# DataFunSummit

# 大数据计算

## 架构峰会

## 实时计算 论坛

# 目录
CONTENTS

京东 | **DataFunSummit**

# 01
# 问题

Subject

实时维度建模的过程中，有很多技术场景，这里我们就只就其中的部分难点场景进行阐述

# 问题

在这里我们只选择了里面的两类代表性问题进行展开阐述

➤ 问题一： 实时多流全量关联的问题

```
select * from A full join B on A.name = B.name;
```

➤ 问题二： 实时流全量分组计算问题

```
select id,name,val,row_number() over (partition by name order by val) as rn from A ;
select name, min(val) from A group by k;
```

上面提到的两类问题 可以直接使用Flink SQL去做简单的SQL完成吗

➤ 流是从 [now, +∞) 但是我们需要历史全量

➤ 状态存放在内存中是有大小限制的

➤ 状态存放在rocksdb性能不能满足

# 问题一：实时多流全量关联的问题

[0, now)

（1）

（2）

（3）

| A | | |
|---|---|---|
| id | name | val |
| 100001 | aaa | 23 |
| 100002 | aaa | 55 |
| 100003 | bbb | 21 |
| 100004 | bbb | 33 |
| 100005 | bbb | 66 |
| 100006 | ccc | 43 |

| B | | |
|---|---|---|
| id | name | val |
| 100001 | aaa | pp1 |
| 100002 | aaa | pp2 |
| 100003 | bbb | pp3 |
| 100004 | bbb | pp4 |
| 100005 | ddd | pp5 |

| A | | | | | B | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100001 | aaa | 23 | | | 100001 | aaa | pp1 | | |
| 100002 | aaa | 55 | | | 100002 | aaa | pp2 | | |

| A | | | | | B | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100003 | bbb | 21 | | | 100003 | bbb | pp3 | | |
| 100004 | bbb | 33 | | | 100004 | bbb | pp4 | | |
| 100005 | bbb | 66 | | | | | | | |

| A | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100006 | ccc | 43 | | | | | | | |

| | | | | | B | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | 100005 | ddd | pp5 | | |

| A full outer join B | | | | | |
|---|---|---|---|---|---|
| Aid | Aname | Aval | Bid | Bname | Bval |
| 100001 | aaa | 23 | 100001 | aaa | pp1 |
| 100002 | aaa | 55 | 100001 | aaa | pp1 |
| 100001 | aaa | 23 | 100002 | aaa | pp2 |
| 100002 | aaa | 55 | 100002 | aaa | pp2 |
| 100003 | bbb | 21 | 100003 | bbb | pp3 |
| 100004 | bbb | 33 | 100003 | bbb | pp3 |
| 100005 | bbb | 66 | 100003 | bbb | pp3 |
| 100003 | bbb | 21 | 100004 | bbb | pp4 |
| 100004 | bbb | 33 | 100004 | bbb | pp4 |
| 100005 | bbb | 66 | 100004 | bbb | pp4 |
| 100006 | ccc | 43 | | | |
| | | | 100005 | ddd | pp5 |

[now, +∞)

（1）　　　　　　　　　　　　　　（2）　　　　　　　　　　　　　　（3）

**A**

| | id | name | val | mid |
|---|---|---|---|---|
| src | 100001 | aaa | 23 | |
| cur | 100001 | bbb | 23 | 1 |

| A | | | | | | B | | |
|---|---|---|---|---|---|---|---|---|
| 100002 | aaa | 55 | | | | 100001 | aaa | pp1 |
| 100001 | aaa | 23 | 1 | D | | 100002 | aaa | pp2 |

| A | | | | | | B | | |
|---|---|---|---|---|---|---|---|---|
| 100003 | bbb | 21 | | | | 100003 | bbb | pp3 |
| 100004 | bbb | 33 | | | | 100004 | bbb | pp4 |
| 100005 | bbb | 66 | | | | | | |
| 100001 | bbb | 23 | 1 | I | | | | |

