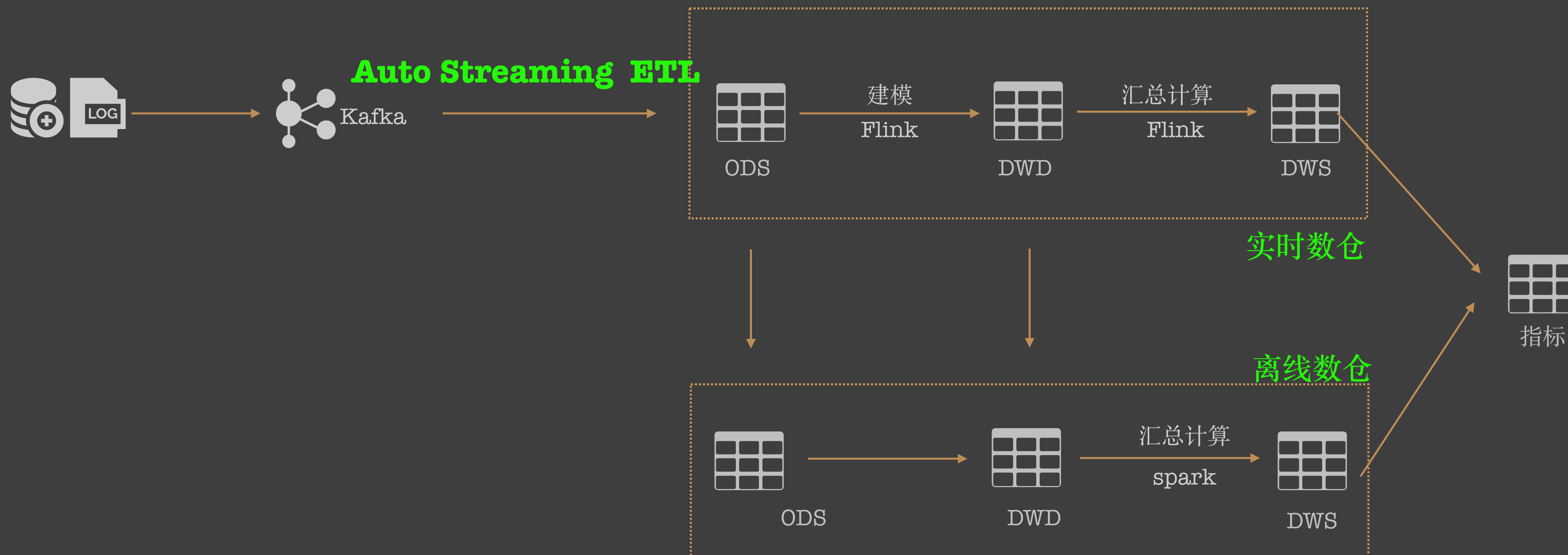
Spark-Compaction任务

eqDeleteFile -> posDeleteFile的转换
posDeleteFile的合并
dataFile的合并优化，重排后写入

Spark-周期快照任务

从iceberg镜像表中合并制作T+1快照

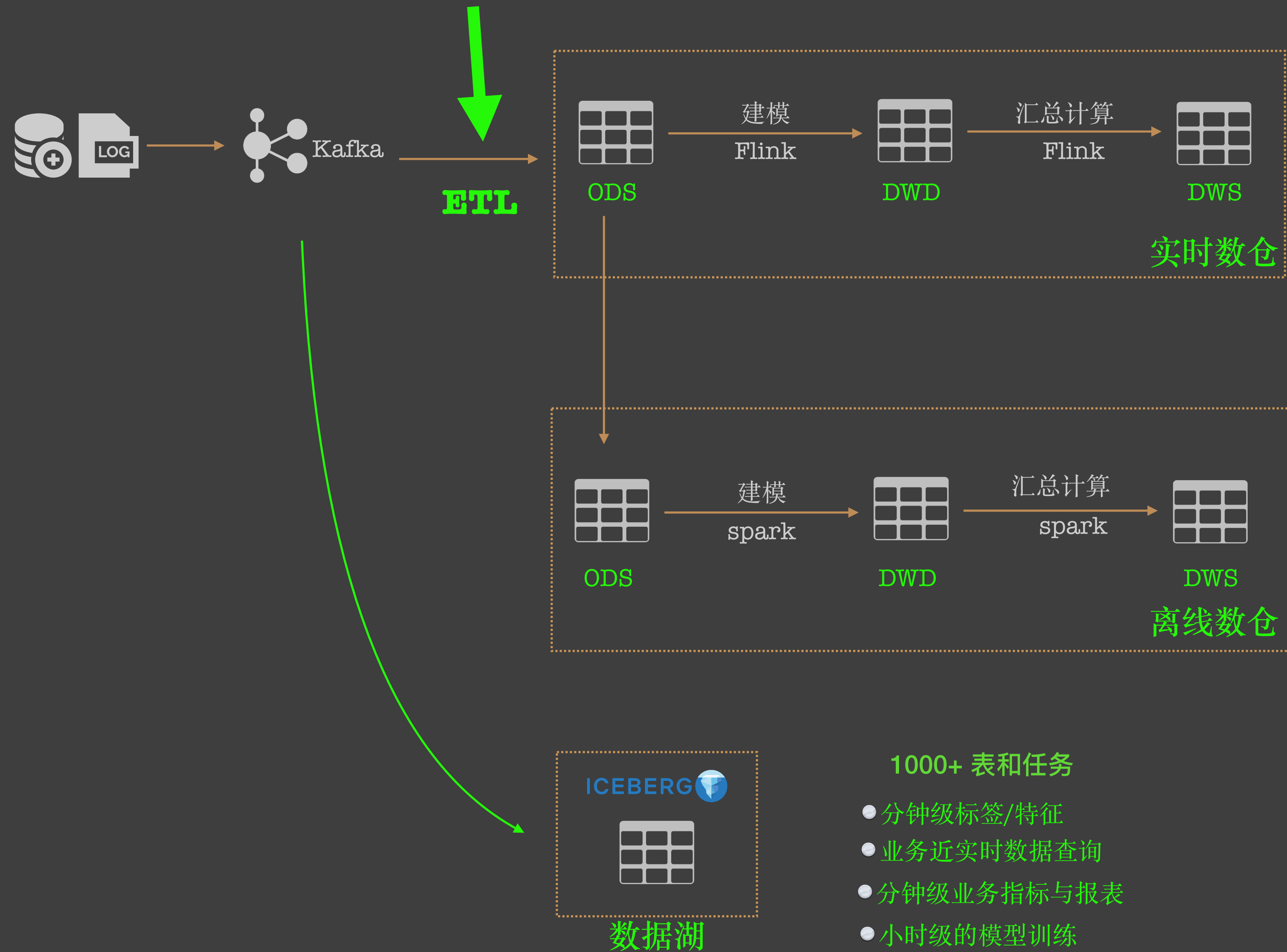


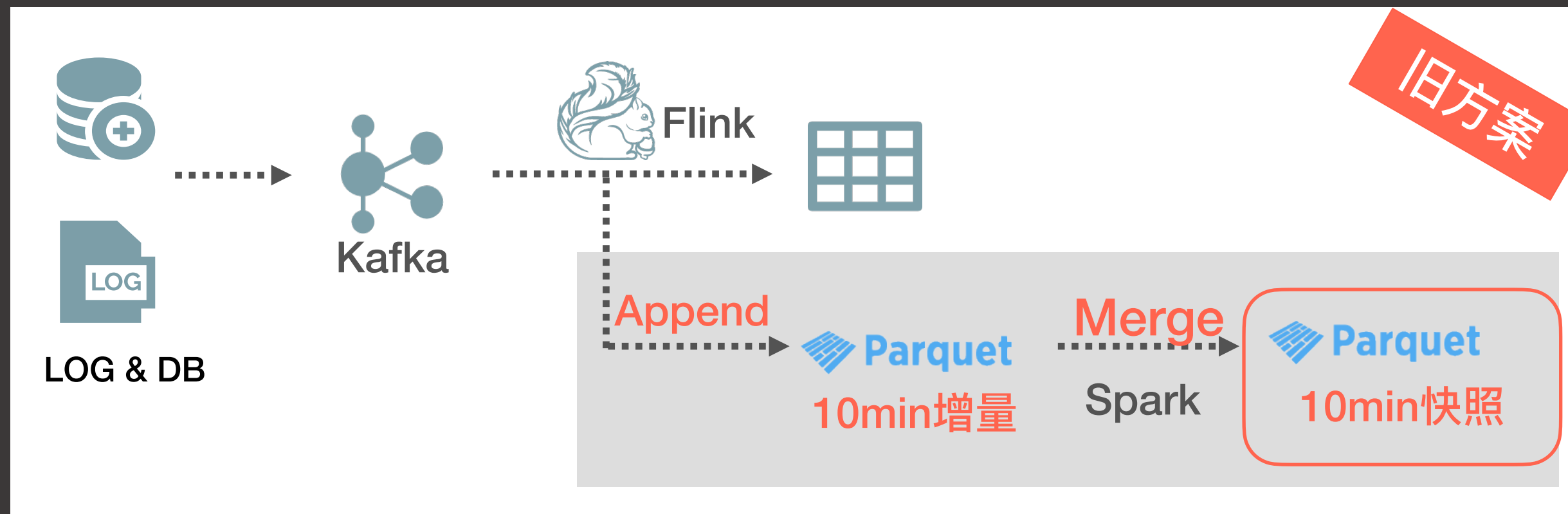


