

第十章 Flink习题解答



毕倪飞、杨溢

華東師範大學



大纲

2

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

3

- 要求：假设每个学生都选修了10学分的课程，现请根据成绩单和公式计算所有学生的平均学分绩点
- 输入：每行为一个学生选修某课程的成绩信息，包括学生id、课程学分、课程成绩
- 输出：学生id及其总学分和平均学分绩点

$$\text{平均学分绩点} = \frac{\sum (\text{课程学分} * \text{绩点})}{\sum \text{课程学分}}$$

输入	输出
t001 4 3.5	(t001,10,3.4) (t002,10,3.7) (t003,10,3.6)
t003 2 4.0	
t003 4 3.5	
<u>t001 2 3.0</u>	
t002 3 3.0	
t002 3 4.0	
<u>t001 4 3.5</u>	
<u>t002 4 4.0</u>	
t003 3 4.0	
t003 1 2.0	



解决方案

4

□ map

- ✚ 将每行映射为(学生id, 课程学分, 课程成绩)元组

□ keyBy

- ✚ 以**学生id**为key对元组分组, 使每个学生各自的成績信息被划分到同一个分组内

□ reduce

- ✚ 根据平均学分绩点公式, 计算得到(学生id, **已修学分, 平均学分绩点**) 元组

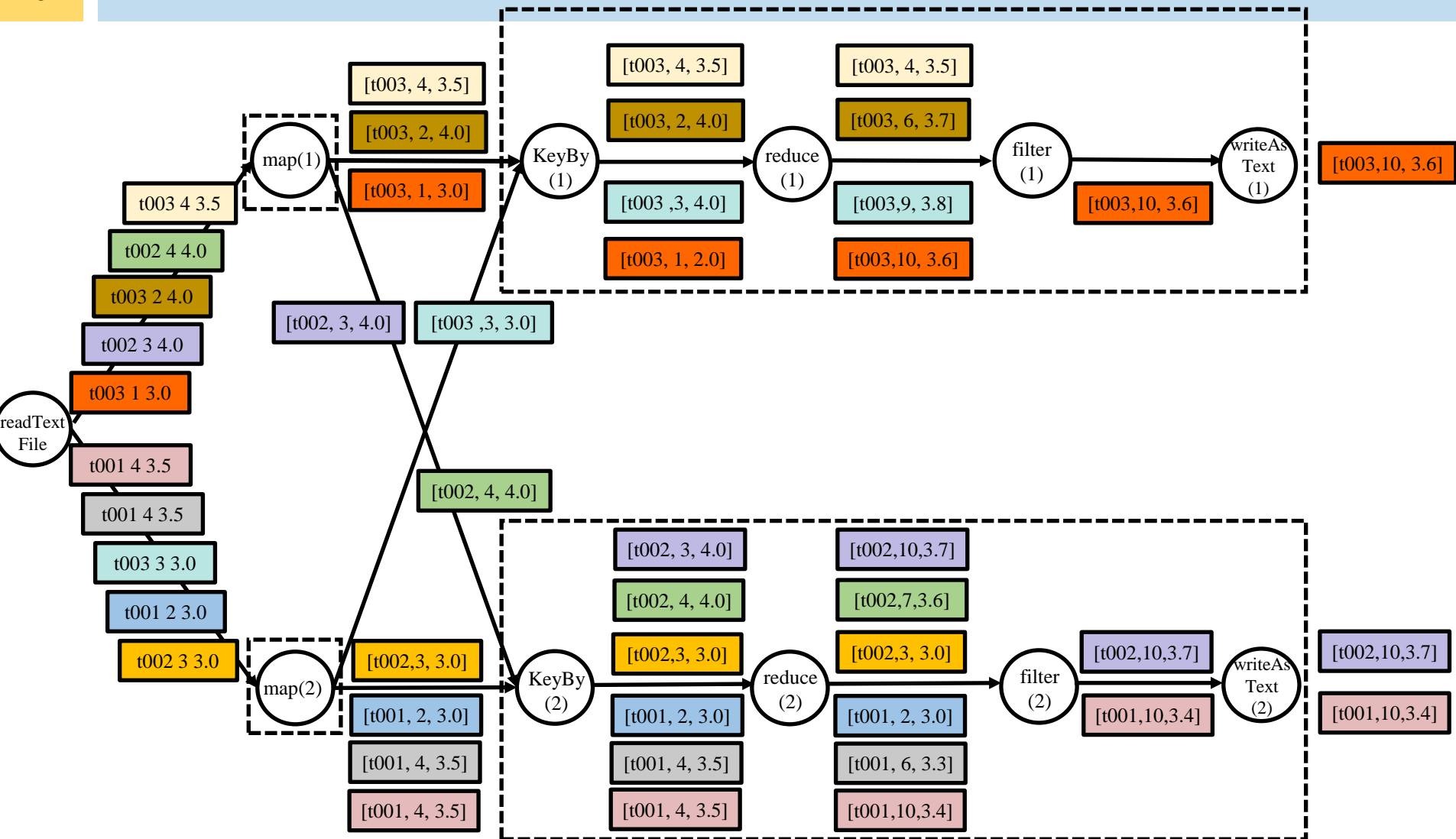
□ filter

- ✚ 筛选出已修学分为10分的元组



计算平均学分绩点的运行过程

5



编写calculate方法

6

```
1 @Override
2 public DataStream<Tuple3<String, Integer, Float>> calculate(DataStream<String> text) {
3     // 将成绩信息映射为映射为[学生id, 课程学分, 课程成绩]
4     DataStream<Tuple3<String, Integer, Float>> points = text.map(
5         new MapFunction<String, Tuple3<String, Integer, Float>>() {
6             @Override
7             public Tuple3<String, Integer, Float> map(String value) throws Exception {
8                 String infos[] = value.split(" ");
9                 float point = Float.parseFloat(infos[2]);
10                if (point >= 0 && point <= 5) {
11                    return new Tuple3<>(
12                        infos[0], Integer.parseInt(infos[1]), point);
13                }
14                // 成绩若为-1,表示该成绩异常,此时该科成绩以0分计。
15                return new Tuple3<>(
16                    infos[0], Integer.parseInt(infos[1]), point + 1);
17            }
18        });
19    // 将学生成绩信息按学生id进行分组,分组后计算每个学生总学分及平均学分绩点
20    DataStream<Tuple3<String, Integer, Float>> pointsGroupByKey = points.keyBy(0).reduce(
21        new ReduceFunction<Tuple3<String, Integer, Float>>() {
22            @Override
23            public Tuple3<String, Integer, Float> reduce(Tuple3<String, Integer, Float> preRecord,
24                Tuple3<String, Integer, Float> currRecord) throws Exception {
25                return Tuple3.of(preRecord.f0, preRecord.f1 + currRecord.f1,
26                    (float) (Math.round(
27                        ((preRecord.f1 * preRecord.f2 + currRecord.f1 * currRecord.f2) / (preRecord.f1
28                            + currRecord.f1)) * 10)) / 10);
29            }
30        });
31    // 筛选总学分为10的成绩记录
32    DataStream<Tuple3<String, Integer, Float>> result = pointsGroupByKey
33        .filter((value) -> value.f1 == 10);
34    // 返回计算结果
35    return result;
36 }
```