| A | | | | | |
|---|---|---|---|---|---|
| 100006 | ccc | 43 | | | |

| | | | | B | | |
|---|---|---|---|---|---|---|
| | | | | 100005 | ddd | pp5 |

| A full outer join B | | | | | | | |
|---|---|---|---|---|---|---|---|
| Aid | Aname | Aval | Bid | Bname | Bval | mid | opt |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | | |
| 100002 | aaa | 55 | 100001 | aaa | pp1 | | |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | | |
| 100002 | aaa | 55 | 100002 | aaa | pp2 | | |
| 100003 | bbb | 21 | 100003 | bbb | pp3 | | |
| 100004 | bbb | 33 | 100003 | bbb | pp3 | | |
| 100005 | bbb | 66 | 100003 | bbb | pp3 | | |
| 100003 | bbb | 21 | 100004 | bbb | pp4 | | |
| 100004 | bbb | 33 | 100004 | bbb | pp4 | | |
| 100005 | bbb | 66 | 100004 | bbb | pp4 | | |
| 100006 | ccc | 43 | | | | | |
| | | | 100005 | ddd | pp5 | | |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 1,0 | D |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | 1,0 | D |
| 100001 | bbb | 23 | 100003 | bbb | pp3 | 1,0 | I |
| 100001 | bbb | 23 | 100004 | bbb | pp4 | 1,0 | I |

# 问题一：实时多流全量关联的问题

（1）

| A | | | |
|---|---|---|---|
| | id | name | val | mid |
| src | 100001 | aaa | 23 | 1 |
| cur | 100001 | bbb | 23 | |

| A | | | |
|---|---|---|---|
| | id | name | val | mid |
| src | 100001 | bbb | 23 | 2 |
| cur | 100001 | aaa | 23 | |

（2）

| A | | | | | | B | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100002 | aaa | 55 | | | | 100001 | aaa | pp1 | |
| 100001 | aaa | 23 | 2 | | | 100002 | aaa | pp2 | |

| A | | | | | | B | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100003 | bbb | 21 | | | | 100003 | bbb | pp3 | |
| 100004 | bbb | 33 | | | | 100004 | bbb | pp4 | |
| 100005 | bbb | 66 | | | | | | | |
| 100001 | bbb | 23 | 2 | | | | | | |

| A | | | |
|---|---|---|---|
| 100006 | ccc | 43 | |

| | | | | B | | |
|---|---|---|---|---|---|---|
| | | | | 100005 | ddd | pp5 |

[now, +∞)　　（3）

| A full outer join B | | | | | | | |
|---|---|---|---|---|---|---|---|
| Aid | Aname | Aval | Bid | Bname | Bval | mid | opt |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | | |
| 100002 | aaa | 55 | 100001 | aaa | pp1 | | |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | | |
| 100002 | aaa | 55 | 100002 | aaa | pp2 | | |
| 100003 | bbb | 21 | 100003 | bbb | pp3 | | |
| 100004 | bbb | 33 | 100003 | bbb | pp3 | | |
| 100005 | bbb | 66 | 100003 | bbb | pp3 | | |
| 100003 | bbb | 21 | 100004 | bbb | pp4 | | |
| 100004 | bbb | 33 | 100004 | bbb | pp4 | | |
| 100005 | bbb | 66 | 100004 | bbb | pp4 | | |
| 100006 | ccc | 43 | | | | | |
| | | | 100005 | ddd | pp5 | | |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 1,0 | D |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | 1,0 | D |
| 100001 | bbb | 23 | 100003 | bbb | pp3 | 1,0 | I |
| 100001 | bbb | 23 | 100004 | bbb | pp4 | 1,0 | I |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 2,0 | I |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | 2,0 | I |
| 100001 | bbb | 23 | 100003 | bbb | pp3 | 2,0 | D |
| 100001 | bbb | 23 | 100004 | bbb | pp4 | 2,0 | D |

# 问题一：实时多流全量关联的问题

（1）

| A | id | name | val | mid |
|---|---|---|---|---|
| src | 100001 | aaa | 23 | 1 |
| cur | 100001 | bbb | 23 | 1 |

| A | id | name | val | mid |
|---|---|---|---|---|
| src | 100001 | bbb | 23 | 2 |
| cur | 100001 | aaa | 23 | 2 |

| B | id | name | val | mid |
|---|---|---|---|---|
| src | 100001 | aaa | pp1 | 1 |
| cur | 100001 | aaa | pp3 | 1 |

