

云上屈原技术方案

大数据驾驶舱

一. 架构设计

1. 数据规模&时效

第六章 系统建设其他需求

5.1 信息量指标

用户量

本项目建设系统主要面向省、市、县区全体教育体系以及文旅单位，同时面向群众提供信息上报和公开入口，具体用户量估算如下：

序号	用户类别	用户范围	用户量	3~5 年
1	综合素质教育	湖北省教师学生群体	50 万人	500 万人
5	研学用户	全国	10 万人	100 万人
6	其他用户	全国	20 万人	200 万人

- 1. 3-5年：800万用户。
- 2. 数据量预估：
 - a. 每日1/10的在线用户估算，80万用户在线。日活均值：40万用户在线。
 - b. 每人每天平均100条行为日志/数据，总量4000万。
 - c. 每条日志1K大小，数据量大小为40G。
 - d. 数据量大小：
 - i. 一年：40*360 = 14T；三年：40*365*3 = 43800G = 42T
 - ii. 三年数据量：2个副本&1/10的压缩 = 8.4T
- 3. 更新周期：最短1h，所有需求均需要当天内指标统计。

2. 框架版本

- 1. Apache/CDH/HDP
 - a. Apache：免费，但运维麻烦。需运维自行调研组件，配置组件兼容性。
 - b. CDH：收费。一个节点一万美金/年。

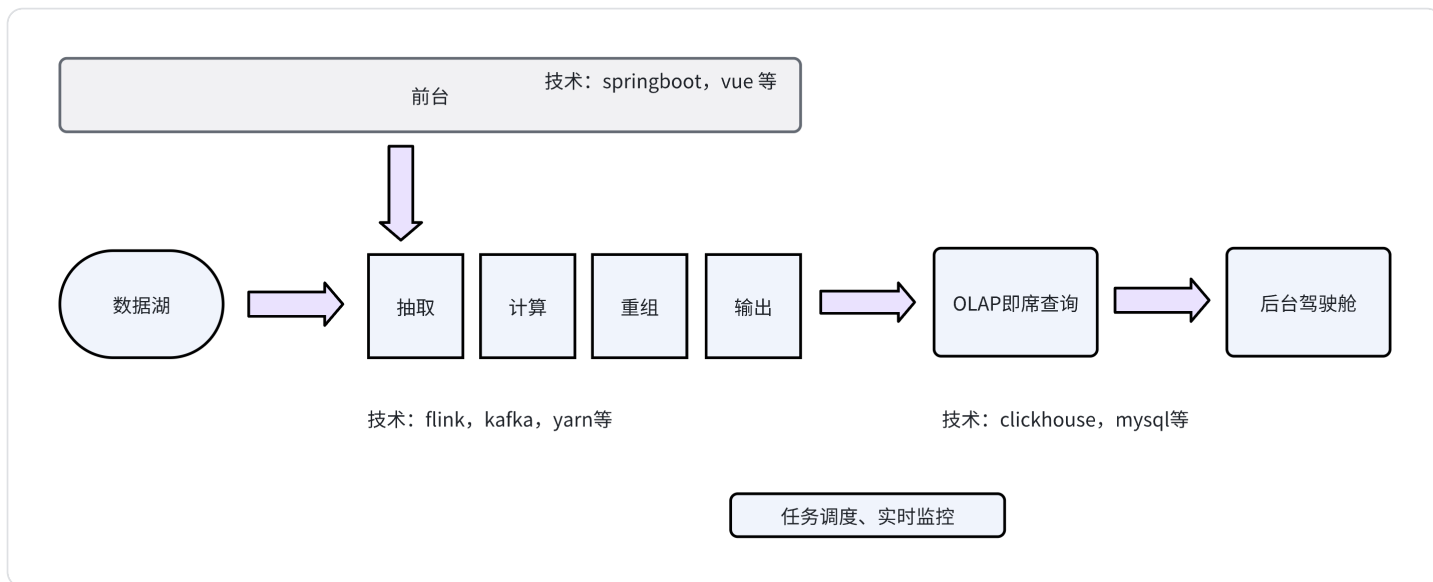
c. HDP：开源，可二次开发。没有CDH稳定，国内使用少。

2. 云服务：

a. 阿里云EMR、MaxCompute、DataWorks；亚马逊云EMR；腾讯云EMR；华为云EMR

b. 运维简单，收费，国内用户逐年增长。

3. 技术架构



4. 技术选型

计算引擎-Flink

需求：T+0当天内，实时。

Flink：流式，状态后端，准确性和时效性都高。

SparkStreaming：微批次，时效性、计算速率不如flink。

Storm：准确性欠缺。

数据存储-ClickHouse

需求：实时聚合查询

Clickhouse：DWS层到ads层需重新分组再聚合，clickhouse是列式存储做聚合操作有优势，且非常擅长单宽表的聚合分析。

Hbase：k-v存储，更适合单条数据的k-v查询。

消息队列-Kafka

Kafka：解耦、削峰、适用于大数据量。

MQ：吞吐量不如Kafka，不适合大数据场景。

数据同步-flinkCDC

选型原因：

- a. 基于Debezium实现，实时同步，且低延迟。且改进了debezium单机的缺陷，可对接分布式系统。
- b. 使用flinkCDC，可直接把维度配置信息同步至flink，不需要再经过kafka。

	基于查询的CDC	基于Binlog的CDC
开源产品	Sqoop、Kafka JDBC Source	Canal、Maxwell、 Debezium
执行模式	Batch	Streaming
是否可以捕获所有数据变化	否	是
延迟性	高延迟	低延迟
是否增加数据库压力	是	否