大纲

7

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

8

- 要求：10km路程的路段内，平均速度 $\leq 60\text{km/h}$ ，途中两测速点瞬时速度 $\leq 60\text{km/h}$
- 输入：每行文本包含两个字符串，第一个代表车牌号，第二个代表时刻或速度
- 输出：超速的车牌号

输入	输出
00001 10:00:00	<u>00012</u>
00001 58	
00001 59	
00001 10:12:00	
00012 10:15:00	
00012 55	
<u>00012 66</u>	
00012 10:27:00	

解决方案

9

□ map

- ✚ 将每一行映射为[车牌号,时刻/速度]键值对

□ keyBy

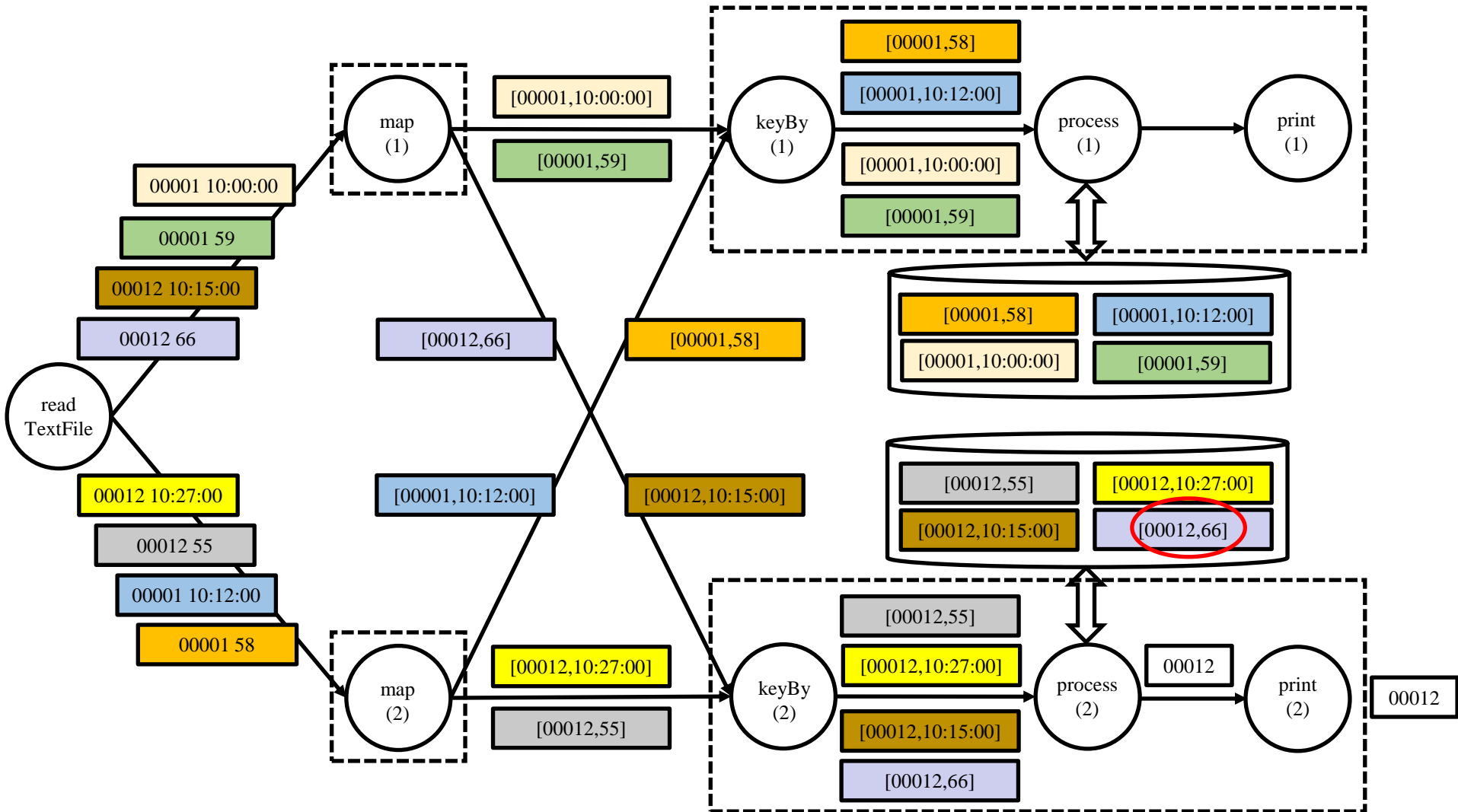
- ✚ 将[车牌号,时刻/速度]键值对按车牌号进行分组

□ process

- ✚ 将收到的[车牌号,时刻/速度]键值对中的**时刻/速度值存入状态**，并判断状态中是否存满4个值
- ✚ 若未存满4个值，则不做任何处理
- ✚ 若已存满4个值，则计算平均速度，并根据平均速度和瞬时速度判断是否超速

区间测速的运行过程

10



编写avgSpeed方法

11

```
1 @Override
2 public float averageSpeed(String[] times) {
3     if (times[0].compareTo(times[1]) > 0) {
4         float hours = ((timeToSeconds(times[0]) - timeToSeconds(times[1])) * 1.0f / 3600);
5         return 10 / hours;
6     } else {
7         float hours = (timeToSeconds(times[1]) - timeToSeconds(times[0])) * 1.0f / 3600;
8         return 10 / hours;
9     }
10 }
11
12 /**
13  * 将一个时刻转换为秒
14  *
15  * @param time 代表某一个时刻
16  * @return 时刻中含有多少秒
17  */
18 public Integer timeToSeconds(String time) {
19     String[] splits = time.split(":");
20     return Integer.parseInt(splits[0]) * 3600
21         + Integer.parseInt(splits[1]) * 60
22         + Integer.parseInt(splits[2]);
23 }
```



