

# Flink超神文档

---

## Flink超神文档

### Flink初次见面

- 什么是Flink?
- 什么是Unbounded streams?
- 什么是Bounded streams?
- 什么是stateful computations?
- Flink使用用户
- Flink的特点和优势

### Flink安装&部署

- Flink基本架构
- Standalone集群安装&测试
  - 集群角色划分
  - 安装步骤
  - 提交Job到standalone集群
- Standalone HA集群安装&测试
  - 集群角色划分
  - 安装步骤

### Flink on Yarn

- 运行流程
- Flink on Yarn两种运行模式
- 配置两种运行模式
  - yarn seesion模式配置
  - Run a Flink job on YARN模式配置

### Flink on YARN HA集群安装&测试

- 安装步骤
- HA集群测试
  - yarn-session模式测试
  - Run a Flink job on YARN模式测试

### Flink API详解&实操

- Flink API介绍
- Dataflows数据流图
- 配置开发环境
- WordCount流批计算程序
- WordCount Dataflows 算子链
- Flink任务调度规则
- Flink并行度设置方式
- Dataflows DataSource数据源
  - File Source
  - Collection Source
  - Socket Source
  - Kafka Source
  - Custom Source
- Dataflows Transformations
  - Map
  - FlatMap
  - Filter
  - KeyBy
  - Reduce
  - Aggregations
  - Union 真合并
  - Connect 假合并
  - CoMap, CoFlatMap

- Split
- Select
- side output侧输出流
- Iterate (比较重要)
- 函数类和富函数类
- 底层API(ProcessFunctionAPI)
- 总结
- Dataflows分区策略
  - shuffle
  - rebalance
  - rescale
  - broadcast
  - global
  - forward
  - keyBy
  - PartitionCustom
- Dataflows Sink
  - Redis Sink
  - Kafka Sink
  - MySQL Sink (幂等性)
  - Socket Sink
  - File Sink
  - HBase Sink
- Flink State状态
- CheckPoint
  - CheckPoint原理
  - SavePoint原理
- StateBackend状态后端
  - MemoryStateBackend
  - FsStateBackend
  - RocksDBStateBackend
  - 集群级配置StateBackend
- Flink Window操作
  - Window窗口分类
  - 窗口聚合函数
    - 增量聚合函数
    - 全量聚合函数
- Flink Time时间语义
- Flink Time Watermark(水印)
  - AllowedLateness
- Flink关联维表实战
- TableAPI和Flink SQL
  - 开发环境构建
  - TableEnvironment
  - Table API
    - 1)创建Table
    - 2) **修改Table中字段名**
    - 3) **查询和过滤**
    - 4) **分组聚合**
    - 5) **UDF自定义的函数**
    - 6) **Window**
- Flink的复杂事件处理CEP
  - 1. **CEP相关概念**
    - 1) **配置依赖**
    - 2) **事件定义**
  - 2. **Pattern API**
    - 1) **模式定义**
    - 2) **模式检测**

### 3) 选择结果

#### Flink内存优化

##### 1) JobManager配置

##### 2) TaskManager配置

#### 3. Flink的网络缓存优化

##### 1) 设定Network Buffer内存数量\*\* (过时了) \*\*

##### 2) 设定Network内存比例\*\* (推荐) \*\*

#### 常见面试问题汇总

1. 面试题一：应用架构
2. 面试题二：压测和监控
3. 面试题三：为什么用Flink
4. 面试题四：checkpoint的存储
5. 面试题五：exactly-once的保证
6. 面试题六：状态机制
7. 面试题七：海量key去重
8. 面试题八：checkpoint与spark比较
9. 面试题九：watermark机制
10. 面试题十：exactly-once如何实现
11. 面试题十一：CEP
12. 面试题十二：三种时间语义
13. 面试题十三：数据高峰的处理

Flink源码(GitHub):

- [git@github.com:bjmashibing/Flink-Study.git](https://github.com/bjmashibing/Flink-Study.git)
- <https://github.com/bjmashibing/Flink-Study>

## Flink初次见面

---

### 什么是Flink?

Apache Flink is a framework and distributed processing engine for **stateful computations** over **unbounded** and **bounded** data streams. Flink has been designed to run in all common cluster environments, perform computations at in-memory speed and at any scale

