

Apache Doris 实时数仓应用实践

李荣谦

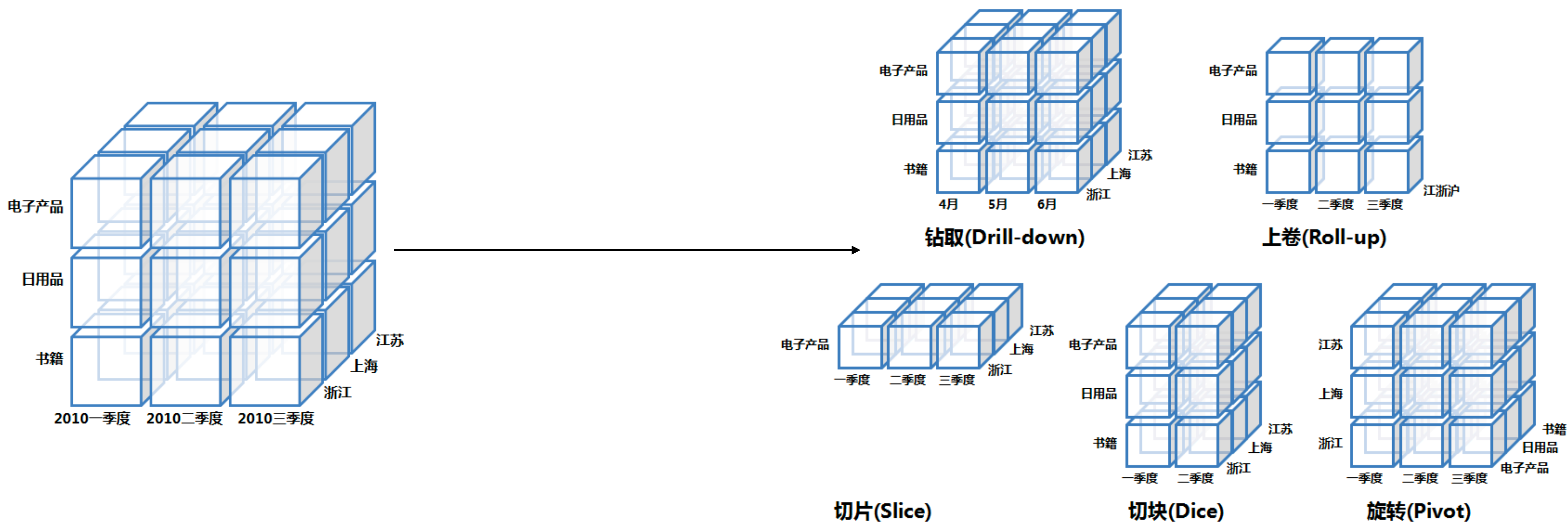
目录

- 何时需要实时数仓
- 什么是多维分析
- 实时数仓优势
- 实时数仓对比
- Doris简介
- Doris适用场景
- Doris 在精品课的应用
- 目前业界落地场景

何时需要实时数仓

- 实时
 - 需要实时导入数据
 - 需要较高的查询性能
- 数据仓库
 - 分析型查询较多
 - 数据量较大
 - 只查询部分列
 - 多维分析

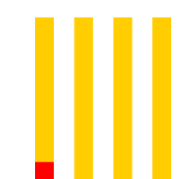
什么是多维分析



实时数仓优势



DORIS



ClickHouse

VS



大规模数据需要分库分表
分析型sql较慢
无法横向扩展



不能实时写入
update/delete效率非常低
查询耗时较高



elasticsearch

不能Join
数据压缩不足

实时数仓对比

- ClickHouse
- TiDB
- Druid
- Kylin
- Doris

实时数仓对比

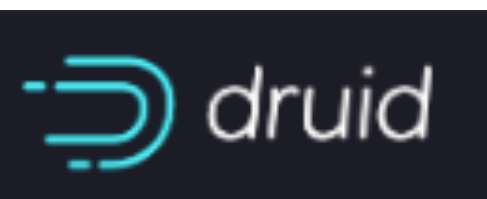
ClickHouse

- 优点
 - 支持80%的 SQL
 - 向量化引擎
 - 数据压缩空间大
 - 支持 MySQL/Kafka/JDBC/ODBC 远程映射
- 缺点
 - 对更新支持不好
 - 集群加减节点，不能自动 Rebalance 数据

实时数仓对比

TiDB

- 优点
 - 兼容 MySQL 协议
 - 低成本数据复制(TiKV Raft 同步 TiFlash)
 - 完善的DML支持
- 缺点
 - TiFlash 不开源，TiKV为行存
 - OLAP支持较少
 - 目前实时数仓领域落地较少



实时数仓对比

Druid

- 优点
 - 可以基于时间分区，摄取数据时预聚合
 - 支持实时数据写入
- 缺点
 - Druid SQL不是标准SQL(DDL,DML,JOIN 等不支持)
 - 不支持查询明细数据(只能查询到聚合粒度的数据)



实时数仓对比

Kylin

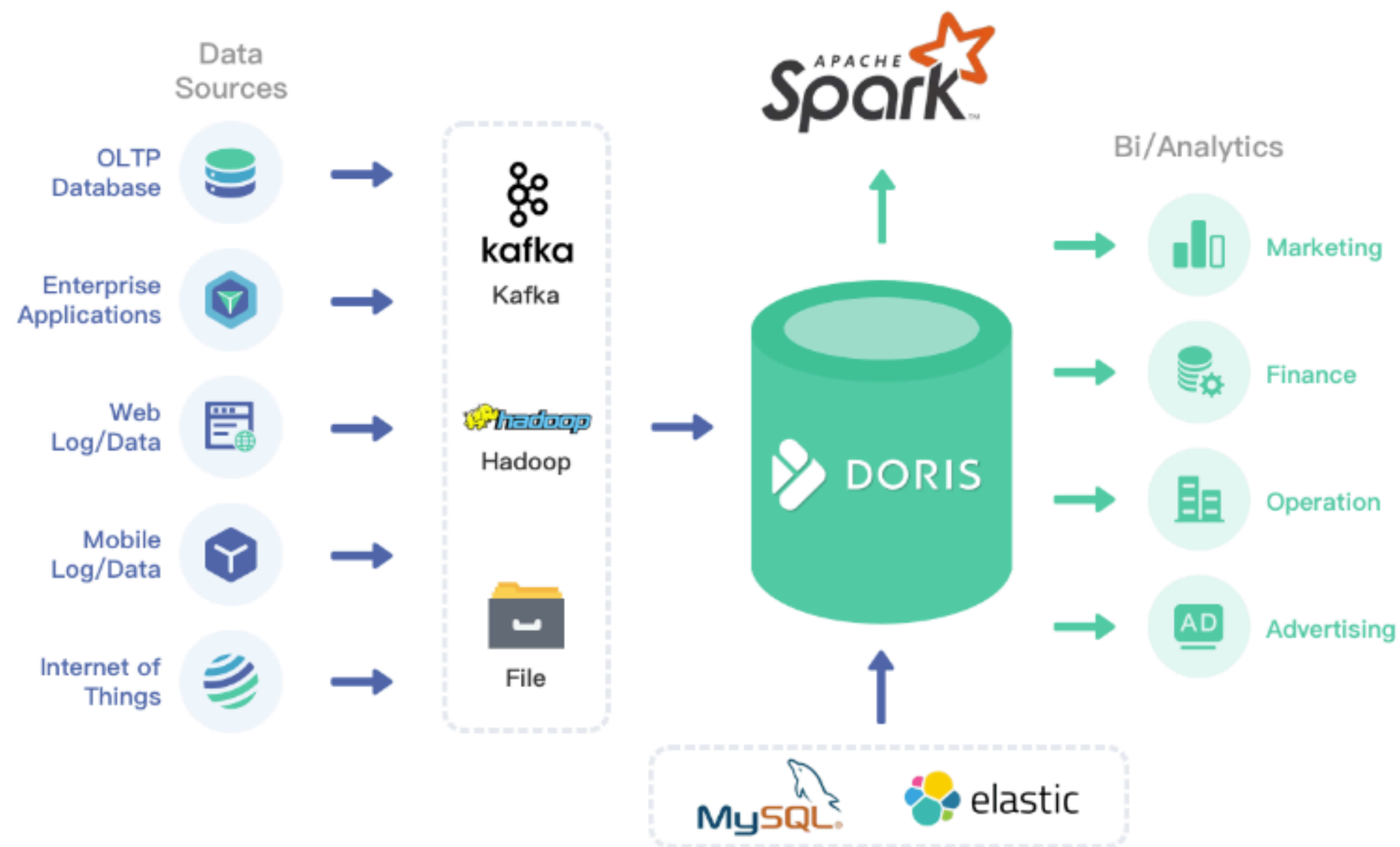
- 优点
 - 预聚合后的数据存储在 Hbase，查询较快
 - 支持标准 SQL 查询
- 缺点
 - 明细查询支持较弱
 - 对 Cube Schema 修改较复杂，周期很长
 - 依赖较多，需要 Hadoop/Hive/Hbase
 - 对多维数据分析支持不好