DM-指标自动化构建

指标定义即开发

标准化-> 自动化





近实时数据方案

缺陷

短周期快照10/30/60分钟级别的快照数据

1. 实时性无法满足，实时计算门槛高，按需开发成本高
2. 以空间换时间，浪费大量的存储资源
3. 没有ACID，更新期间并发访问不可用，可用率60%



Iceberg存储方案



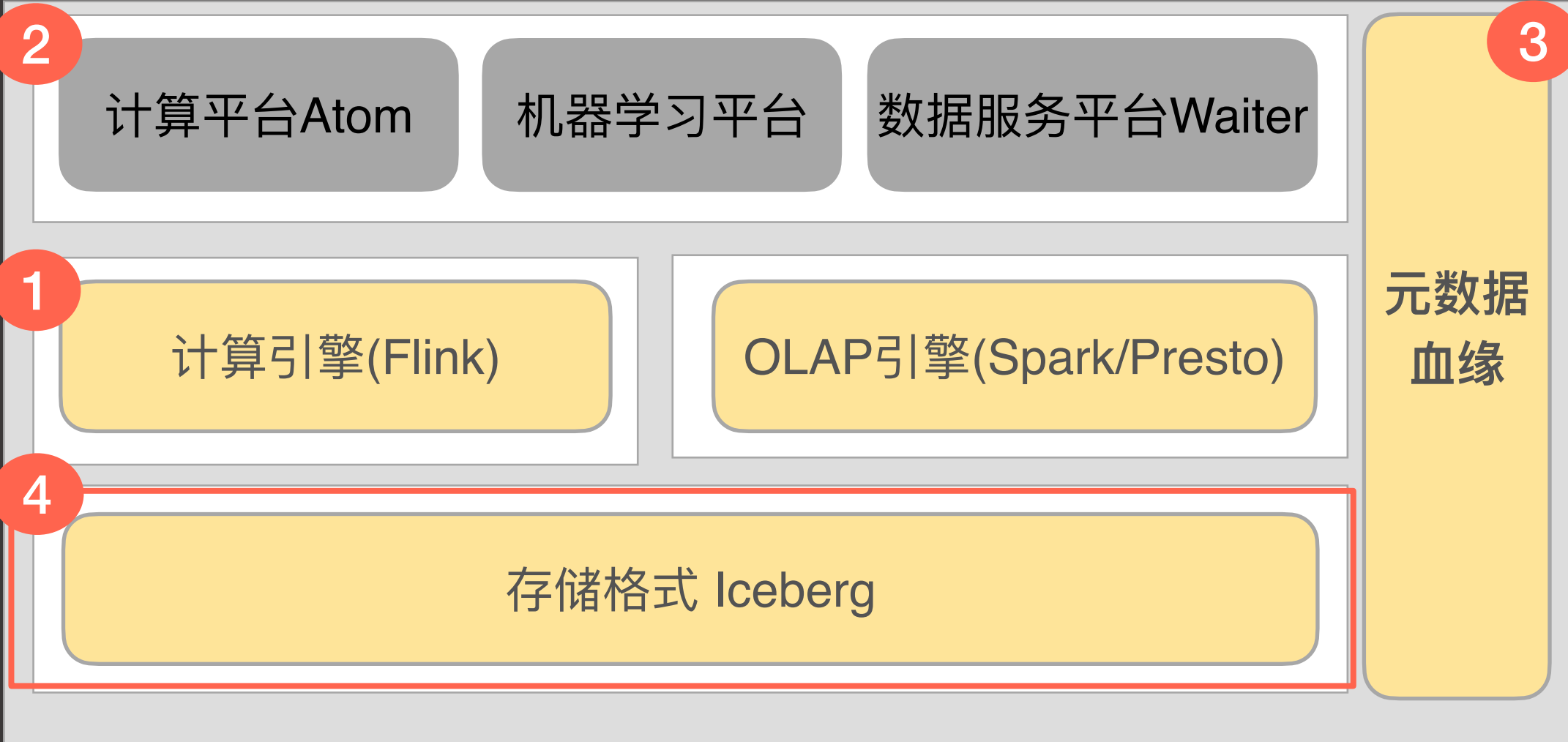
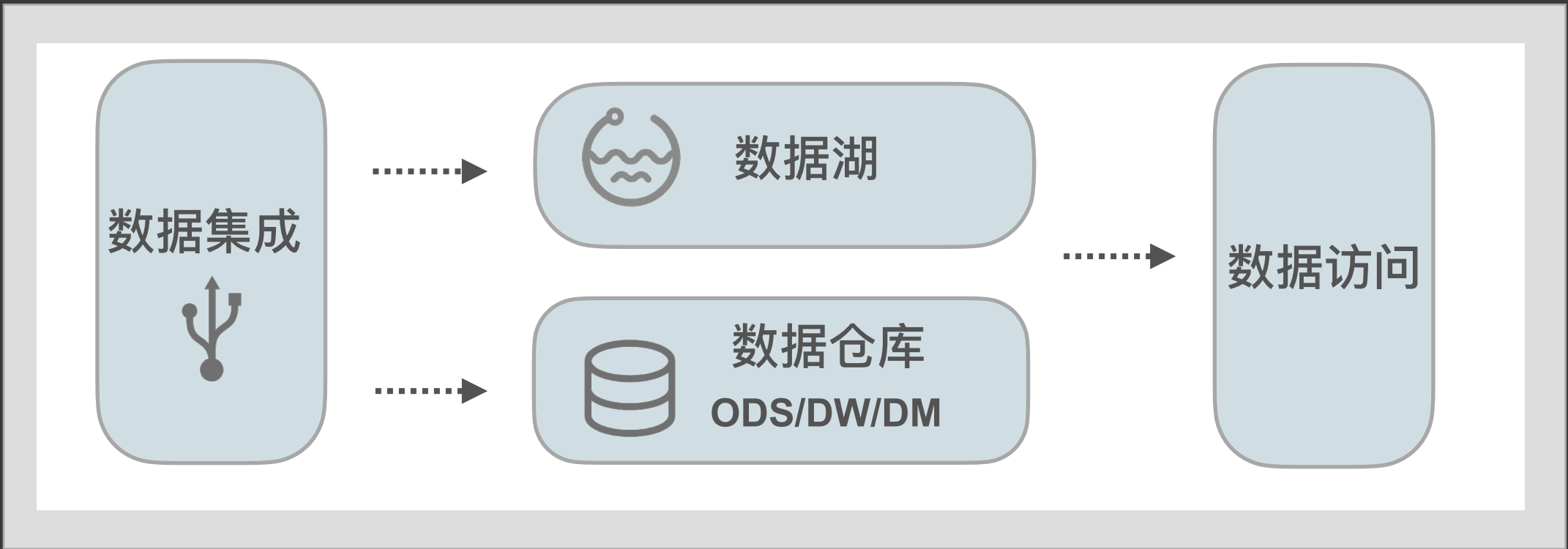
数据延迟