Flink的世界观是数据流

对 Flink 而言，其所要处理的主要场景就是流数据，批数据只是流数据的一个极限特例而已，所以 Flink 也是一款真正的流批统一的计算引擎

---

### 什么是Unbounded streams?

无界流 有定义流的开始，但没有定义流的结束。它们会无休止地产生数据。无界流的数据必须持续处理，即数据被摄取后需要立刻处理。我们不能等到所有数据都到达再处理，因为输入是无限的，在任何时候输入都不会完成。处理无界数据通常要求以特定顺序摄取事件，例如事件发生的顺序，以便能够推断结果的完整性

### 什么是Bounded streams?

有界流 有定义流的开始，也有定义流的结束。有界流可以在摄取所有数据后再进行计算。有界流所有数据可以被排序，所以并不需要有序摄取。有界流处理通常被称为批处理

一图秒懂：无界流与有界流

## 什么是stateful computations?

**有状态的计算：**每次进行数据计算的时候基于之前数据的计算结果（状态）做计算，并且每次计算结果都会保存到存储介质中，计算关联上下文context

基于有状态的计算不需要将历史数据重新计算，提高了计算效率

**无状态的计算：**每次进行数据计算只是考虑当前数据，不会使用之前数据的计算结果

## Flink使用用户

自 2019 年 1 月起，阿里巴巴逐步将内部维护的 Blink 回馈给 Flink 开源社区，目前贡献代码数量已超过 100 万行。国内包括腾讯、百度、字节跳动等公司，国外包括 Uber、Netflix 等公司都是 Flink 的使用者

## Flink的特点和优势

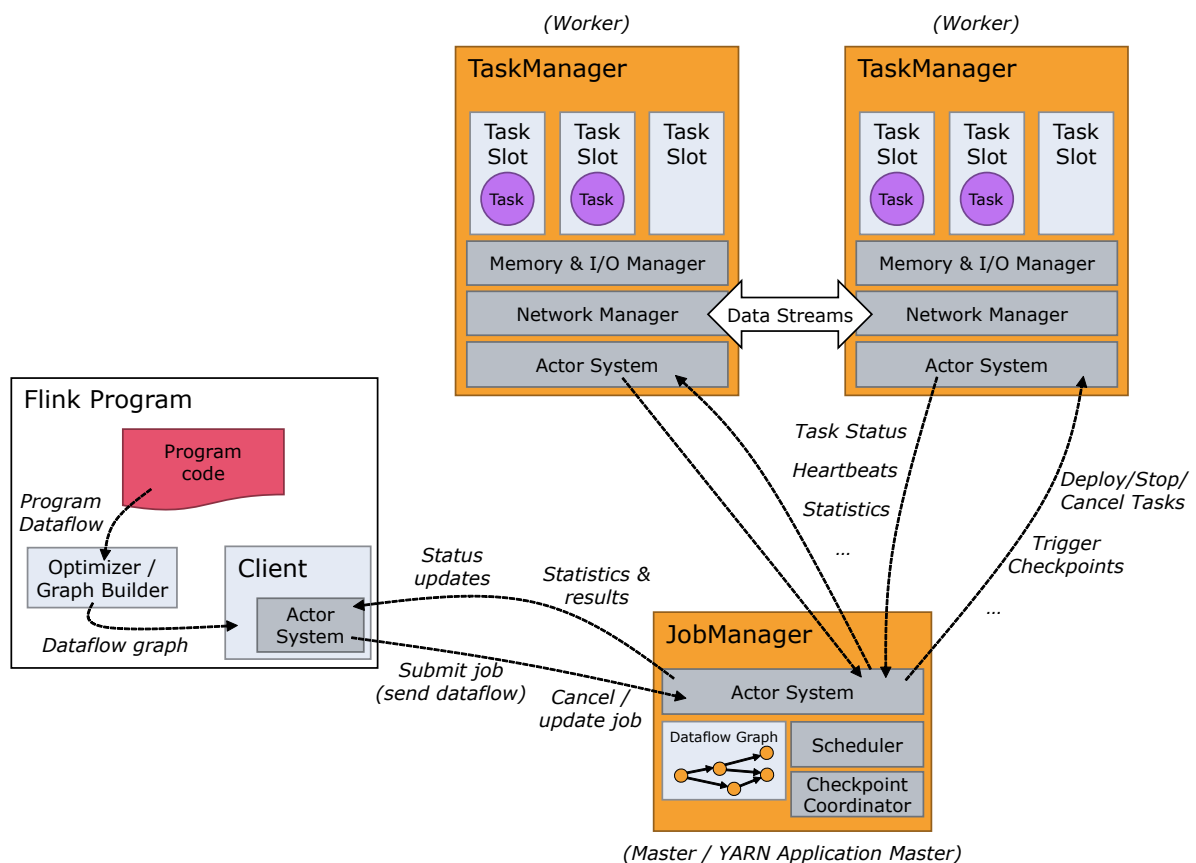
- 1、同时支持高吞吐、低延迟、高性能
  - 2、支持事件时间（Event Time）概念，结合Watermark处理乱序数据
  - 3、支持有状态计算，并且支持多种状态 内存、文件、RocksDB
  - 4、支持高度灵活的窗口（Window）操作 time、count、session
  - 5、基于轻量级分布式快照（CheckPoint）实现的容错 保证exactly-once语义
  - 6、基于JVM实现独立的内存管理
  - 7、Save Points（保存点）
- 