实时数仓对比

Doris

- 优点
 - 兼容 MySQL 协议
 - 支持物化视图实现预聚合
 - 支持基于时间分区，支持数据冷热分离
 - 支持将 Elasticsearch/Mysql 作为外部表，实现跨源 join 查询
 - 无额外依赖，运维简单
- 缺点
 - 开源较晚，还在孵化中
 - 相关生态不完善

Doris 简介



Apache Doris

Apache Doris是一个现代化的MPP分析型数据库产品。仅需亚秒级响应时间即可获得查询结果，有效地支持实时数据分析。Apache Doris的分布式架构非常简洁，易于运维，并且可以支持10PB以上的超大数据集。

Apache Doris可以满足多种数据分析需求，例如固定历史报表，实时数据分析，交互式数据分析和探索式数据分析等。令您的数据分析工作更加简单高效！

了解更多

参考链接: <http://doris.incubator.apache.org/master/zh-CN/>

Doris 简介

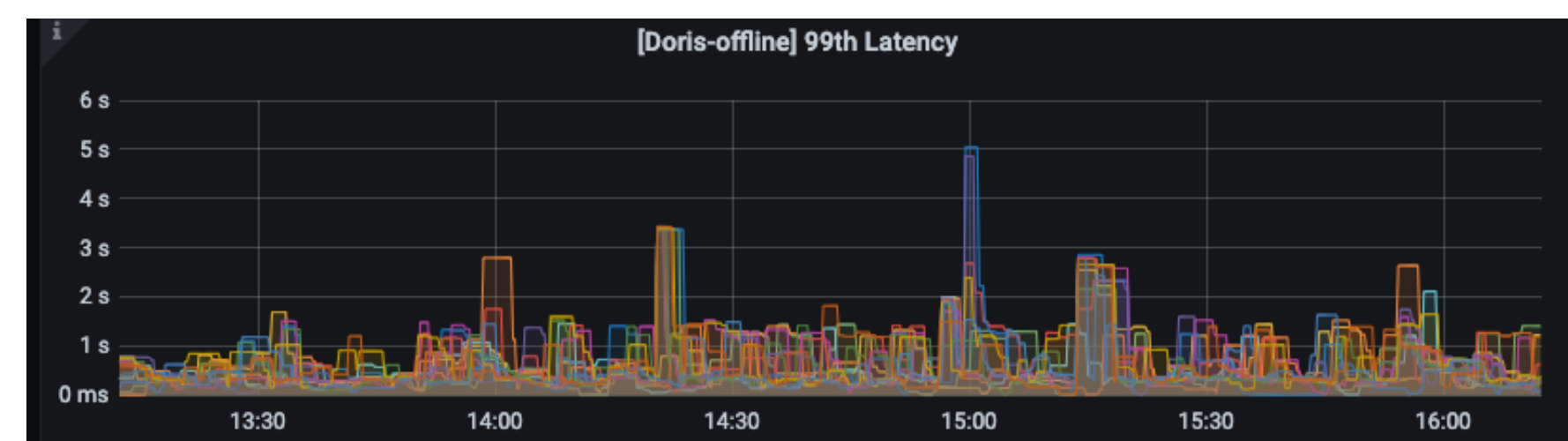
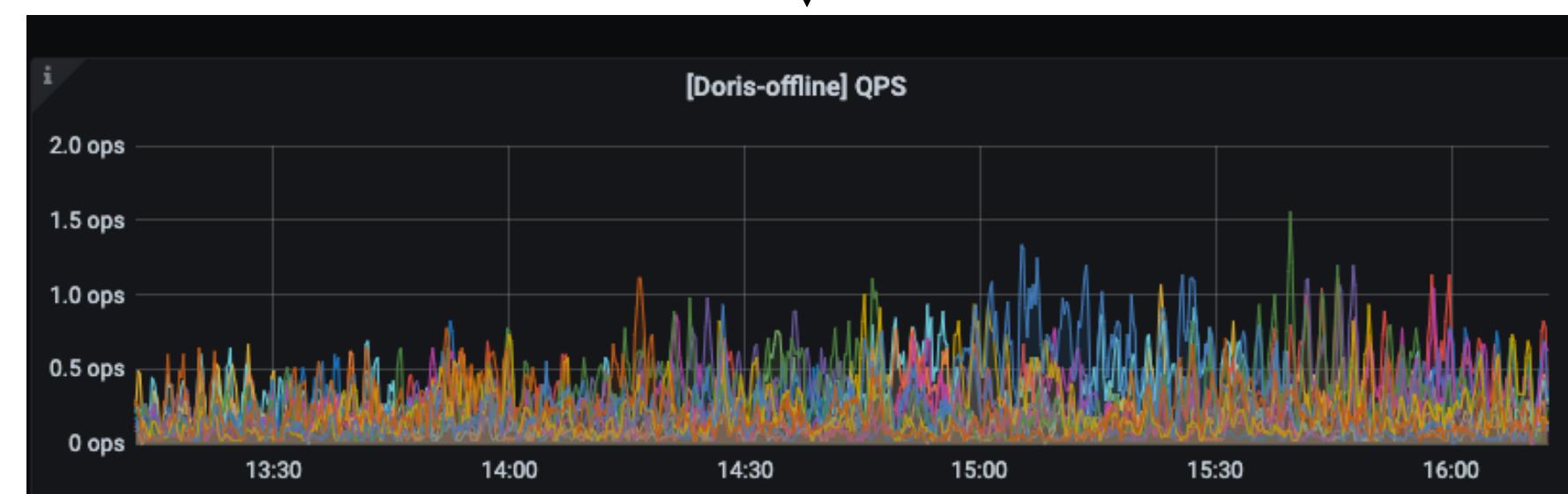