可用率



节省资源



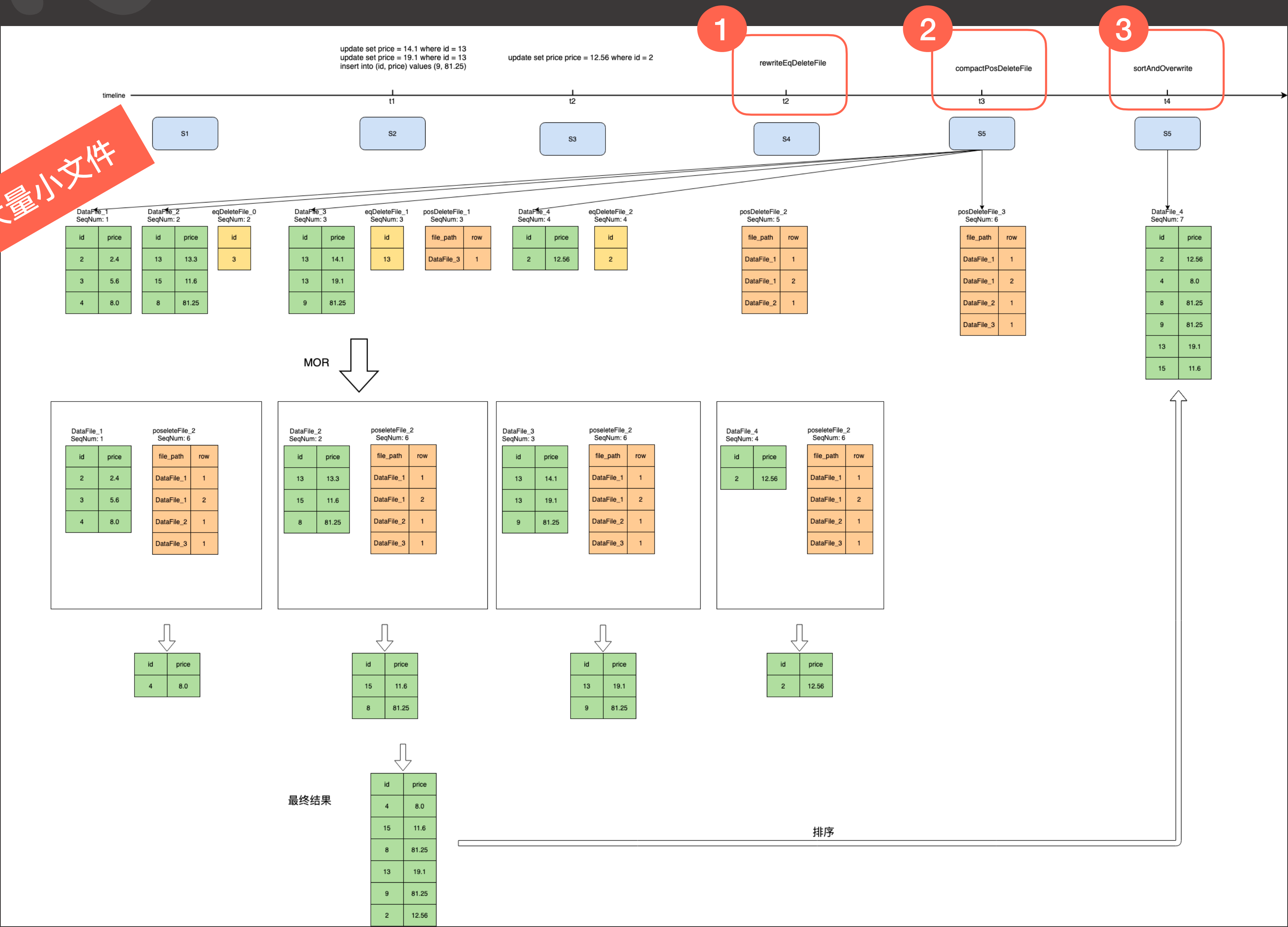
基础设施的建设

0故障

- 1 计算引擎层: Flink/Spark/Presto社区版本支持
- 2 统一元数据及血缘: Iceberg元数据的管理和查询
血缘的接入-提供灰度自动切换的保障
- 3 开发平台层: 计算平台/机器学习平台/数据服务平台
- 4 Iceberg的优化 Compaction的优化, 性能提升10倍
- 5 DataCheckTool保证数据100%一致

数据仓库 与 数据湖 平台架构上并无本质区别, 并不引入额外的成本

大量小文件



问题及原因

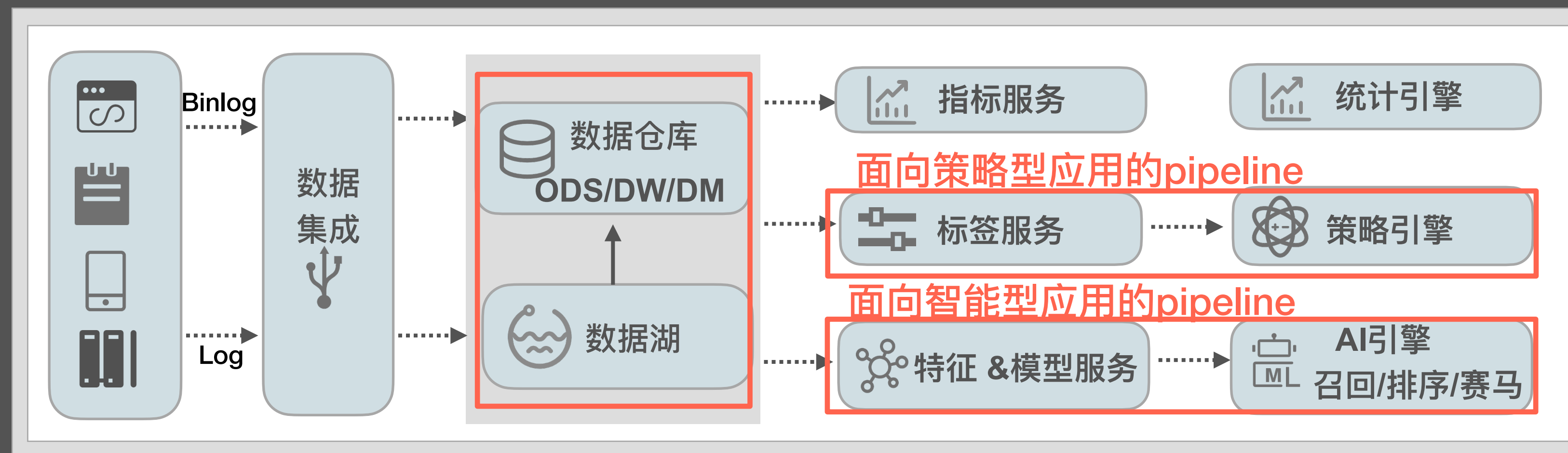
Iceberg表读性能差，文件多影响MOR性能

主要优化点

- 1. rewrite EqDeleteFile ->PosDeleteFile
- 2. compact PosDeleteFile
- 3. order by index & overwrite

