

DMP 项目总结V1.0

业务场景

根据用户id查询该用户对应的标签

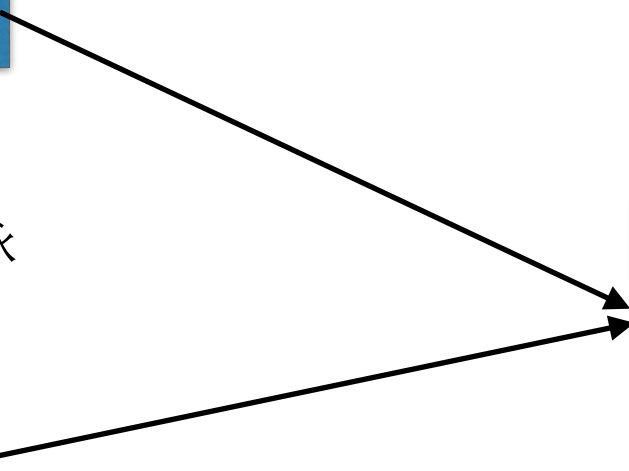
用户A: 男,电影,游戏,北京,游泳

根据标签查询含有该标签的用户

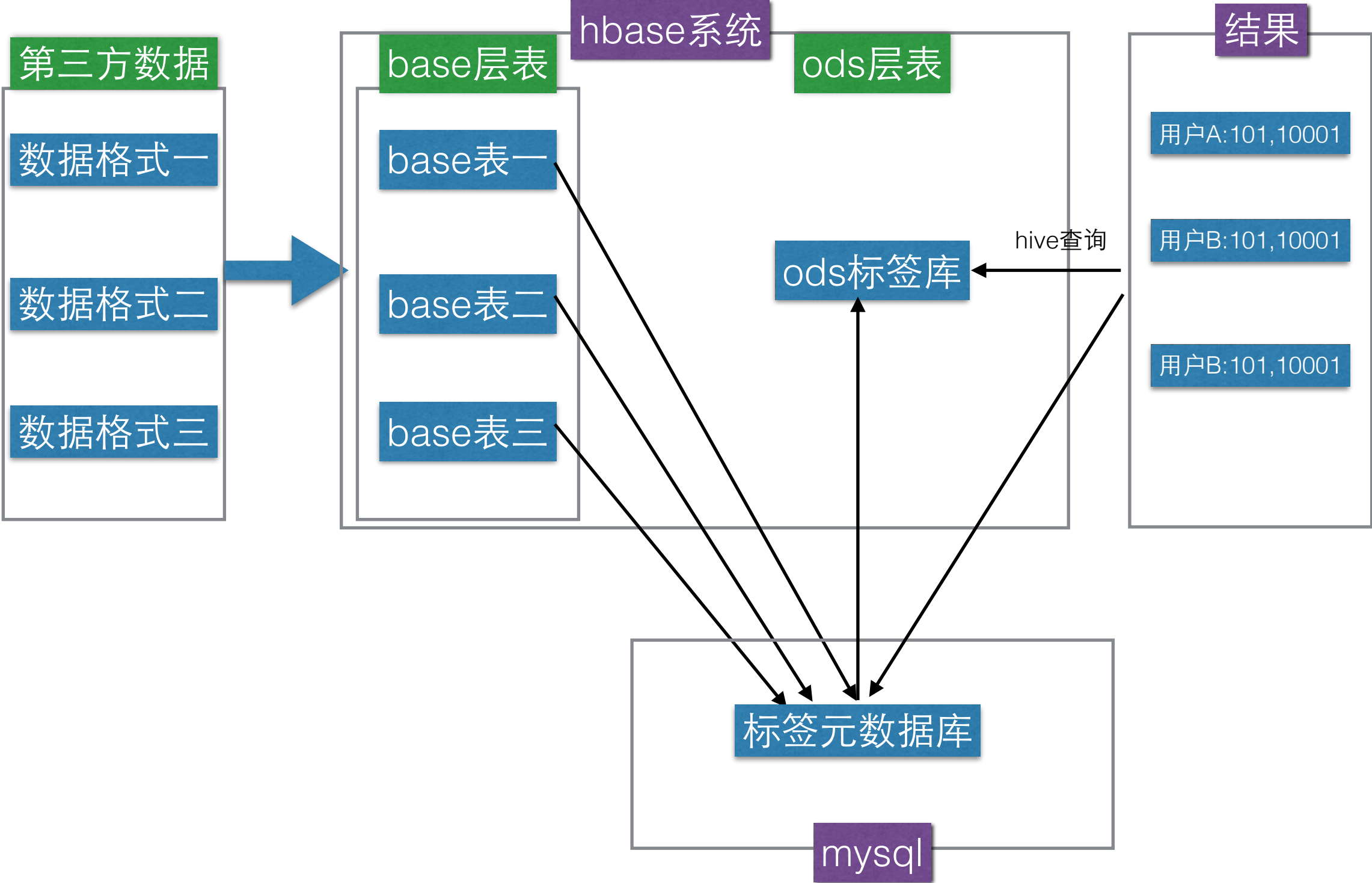
游泳:用户A, 用户B

准实时

hive查询



架构设计图



详细设计

数据清洗模块

第三方过来的数据，根据不同的数据源建不同的base表，这部分作为dmp的基础数据，保持完整性，同步到hbase表中，第三方的数