自定义HashMap来维护状态

12

```
1 // 通过一个Map来维护每个Key所对应的状态
2 Map<String, List<String>> customState;
3
4 @Override
5 public void processElement(Tuple2<String, String> t, Context context, Collector<String> collector)
6     throws Exception {
7     if (customState == null) {
8         customState = new HashMap<>();
9     }
10
11 // 取出当前Key对应的状态
12 List<String> curState = customState.get(t.f0);
13 if (curState == null) {
14     curState = new ArrayList<>();
15 }
16 // 若状态中已包含3个值,则表示加上当前收到的数据后,就已收集到当前车牌号所对应的所有数据
17 if (curState.size() == 3) {
18     // 加上当前收到的数据
19     curState.add(t.f1);
20     String[] times = new String[2];
21     int[] speeds = new int[2];
22     int cntTimes = 0, cntSpeeds = 0;
23     // 分离出入时刻数据与速度数据
24     for (String s : curState) {
25         if (s.contains(":")) {
26             times[cntTimes++] = s;
27         } else {
28             speeds[cntSpeeds++] = Integer.parseInt(s);
29         }
30     }
31 // 求区间平均速度
32 float avgSpeed = averageSpeed(times);
33 if (speeds[0] > 60 || speeds[1] > 60 || avgSpeed > 60) {
34     collector.collect(t.f0);
35 }
36 } else { // 否则继续将数据收集进状态
37     curState.add(t.f1);
38     customState.put(t.f0, curState);
39 }
40 }
41
```



通过ValueState来维度状态

13

□ 编写open方法

```
1 private transient ValueState<List<String>> state;
2
3 @Override
4 public void open(Configuration parameters) throws Exception {
5     ValueStateDescriptor<List<String>> stateDescriptor = new ValueStateDescriptor<>("s",
6         Types.LIST(Types.STRING));
7     state = getRuntimeContext().getState(stateDescriptor);
8 }
```

通过ValueState来维度状态

14

□ 编写processElement方法

```
1 @Override
2 public void processElement(Tuple2<String, String> t, Context context, Collector<String> collector)
3     throws Exception {
4     if (state.value() == null) {
5         state.update(new ArrayList<>());
6     }
7
8     // 取出状态
9     List<String> curState = state.value();
10    // 若状态中已包含3个值, 则表示加上当前收到的数据后, 就已收集到当前车牌号所对应的所有数据
11    if (curState.size() == 3) {
12        // 加上当前收到的数据
13        curState.add(t.f1);
14        String[] times = new String[2];
15        int[] speeds = new int[2];
16        int cntTimes = 0, cntSpeeds = 0;
17        // 分离出入时刻数据与速度数据
18        for (String s : curState) {
19            if (s.contains(":")) {
20                times[cntTimes++] = s;
21            } else {
22                speeds[cntSpeeds++] = Integer.parseInt(s);
23            }
24        }
25        // 求区间平均速度
26        float avgSpeed = averageSpeed(times);
27        if (speeds[0] > 60 || speeds[1] > 60 || avgSpeed > 60) {
28            collector.collect(t.f0);
29        }
30    } else { // 否则继续将数据收集进状态
31        curState.add(t.f1);
32        state.update(curState);
33    }
34 }
```

大纲

15

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

16

- 输入：包含操作事件和查询事件的流
 - ✚ 格式: (是否查询事件, 泳池编号, 时间戳, 速度)
- 输出：每个查询事件对应的泳池的水量
 - ✚ 格式:(泳池编号, 时间戳, 水量)

输入	输出
<u>false,1,0,10</u>	<u>1,5,50</u>
false,2,0,20	2,5,100
<u>true,1,5,0</u>	1,20,1100
true,2,5,0	2,20,200
false,1,10,100	
false,2,10,0	
true,1,20,0	
true,2,20,0	