平均读取性能提升10倍

中台



平台



数据治理

指标治理

标签治理

表治理

任务治理

基础服务

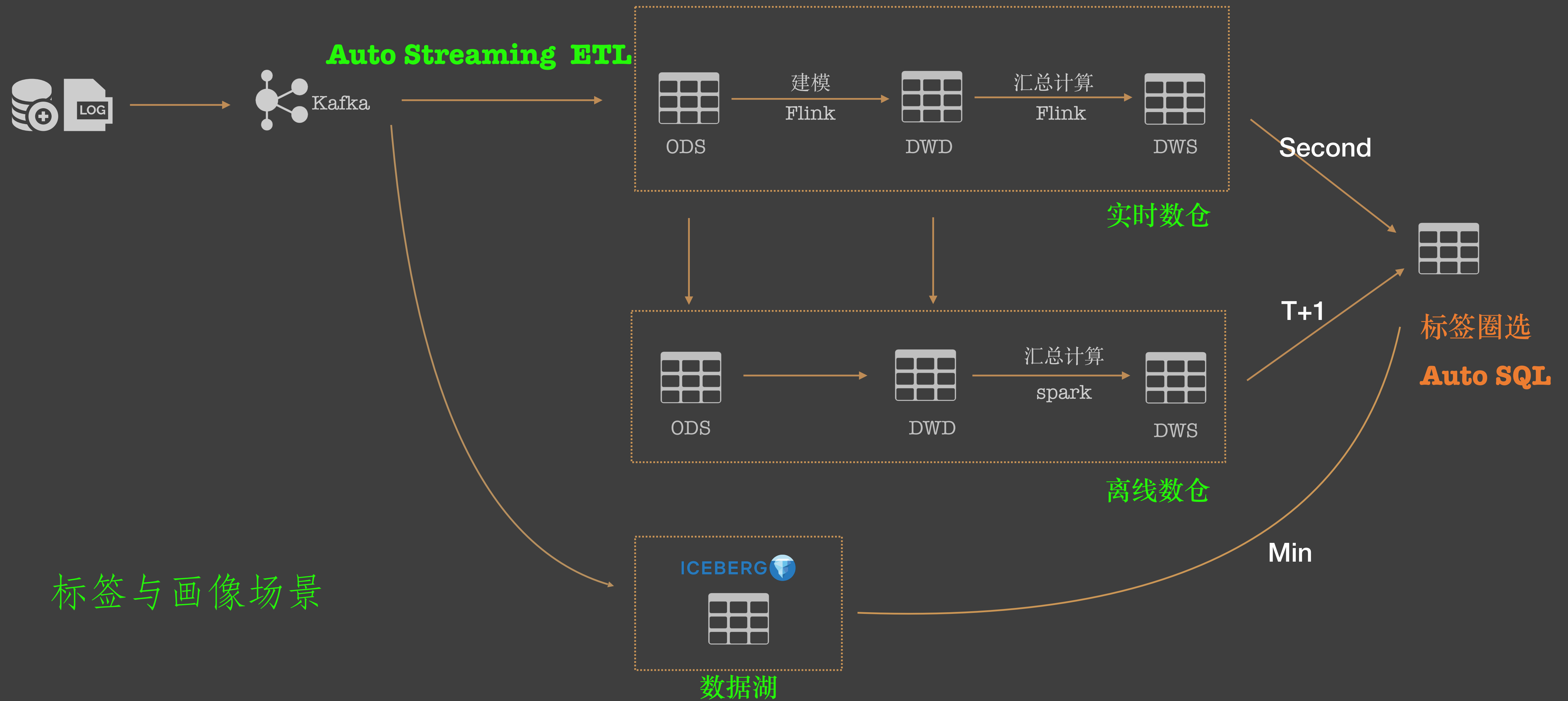
统一元数据

数据血缘

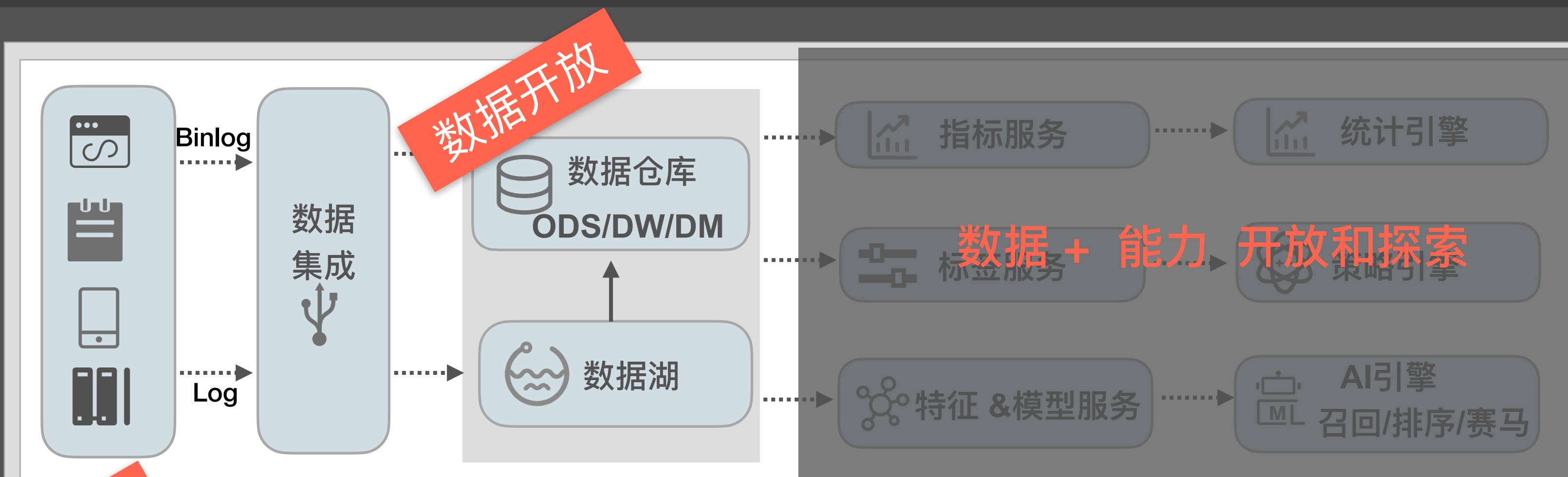
监控服务

特征和样本近实时存储 - 支持模型实时训练

标签服务 - 支撑近实时的用户圈选运营活动场景



中台



能力开放

平台



数据治理

指标治理

标签治理

表治理

任务治理

基础服务

统一元数据

数据血缘

监控服务

数据湖开放

如何使用湖和仓？平台提供数据集成-数据访问能力