（2）

| A | | | | | B | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100002 | aaa | 55 | | | 100001 | aaa | pp1 | 1 | |
| 100001 | aaa | 23 | 2 | | 100002 | aaa | pp2 | | |

| A | | | | | B | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 100003 | bbb | 21 | | | 100003 | bbb | pp3 | | |
| 100004 | bbb | 33 | | | 100004 | bbb | pp4 | | |
| 100005 | bbb | 66 | | | 100001 | bbb | pp1 | 1 | |

| A | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 100006 | ccc | 43 | | | | | | |

| | | | | | B | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | 100005 | ddd | pp5 | |

[now, +∞)

（3）

| A full outer join B | | | | | | | |
|---|---|---|---|---|---|---|---|
| Aid | Aname | Aval | Bid | Bname | Bval | mid | opt |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | | |
| 100002 | aaa | 55 | 100001 | aaa | pp1 | | |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | | |
| 100002 | aaa | 55 | 100002 | aaa | pp2 | | |
| 100003 | bbb | 21 | 100003 | bbb | pp3 | | |
| 100004 | bbb | 33 | 100003 | bbb | pp3 | | |
| 100005 | bbb | 66 | 100003 | bbb | pp3 | | |
| 100003 | bbb | 21 | 100004 | bbb | pp4 | | |
| 100004 | bbb | 33 | 100004 | bbb | pp4 | | |
| 100005 | bbb | 66 | 100004 | bbb | pp4 | | |
| 100006 | ccc | 43 | | | | | |
| | | | 100005 | ddd | pp5 | | |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 1,0 | D |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | 1,0 | D |
| 100001 | bbb | 23 | 100003 | bbb | pp3 | 1,0 | I |
| 100001 | bbb | 23 | 100004 | bbb | pp4 | 1,0 | I |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 2,0 | I |
| 100001 | aaa | 23 | 100002 | aaa | pp2 | 2,0 | I |
| 100001 | bbb | 23 | 100003 | bbb | pp3 | 2,0 | D |
| 100001 | bbb | 23 | 100004 | bbb | pp4 | 2,0 | D |
| 100001 | aaa | 23 | 100001 | aaa | pp1 | 2,1 | D |
| 100002 | aaa | 55 | 100001 | aaa | pp1 | 0,1 | D |
| 100003 | bbb | 21 | 100001 | bbb | pp1 | 0,1 | I |
| 100004 | bbb | 33 | 100001 | bbb | pp1 | 0,1 | I |
| 100005 | bbb | 66 | 100001 | bbb | pp1 | 0,1 | I |

# 问题二：实时流全量分组计算问题

| A | | | | | |
|---|---|---|---|---|---|
| id | name | val | rn | mid | opt |
| 100001 | aaa | 23 | 1 | | |
| 100002 | aaa | 55 | 2 | | |
| 100003 | bbb | 21 | 1 | | |
| 100004 | bbb | 33 | 2 | | |
| 100005 | bbb | 66 | 3 | | |
| 100006 | ccc | 43 | 1 | | |
| 100003 | bbb | 21 | 1 | 1 | D |
| 100004 | bbb | 33 | 2 | 1 | D |
| 100005 | bbb | 66 | 3 | 1 | D |
| 100004 | bbb | 33 | 1 | 1 | I |
| 100005 | bbb | 66 | 2 | 1 | I |
| 100001 | aaa | 23 | 1 | 1 | D |
| 100002 | aaa | 55 | 2 | 1 | D |
| 100003 | aaa | 21 | 1 | 1 | I |
| 100001 | aaa | 23 | 2 | 1 | I |
| 100002 | aaa | 55 | 3 | 1 | I |

| A | | | | |
|---|---|---|---|---|
| 100001 | aaa | 23 | | |
| 100002 | aaa | 55 | | |
| 100003 | aaa | 21 | | |

| A | | | | |
|---|---|---|---|---|
| 100003 | bbb | 21 | | |
| 100004 | bbb | 33 | | |
| 100005 | bbb | 66 | | |

| A | | | | |
|---|---|---|---|---|
| 100006 | ccc | 43 | | |

| A | | | | |
|---|---|---|---|---|
| | id | name | val | mid |
| src | 100003 | bbb | 21 | 1 |
| cur | 100003 | aaa | 21 | |