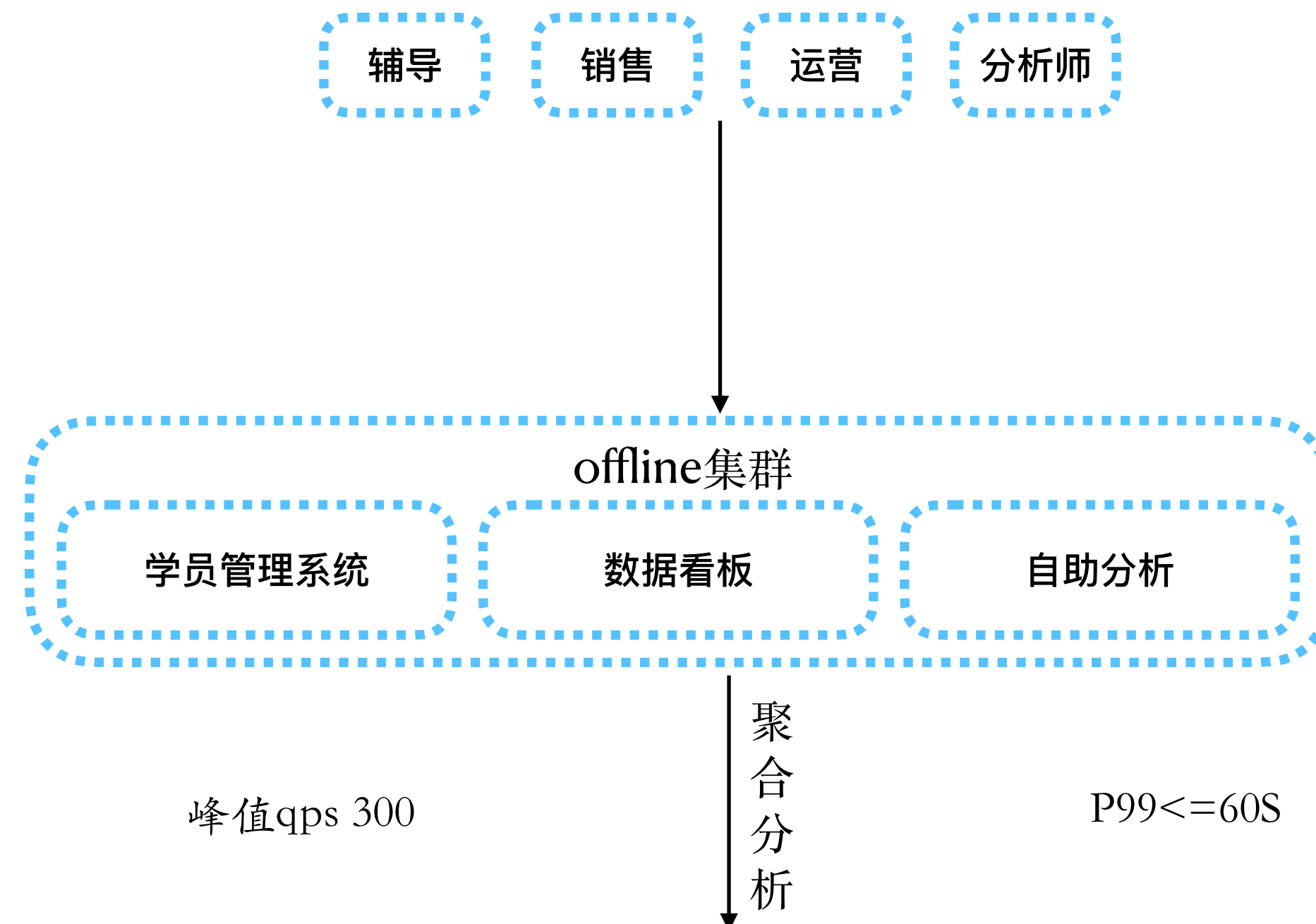
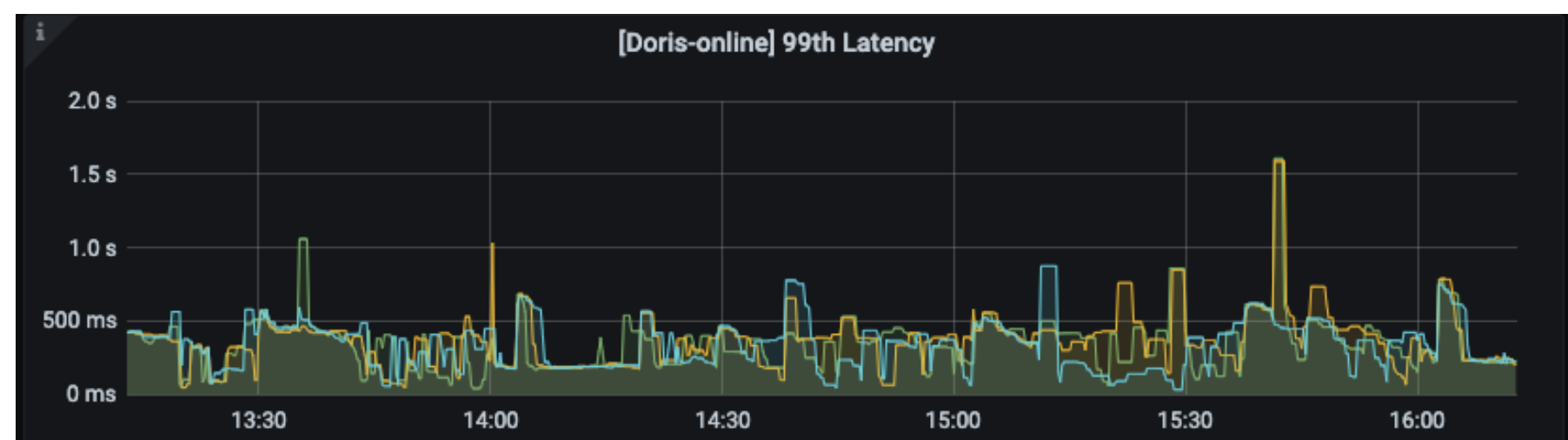
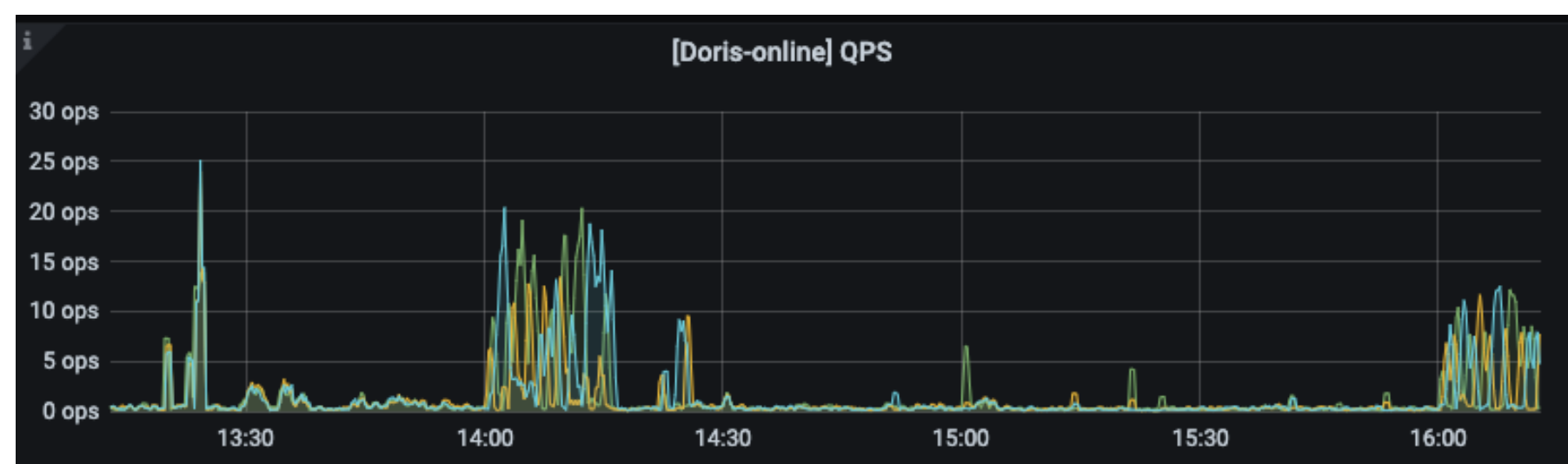
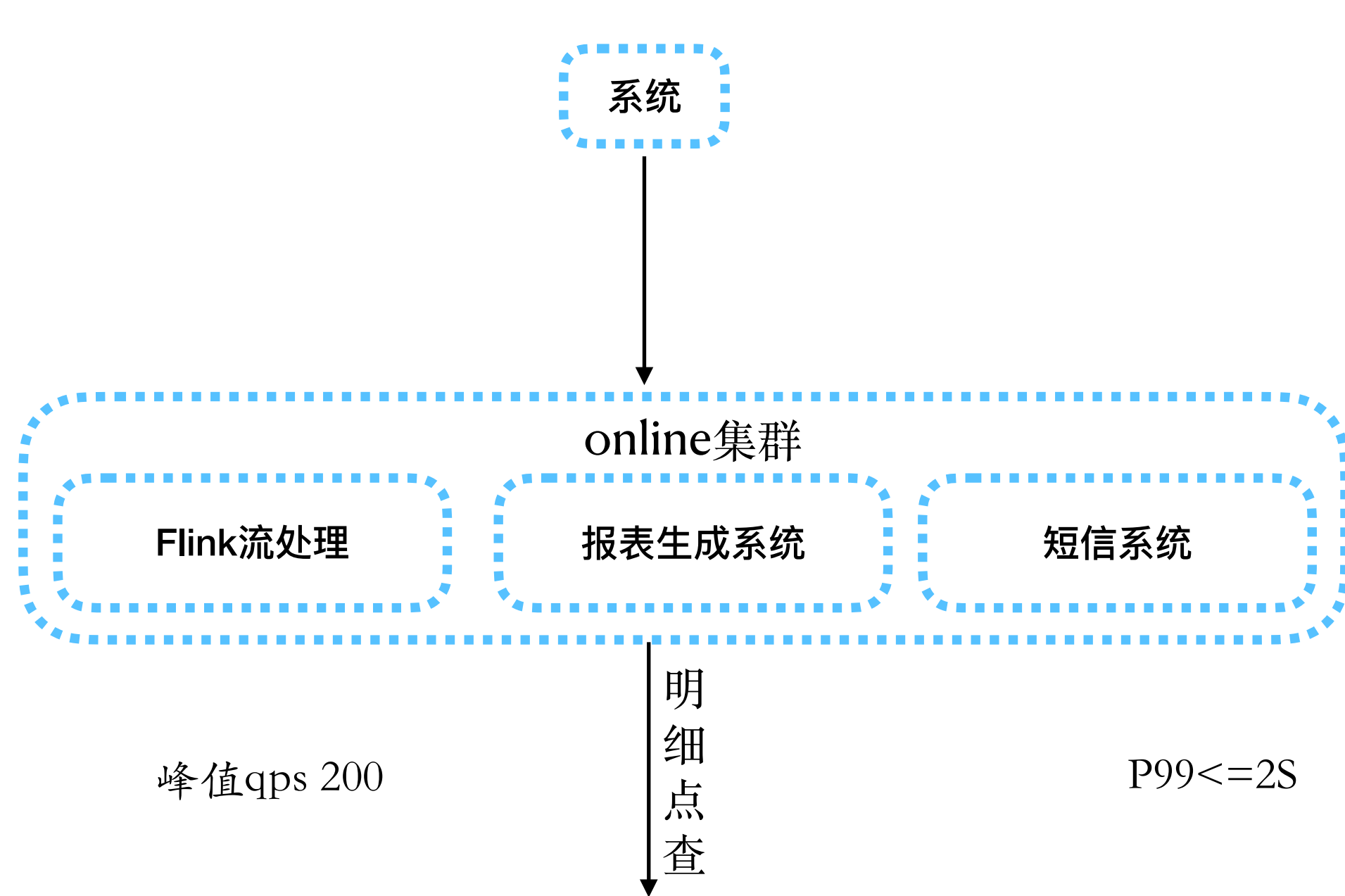
- 数据模型
 - Aggregate
 - Unique
 - Duplicate
- 导入方式
 - RoutineLoad
 - StreamLoad
 - BrokerLoad