近实时数据查询 - 支撑供应链/商品中心/财务等业务实时数据查询

1

数据驱动技术体系

数据分析-> 数据决策

2

数据及算法中台

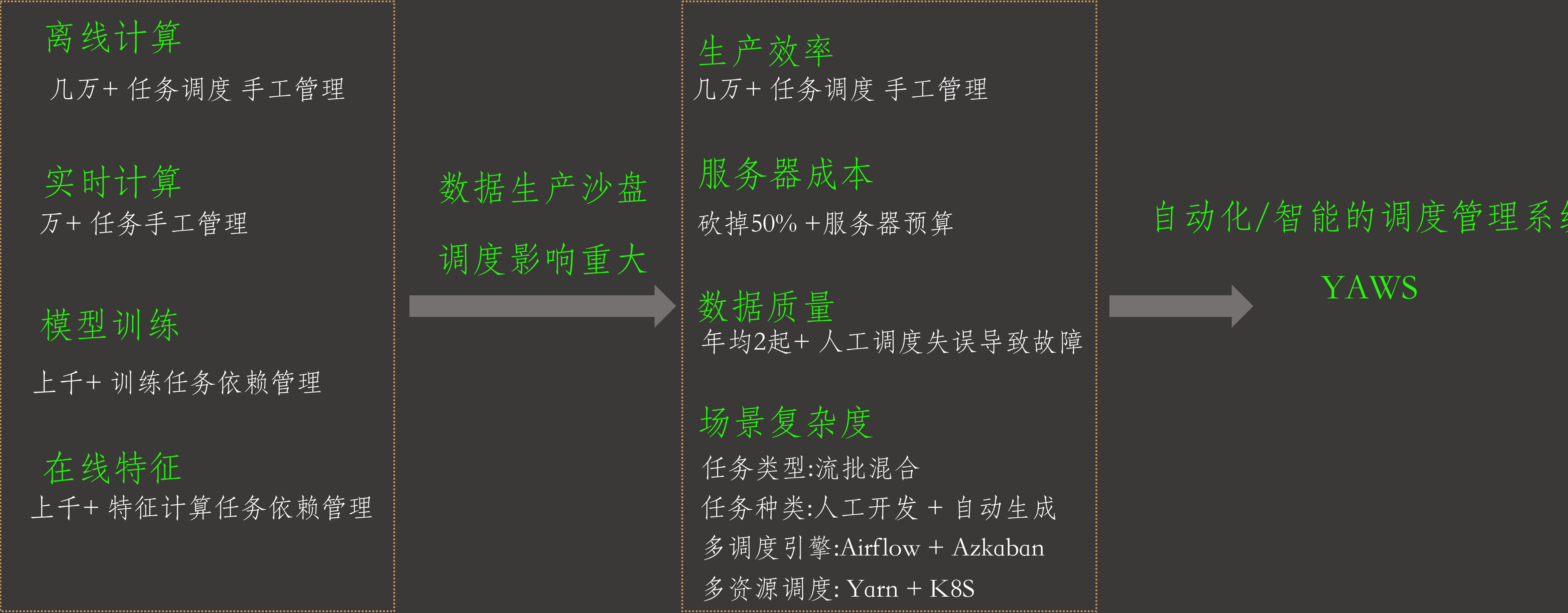
DataLake
AutoDW

3

数据及算法平台

智能任务调度
Cloud Native





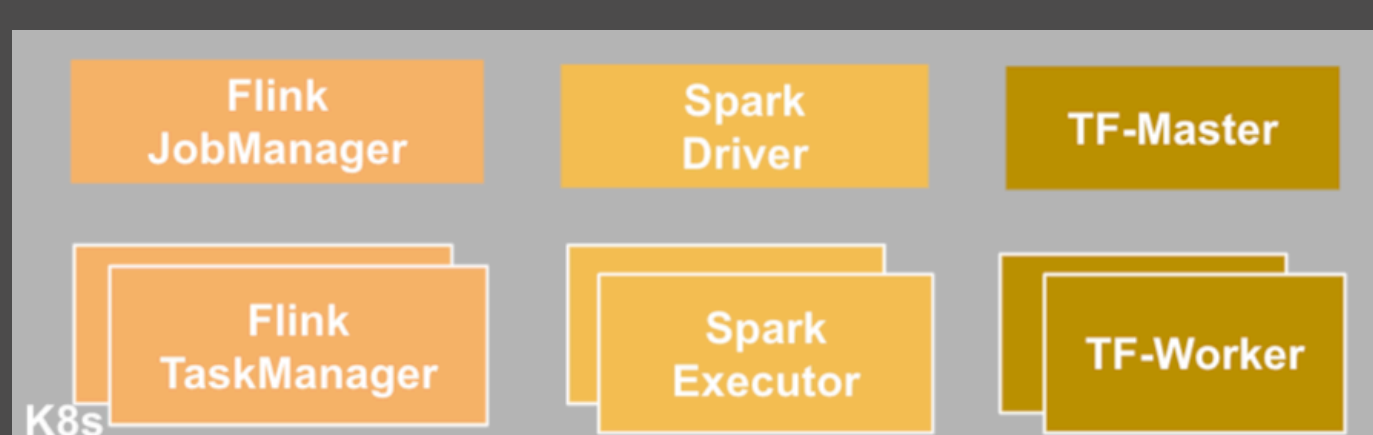
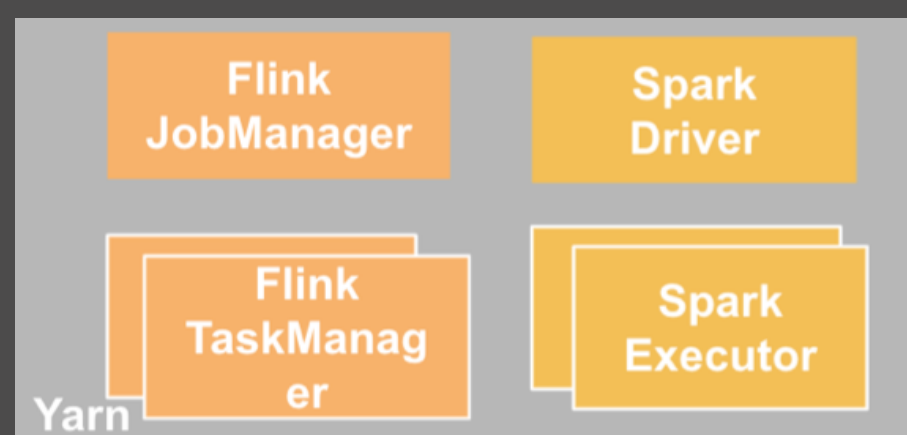
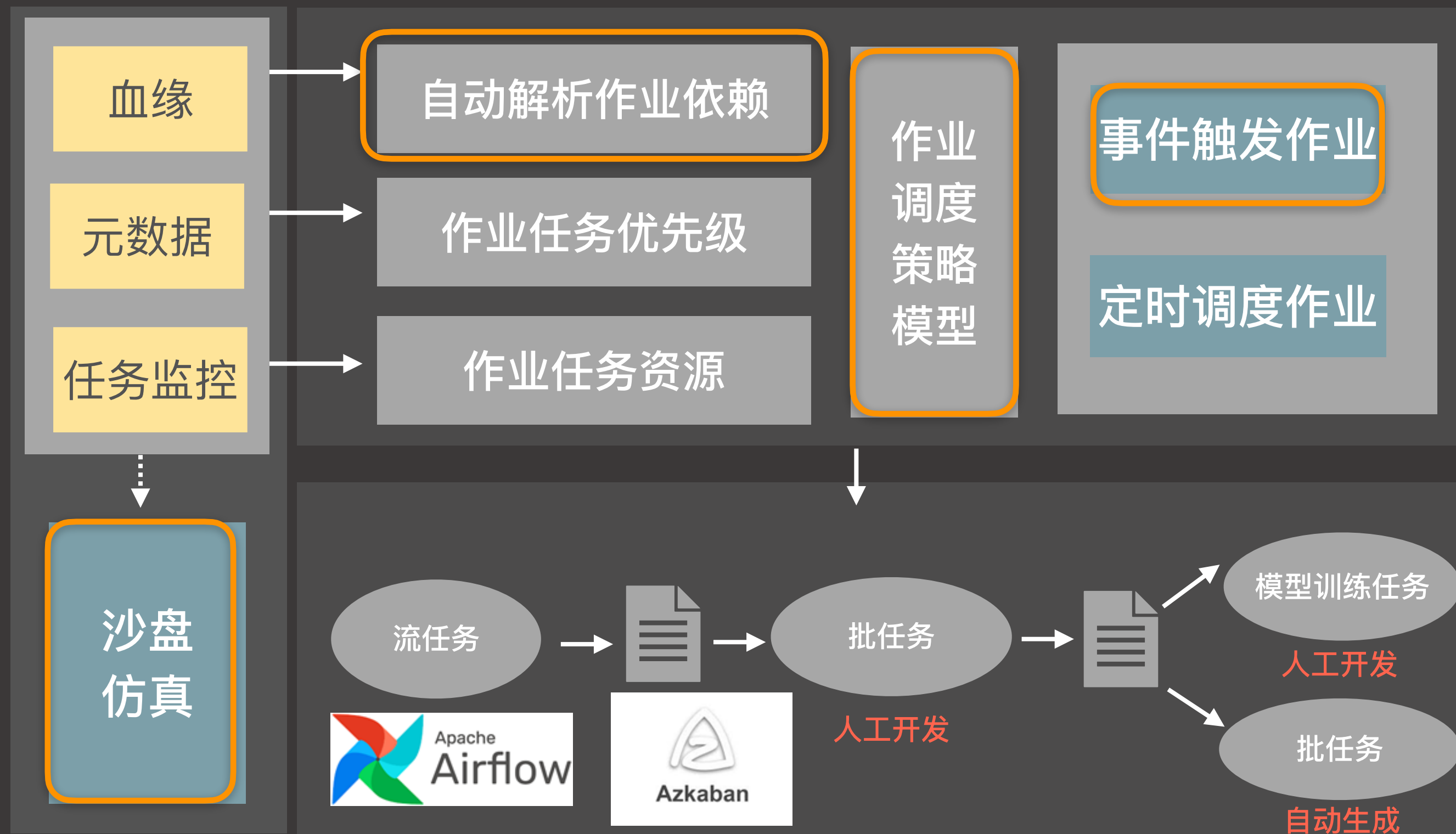
现状

