



# T3出行湖仓一体架构 下的统一指标平台

---

郑平贺



目录  
CONTENTS

- 01 为什么需要统一指标
- 02 湖仓一体下的统一指标平台架构
- 03 统一指标平台的实现
- 04 未来展望



# 为什么需要统一指标

---

# 什么是指标



## 指标

- 用来量化事物的一个工具，用数字来帮助我们描述一些抽象的事件
- 一组能反映某一业务在单位时间内的规模、程度、比例的数字



T3出行  
只为更好



## 为什么需要统一指标



T3出行  
只为更好

| DataFun.

## 简介：统一指标平台

指标一体化  
平台



定位

核心指标分  
析、拆解、  
异常运营、  
归因、预测、  
估算



功能

运营决策、  
全公司



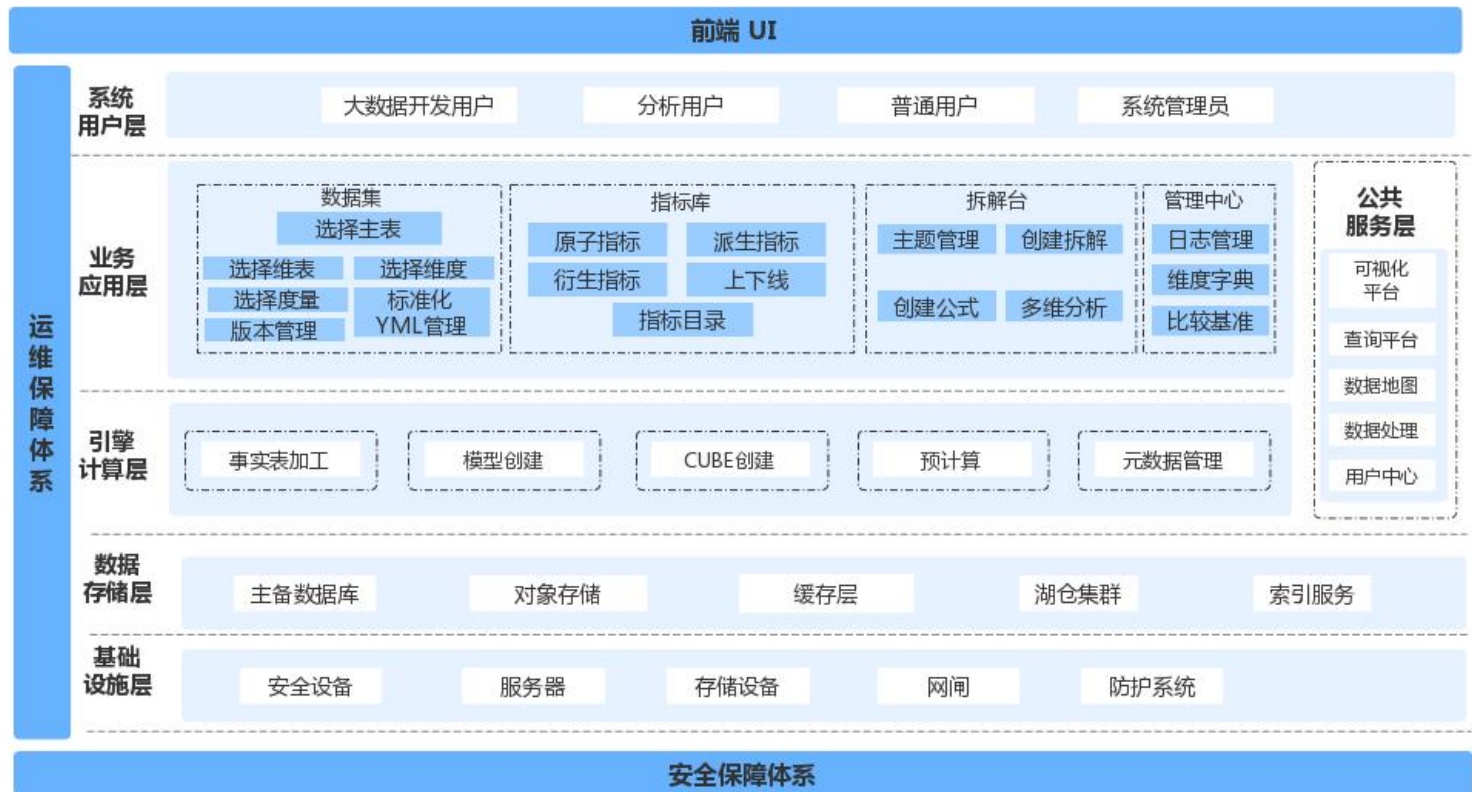
受众



## 湖仓一体下的统一指标平台架构

---

# 系统架构



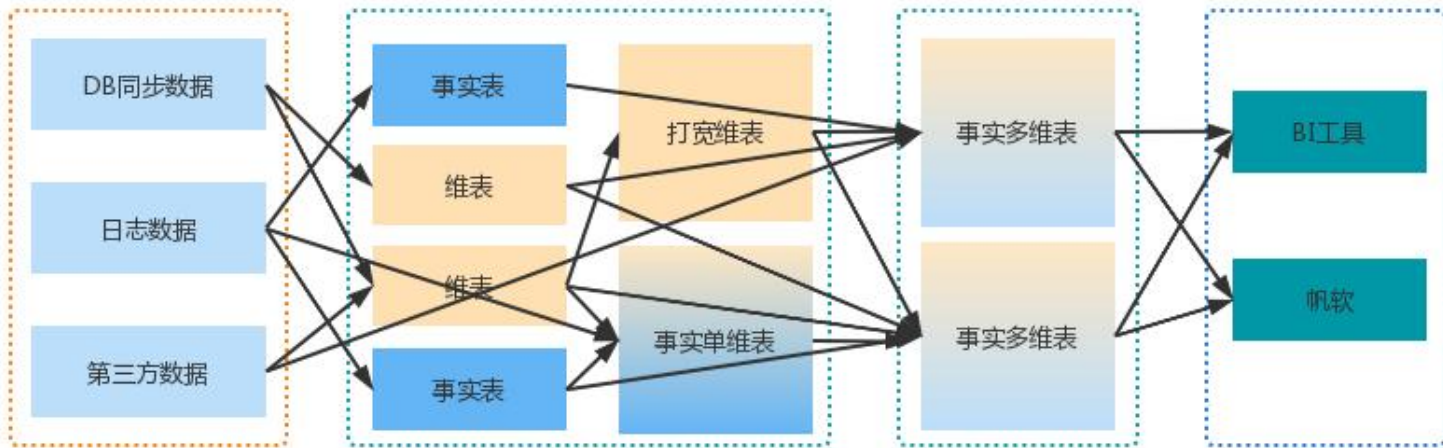
T3出行  
只为更好



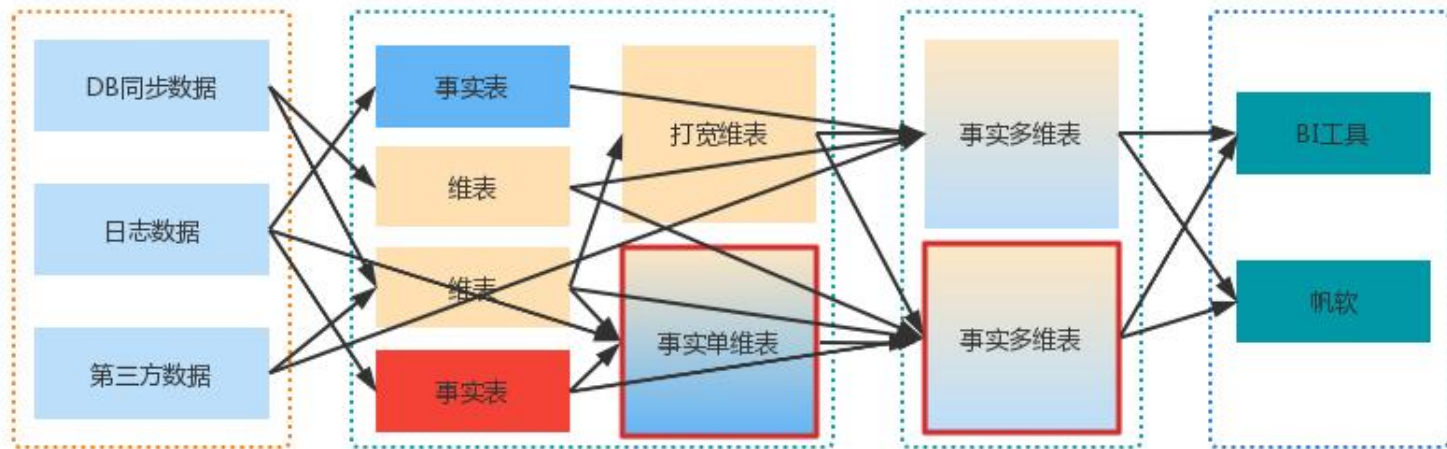


# 旧架构什么问题？

”