解决方案

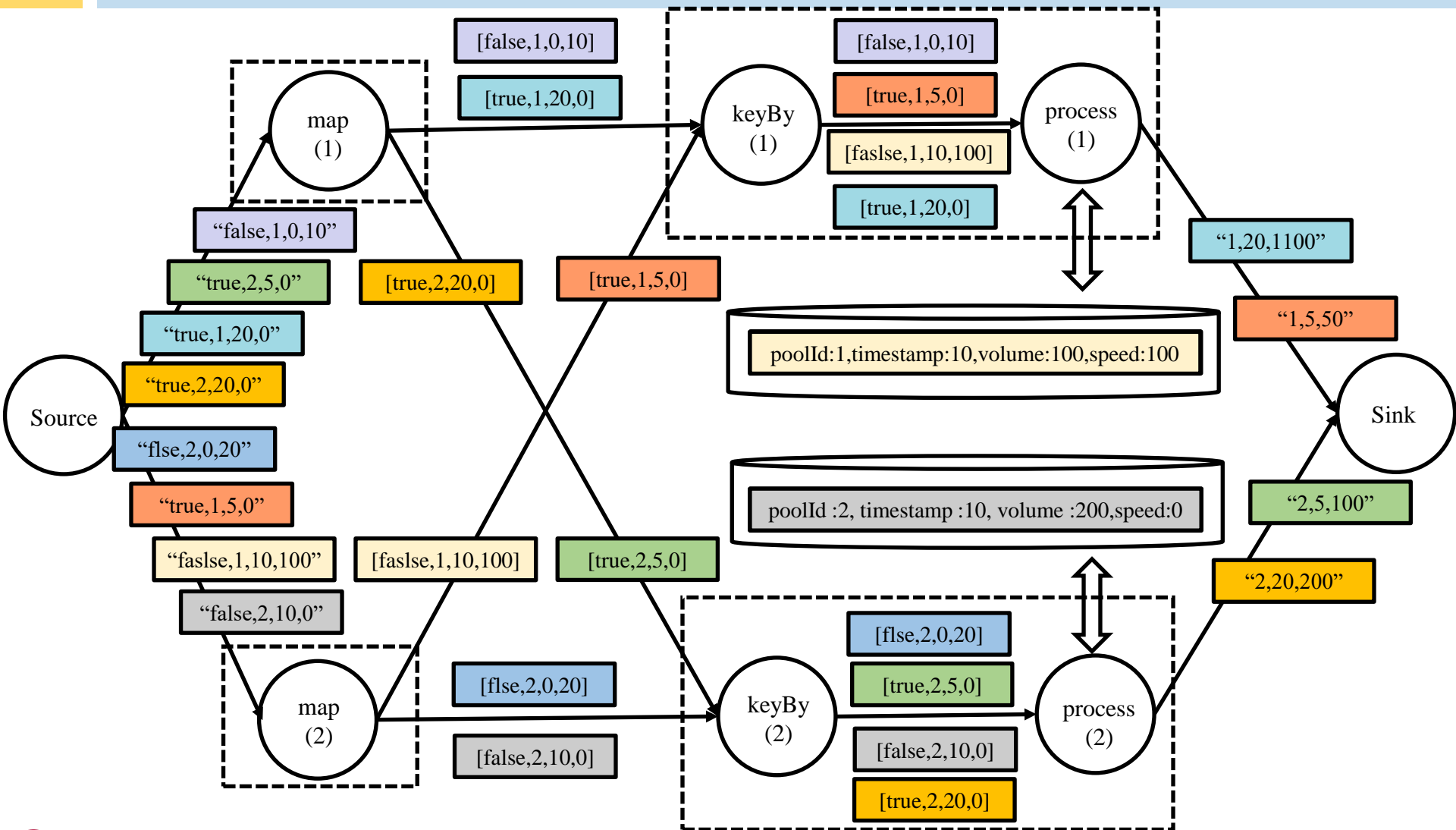
17

- 解析字符串, 得到(是否查询事件, 泳池编号, 时间戳, 速度)元组 `map`
- 按泳池编号进行分组 `keyBy`
- 处理事件 `process`
 - ✚ 计算当前水量 (注意溢出和负数)
$$\text{水量} = \text{上次水量} + \text{上次速度} * (\text{时间戳} - \text{上次时间戳})$$
 - ✚ 若为修改事件, 则用时间戳、**水量**和速度来更新时间戳状态、水量状态和速度状态
 - ✚ 若为查询事件, 则输出结果 (泳池编号, 时间戳, 水量)



WaterProblem的运行过程

18



编写execute方法

19

```
1  @Override
2  public DataStream<String> execute(DataStream<String> dataStream) {
3      return dataStream
4          // 把数据格式从String转换到结构体
5          .map(line -> {
6              String[] items = line.split(",");
7              boolean isquery = "true".equals(items[0]);
8              int poolId = Integer.parseInt(items[1]);
9              int time = Integer.parseInt(items[2]);
10             int speed = Integer.parseInt(items[3]);
11             return new InputEvent(isquery, poolId, time, speed);
12         })
13         // 按泳池编号进行分区
14         .keyBy(x -> x.poolId)
15         // 处理输入事件、维护状态、输出结果
16         .process(new ProcessEventFunction());
17     }
18
19     public static class InputEvent {
20         boolean isQuery; // 查询事件为true, 修改事件为false
21         int poolId; // 泳池编号
22         int time; // 事件戳
23         int speed; // 速度
24         public InputEvent(boolean isQuery, int poolId, int time, int speed) {
25             this.poolId = poolId;
26             this.time = time;
27             this.speed = speed;
28             this.isQuery = isQuery;
29         }
30     }
```



编写ProcessEventFunction类

20

```
1 private static class ProcessEventFunction extends
2     KeyedProcessFunction<Integer, InputEvent, String> implements Serializable {
3     // 记录上次修改事件发生时的水量
4     private transient ValueState<Long> volume;
5     // 记录上次修改事件的速度
6     private transient ValueState<Integer> speedState;
7     // 记录上次修改事件的时间戳
8     private transient ValueState<Integer> timeState;
9
10    public ProcessEventFunction() {}
11
12    @Override
13    public void open(Configuration parameters) throws Exception {
14        // 在open函数中注册状态
15        ValueStateDescriptor<Long> volumeDescriptor = new ValueStateDescriptor<>("v", Types.LONG);
16        volume = getRuntimeContext().getState(volumeDescriptor);
17        ValueStateDescriptor<Integer> inSpeedDescriptor = new ValueStateDescriptor<>("is", Types.INT);
18        speedState = getRuntimeContext().getState(inSpeedDescriptor);
19        ValueStateDescriptor<Integer> timeDescriptor = new ValueStateDescriptor<>("t", Types.INT);
20        timeState = getRuntimeContext().getState(timeDescriptor);
21    }
22
23    void initState() throws IOException {
24        if (volume.value() == null) volume.update(0L);
25        if (speedState.value() == null) speedState.update(0);
26        if (timeState.value() == null) timeState.update(0);
27    }
28
29    @Override
30    public void processElement(InputEvent inputEvent, Context context,
31        Collector<String> collector) throws Exception {
32        .....
33    }
34 }
```



编写processElement方法

21

```
1  @Override
2  public void processElement(InputEvent inputEvent, Context context,
3      Collector<String> collector) throws Exception {
4      // 对状态进行初始化
5      initState();
6
7      // 计算当前水量，因为水量可能会超出int的取值范围，所以要使用long
8      long currentVolume = volume.value() +
9          (long) speedState.value() * (inputEvent.time - timeState.value());
10
11     // 水量不能是负的
12     currentVolume = Math.max(0, currentVolume);
13
14     if (inputEvent.isQuery) {
15         // 如果是查询事件，就输出查询结果
16         collector.collect(inputEvent.poolId + "," + inputEvent.time + "," + currentVolume);
17     } else {
18         // 如果是修改事件，就更新各个状态
19         volume.update(currentVolume);
20         speedState.update(inputEvent.speed);
21         timeState.update(inputEvent.time);
22     }
23 }
```

注意溢出和负数

大纲

22

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

23

- 要求：根据别名分配规则以及顺序的别名申请或注销请求序列，**计算各请求的状态**（成功或失败）
- ✚ 别名长度的取值范围是 $[5, 11]$ ，且别名只能由英文字母、数字或下划线（ $_$ ）组成
- ✚ 同一学院内的别名不能发生重复
- ✚ 不同学院间的别名可以发生重复

