



大数据计算

架构峰会

大数据计算实践

2021.06.19 (周六) 14:50~15:30





好未来 实时平台建设实践

毛祥溢 数据平台开发专家





- 公司介绍
- 批流融合方案及技术架构
- 实时平台及技术架构
- 续报解决方案
- 分钟级数仓解决方案
- 总结回顾

毛祥溢 数据平台开发专家



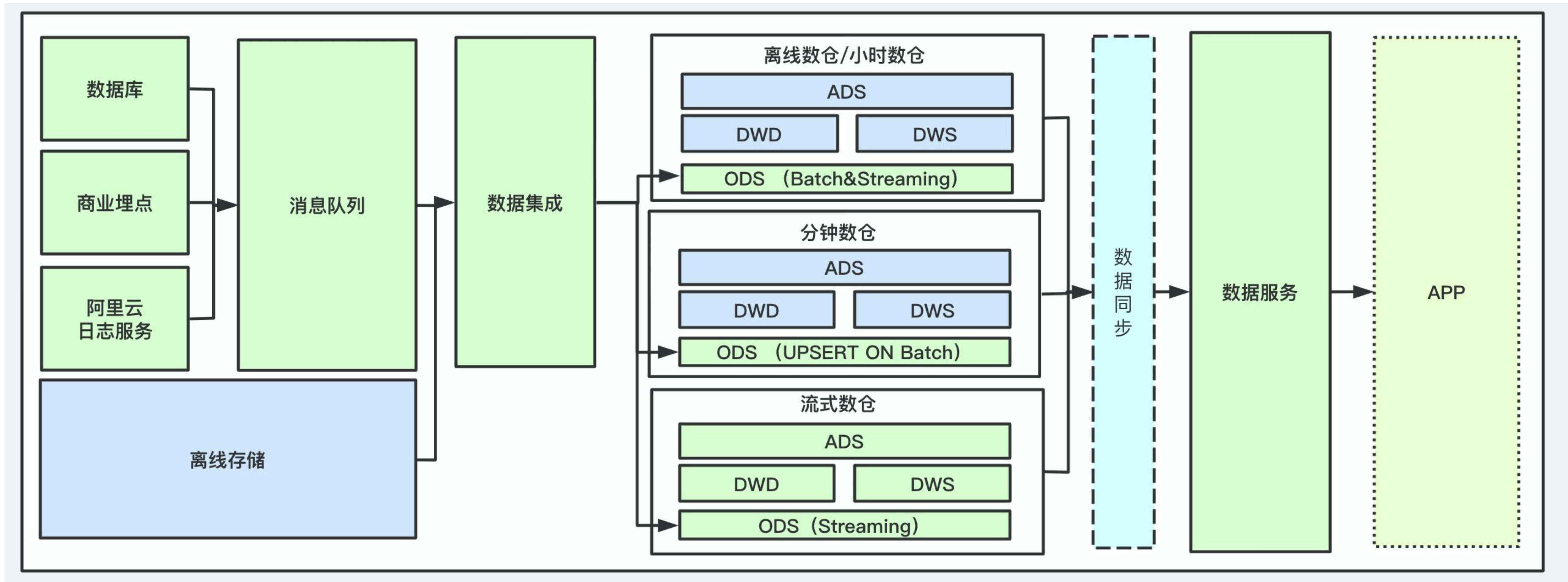
好未来-爱和科技让教育更美好



DataFunSummit

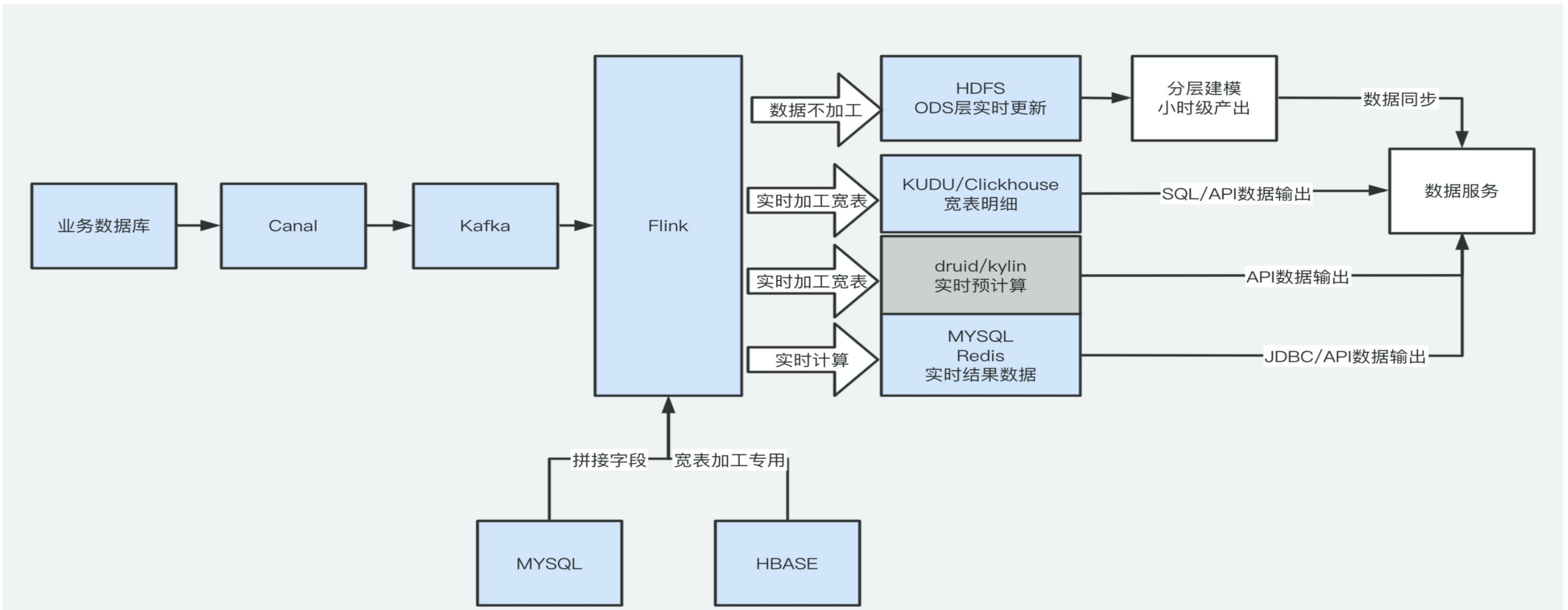
好未来-批流融合加速数据分析

- 重点放在数据分析的链路上，做到全链路数据分析实时化
- 提供3种不同层级的实时化方案，适配不同的场景和用户
- 通过Flink实现不同存储引擎的ODS层实时化



好未来-批流融合技术架构

- 流式计算框架擅长解决实时ETL的工作，并不大适合很重的分析场景，**实时OLAP引擎会是现在要突破的方向**
- ODS层实时化技术方案推荐 Flink + DataX，维表关联直连MYSQL、HBASE，必要时才做内存缓存
- 中小数据量**优先考虑Flink+实时OLAP引擎**的解决方案，开源方案 Flink+KUDU/数据湖、**商业方案Flink+Hologres**