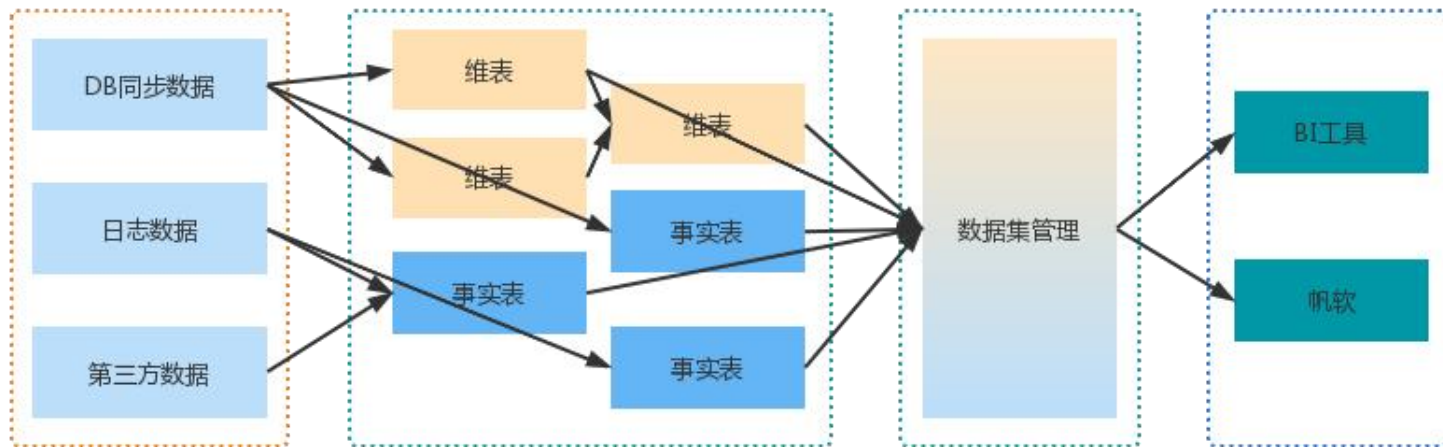
## 数据处理架构-网状模式（旧）



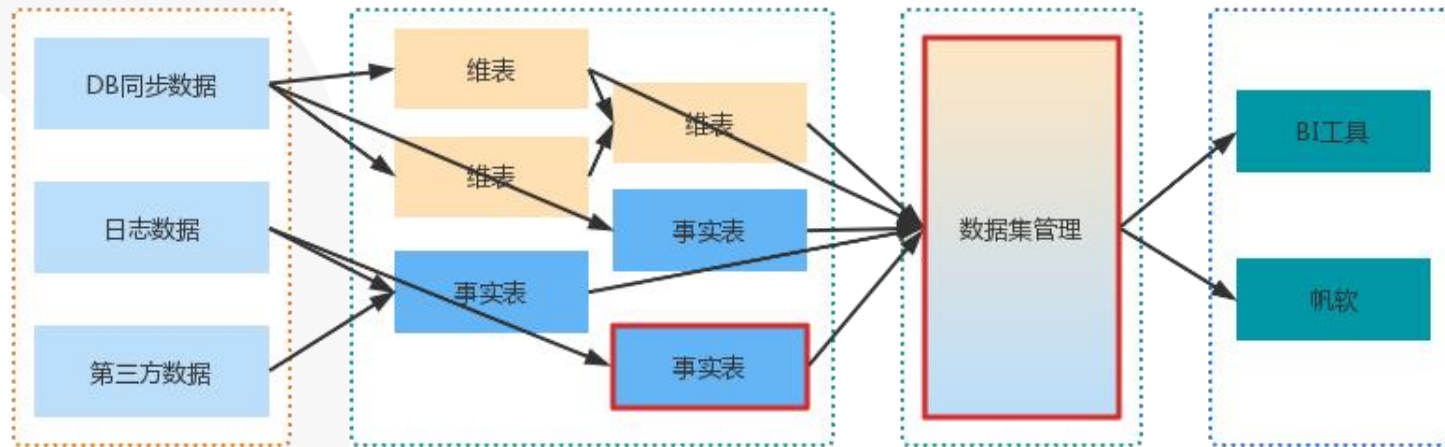
- 每天基于事实表和维表生成的打宽表，但却不知道存量表中是否已经存在，而且复杂性变得越来越复杂，当上游数据发生变更时，并不能保证下游作业完全被更新。
- 业务口径没办法统一，无法确认哪个部门提供的数据是准确的。



# 新架构有什么优势？



## 数据处理架构-中心化模式（新）



- 重建业务模型，生成精简的、规范化的、经过认证的事实表和维表。这样可以统一业务口径，减少表业务逻辑冗余
- 增加数据集管理，可经过确认的事实表和维表生成统一度量的聚合表，并结合业务变更，自动维护数据集版本，将变更数据推送给下游涉及表。