推演

价值

策略

统一任务管理



Cloud Native 架构

批/流/AI跨多平台任务提交与管理

支持Yarn/K8S调度，屏蔽基础设施细节

统一的API服务，支撑上层应用计算任务

eg. 日志平台/标签服务/指标服务等

调度策略

全链路数据血缘调度，不依赖人工配置，准确率99.9%

事件触发作业，解决流批混合调度，数据准确率100%

沙盘仿真模块，量化策略效果，准确率99%

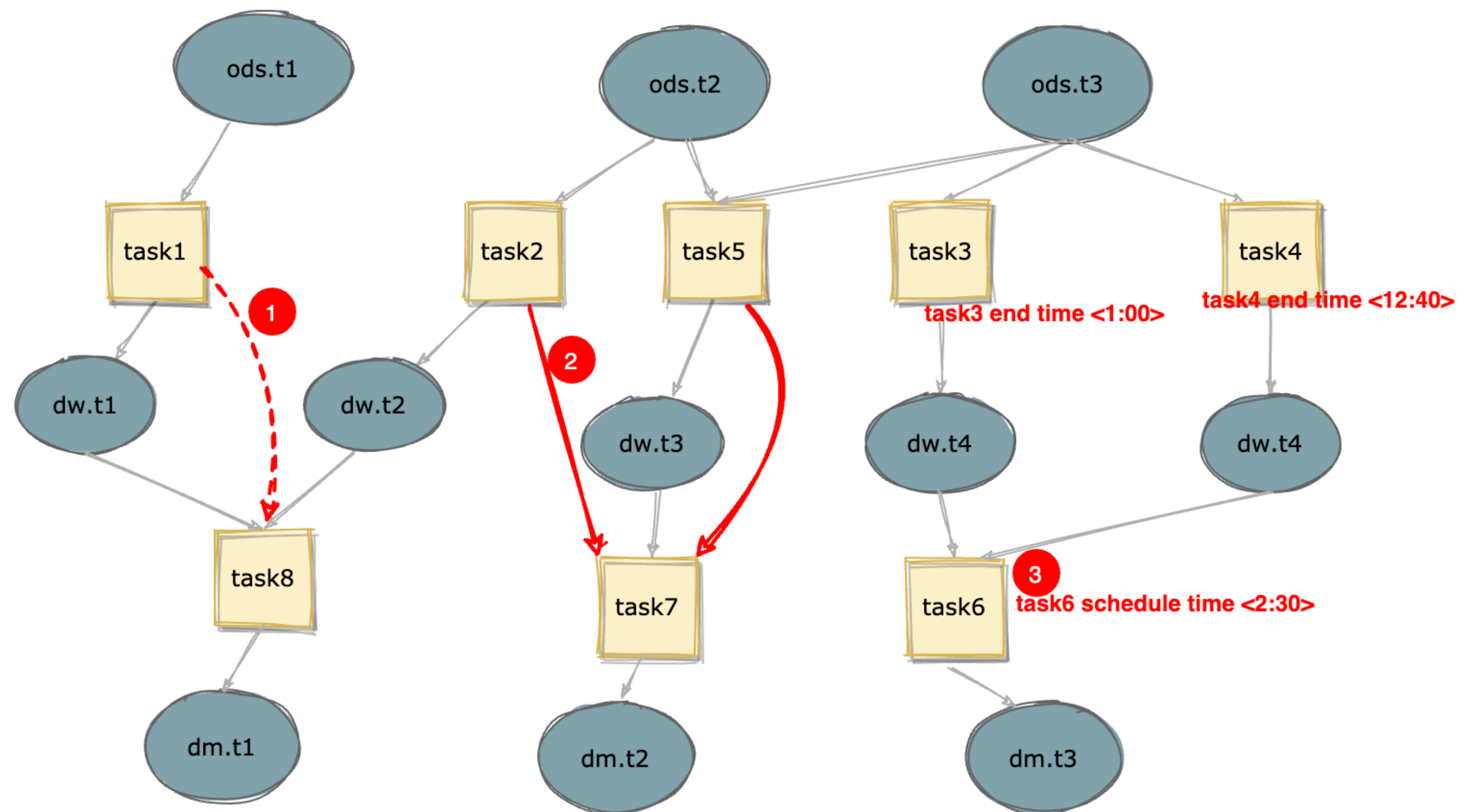
智能调度策略，血缘+任务优先级+任务资源的策略模型