数据同步策略：因用户表数据量不小，采用增量同步。

mysql：开启binlog。

应用数据-SpringBoot接口发布&MySQL

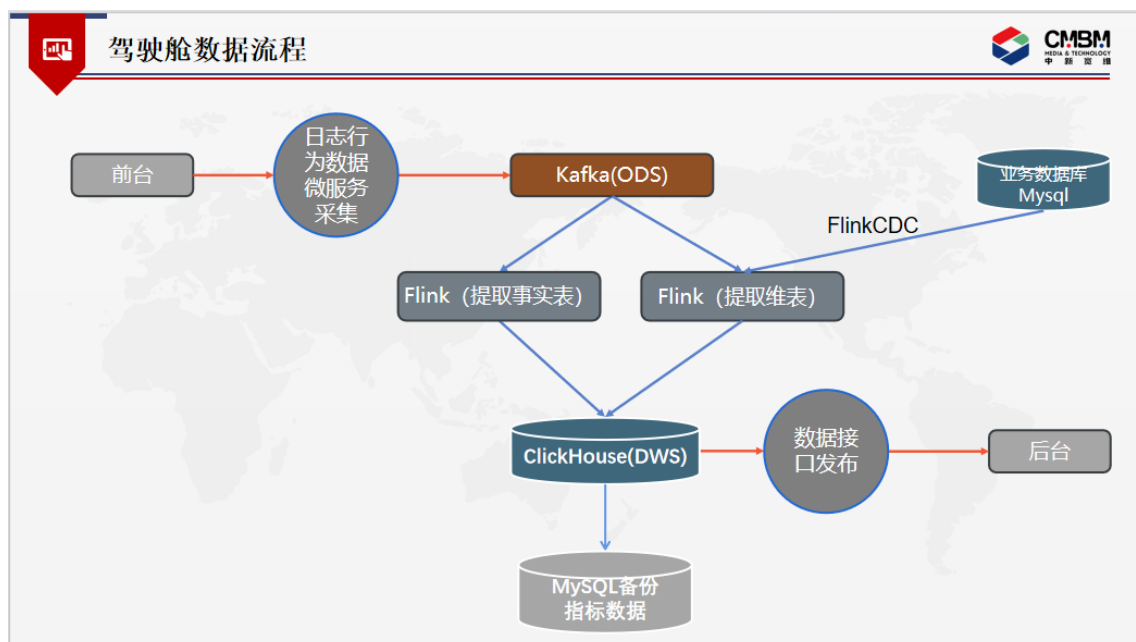
SpringBoot:使用SpringBoot微服务查询OLAP库进行聚合分析，计算结果发布成数据接口至驾驶舱。

MySQL：指标数据备份。

日志数据采集-SpringBoot

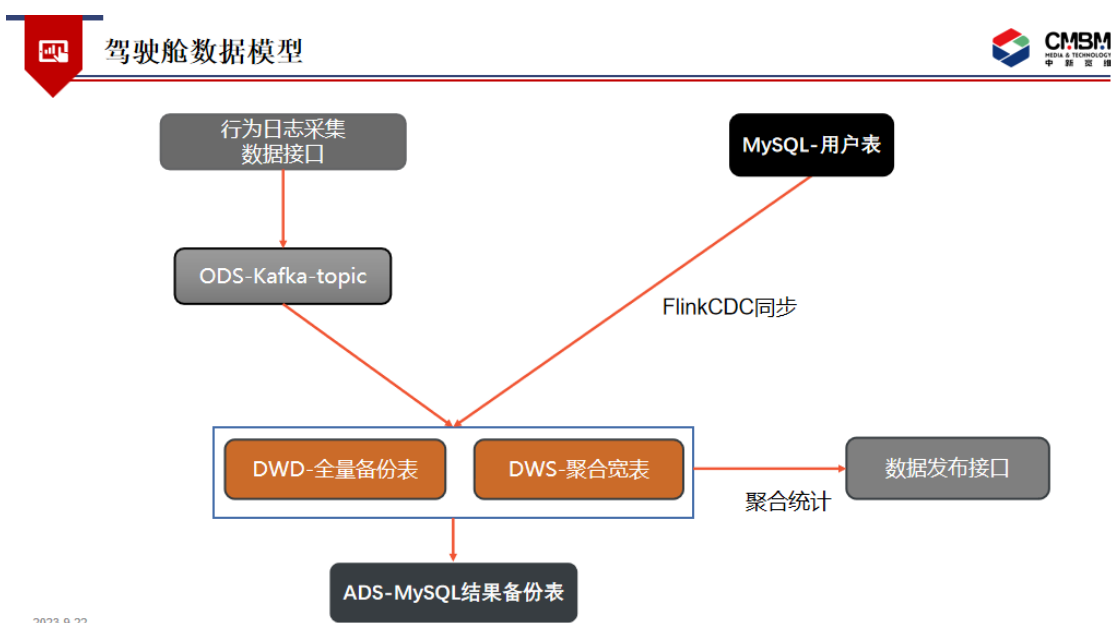
前后端通用、成熟交互方式。

5. 数据流程



二. 数据体系

分层模型



2023-9-22

维度模型

1. 事实：对应业务过程（一个个不可拆分的行为事件）
 - a. 行为日志表，分主题。
2. 维度：对应业务过程发生时所处的环境。
 - a. 用户
 - b. 地区
 - c. 日期

DWS

ClickHouse：日志宽表

ODS

Kafka：行为日志topic

1 数据字段：浏览行为，行为时间，

2 样例数据：

3 {}

MySQL：用户业务表

字段名	字段说明	数据类型	缺省值	样例数据	提取规则
	用户ID				
	用户名				
	性别				
	年龄				
	地区				
	注册时间				

数据采集API接口

接口名称	接口路径	请求方式	样例数据

ADS

MySQL：指标结果表

发布API接口

接口名称	接口路径	请求方式	样例数据