Doris适用场景

- MPP，现场计算
- 查询加速，同步 HDFS 数据(Broker Load)
- 屏蔽 Elasticsearch 查询
- Elasticsearch 多表 Join
- Rollup
- Mysql/Elasticsearch/Doris 多数据源 Join
- 多流生成宽表
- 准实时生成分层表/宽表

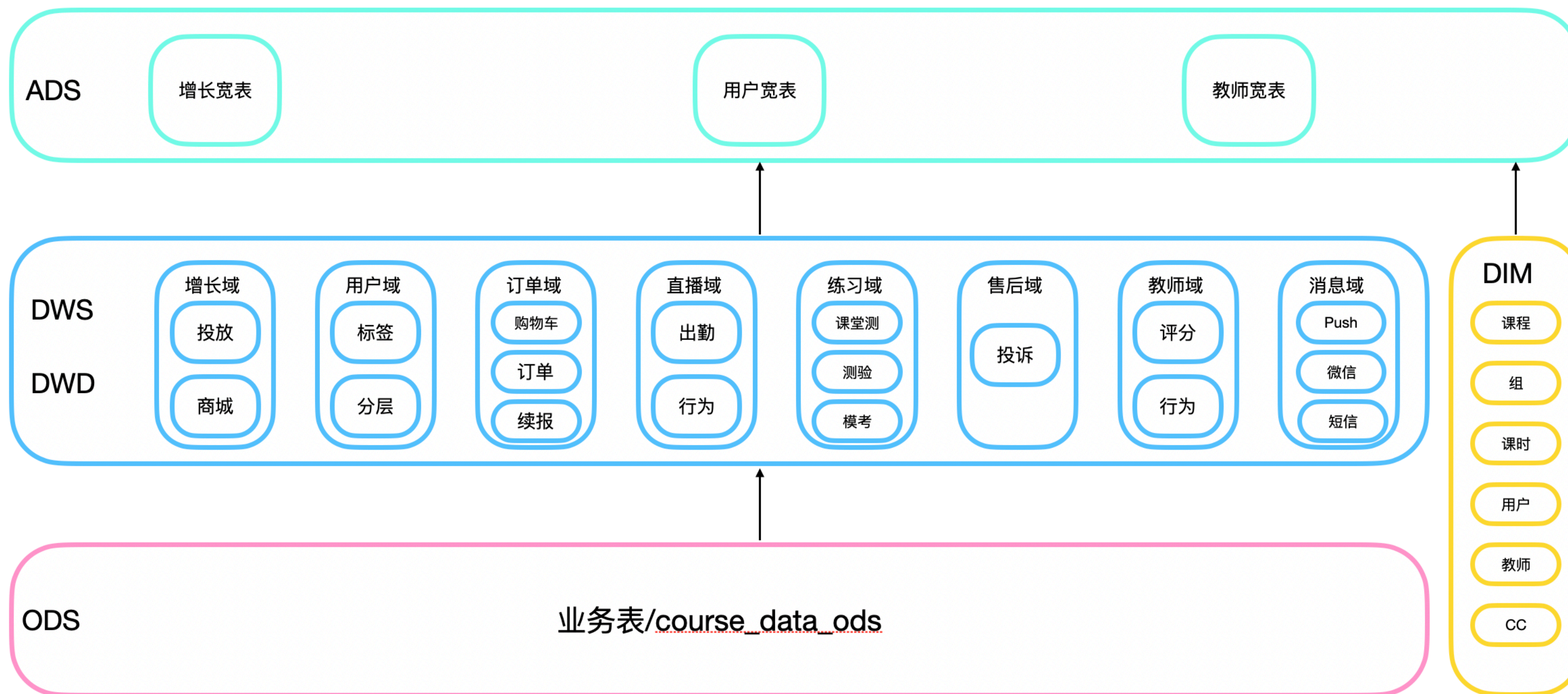
Doris在精品课的应用

集群情况



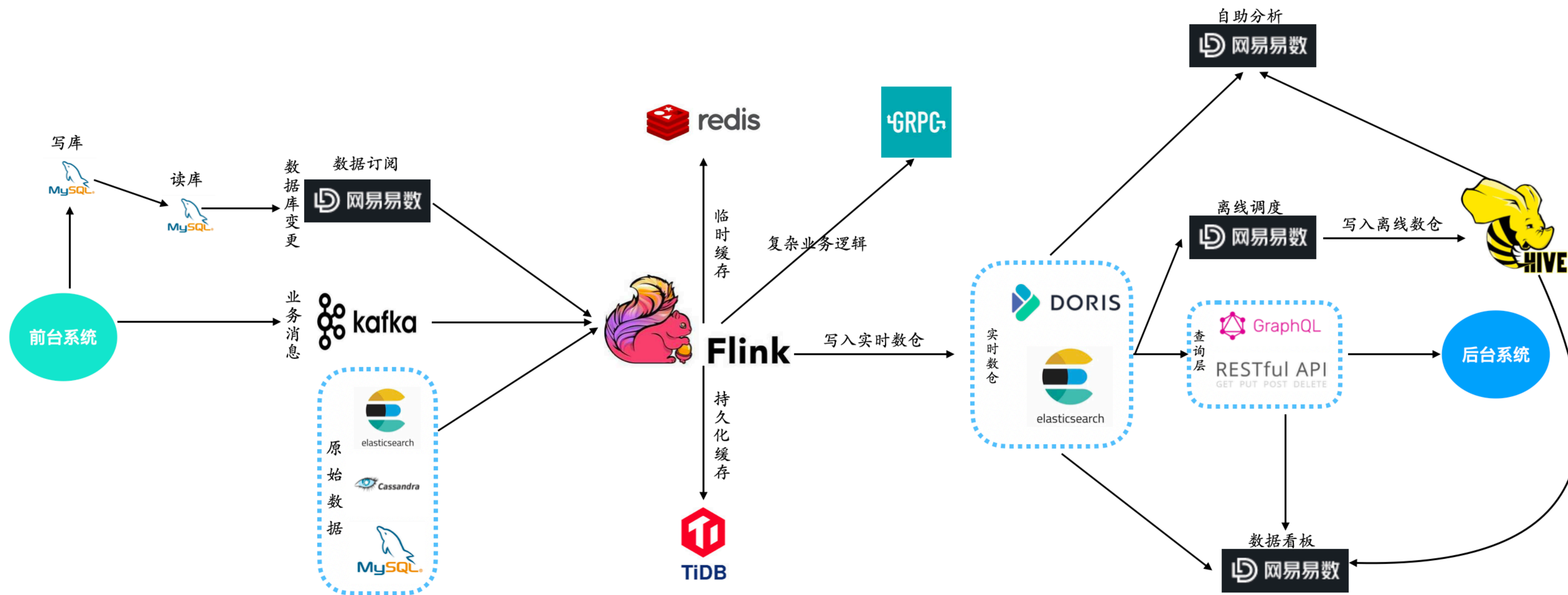
Doris在精品课的应用

数据分层



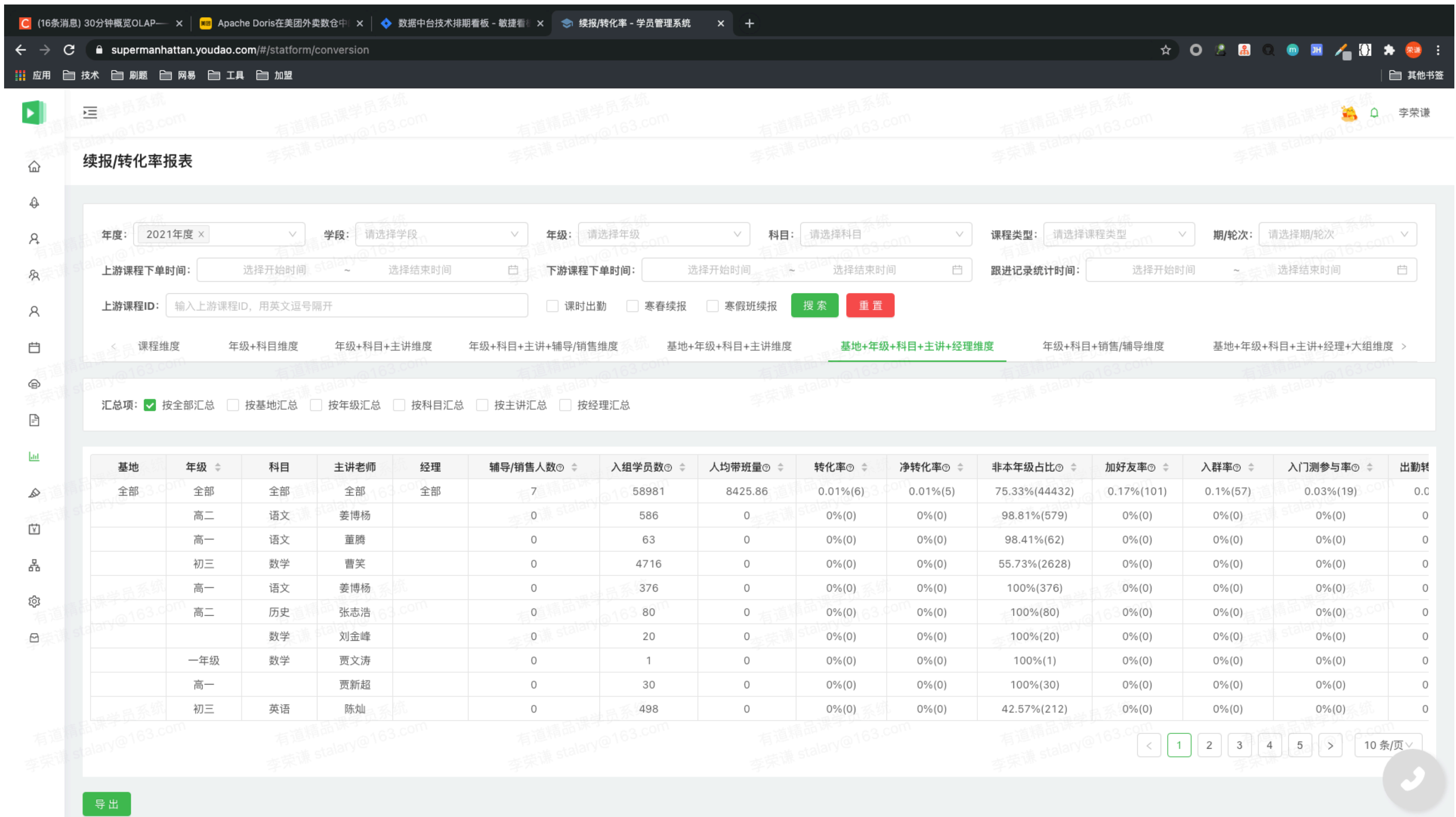
Doris在精品课的应用

数据中台实时流



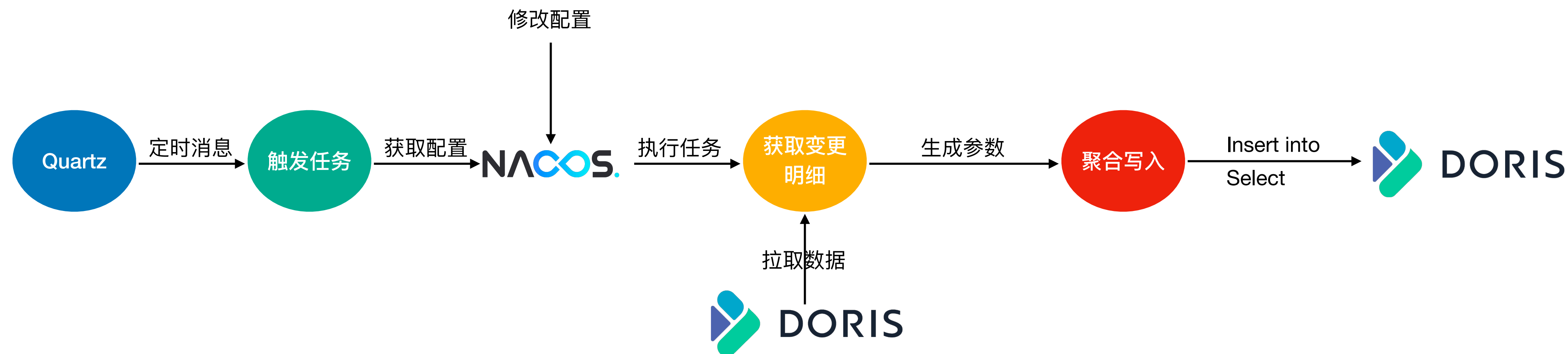
Doris在精品课的应用

多维分析



Doris在精品课的应用

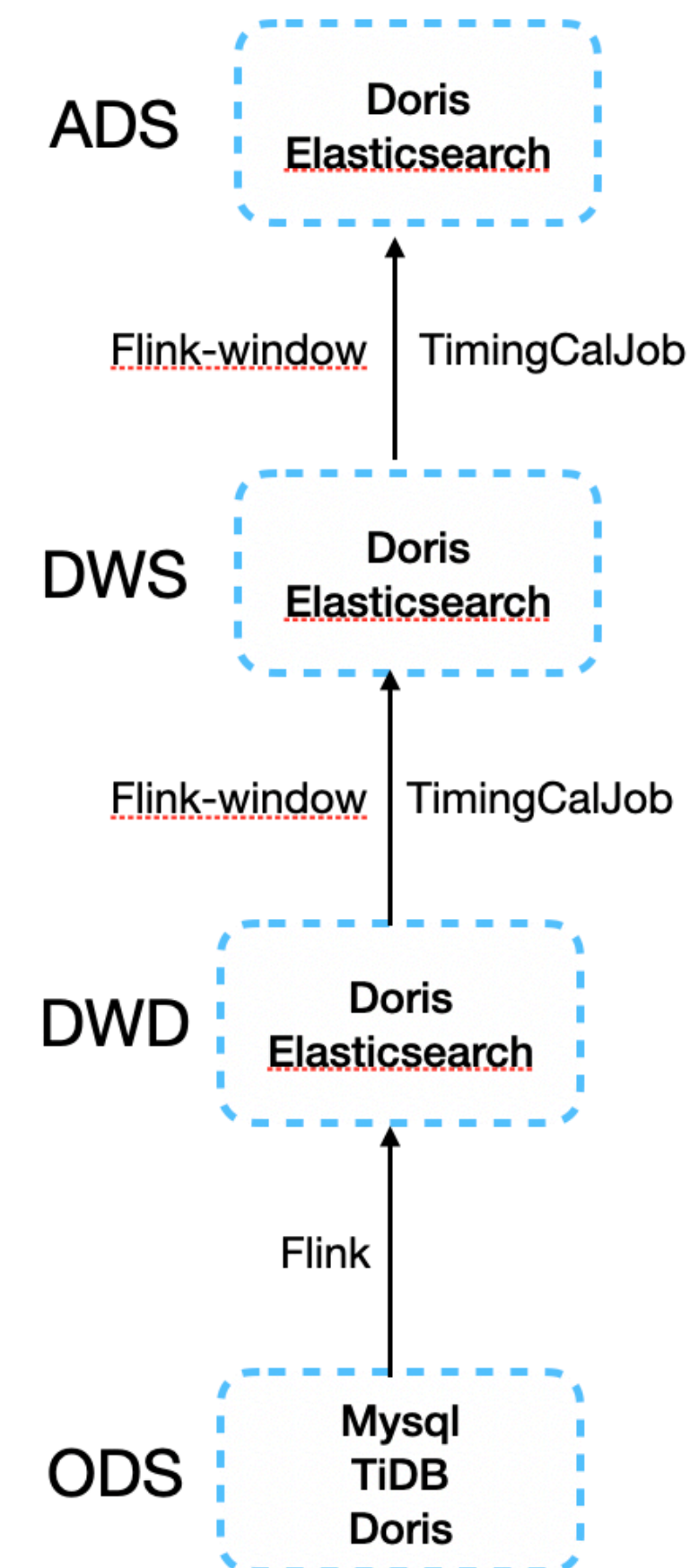
分层数据生成



配置格式: ☐ TEXT ☐ JSON ☐ XML ☒ YAML ☐ HTML ☐ Properties

配置内容 ? :

```
1 timing-cal:
2 jobs[0]:
3   calName: Cal_offline_dws_renew_ws_up
4   cluster: offline
5   listenRules[0]:
6     tableName: dwd_renew_more
7     fieldName: userId
8     targetTable: dws_renew_ws_up
9     whereFieldName: userId
10    whereFieldType: String
11    sqlList:
12      - Insert into dws_renew_ws_up(upStreamCourseId, userId, cntDownCourseId, isCanExpWinterSpring,
isRenewSpringAndWinter, isRenewNewSpringAndWinter, isRenewBeforeSpringAndWinter, isShopSpringAndWinter,
isUnRenewSpringAndWinter, cntContinueSpring, cntContinueNewSpring, cntContinueBeforeSpring, cntShopSameSubjectSpring,
cntCanExpSpring, cntShopDiffSubjectSpring, updateStamp) select drm.upStreamCourseId, drm.userId, count(1) as
cntDownCourseId, IF( count(if((drm.status = 1 or drm.status = 3) and drm.classType = '春季班', 1, null)) > 0 and count(if
((drm.status = 1 or drm.status = 3) and drm.classType = '寒假班', 1, null)) > 0 and count(if((drm.status = 1 or
```