调度执行

调度引擎-Airflow 和 Azkaban

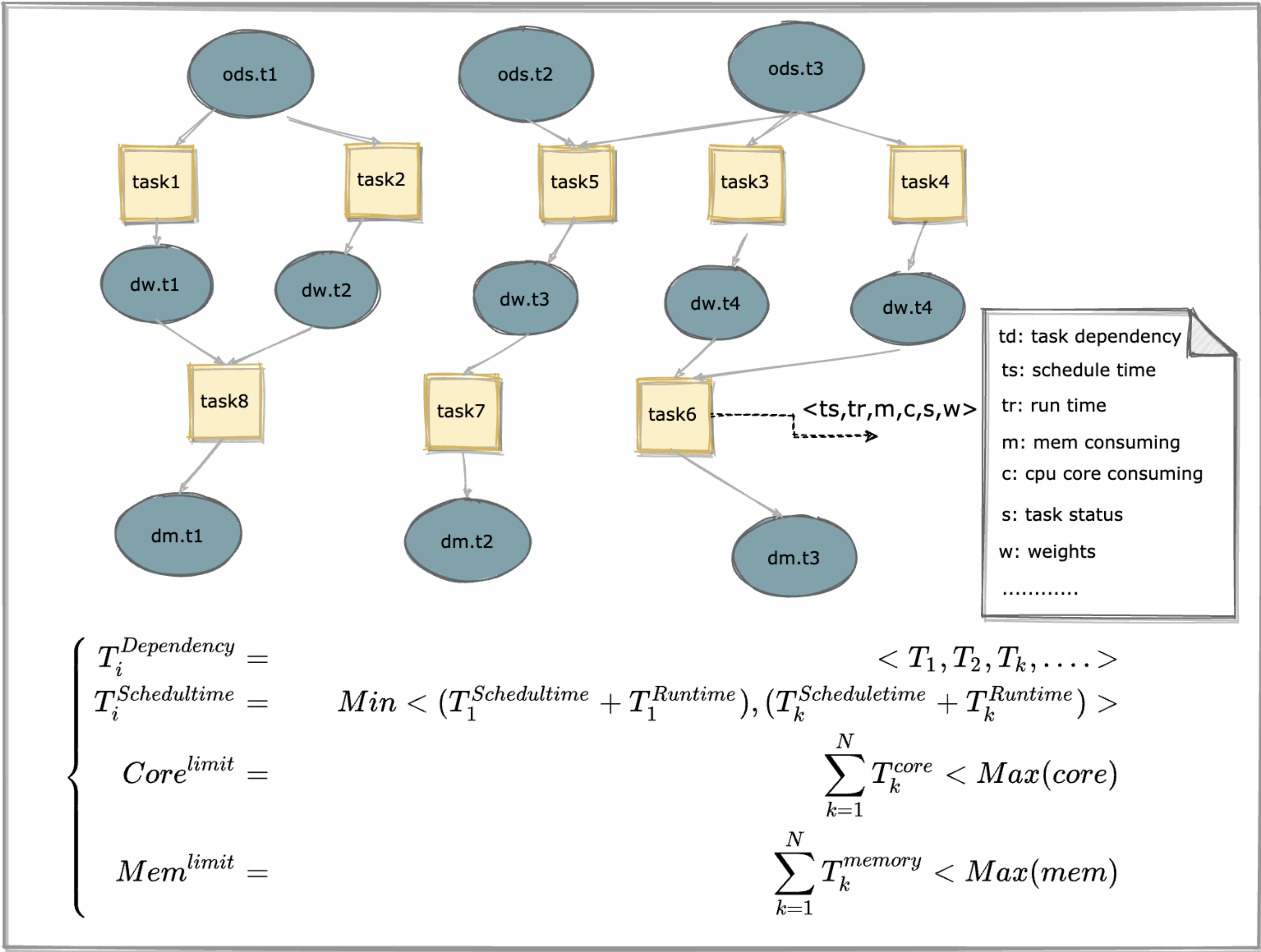
K8S- 批处理/流计算/算法训练任务资源混合调度

Data & AI 混合调度，AI-Pipeline 端到端的流调度管理

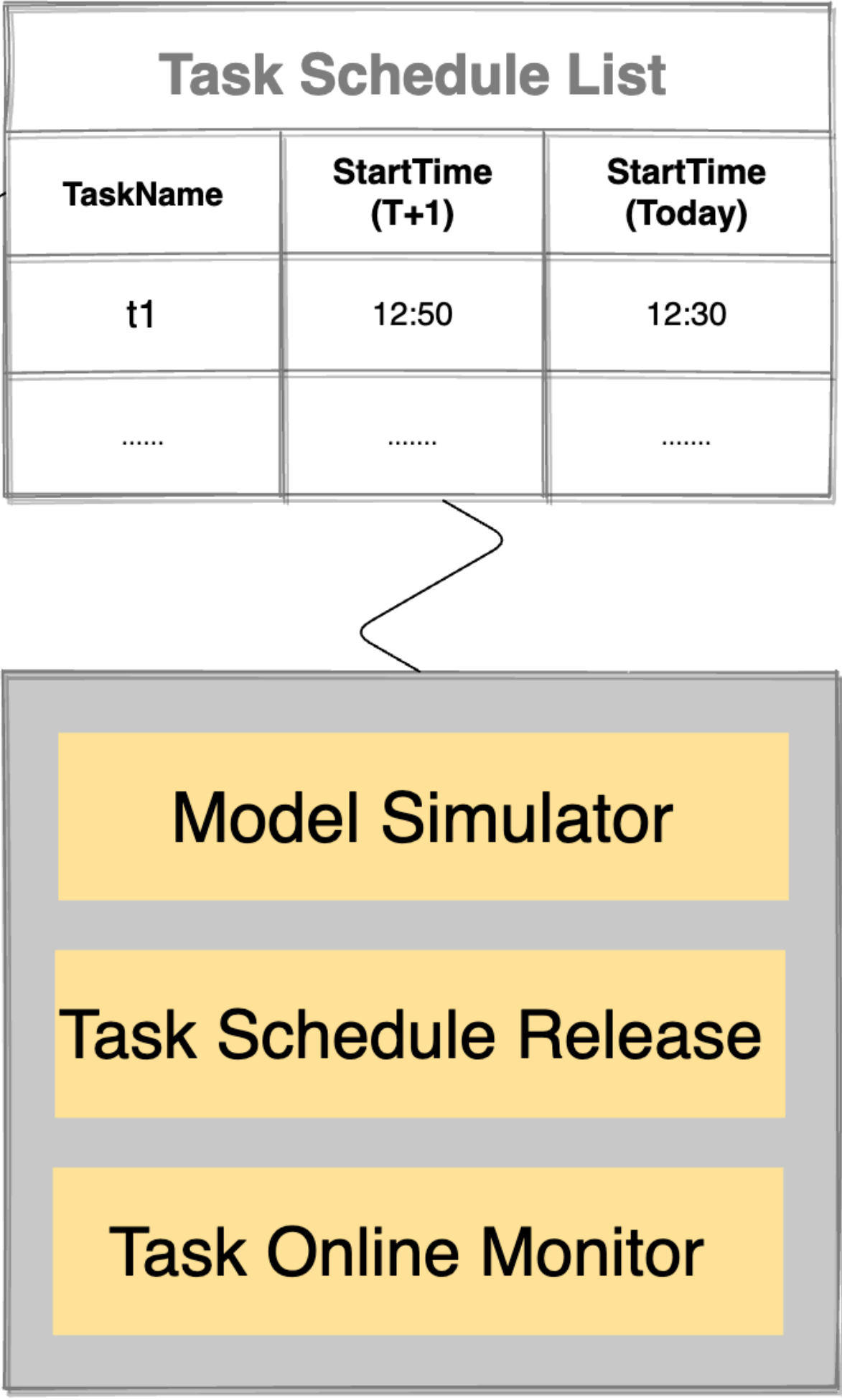


原来存在问题:

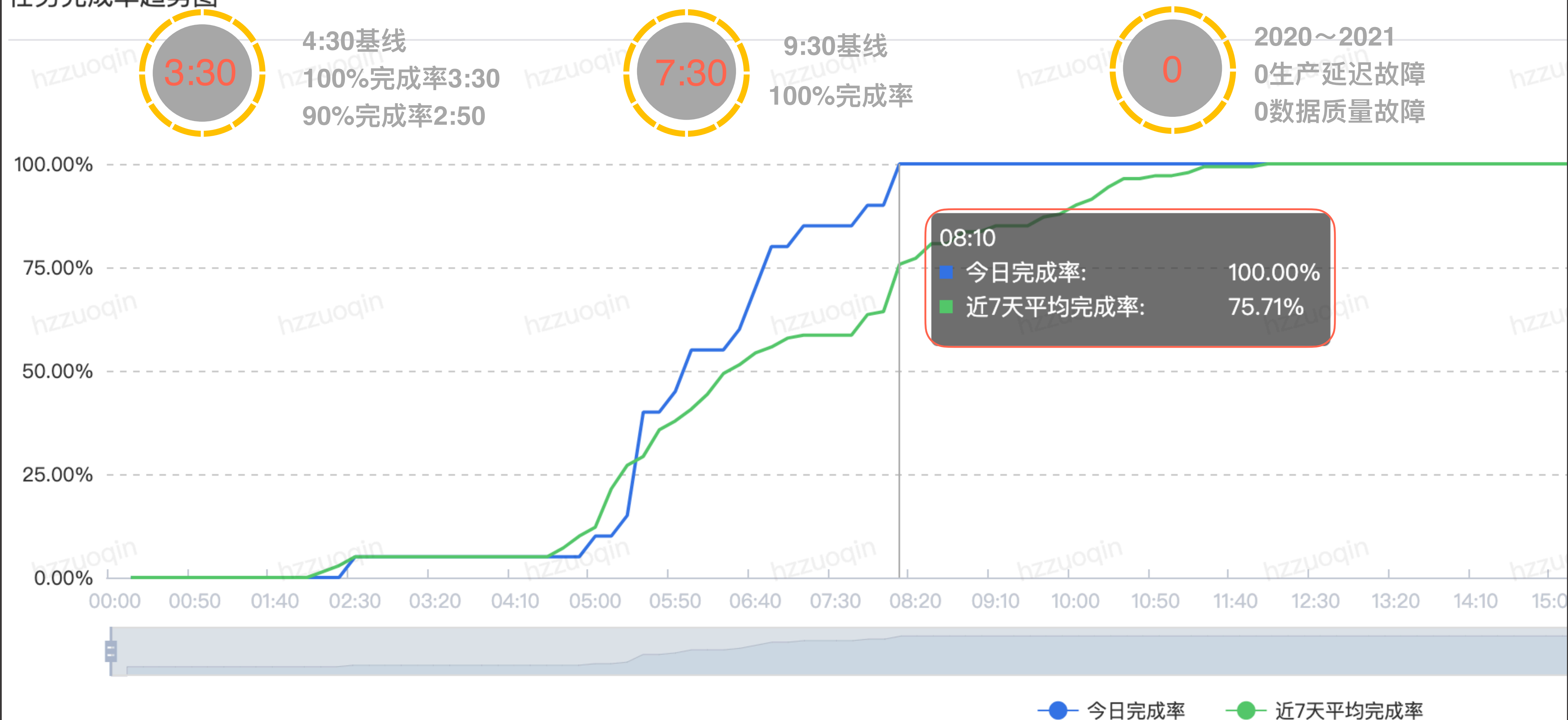
- 1 缺依赖导致数据准确性问题
- 2 多依赖导致下游不必要的等待
- 3 调度时间设置不优化