题目描述

24

- 输入：由请求类型type、用户工号id、用户所属学院depart以及请求的别名组成
- 输出：每个请求的完成状态，“SUCCESS”表示请求成功执行，“FAILURE”表示请求执行失败

输入	Request{type=APPLY, id=1, depart=Department{1,1}, alias="*invalid"} Request{type=APPLY, id=2, depart=Department{1,3}, alias="a_name"} Request{type=APPLY, id=3, depart=Department{1,3}, alias="b_name"} Request{type=APPLY, id=3, depart=Department{1,3}, alias="a_name"} Request{type=REVOKE, id=3, depart=Department{1,3}, alias="a_name"} Request{type=REVOKE, id=2, depart=Department{1,3}, alias="a_name"}
输出	FAILURE SUCCESS SUCCESS FAILURE FAILURE SUCCESS



解决方案

25

□ keyBy

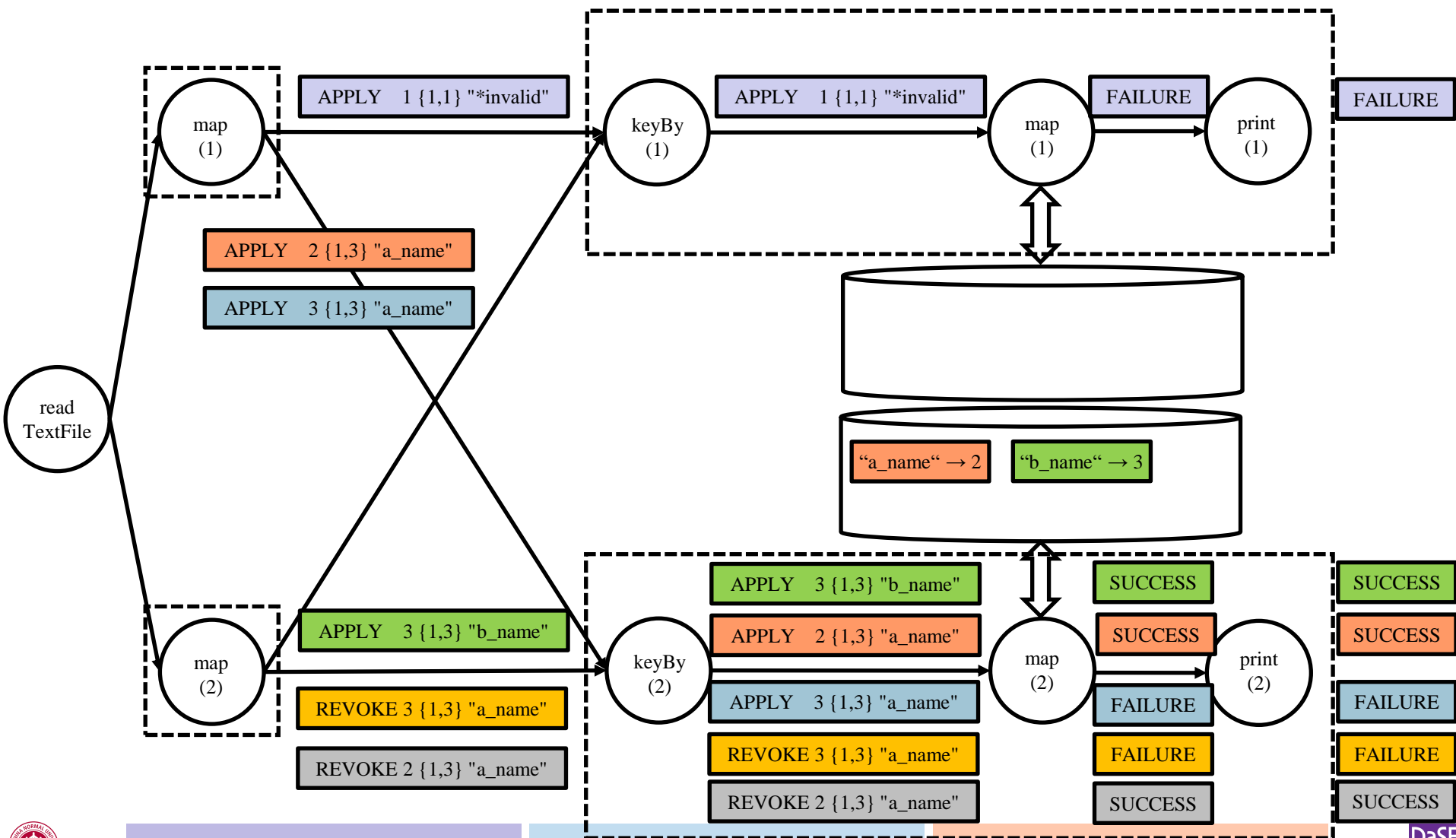
- ✚ 按学院对请求进行分组，需要自定义哈希函数

□ map: 使用MapState计算请求完成状态

- ✚ 若别名格式不符合要求，直接判定失败
- ✚ 若是申请别名请求：
 - 检测状态中是否存在别名，若存在则判定失败，否则判定成功，并向MapState中存入[别名,工号]
- ✚ 若是注销别名请求：
 - 检测状态中是否存在与当前请求对应的[别名,工号]，若是则判定成功并从状态中删除当前别名，否则判定失败

运行过程

26



按部门对请求进行分组

27

```
public class EmailAssignmentImpl extends EmailAssignment {

    @Override
    public DataStream<String> processRequest(DataStream<Request> requests) {
        return requests
            .keyBy(r -> hash(r.getDepart()))
            // ...
    }

    /**
     * 生成部门的哈希值
     * @param depart 部门
     * @return 对应哈希值
     */
    private int hash(Department depart) {
        return 31 * depart.getFirstLevelCode() + depart.getSecondLevelCode();
    }
}
```



获取MapState

28

```
public class EmailAssignmentImpl extends EmailAssignment {

    @Override
    public DataStream<String> processRequest(DataStream<Request> requests) {
        return requests
            .keyBy(r -> hash(r.getDepart()))
            .map(new RichMapFunction<Request, String>() {
                // 别名 -> 工号
                private MapState<String, Integer> aliasMap;
                private final Pattern aliasPattern = Pattern.compile("^[_a-zA-Z0-9]{5,11}$");

                @Override
                public void open(Configuration parameters) throws Exception {
                    super.open(parameters);
                    aliasMap = getRuntimeContext().getMapState(
                        new MapStateDescriptor<>("aliasMap", String.class, Integer.class));
                }

                @Override
                public String map(Request request) throws Exception {
                    // ...
                }
            });
    }
    // ...
}
```