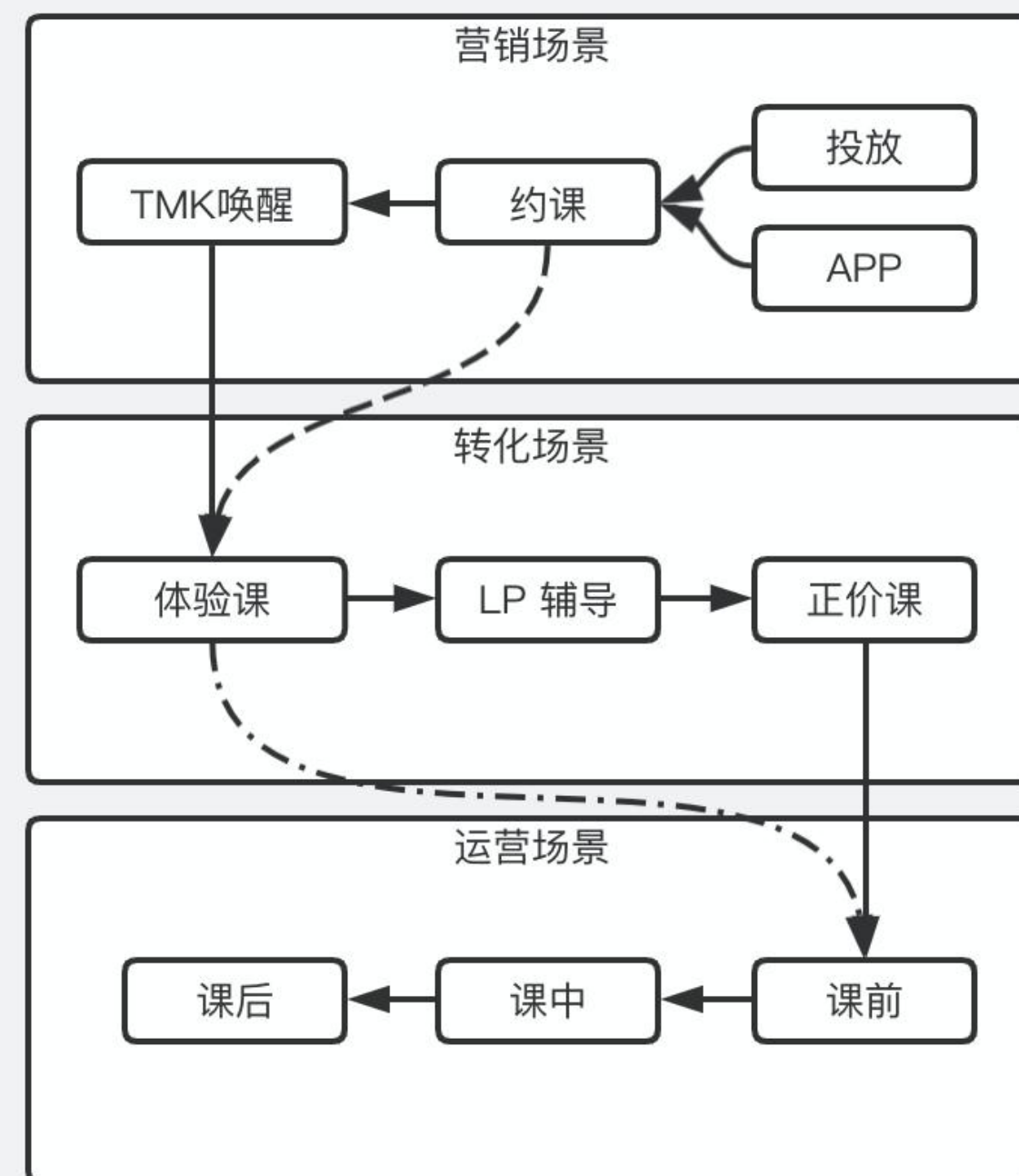


好未来-实时平台-功能介绍

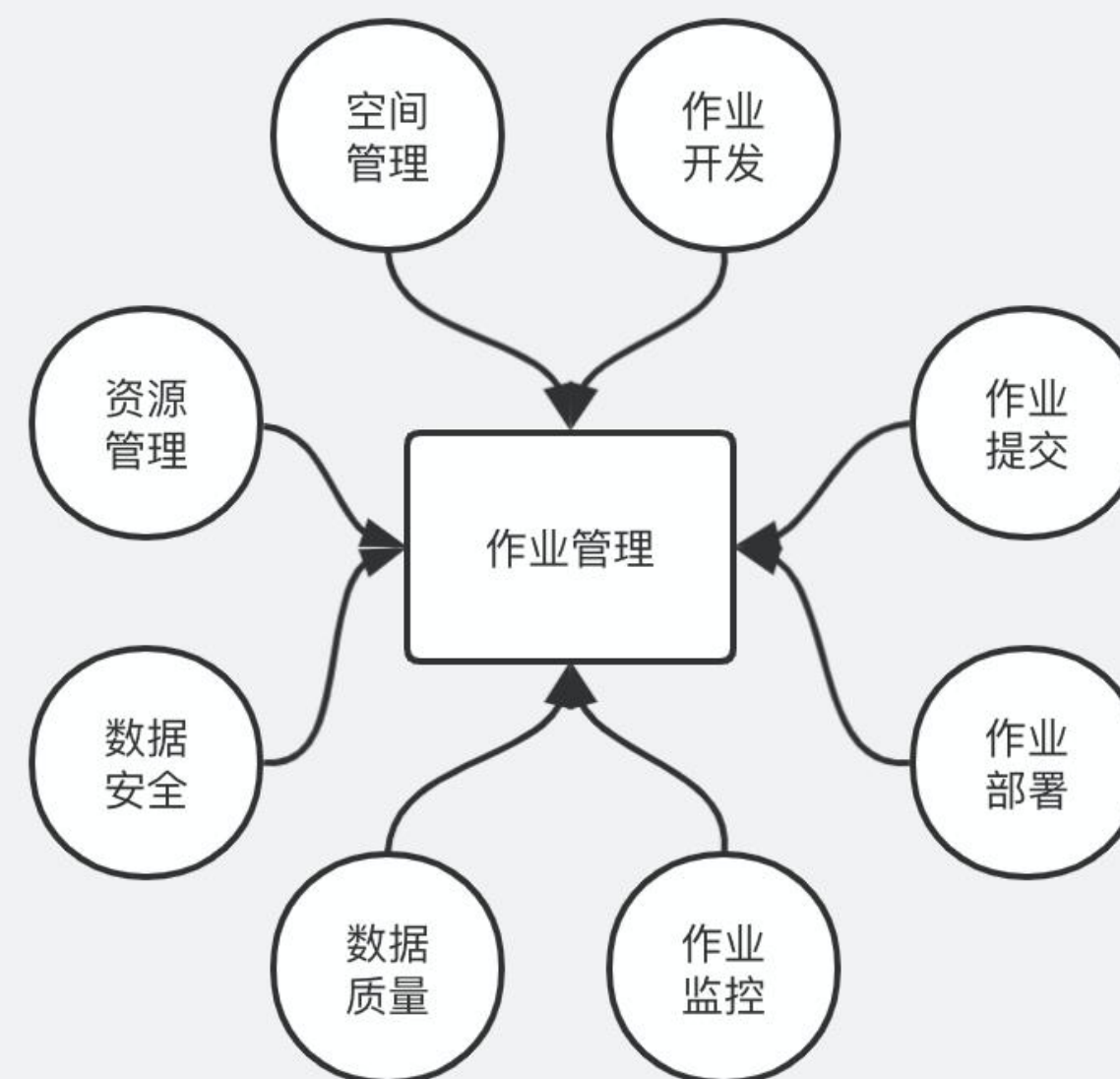


- 实时平台不应该局限于作业管理工具，还可能有**实时场景的通用解决方案**、通用的技术架构以及**计算链路的底层优化**。
- 开发管理包括：空间管理、作业开发、作业提交、作业部署、作业监控、数据质量、数据安全、资源管理等。
- 解决方案包括：打磨一套贯穿营销场景、转化场景、运营场景的实时分析解决方案，做到通用、可复制、门槛低。
- 底层优化包括：Canal、Kafka、Yarn、K8S、HBase、Kudu、Hologres等全链路的优化，当前正在解决大表实时同步的问题。