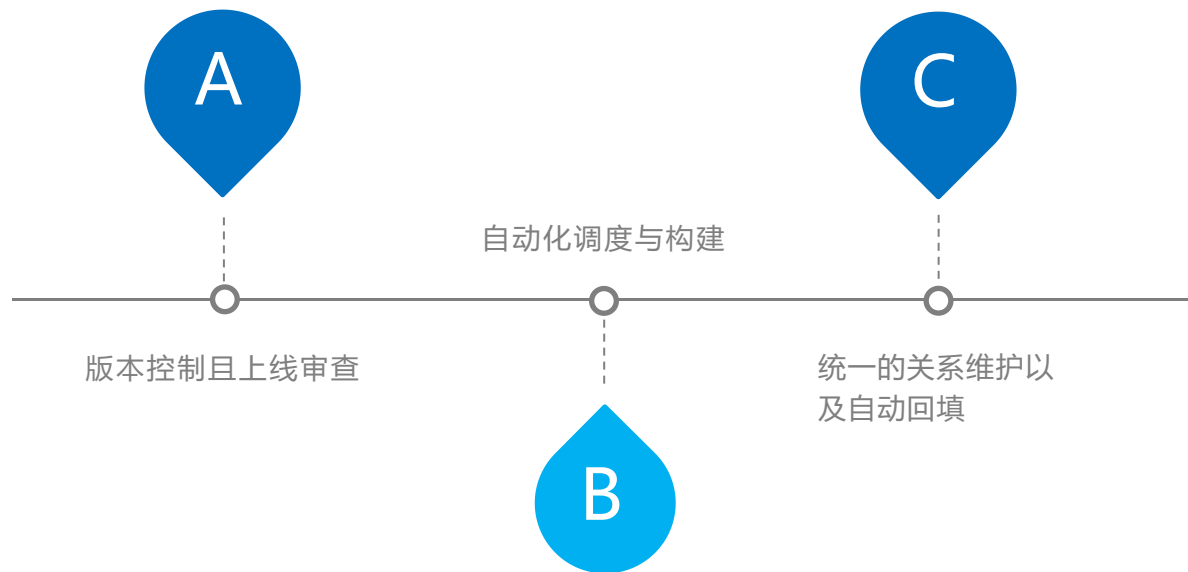
# 数据处理架构

两种架构 模式对比	功能点	网状模式	中心化模式
	规范化程度	低	高
	可观察性	低	高
	一致性	低	高
	上下游耦合性	高	低
	特殊场景适配度	高	低
	被权威认证	无	有



## 统一指标平台的实现

为什么需要语义?



# 数据集语义

```
#数据集定义, 必填
data_set:
  #kylin工程名, 必填
  project: phol
  #数据集名称, 必填
  name: ds_dev_1
  #数据集描述, 必填
  description: 测试开发
  #数据集定义人, 必填
  owner: admin
  #主表, 必填
  main_table: T3CKDW.DWS_DRI_ONLINE_HOUR_CUBE_DS
  #关联表
  correlation_tables:
    #关联表/sql
    - sql_table: T3CKDW.DIM_CITY_CUBE
      as_name: DIM_CITY_CUBE
      from_table: T3CKDW.DWS_DRI_ONLINE_HOUR_CUBE_DS
      #方式left_inner
      join_type: left
      join_condition: DWS_DRI_ONLINE_HOUR_CUBE_DS.CITY_ID = DIM_CITY_CUBE.COMBINE_CODE

#度量定义
measures:
  #度量名称, 必填
  - name: SUM_CHARGING_DURATION
    #度量描述, 必填
    description:
    #聚合模式, 必填
    agg: SUM
    #度量操作类型, 必填column或者constant
    expr_type: column
    #度量操作表达式, 必填
    expr_value: CHARGING_DURATION
    #度量返回, 必填
    return_type: double
  - name: SUM_PICKUP_DURATION
    #度量描述
    description:
    #聚合模式
    agg: SUM
```

```
#度量操作类型
expr_type: column
#度量操作表达式
expr_value: PICKUP_DURATION
return_type: double
- name: SUM_SERVICE_DURATION
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: SERVICE_DURATION
  return_type: double
- name: SUM_ONLINE_DURATION
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: ONLINE_DURATION
  return_type: double
- name: SUM_CRUISE_DURATION
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: CRUISE_DURATION
  return_type: double
- name: SUM_HASH_DRIVER_ID
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: HASH_DRIVER_ID
```

```
- name: SUM_CRUISE_DURATION
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: CRUISE_DURATION
  return_type: double
- name: SUM_HASH_DRIVER_ID
  #度量描述
  description:
  #聚合模式
  agg: SUM
  #度量操作类型
  expr_type: column
  #度量操作表达式
  expr_value: HASH_DRIVER_ID
  return_type: bigint
- name: count_total
  description:
  agg: COUNT
  #度量操作类型
  expr_type: constant
  #度量操作表达式
  expr_value: 1
  return_type: bigint
```

```
#关联定义
correlation:
  #维度字段
  - name: DS
    #维度类型: 时间
    type: string
    table: DWS_DRI_ONLINE_HOUR_CUBE_DS
    #是否为分区
    is_partition: true
    #类型参数
    type_params:
      #时间格式
      time_format: yyyy-MM-dd
      #时间粒度
      time_granularity: day
```