计算请求完成状态

29

```
@Override
public String map(Request request) throws Exception {
    String alias = request.getAlias();
    if (!aliasPattern.matcher(alias).matches()) {
        return "FAILURE";
    }

    if (request.getType() == RequestType.APPLY) {
        if (aliasMap.containsKey(alias)) {
            return "FAILURE";
        }
        aliasMap.put(alias, request.getId());
        return "SUCCESS";
    } else { // request.getType() == RequestType.REVOKE
        if (aliasMap.containsKey(alias) && aliasMap.get(alias) == request.getId()) {
            aliasMap.remove(alias);
            return "SUCCESS";
        }
        return "FAILURE";
    }
}
```



大纲

30

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

31

- 要求：某程序会持续地产生一些位于 $[0, 9]$ 区间的数。这些数均匀地分布在 $[0, 4]$ 和 $[5, 9]$ 两个区间内。要求将这些数均匀地划分到两个任务中，并将其中的整数转换为英文单词。
- 输入：每行为一个数
- 输出：整数转换后的结果

输入	输出	
	任务一输出	任务二输出
0	One Three Zero Two Four	Five
1		Seven
2		Nine
3		Six
4	One Three Zero Two Four	Eight
5		
6		
7		
8		
9		



解决方案

32

□ 均匀划分数

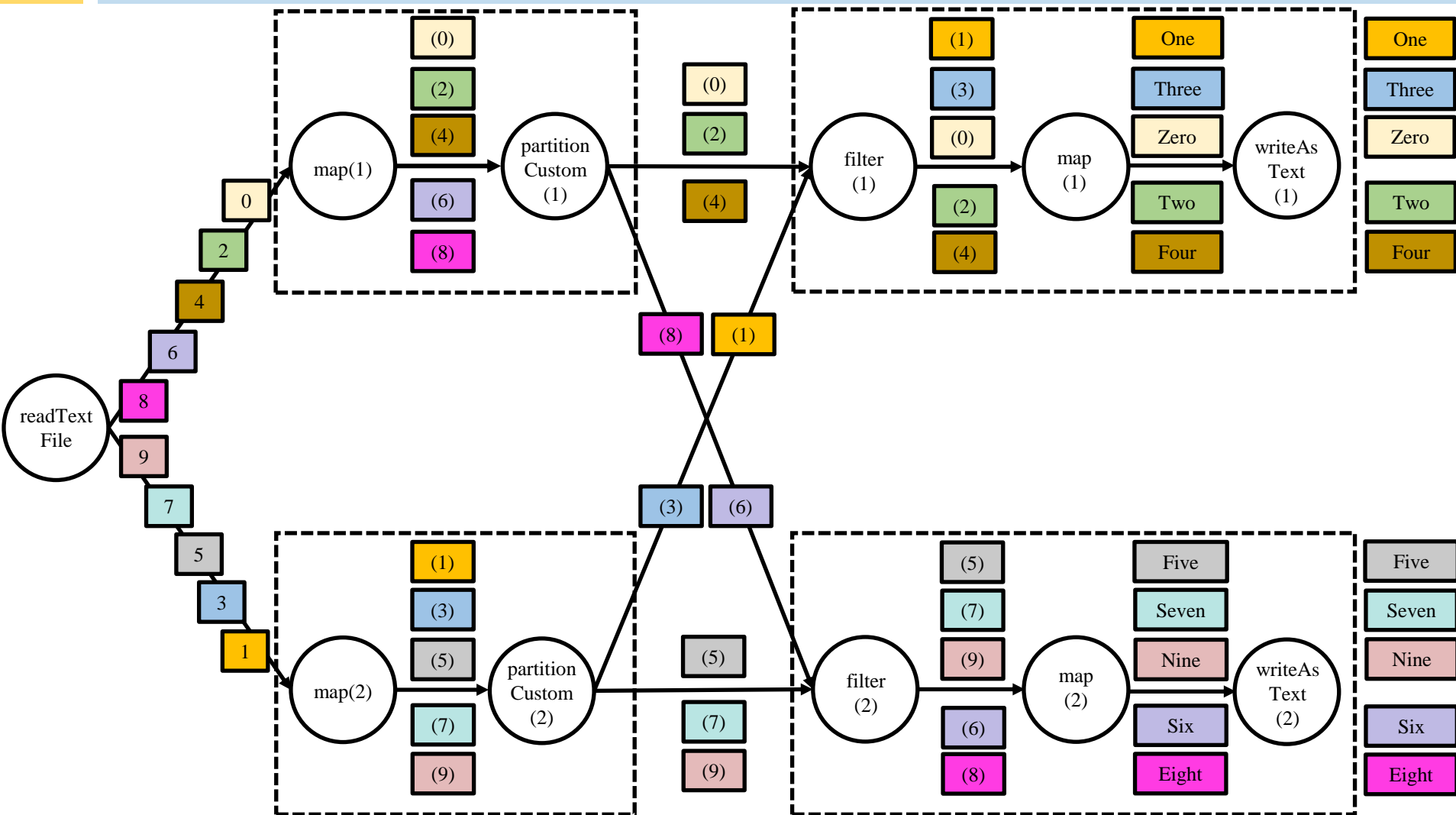
- ✚ 通过自定义 **Partitioner** 划分数

□ 将整数转换为英文单词

- ✚ 通过 **filter** 过滤非整数
- ✚ 通过 **map** 将整数转换为英文单词

整数转换的运行过程

33



编写partition方法

34

```
public class DigitalPartitionerImpl extends DigitalPartitioner {  
  
    @Override  
    public int partition(String key, int numPartitions) {  
        double parsedKey = Double.parseDouble(key);  
        // 条件表达式, [0, 4]的数字划分给任务一, 而[5, 9]的数字划分给任务二  
        return parsedKey >= 0 && parsedKey <= 4 ? 0 : 1;  
    }  
}
```