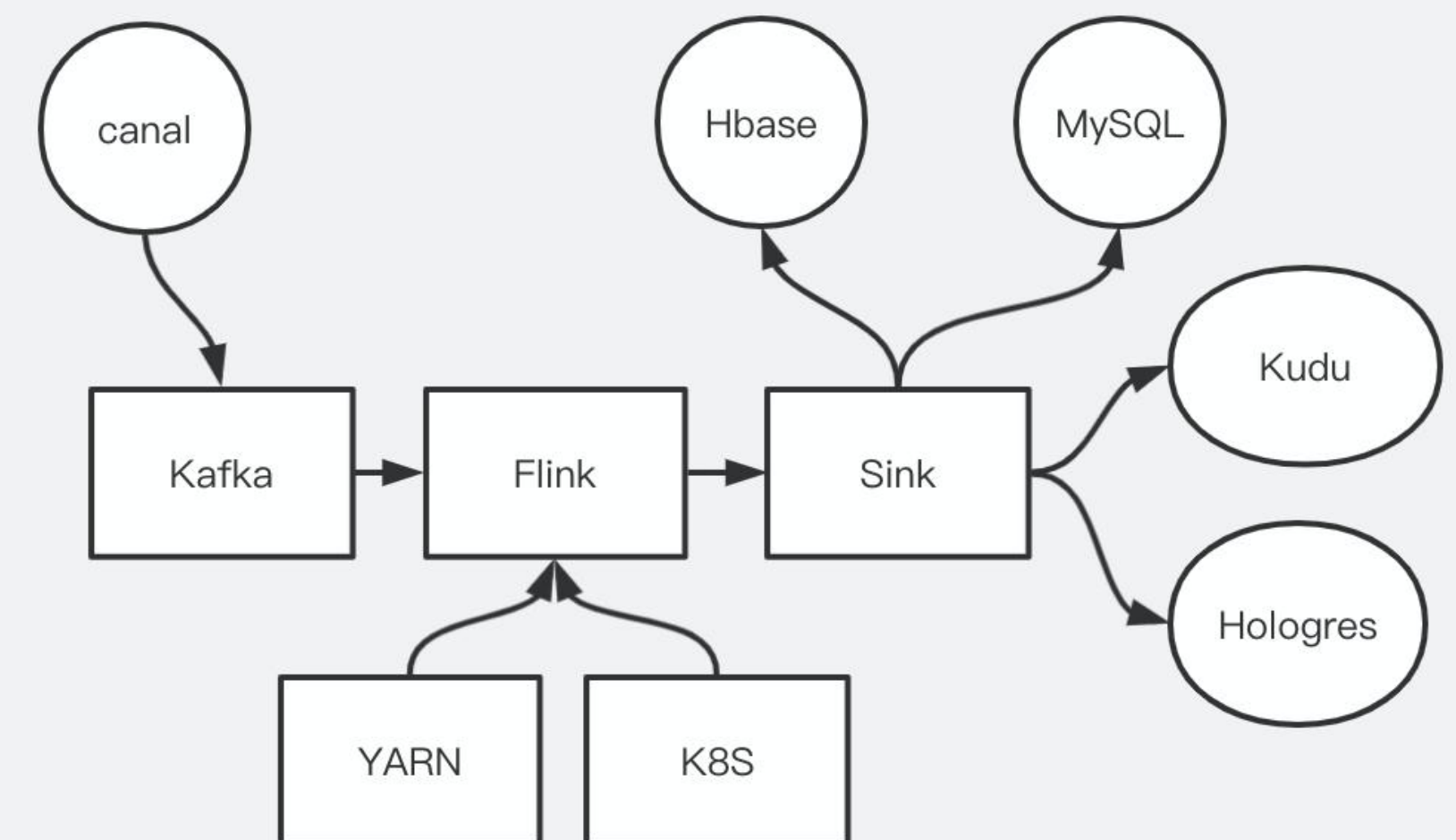
基于场景的解决方案



基于开发的作业管理



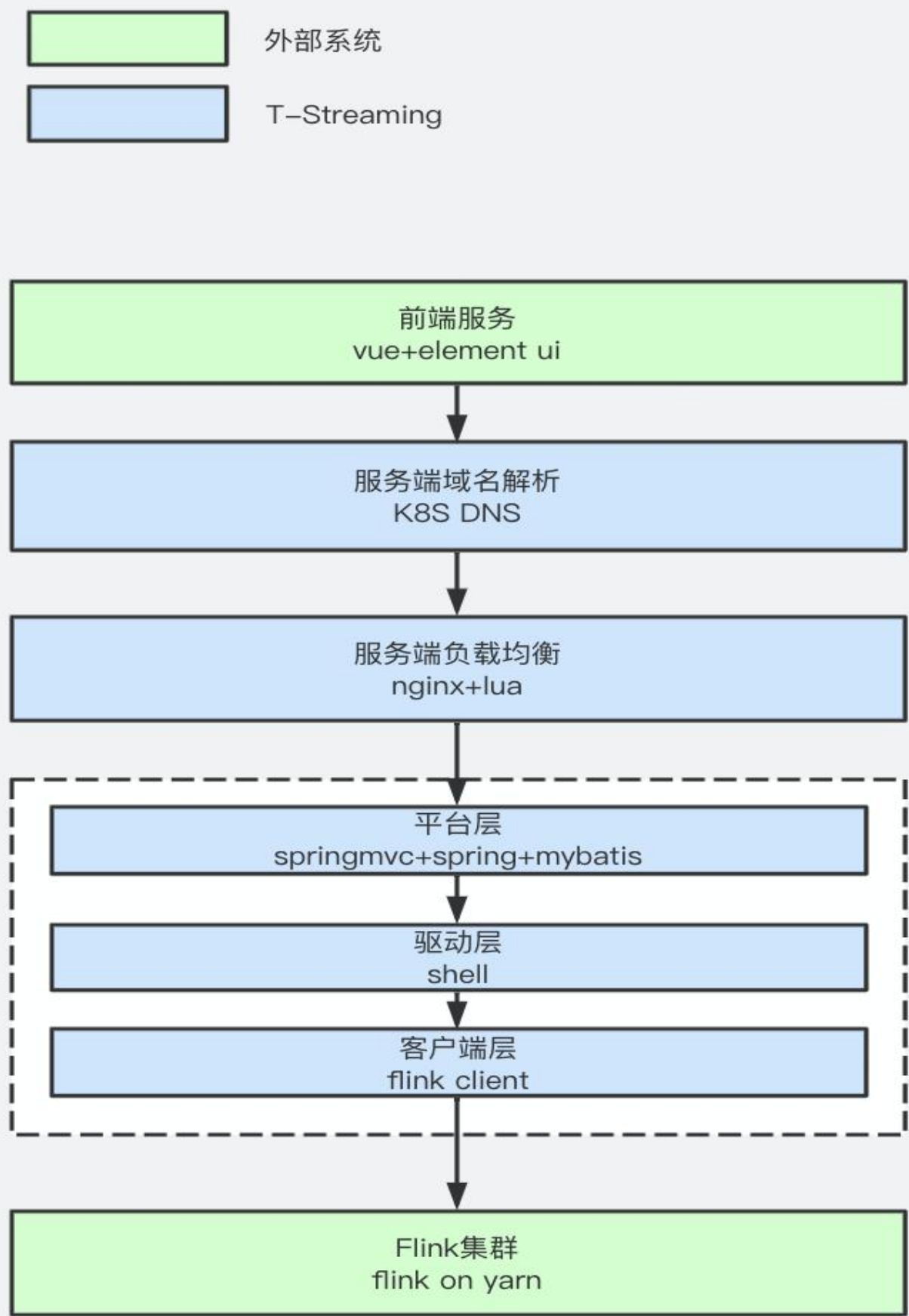
基于运维的底层优化



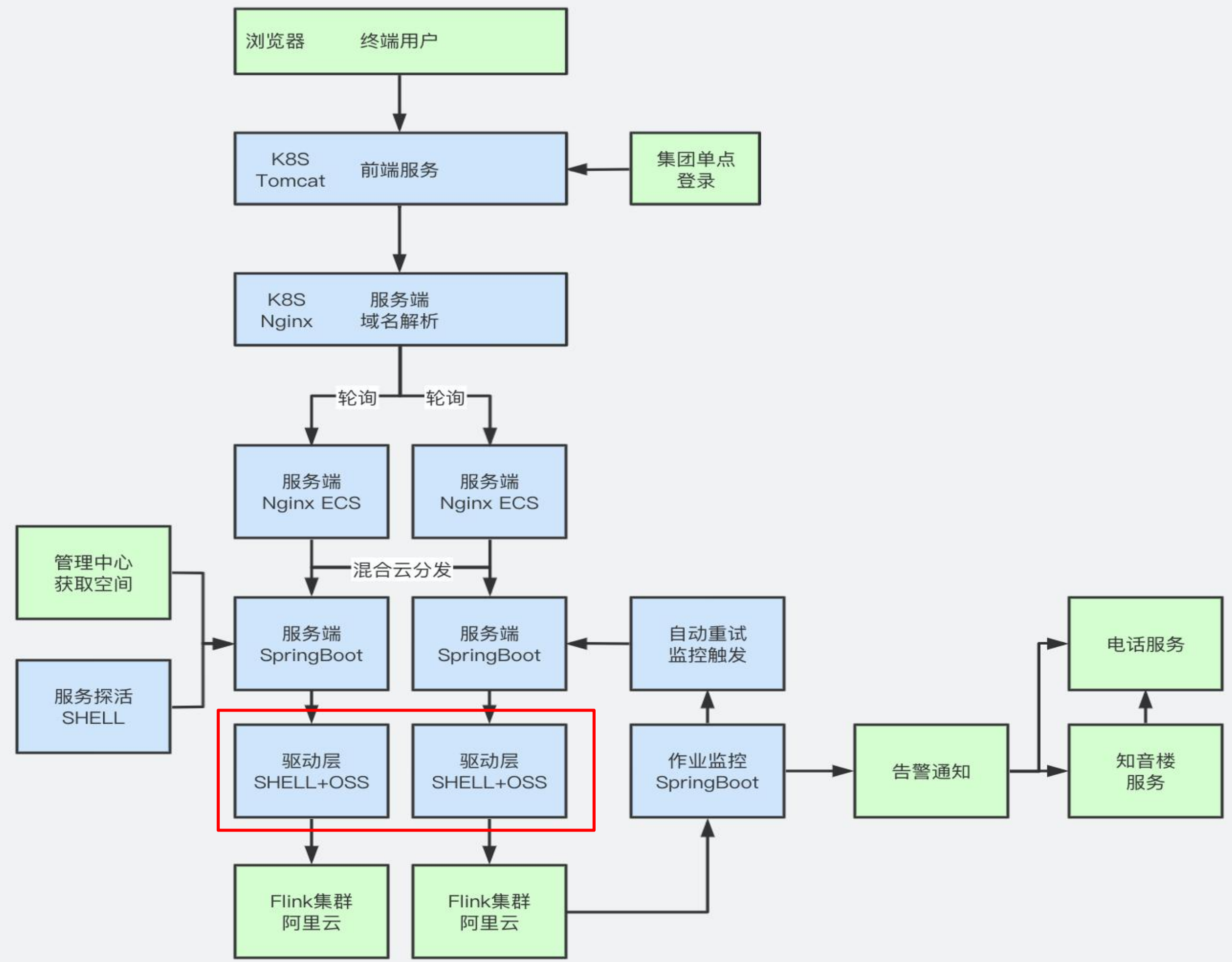
好未来-实时平台-技术架构



T-Streaming技术架构图



概览图