## Flink安装&部署

---

### Flink基本架构

Flink系统架构中包含了两个角色，分别是JobManager和TaskManager，是一个典型的Master-Slave架构。JobManager相当于是Master，TaskManager相当于是Slave



## JobManager (JVM进程) 作用

JobManager负责整个集群的资源管理与任务管理，在一个集群中只能由一个正在工作（active）的JobManager，如果HA集群，那么其他JobManager一定是standby状态

### (1) 资源调度

- 集群启动，TaskManager会将当前节点的资源信息注册给JobManager，所有TaskManager全部注册完毕，集群启动成功，此时JobManager就掌握整个集群的资源情况
- client提交Application给JobManager，JobManager会根据集群中的资源情况，为当前的Application分配TaskSlot资源

### (2) 任务调度

- 根据各个TaskManager节点上的资源分发task到TaskSlot中运行
- Job执行过程中，JobManager会根据设置的触发策略触发checkpoint，通知TaskManager开始checkpoint
- 任务执行完毕，JobManager会将Job执行的信息反馈给client，并且释放TaskManager资源

## TaskManager (JVM进程) 作用

- 负责当前节点上的任务运行及当前节点上的资源管理，TaskManager资源通过TaskSlot进行了划分，每个TaskSlot代表的是一份固定资源。例如，具有三个 slots 的 TaskManager 会将其管理的内存资源分成三等份给每个 slot。划分资源意味着 subtask 之间不会竞争内存资源，但是也意味着它们只拥有固定的资源。注意这里并没有 CPU 隔离，当前 slots 之间只是划分任务的内存资源
- 负责TaskManager之间的数据交换

## client客户端

负责将当前的任务提交给JobManager，提交任务的常用方式：命令提交、web页面提交。获取任务的执行信息

## Standalone集群安装&测试

Standalone是独立部署模式，它不依赖其他平台，不依赖任何的资源调度框架

Standalone集群是由JobManager、TaskManager两个JVM进程组成

### 集群角色划分

node01	node02	node03	node04
JobManager	TaskManager	TaskManager	TaskManager

### 安装步骤

1. 官网下载Flink安装包

Apache Flink® 1.10.0 is our latest stable release.现在最稳定的是1.10.0，不建议采用这个版本，刚从1.9升级到1.10，会存在一些bug，不建议采用小版本号0的安装包，所以我们建议使用1.9.2版本

下载链接:[https://mirrors.tuna.tsinghua.edu.cn/apache/flink/flink-1.9.2/flink-1.9.2-bin-scala\\_2.11.tgz](https://mirrors.tuna.tsinghua.edu.cn/apache/flink/flink-1.9.2/flink-1.9.2-bin-scala_2.11.tgz)