编写digitalConversion方法

35

```
public class DigitalConversionImpl extends DigitalConversion {

    @Override
    public DataStream<String> digitalConversion(DataStream<Tuple1<String>> digitals) {
        // 过滤掉小数
        DataStream<Tuple1<String>> integer = digitals.filter((FilterFunction<Tuple1<String>>) value
            -> value.getField(0).toString().matches("[0-9]+")).returns(Types.TUPLE(Types.STRING));

        // 将整数转换为对应的单词
        DataStream<String> word = integer.map((MapFunction<Tuple1<String>, String>) value -> {
            // 解析出整数
            int index = Integer.parseInt(value.getField(0));
            // 通过枚举类获取整数对应的英文单词
            return DigitalWord.values()[index].getWord();
        }).returns(Types.STRING);

        return word;
    }
}
```

大纲

36

- 计算平均学分绩点
- 区间测速
- 泳池水量
- 邮箱别名分配
- 整数转换
- 最大公约数



题目描述

37

- 要求：给定两个正整数，利用辗转相除法以迭代方式求两者的最大公约数（GCD）

$$gcd(a, b) = \begin{cases} gcd(b, a \% b), & (a \% b) \neq 0 \\ b, & (a \% b) = 0 \end{cases}$$

- 输入：字母标识符以及两个正整数a, b组成的计算请求
- 输出：以格式(字母标识符,除数,余数)输出每个计算请求每轮迭代的除数和余数

输入	输出
A 5 2 B 18 6	(A,2,1) (B,6,0) (A,1,0)

辗转相除法

38

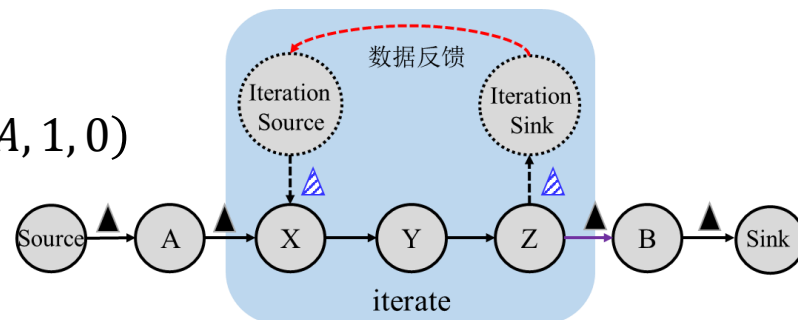
□ 算法执行过程

计算请求(A, 5, 2)的迭代过程

迭代1 $\gcd(a, b) \rightarrow \gcd(b, a \% b)$
 $\gcd(5, 2) \rightarrow \gcd(2, 1) \rightarrow (A, 2, 1)$

迭代2 $\gcd(a, b) \rightarrow \gcd(b, a \% b)$
 $\gcd(2, 1) \rightarrow \gcd(1, 0) \rightarrow (A, 1, 0)$

$a \% b = 0$ ，不需要再向迭代前端反馈



▲ 输入/输出数据记录

▲ 参与迭代反馈的数据记录

将(A, b, a%b)作为输出向后传递

若 $a \% b \neq 0$ ，则将(A, b, a%b)作为下一轮迭代的输入

解决方案

39

□ 一轮迭代的处理逻辑

✚ 根据(标识符, a , b)计算出(标识符, b , $a \% b$)