T3出行  
只 为 更 好





## 指标语义

```
name: finish_cnt_total
#指标显示名
display_name: 完单总量
#指标描述
description: 完单量指标测试
owner: test
#指标类型: 1表示原子指标, 2表示衍生指标
type: 1
#依赖的数据集
depend_data_set:
  #数据集名称
  name: ord_tra_cash_fact_set
  #度量名
  measure: finish_cnt_total
#指标所关联的维度
dimensions:
```



T3出行  
只为更好

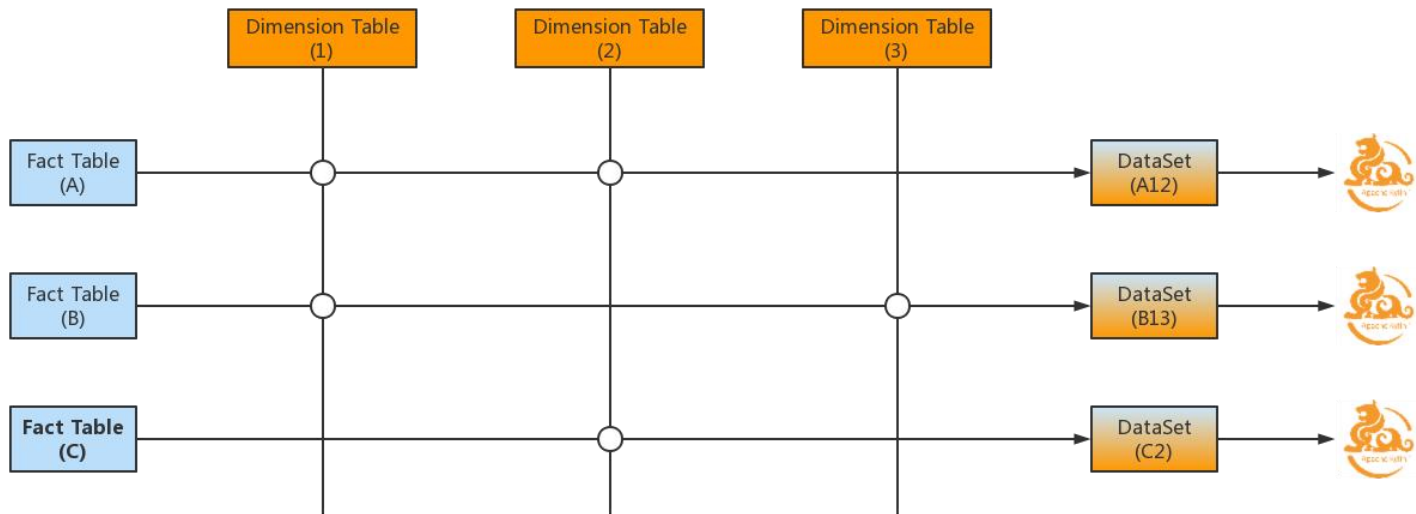


# 数据集伸缩设计

基于星型模型

事实表对应唯一  
数据集

维表可对应多个数  
数据集

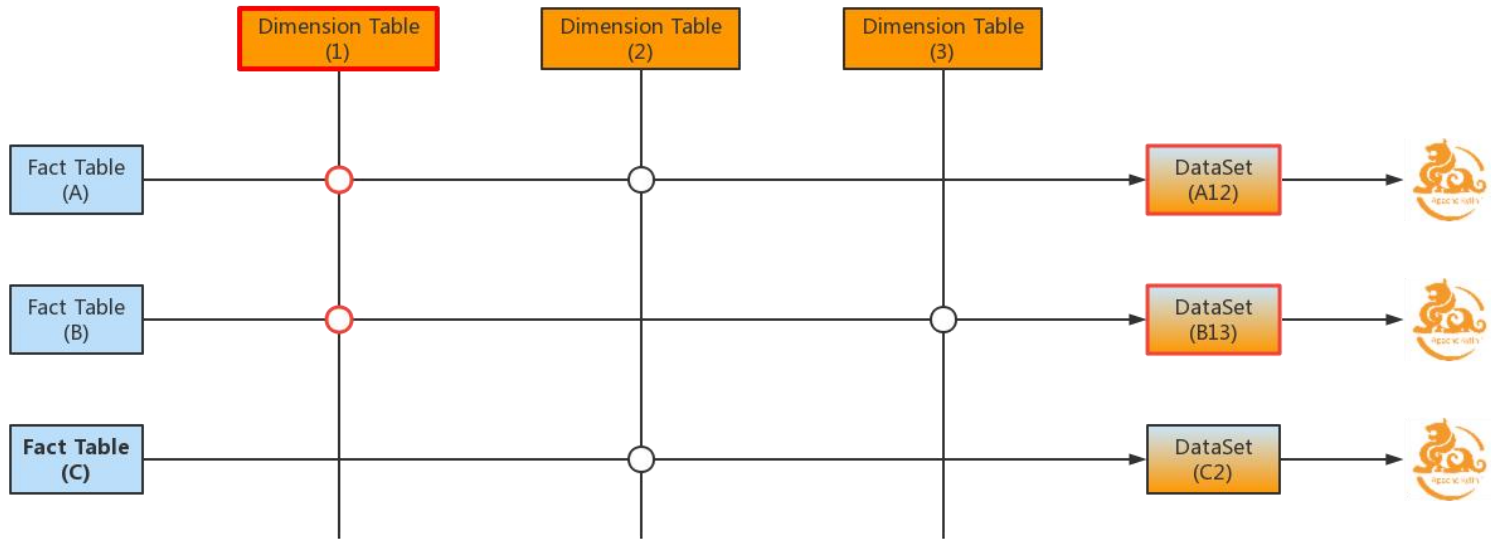


T3出行  
只为更好



## 数据集一致设计

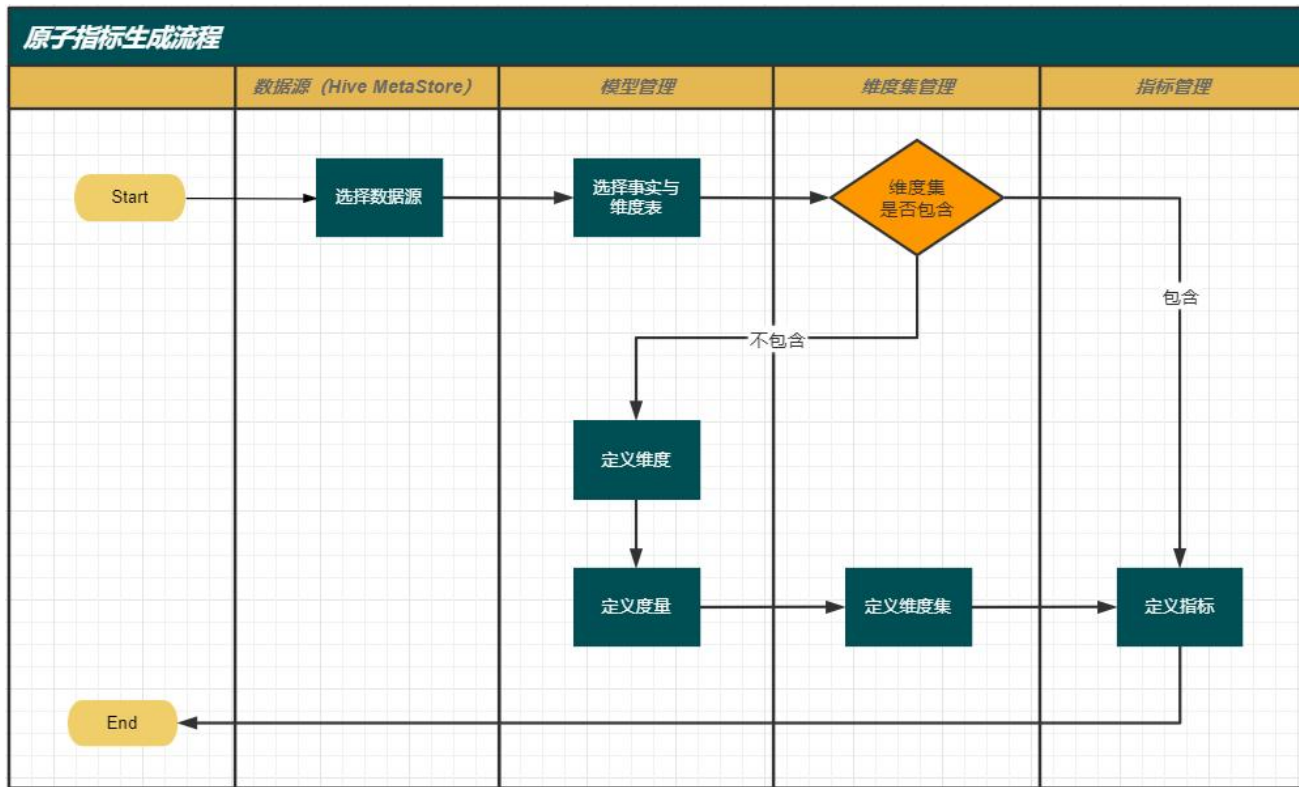
当有维表发生变更时，可自动探索到变更并将变更提醒通知给数据集维护同学，并进行数据集调整，创建数据集版本，形成多版本Cube。



T3出行  
只为更好



## 指标定义流程



T3出行  
只为更好

DataFun.