据有几种格式，就需要建几张base层表，所以数据清洗加载数据到hbase的base层表，较易实现

mysql标签元数据库

针对标签可能扩展为几万、十几万、甚至百万级的量，所以考虑用mysql来存储标签定义，每一个主键id作为一个标签号，例如：

一级标签号从1-999

二级标签用它对应的一级标签id为前缀，例如标签1的下级标签是1001-1999，标签10的下级标签是10001-10999

同理，三级标签用它对应的二级标签id为前缀

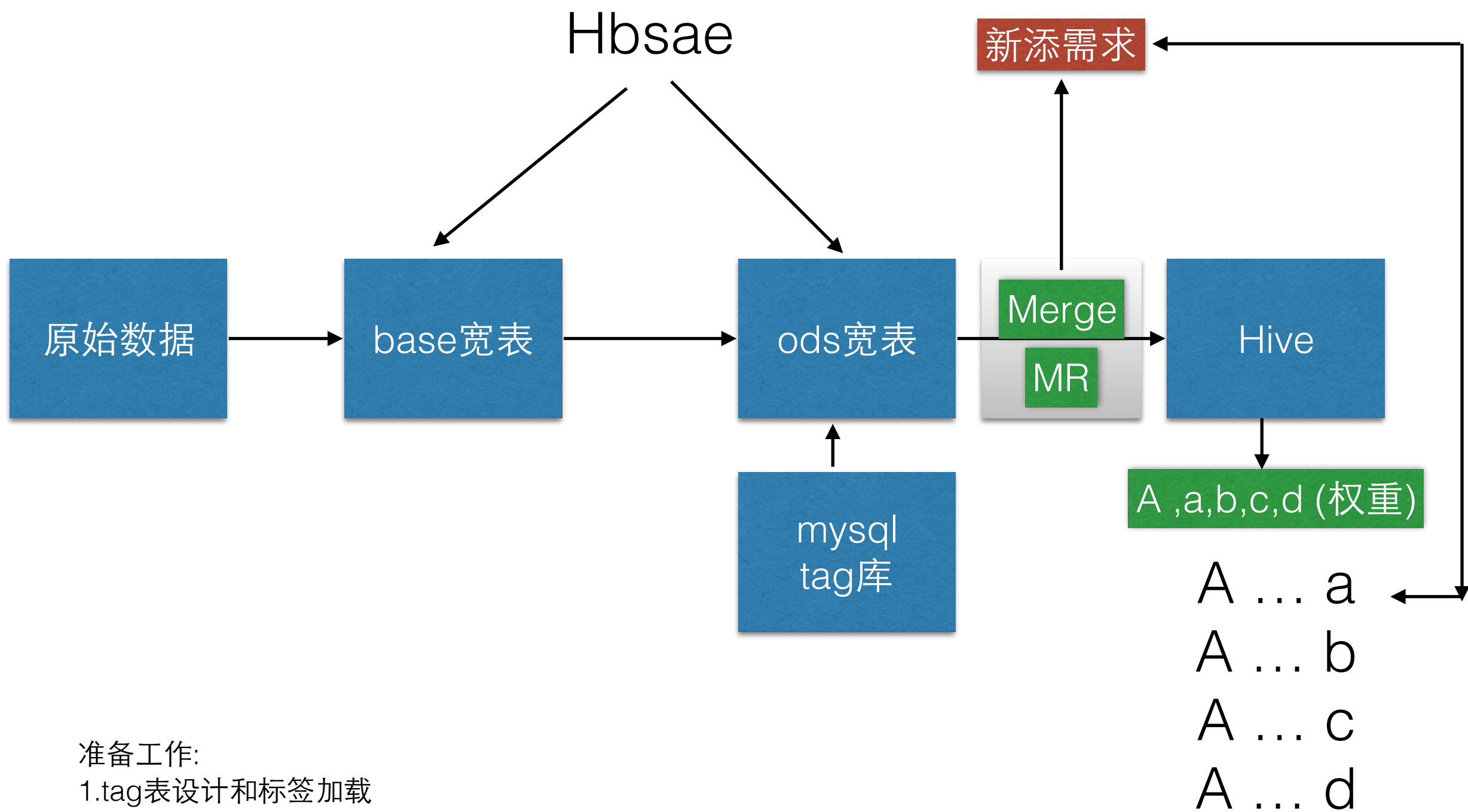
.....