➤ map

✚ 将(标识符, b , $a \% b$)作为输入反馈给迭代前端

➤ filter

✚ 将(标识符, b , $a \% b$)作为输出向后传递

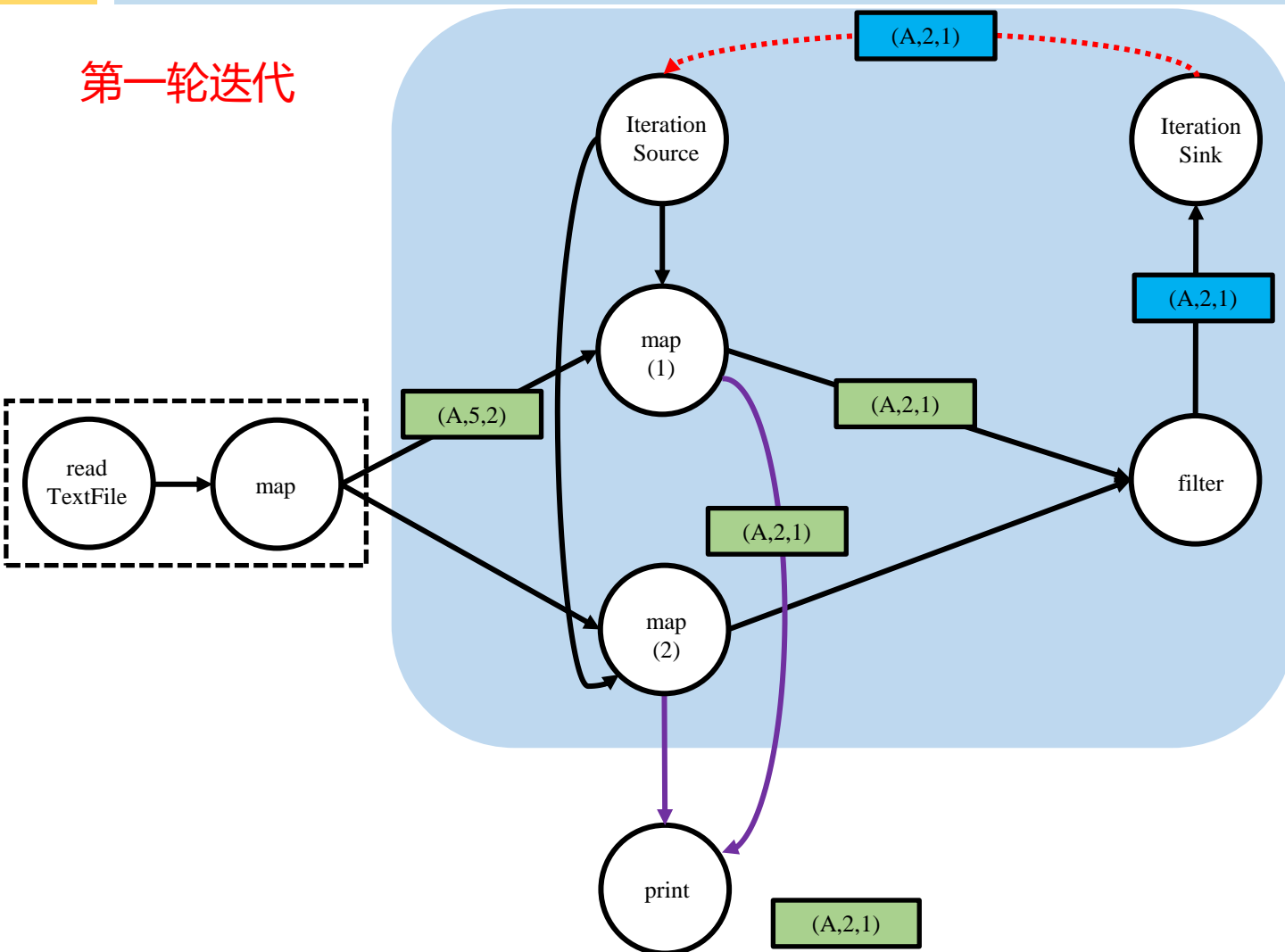
➤ 无需操作

是否需要保证 $a > b$?

最大公约数的运行过程

40

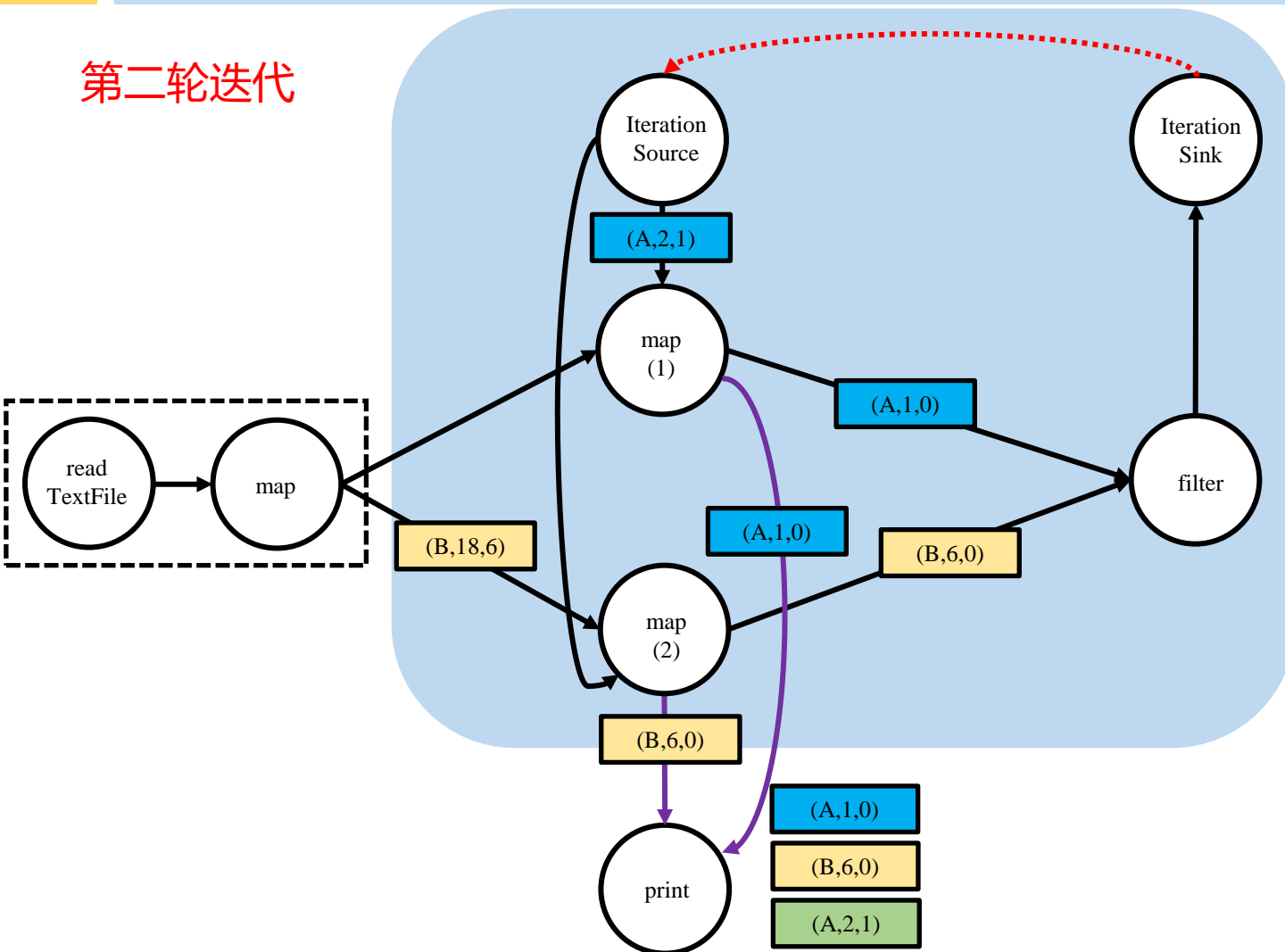
第一轮迭代



最大公约数的运行过程

41

第二轮迭代



编写calGCD方法

42

```
public class GCDImpl extends GCD {
    @Override
    public DataStream<Tuple3<String, Integer, Integer>> calGCD(IterativeStream<Tuple3<String, Integer, Integer>>
iteration) {
    // 实现每一轮迭代计算的逻辑
    DataStream<Tuple3<String, Integer, Integer>> stepStream =
        iteration.map(
            new MapFunction<Tuple3<String, Integer, Integer>, Tuple3<String, Integer, Integer>>() {
                @Override
                public Tuple3<String, Integer, Integer> map(Tuple3<String, Integer, Integer> tuple3)
                    throws Exception {
                    int a = tuple3.f1;
                    int b = tuple3.f2;
                    if (a < b) {
                        int temp = a;
                        a = b;
                        b = temp;
                    }
                    return new Tuple3<>(tuple3.f0, b, a % b);
                }
            }
        );
    // 创建反馈流
    DataStream<Tuple3<String, Integer, Integer>> feedbackStream =
        stepStream.filter(tuple3 -> tuple3.f2 != 0);
    iteration.closeWith(feedbackStream);
    // 创建输出流
    DataStream<Tuple3<String, Integer, Integer>> outputStream = stepStream;
    return outputStream;
}
```

谢谢! Q&A



Apache Flink