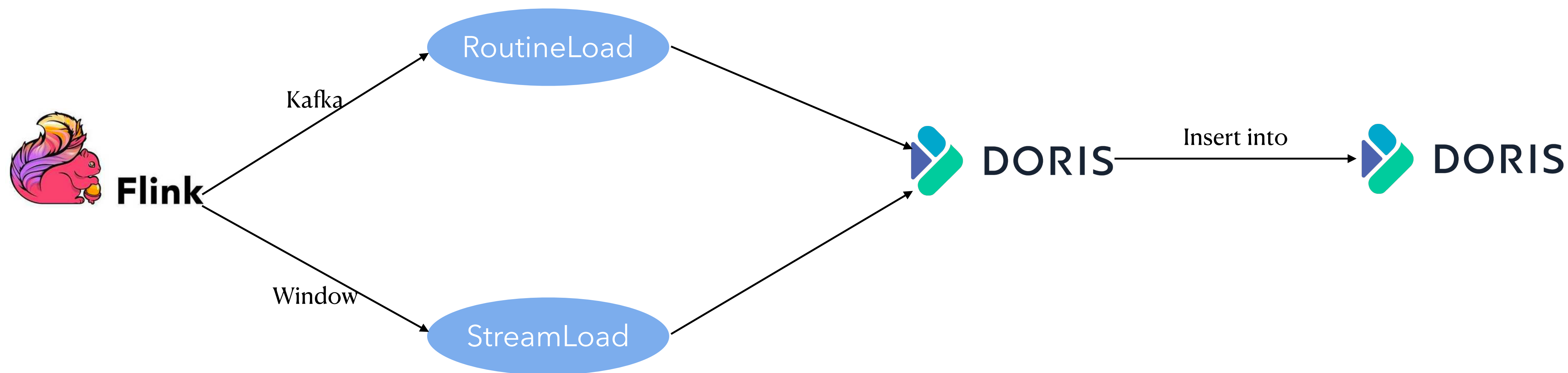
Doris在精品课的应用

使用优化

- Join优化(broadcast/colocation)
- bitmap索引
- 数据分层

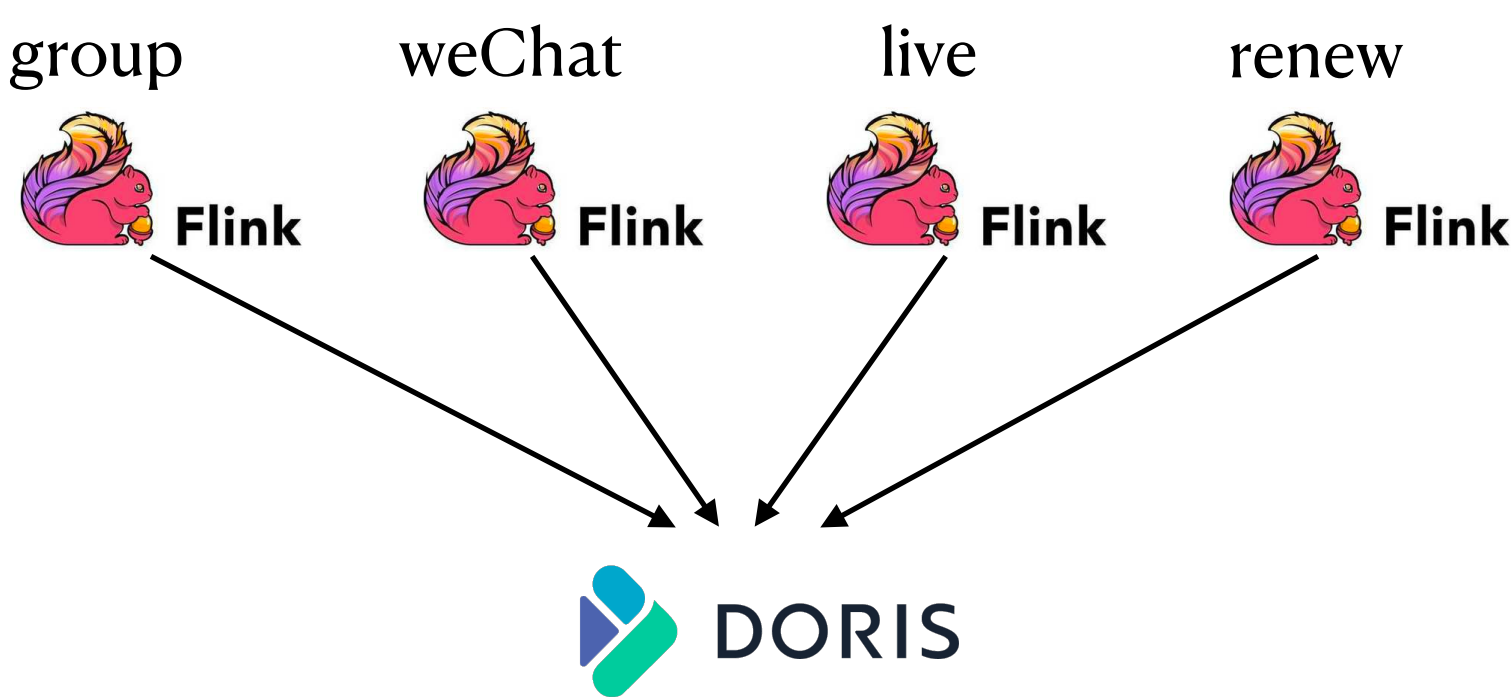
Doris在精品课的应用

数据导入

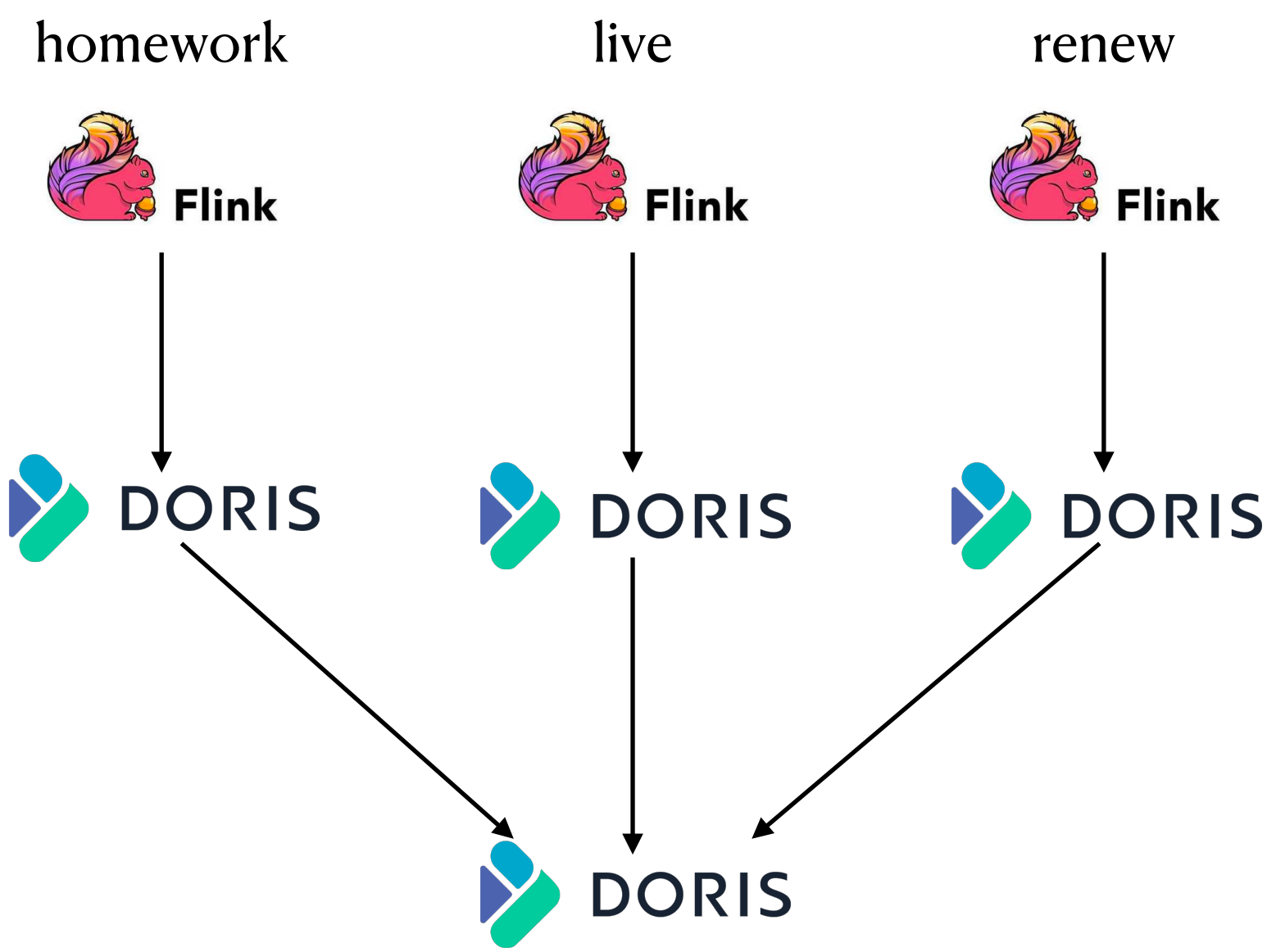


Doris在精品课的应用

宽表生成



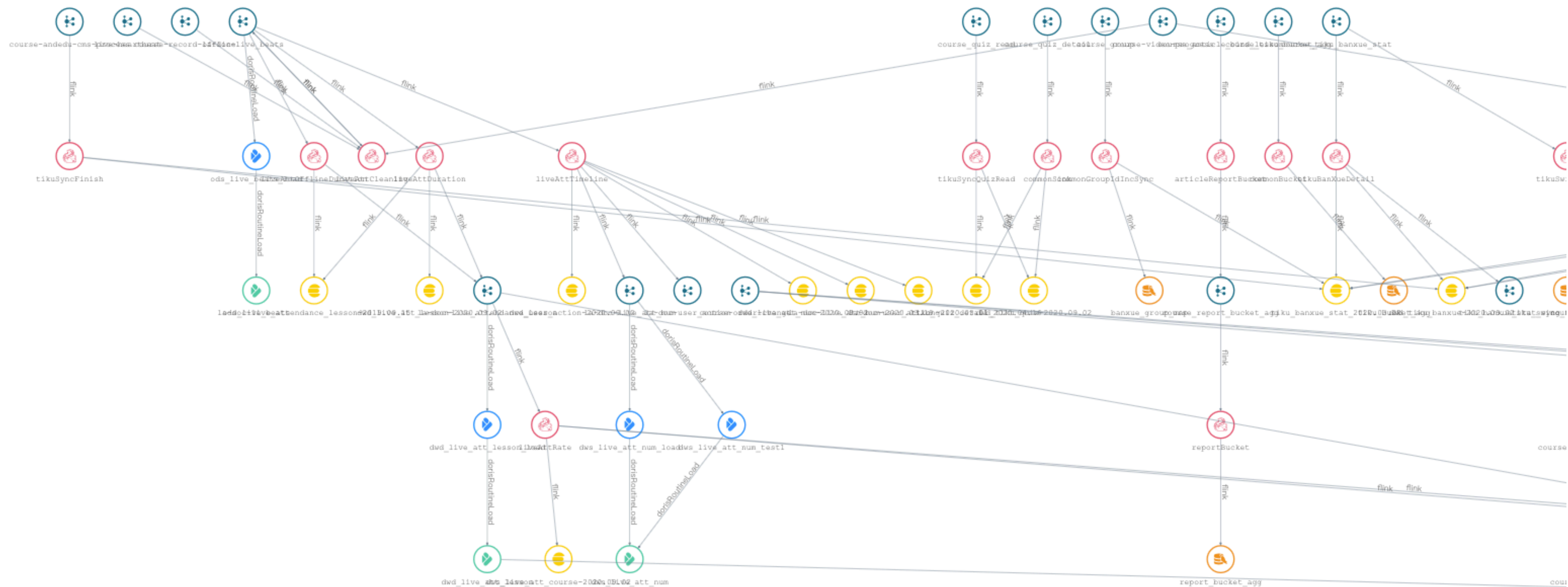
courseld	userId	groupId	weChatFriend	attendance	userRenewType
2276	stalary@163.com	3348	1	0	1



courseld	groupId	homeworkCount	attendanceCount	renewCount
2276	3348	100	80	9

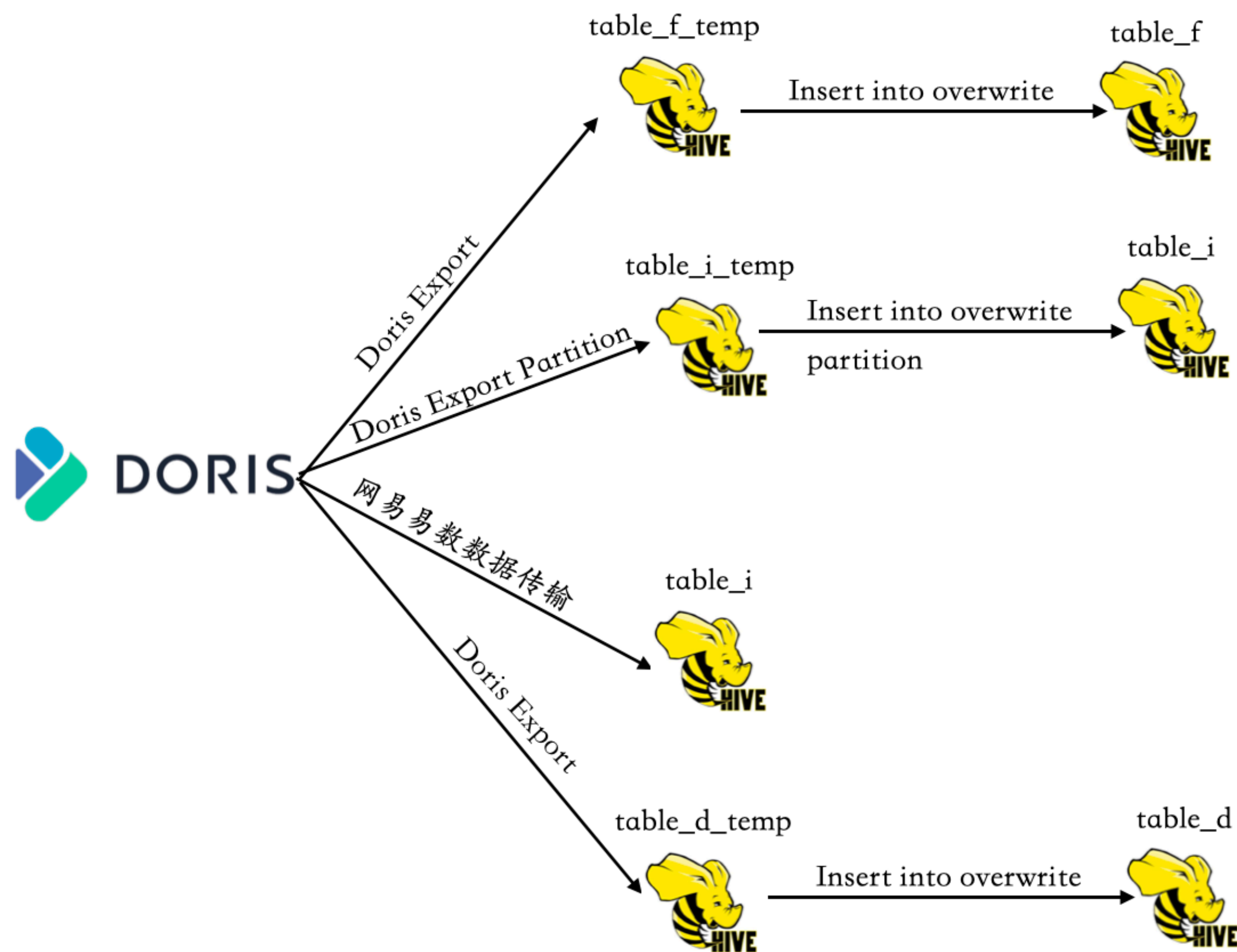
Doris在精品课的应用

数据血缘



Doris在精品课的应用

数据T+1离线同步



目前业界落地场景

Apache Doris 用户



目前业界落地场景



百度小程序用户画像 百度商业日志检索



美团外卖实时数仓



快手商业化
Doris On Es



京东广告平台报表



小米增长分析平台



作业帮实时数仓

参考资料

- 官方文档: <http://doris.incubator.apache.org/master/zh-CN/installing/compilation.html>
- OLAP介绍: <https://blog.csdn.net/xwc35047/article/details/86369465>
- Apache Doris在美团外卖数仓中的应用实践: <https://tech.meituan.com/2020/04/09/doris-in-meituan-waimai.html>

Q&A



评价二维码



Github二维码