也就是说如果一个用户都含有这三种标签，那么应该是：用户A:2,1003,1001001

Hbase ods层表

ods表设计

rowkey(用户标识)	列一(唯一标识)	列二(标签id)
A	A1	2,1003,1001001
B	B1	3,1001

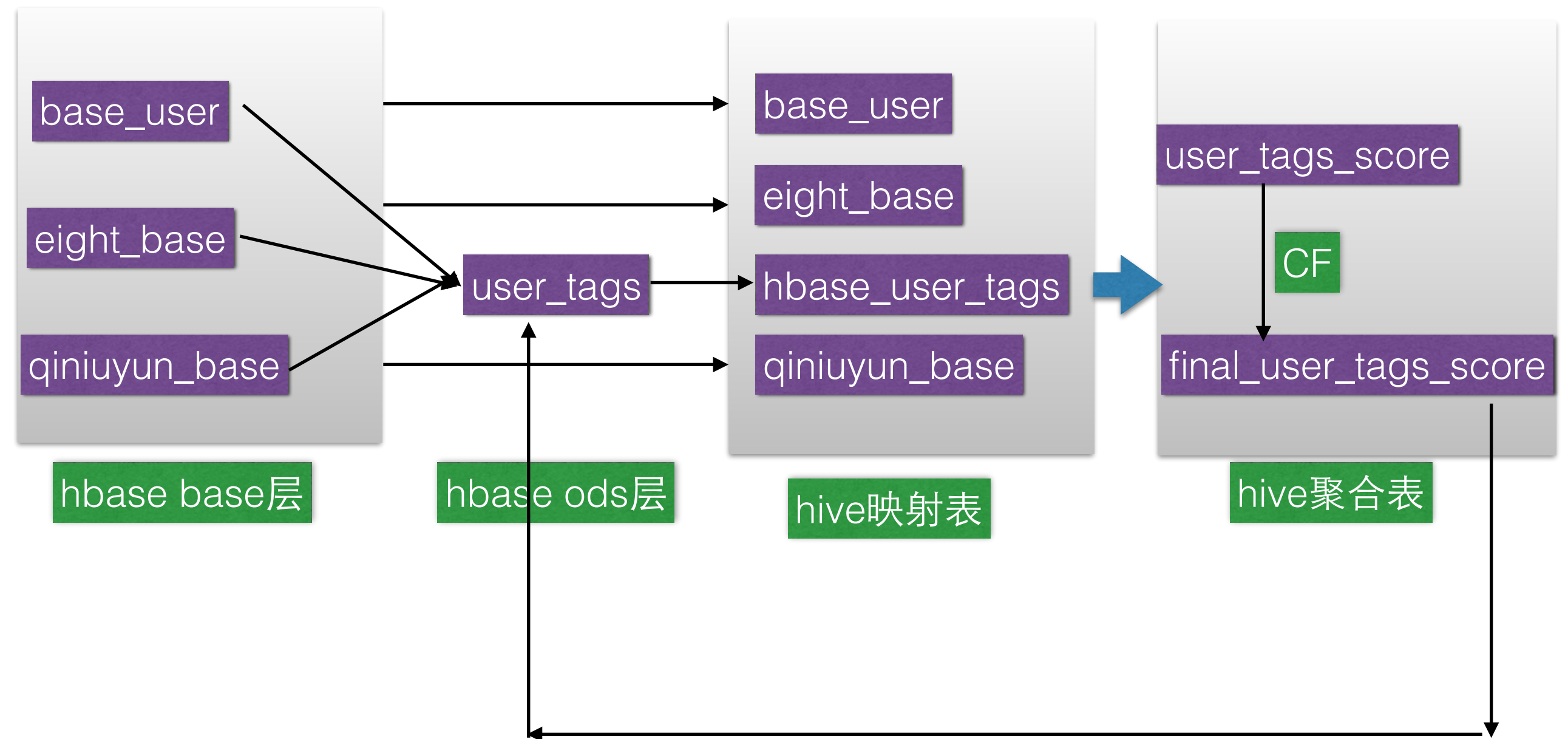


准备工作:

- 1.tag表设计和标签加载
- 2.融合数据用户标识生成规则

时间0609

实践过程



数据融合和映射方案

- 1.根据用户id查询标签
- 2.根据标签查询用户id