DataFunSummit

# 问题二：实时流全量分组计算问题

| A | | | |
|---|---|---|---|
| 100001 | aaa | 23 | |
| 100002 | aaa | 55 | |
| 100003 | aaa | 21 | |

| A | | | | |
|---|---|---|---|---|
| | id | name | val | mid |
| src | 100003 | bbb | 21 | 1 |
| cur | 100003 | aaa | 21 | |

| A | | | |
|---|---|---|---|
| 100003 | bbb | 21 | |
| 100004 | bbb | 33 | |
| 100005 | bbb | 66 | |

| A | | | |
|---|---|---|---|
| 100006 | ccc | 43 | |

| A | | | |
|---|---|---|---|
| name | min(val) | mid | opt |
| aaa | 23 | | |
| bbb | 21 | | |
| ccc | 43 | | |
| aaa | 23 | 1 | D |
| aaa | 21 | 1 | I |
| bbb | 21 | 1 | D |
| bbb | 33 | 1 | I |

# 02
## 难点

Subject

实时维度模型计算过程中，如何获取全量相关历史数据、如何提升**处理**性能、如何降低开发难度及维护成本

# 难点

使用Flink SQL目前主要存在以下几个难点

> 难点一： 如何得到历史数据

> 难点二： 如何提升性能

> 难点三： 如何简化开发难度

这里的说是难点，其实更确切的说，是比较繁琐的地方，可能不是难度高，而是人为
需要思考或者操作的地方比较多的部分。

# 03
# 方案

使用组件化设计，使得大家可以面向应用
编程、计算逻辑使用Flink SQL表达 简化代
码开发成本
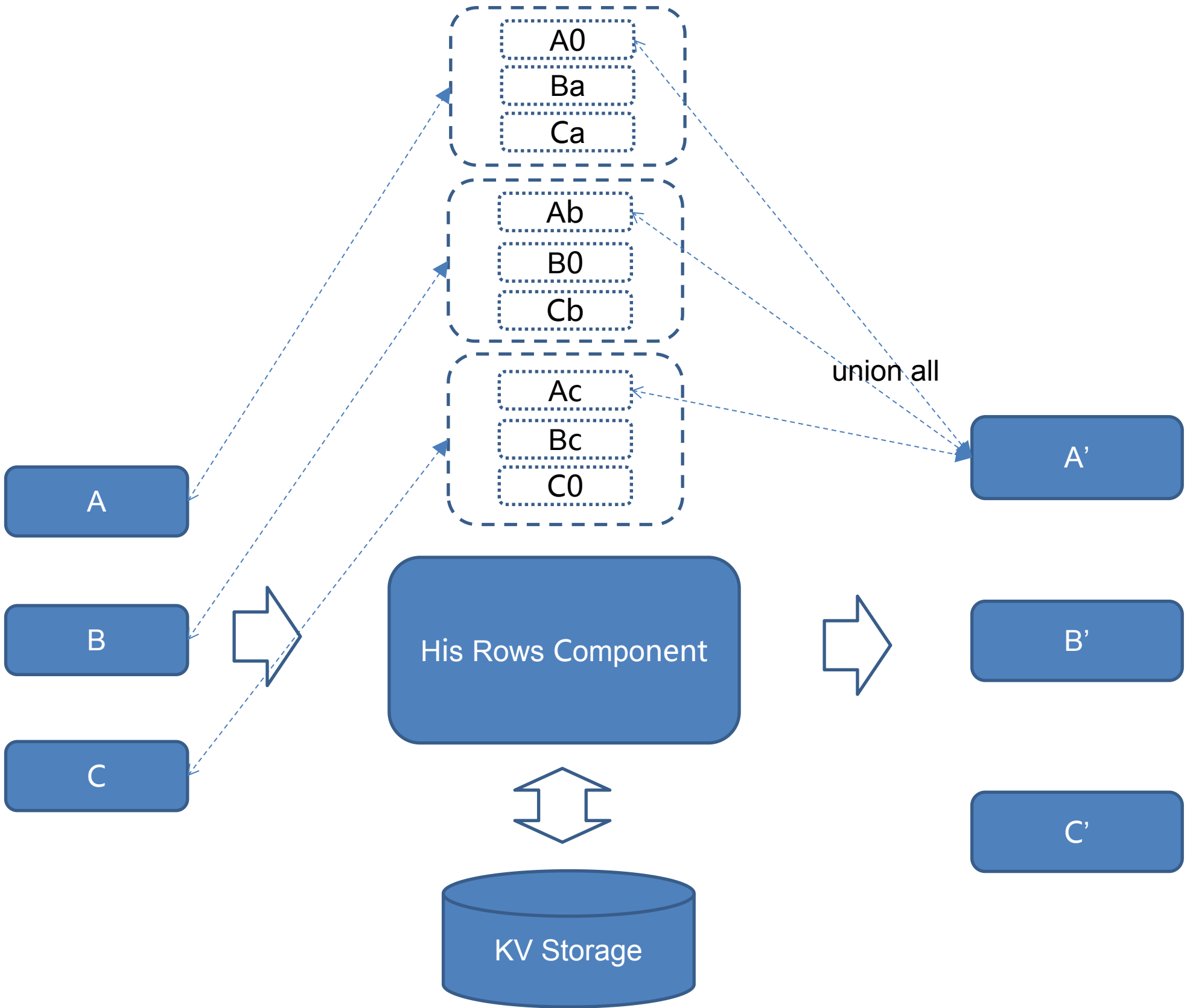
# 方案（RDM Building）