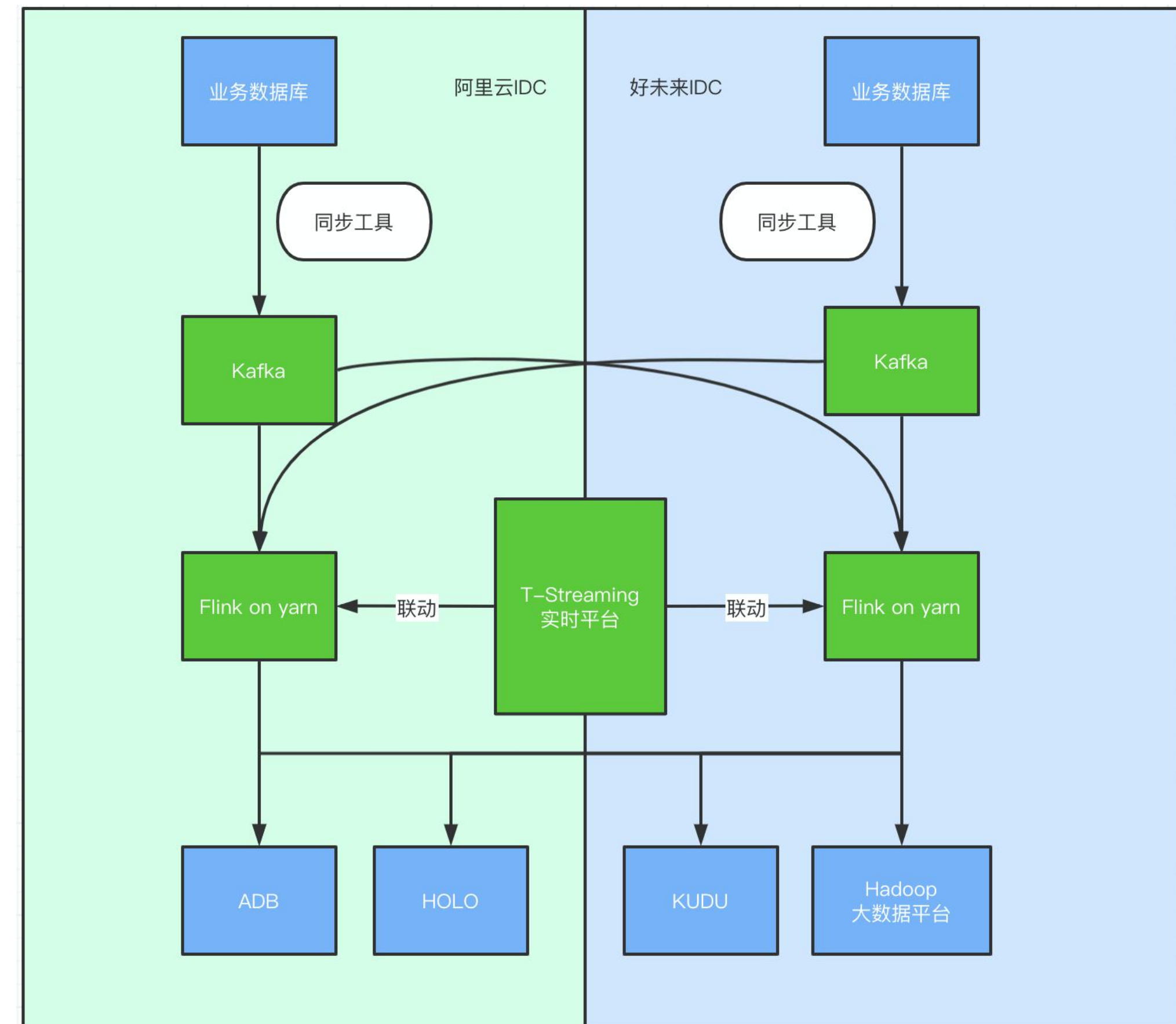
系统依赖图

好未来-实时平台-技术优化



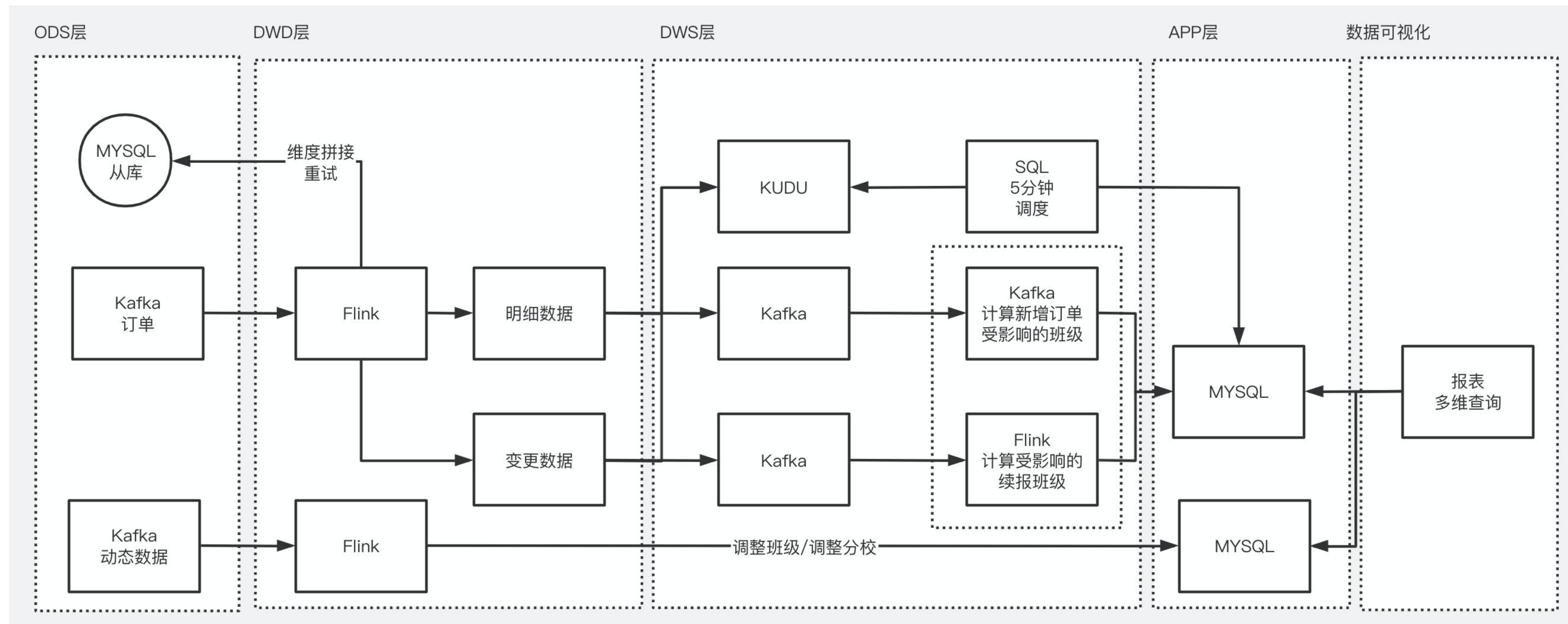
技术优化

- **混合云部署架构**，支持作业在多个集群间切换
打通线上+线下集群kerberos权限认证，实现跨集群读写
源码级修改使用Flink API方式提交作业，减低环境依赖
- 降低数据开发门槛，支持3种开发模式