## 计算引擎选择

### 计算引擎支撑能力

1

支持超大数据集计算和存储

2

支持湖仓生态

3

可多维度字段预计算，提供快速灵活的查询能力

4

支持标准sql语义

5

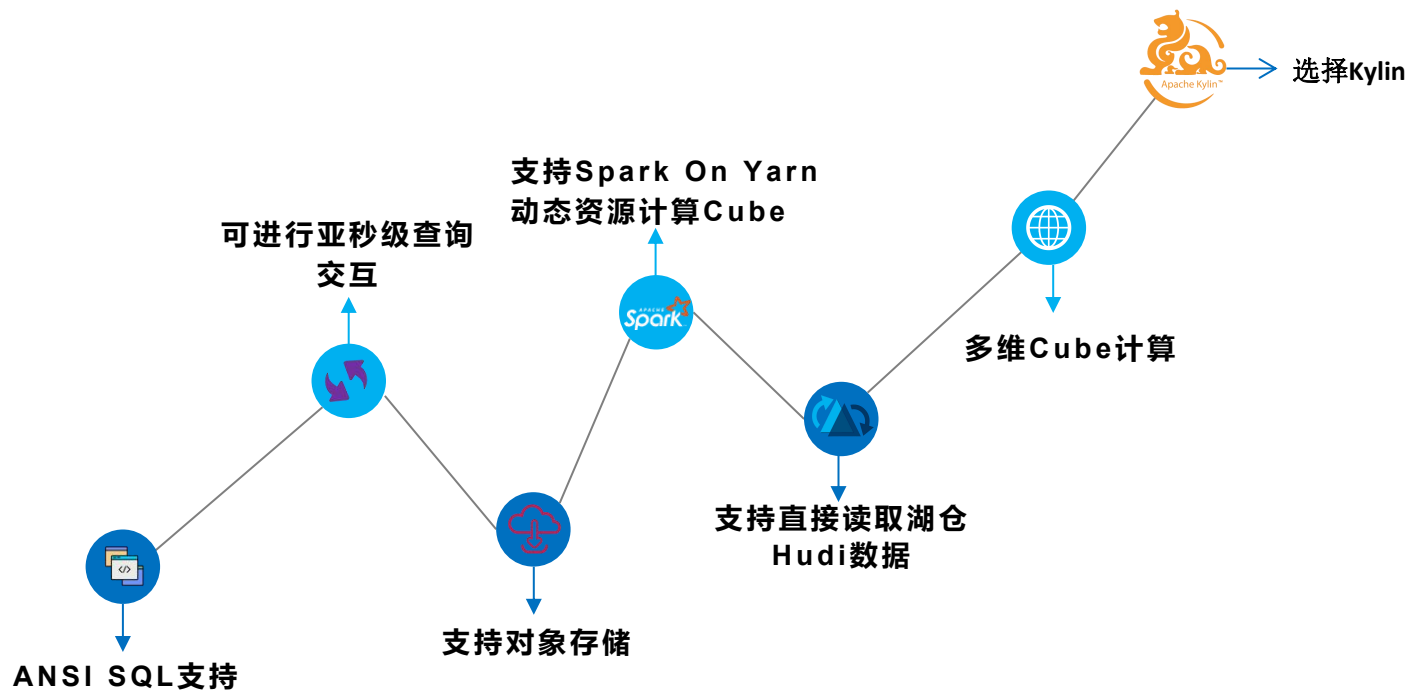
可通过jdbc方便访问和管理数据集



T3出行  
只为更好



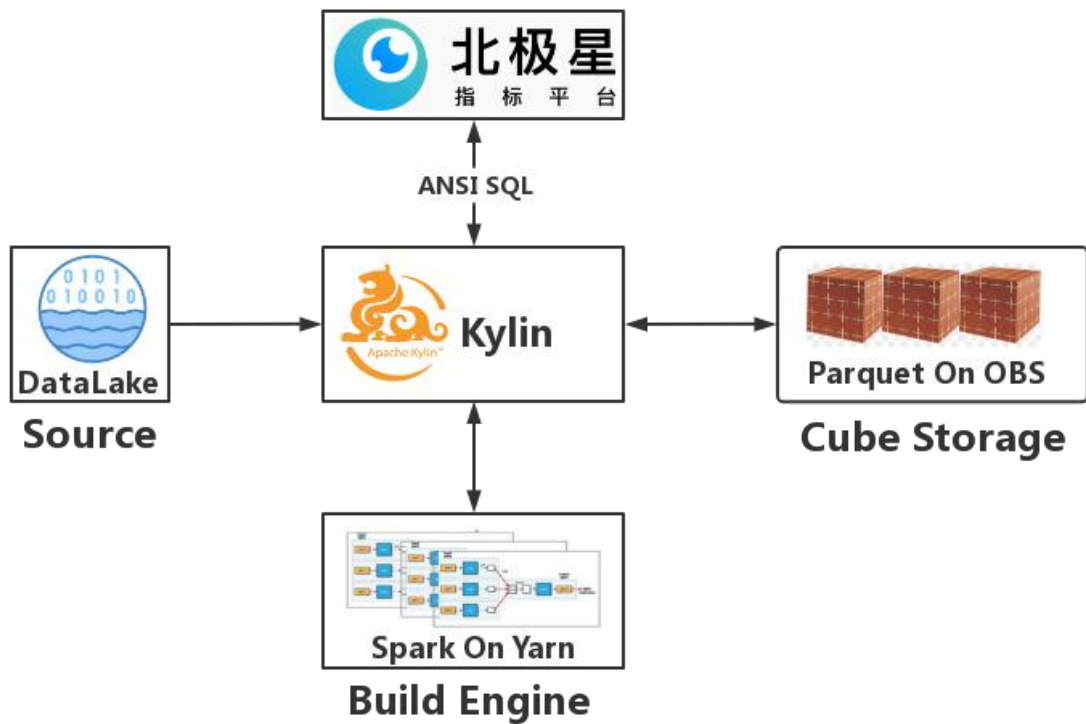
# 计算引擎选择



T3出行  
只为更好

DataFun.

## Kylin引擎数据处理流程



T3出行  
只为更好



## 一站式数据集创建

01

01 界面化编辑语义文件

02

02 语义上线

03

03 解析语义

04

04 创建model

05

05 Cube创建

06

06 配置定时Cube构建



T3出行  
只为更好

| DataFun.



## Kylin引擎优化

存算分离，存储使用独立的对象存储，计算使用独立的MRS集群

A

Spark开启  
`spark.dynamicAllocation.enabled`动态资源

B

支持Hudi表读取生成宽表

C

Build Cube 开启动态  
Shuffle

D

开启Cube build  
planner

E



T3出行  
只为更好