2. 安装包上传到node01节点

3. 解压、修改配置文件

解压：tar -zxf flink-1.9.2-bin-scala\_2.11.tgz

修改flink-conf.yaml配置文件

jobmanager.rpc.address: node01

JobManager地址

jobmanager.rpc.port: 6123

JobManagerRPC通信端口

jobmanager.heap.size: 1024m

JobManager所能使用的堆内存大小

taskmanager.heap.size: 1024m

TaskManager所能使用的堆内存大小

taskmanager.numberOfTaskSlots: 2

TaskManager管理的TaskSlot个数，依据当前物理机的核心数来配置，一般预留出一部分核心（25%）给系统及其他进程使用，一个slot对应一个core。如果core支持超线程，那么slot个数\*2

rest.port: 8081

指定WebUI的访问端口

修改slaves配置文件

node02

node03

node04

4. 同步安装包到其他的节点

同步到node02 scp -r flink-1.9.2 node02: pwd

同步到node03 scp -r flink-1.9.2 node03: pwd

同步到node04 scp -r flink-1.9.2 node04: pwd

## 5. node01配置环境变量

```
vim ~/.bashrc
export FLINK_HOME=/opt/software/flink/flink-1.9.2
export PATH=$PATH:$FLINK_HOME/bin
source ~/.bashrc
```


## 6. 启动standalone集群

启动集群: start-cluster.sh

关闭集群: stop-cluster.sh

## 7. 查看Flink Web UI页面

<http://node01:8081/> 可通过rest.port参数自定义端口

 1586705426320

## 提交Job到standalone集群

常用提交任务的方式有两种，分别是命令提交和Web页面提交

### 1. 命令提交:

flink run -c com.msb.stream.WordCount StudyFlink-1.0-SNAPSHOT.jar

-c 指定主类


-d 独立运行、后台运行

-p 指定并行度

### 2. Web页面提交:

在Web中指定Jar包的位置、主类路径、并行数等

web.submit.enable: true一定是true, 否则不支持Web提交Application

 1586705887854

---

## 3. 启动scala-shell测试

```
start-scala-shell.sh remote <hostname> <portnumber>
```

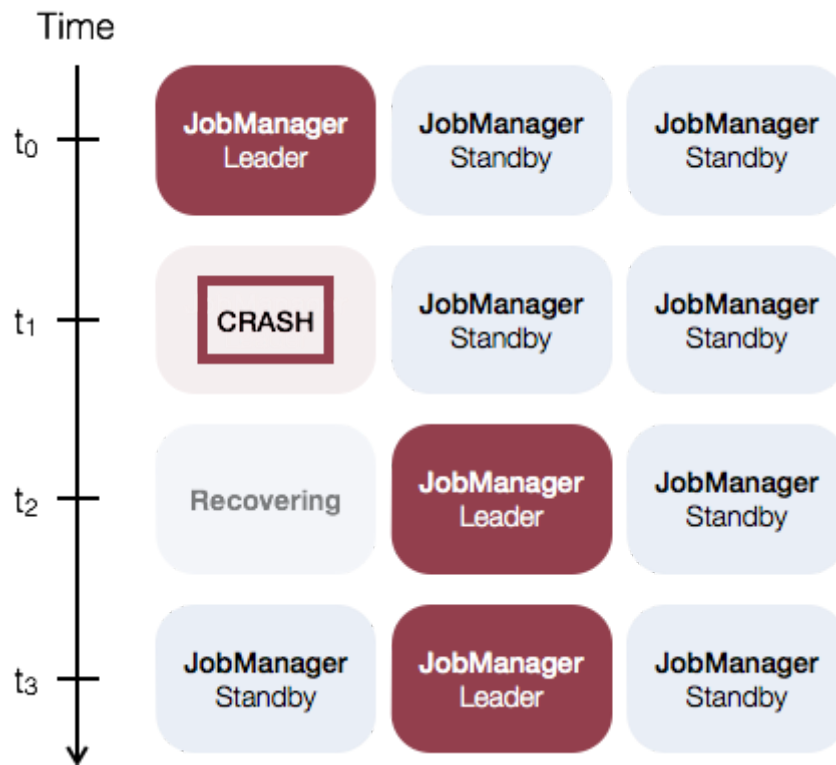
## Standalone HA集群安装&测试

JobManager协调每个flink任务部署,它负责调度和资源管理

默认情况下, 每个flink集群只有一个JobManager, 这将导致一个单点故障(SPOF single-point-of-failure): 如果JobManager挂了, 则不能提交新的任务, 并且运行中的程序也会失败。

使用JobManager HA, 集群可以从JobManager故障中恢复, 从而避免SPOF

Standalone模式(独立模式)下JobManager的高可用性的基本思想是, 任何时候都有一个Active JobManager, 并且多个Standby JobManagers。Standby JobManagers可以在Master JobManager挂掉的情况下接管集群成为Master JobManager。这样保证了没有单点故障, 一旦某一个Standby JobManager接管集群, 程序就可以继续运行。Standby JobManager和Active JobManager实例之间没有明确区别。每个JobManager可以成为Active或Standby节点