自研**TS SQL**，封装原生Flink SQL实现极简操作
内置维表拼接函数、丰富的SLS\HIVE\KUDU模板
- **支持Flink作业自动恢复容错机制**
作业失败重启重试，自动选择checkpoint版本
根据指定策略消费Kafka数据
- **实时作业监控报警服务**、消息自动收敛、自动重试
- 作业在线调试，运行日志归档下载
- Topic字段采集、Sink表自动推测建表



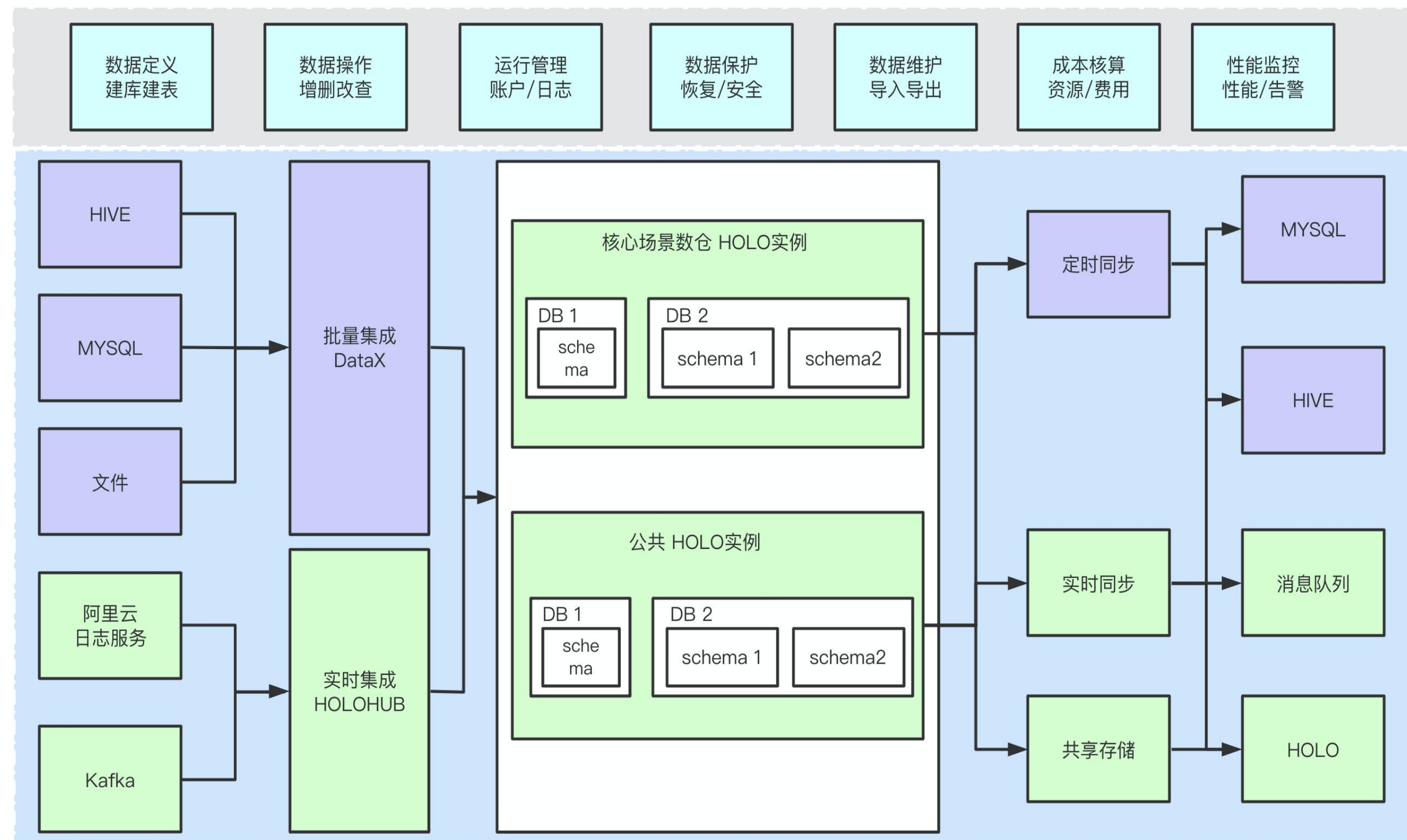
好未来-解决方案-实时续报方案

- 解决方案：打磨一套贯穿营销场景、转化场景、运营场景的实时分析解决方案
- 问题：依赖历史、退单影响历史订单、维度变更频繁
- 解决：实时ETL进行宽表拼接、2条互为容错链路（宽表实时落库使用SQL进行计算、Flink实时统计）
- 技术：Flink分布式维表缓存、KUDU UPSERT数据更新