用户编写组建配置 — Components Config

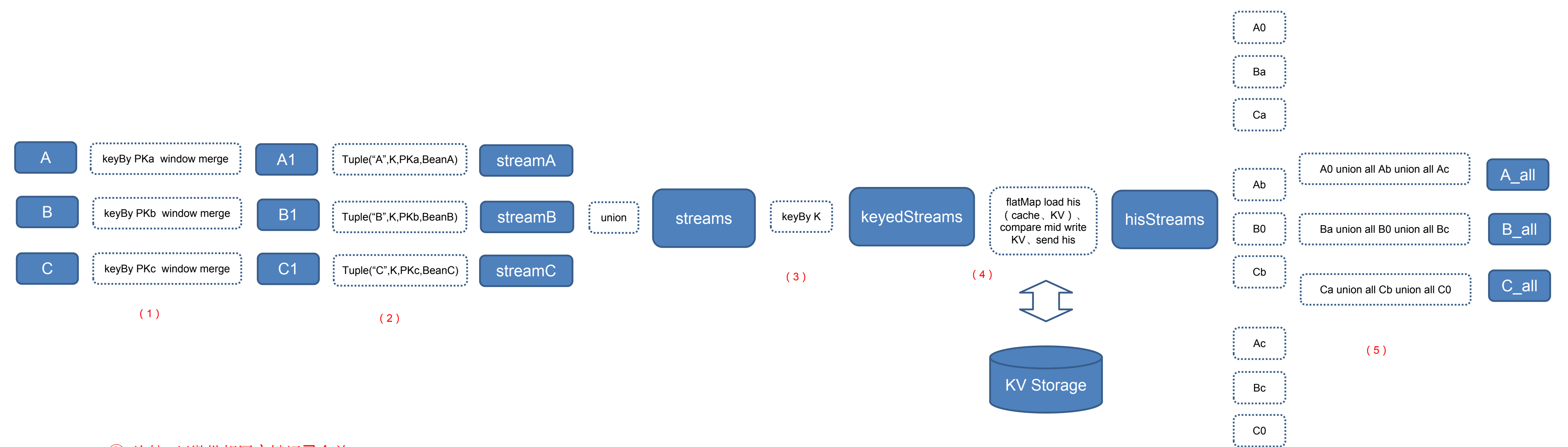构建维度建模组建 — RDDM Component Builder

转化flink执行计划 — RDDM Component Parse

提交执行flink程序 — Flink 、Flink SQL

A0
Ba
Ca

Ab
B0
Cb

Ac
Bc
C0

A
B
C

His Rows Component

KV Storage

union all

A'
B'
C'

# 方案（His Rows Component）



① 比较mid微批相同主键记录合并

② 转化相同流格式处理update retract逻辑

③ 合流统一keyBy使相同存储键的记录分发到相同slot 提升缓存利用率

④ 首先从缓存获取比较mid如果是最新记录则就写入cache及KV存储并向下发送记录

⑤ 从hisStreams流中分拆出来加载出来的数据流 合并得到包含历史数据的A_all、B_all、C_all 用于下一步SQL计算

# 04
# 规划

Subject

增加前端页面、扩展底层对多种实时计算
引擎的支持、将KV存储彻底抽象独立出来
（目前支持hbase、redis这两类KV存储）