如何单独启动JobManager `jobmanager.sh`

如何单独启动TaskManager `taskmanager.sh`

## 集群角色划分

	node01	node02	node03	node04
JobManager	√	√	×	×
TaskManager	×	√	√	√

## 安装步骤

1. 修改配置文件`conf/flink-conf.yaml`

```
high-availability: zookeeper
high-availability.storageDir: hdfs://node01:9000/flink/ha/ 保存JobManager恢复
所需要的所有元数据信息
high-availability.zookeeper.quorum: node01:2181,node02:2181,node03:2181
zookeeper地址
```

2. 修改配置文件`conf/masters`

```
node01:8081
node02:8081
```

3. 同步文件到各个节点

4. 下载支持Hadoop插件并且拷贝到各个节点的安装包的lib目录下



下载地址: <https://repo.maven.apache.org/maven2/org/apache/flink/flink-shaded-hadoop-2-uber/2.6.5-10.0/flink-shaded-hadoop-2-uber-2.6.5-10.0.jar>

- HA集群测试

<http://node01:8081/>

<http://node02:8081/>

两个页面一模一样 存在bug

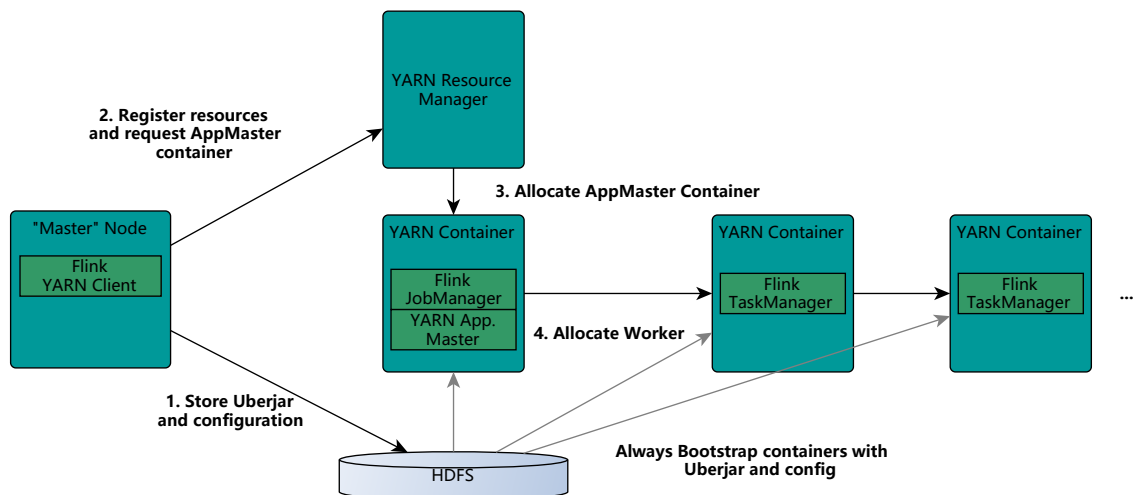
## Flink on Yarn

Flink on Yarn是依托Yarn资源管理器, 现在很多分布式任务都可以支持基于Yarn运行, 这是在企业中使用最多的方式。Why?

(1) 基于Yarn的运行模式可以充分使用集群资源, Spark on Yarn、MapReduce on Yarn、Flink on Yarn等 多套计算框架都可以基于Yarn运行, 充分利用集群资源

(2) 基于Yarn的运行模式降低维护成本

### 运行流程



1. 每当创建一个新flink的yarn session的时候, 客户端会首先检查要请求的资源(containers和memory)是否可用。然后, 将包含flink相关的jar包盒配置上传到HDFS
2. 客户端会向ResourceManager申请一个yarn container 用以启动ApplicationMaster。由于客户端已经将配置和jar文件上传到HDFS, ApplicationMaster将会下载这些jar和配置, 然后启动成功
3. JobManager和AM运行于同一个container
4. AM开始申请启动Flink TaskManager的containers, 这些container会从HDFS上下载jar文件和已修改的配置文件。一旦这些步骤完成, flink就可以接受任务了