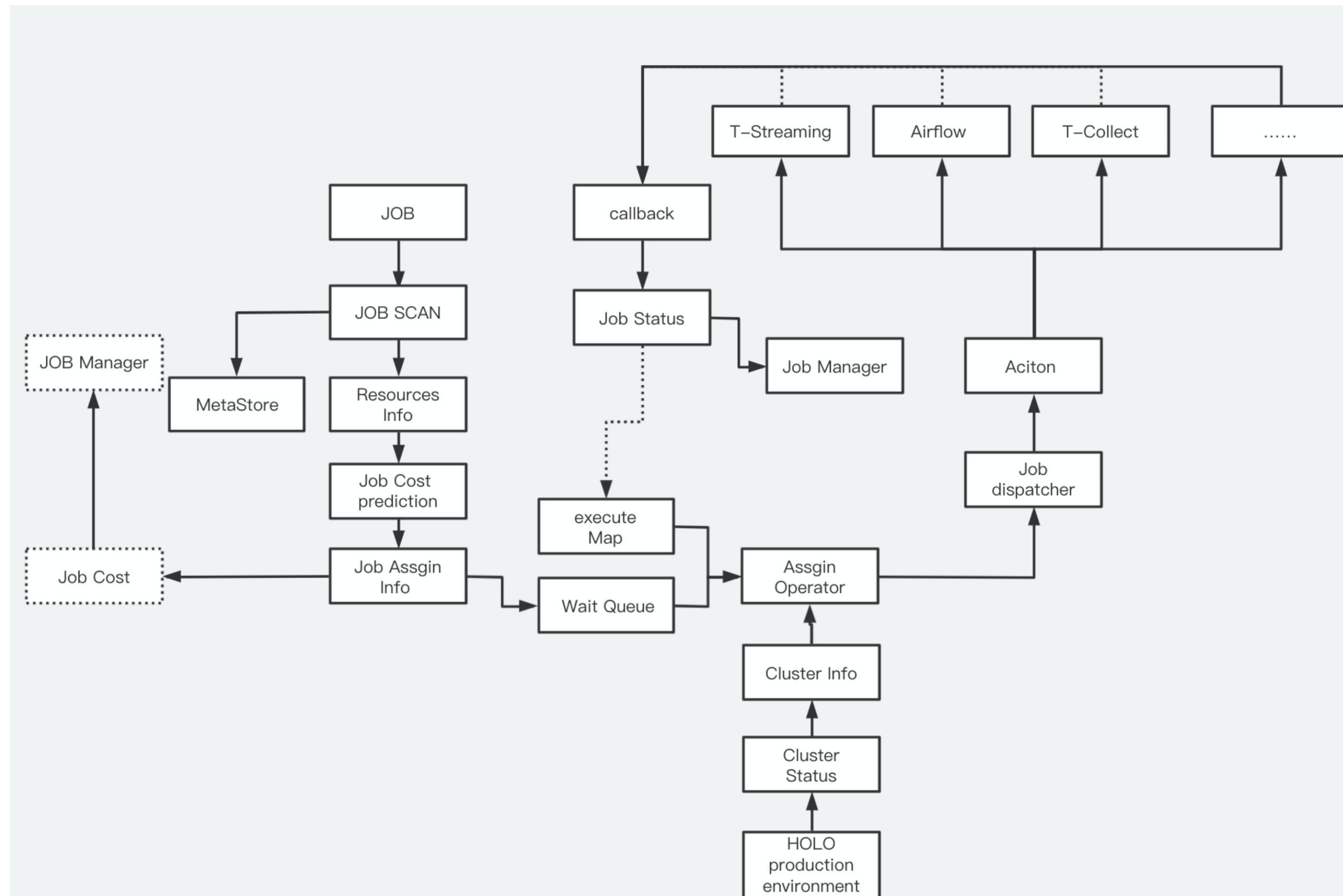
好未来-解决方案-分钟级数仓

- 打磨一套贯穿营销场景、转化场景、运营场景的解决方案的技术支撑。
- 分钟级数仓接入规范，离线数仓>小时级数仓>秒级数仓>分钟级数仓，选场景、控数据，避免数据膨胀。
- JDBC连接方式，数据集成和数据同步比较随意，常有拖库操作，不利于长期、稳定运维。
- 调度任务和即席查询错峰调度，避免任务推挤导致不可用。



好未来-解决方案-分钟级数仓作业调度

- 更好的管控：统一建库建表、统一作业调度（集成、计算、同步）、统一监报告警
- 开发一套外部调度，解决三个问题：主流实时OLAP资源隔离不完善、后期切换OLAP引擎需求、内部管控需求
- 通过SQL SCAN模块分析每个任务的描述信息，做资源消耗评估，生成调度作业对象，根据当前集群运行状态触发调度。



总结

- 批流融合是企业数据分析的趋势
- 通过ODS层实时化提供不同级别的时效性
- 实时平台不应该局限于作业管理工具，可以尝试做通用解决方案和标准技术架构。
- Flink+实时OLAP架构是目前需要突破完善的技术架构。

展望

- 稳定第一、提效第二、降本第三，实现简单可依赖。
- 优化开发生命周期、数据计算上下游链路，加强风控能力。
- 推进标准解决方案，技术架构，持续推广布道。



THANKS!

今天的分享就到这里...

Ending