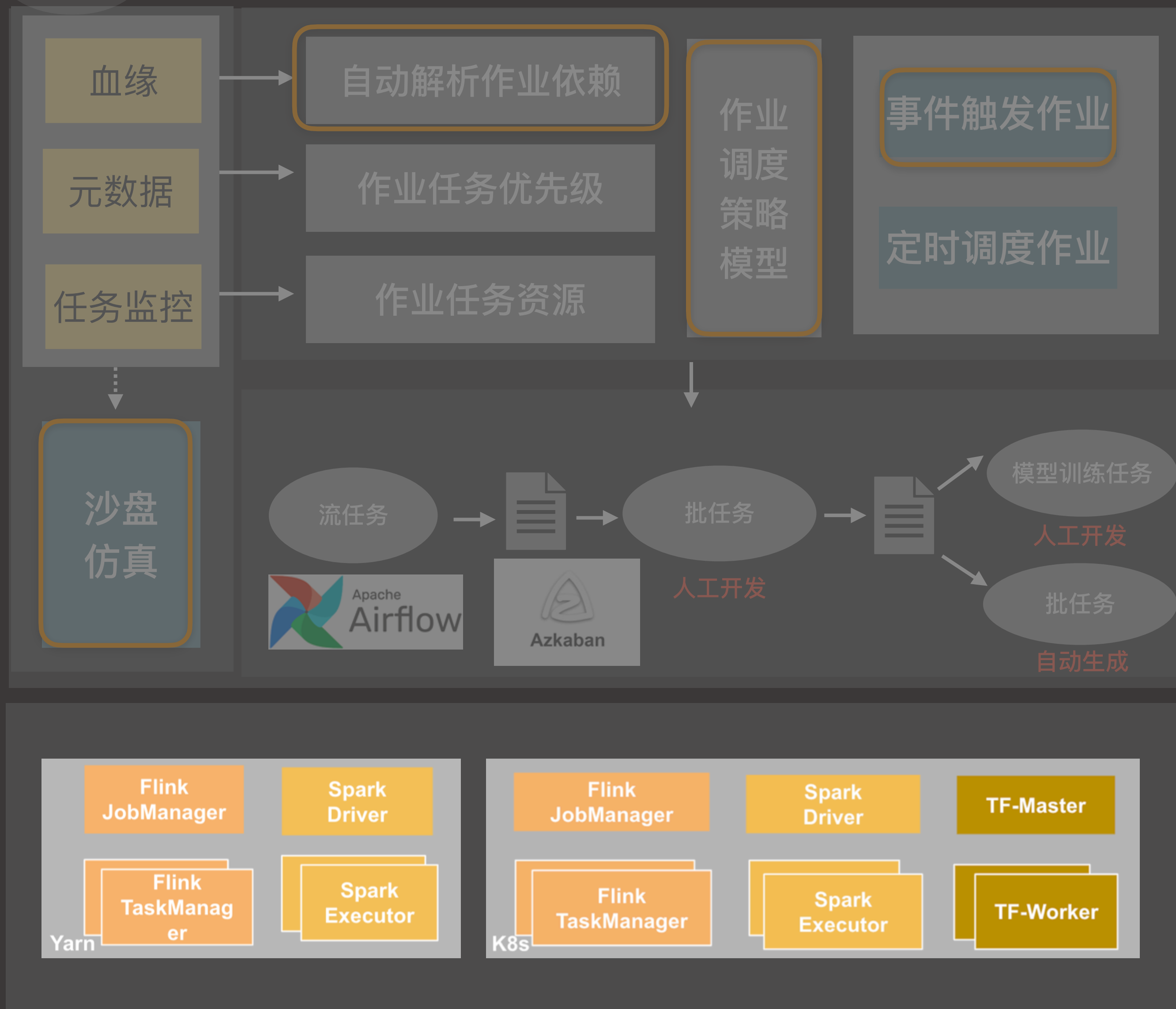


Lineage-Based Model for Schedule



任务完成率趋势图





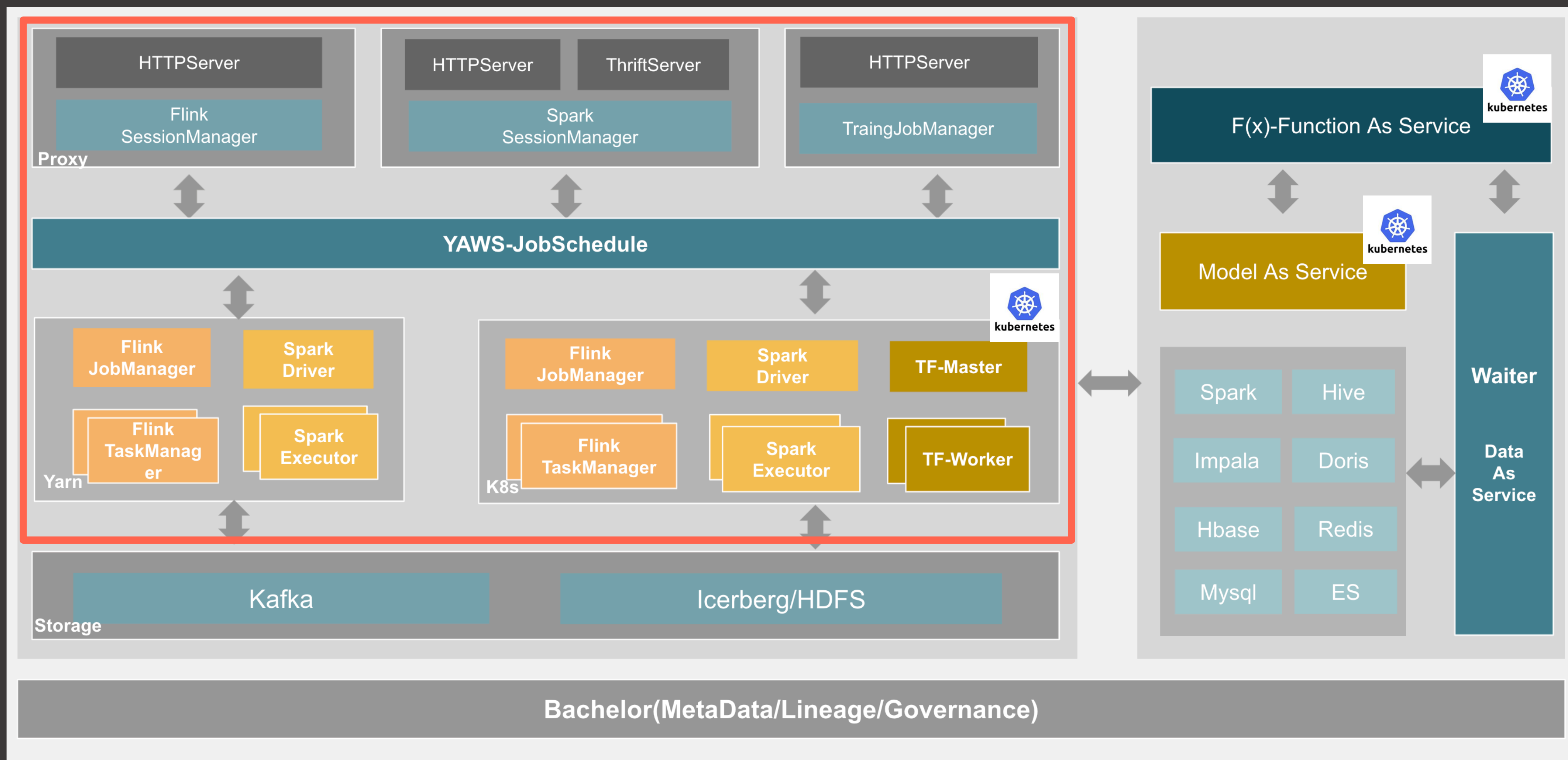
资源混合调度

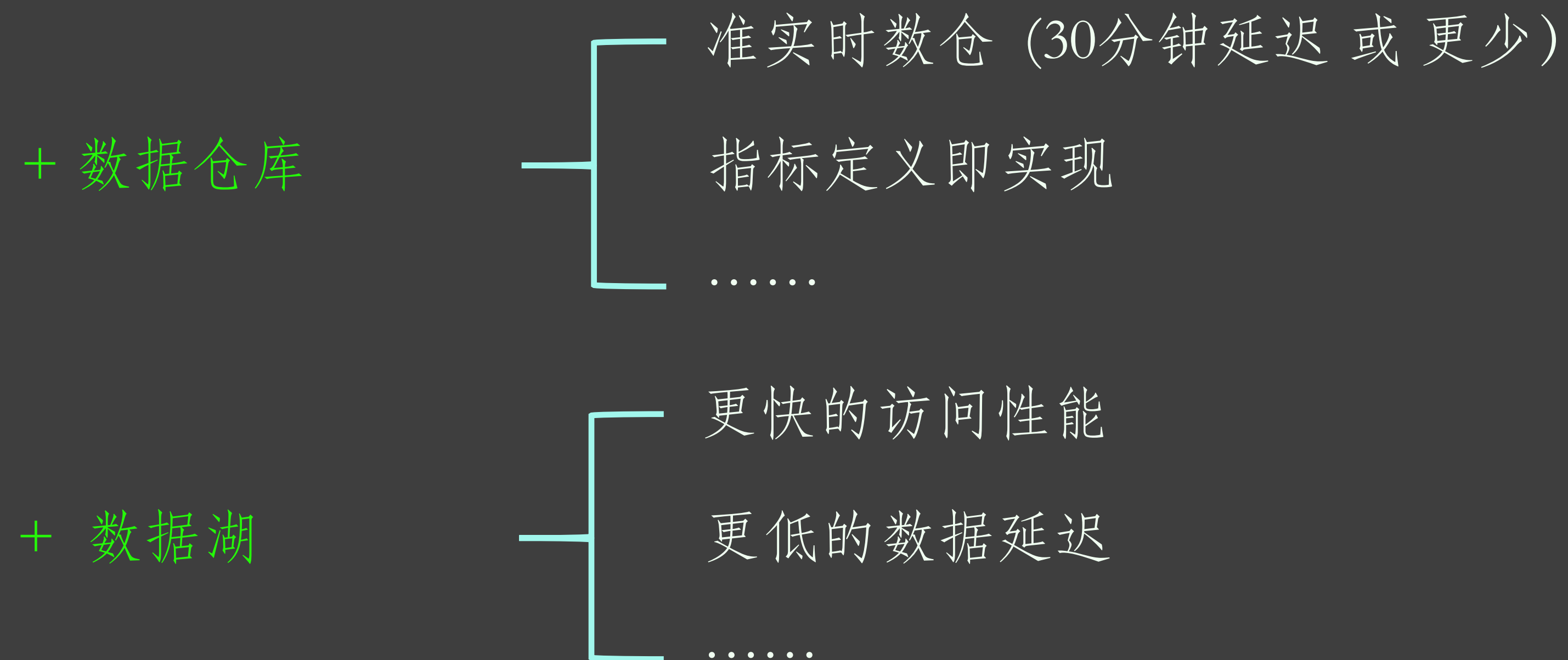
- 实时和离线计算使用不同的Yarn集群
- Yarn运维复杂，缺乏好用的运维工具
- Yarn资源隔离不彻底，任务之间容易相互影响
- 峰值流量场景，集群资源无法快速伸缩



任务资源调度- Cloud Native

- 批处理/流计算/算法训练任务混合部署
- 提供资源统一管理和分配
- 大促期间，计算资源快速扩缩容，抵抗峰值压力





不同的场景 | 不同的数据使用理念

数据湖不会替代数仓，而是长期共存

总结和展望

更可靠的数据生产和访问