## Kylin压测结果

由于业务模型**维度字段**较多且支持客户端灵活查询，所以在Cube维度选择上增大到**16**个，固定**强制维度**为**分区字段**。测试单分区数据量**8000w**。

Label	# 样本	平均值	中位数	最小值	最大值	异常 %	吞吐量	发送 KB/sec
测试接口	50000	421	424	33	1369	0.00%	1178.66	521.42
总体	50000	421	424	33	1369	0.00%	1178.66	521.42



T3出行  
只为更好





## 未来展望

---

## 未来展望

01

### 持续迭代功能

持续迭代指标趋势预测，热点指标关注，预测，估算等功能

02

### 无缝测试上线生产

数据集模型测试配置并无缝切换引擎环境，然后自动完成在不同环境配指上线构建等工作。

03

### 增加实时指标

增加实时指标接入能力，引入实时计算引擎，并支撑离线指标。

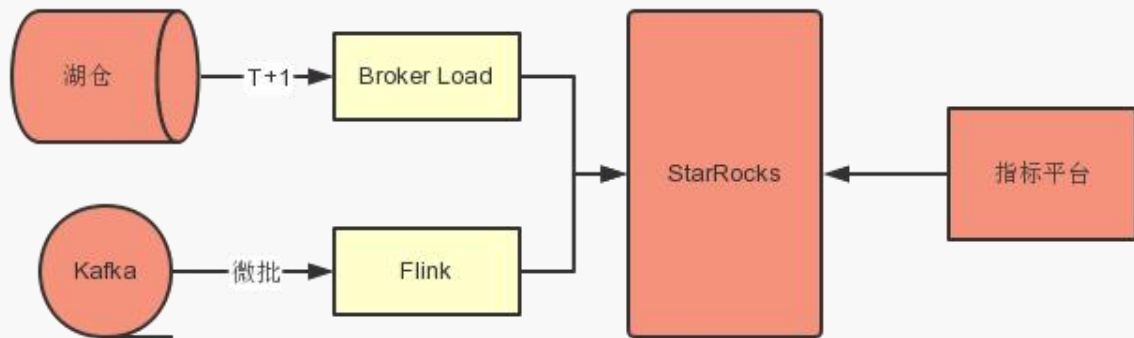
04

### 指标管理独立

为了兼容不同类型指标，包括运维指标，画像指标等.抽取数据集和指标管模块。

## 融合离线实时指标

为了能够同时满足离线和实时指标元数据统一、引擎统一能力，  
我们POC了StarRocks引擎。



## 融合离线实时指标

场景	离线	实时
导入方式	broker load	flink
频次	T-1	10s-15s
并发	-	100

# 非常感谢您的观看

---