### Flink on Yarn两种运行模式

解脱了JobManager的压力 RM做资源管理 JobManager只负责任务管理

- yarn seesion(Start a long-running Flink cluster on YARN)这种方式是在yarn中先启动Flink集群, 然后再提交作业, 这个Flink集群一直停留再yarn中, 一直占据了yarn集群的资源 (只是JobManager会一直占用, 没有Job运行TaskManager并不会运行), 不管有没有任务运行。这种方式能够降低任务的启动时间

- Run a Flink job on YARN 每次提交一个Flink任务的时候，先去yarn中申请资源启动JobManager和TaskManager，然后在当前集群中运行，任务执行完毕，集群关闭。任务之间互相独立，互不影响，可以最大化的使用集群资源，但是每个任务的启动时间变长了

## 配置两种运行模式

### yarn seesion模式配置

- Flink on Yarn依赖Yarn集群和HDFS集群，启动Yarn、HDFS集群 start-all.sh
- 下载支持Hadoop插件并且拷贝到各个节点的安装包的lib目录下

下载地址: <https://repo.maven.apache.org/maven2/org/apache/flink/flink-shaded-hadoop-2-uber/2.6.5-10.0/flink-shaded-hadoop-2-uber-2.6.5-10.0.jar>

- 在yarn中启动Flink集群

```
启动: yarn-session.sh -n 3 -s 3 -nm flink-session -d -q
关闭: yarn application -kill applicationId
```

yarn-session选项:

```
-n,--container <arg>: 在yarn中启动container的个数，实质就是TaskManager的个数
-s,--slots <arg>: 每个TaskManager管理的Slot个数
-nm,--name <arg>: 给当前的yarn-session(Flink集群)起一个名字
-d,--detached: 后台独立模式启动，守护进程
-tm,--taskManagerMemory <arg>: TaskManager的内存大小 单位: MB
-jm,--jobManagerMemory <arg>: JobManager的内存大小 单位: MB
-q,--query: 显示yarn集群可用资源（内存、core）
```

 1586763411234

- 提交Flink Job到yarn-session集群中运行

```
flink run -c com.msb.stream.wordCount -yid application_1586794520478_0007
~/StudyFlink-1.0-SNAPSHOT.jar
```

yid: 指定yarn-session的ApplicationID

不使用yid也可以，因为在启动yarn-session的时候，在tmp临时目录下已经产生了一个隐藏小文件

```
[root@node01 bin]# vim /tmp/.yarn-properties-root
#Generated YARN properties file
#Mon Apr 13 23:39:43 CST 2020
parallelism=9
dynamicPropertiesString=
applicationID=application_1586791887356_0001
```

### Run a Flink job on YARN模式配置

```
flink run -m yarn-cluster -yn 3 -ys 3 -ynm flink-job -c com.msb.stream.WordCount
~/StudyFlink-1.0-SNAPSHOT.jar
```

-yn, --container <arg> 表示分配容器的数量，也就是TaskManager的数量。  
-d, --detached: 设置在后台运行。  
-yjm, --jobManagerMemory<arg>: 设置JobManager的内存，单位是MB。  
-ytm, --taskManagerMemory<arg>: 设置每个TaskManager的内存，单位是MB。  
-ynm, --name: 给当前Flink application在Yarn上指定名称。  
-yq, --query: 显示yarn中可用的资源（内存、cpu核数）  
-yqu, --queue<arg> : 指定yarn资源队列  
-ys, --slots<arg> : 每个TaskManager使用的Slot数量。

## Flink on YARN HA集群安装&测试

无论以什么样的模式提交Application到Yarn中运行，都会启动一个yarn-session(Flink 集群)，依然是由JobManager和TaskManager组成，那么JobManager节点如果宕机，那么整个Flink集群就不会正常运转，所以接下来搭建Flink on YARN HA集群

### 安装步骤

- 修改Hadoop安装包下的yarn-site.xml文件

```
<property>
  <name>yarn.resourcemanager.am.max-attempts</name>
  <value>10</value>
  <description>
    The maximum number of application master execution attempts AppMaster最大
    重试次数
  </description>
</property>
```

- 修改Flink安装包下的flink-conf.yaml文件

```
high-availability: zookeeper
high-availability.storageDir: hdfs://node01:9000/flink/ha/
high-availability.zookeeper.quorum: node01:2181,node02:2181,node03:2181
```

### HA集群测试


两种模式都可以测试，因为不管哪种模式都会启动yarn-session


#### yarn-session模式测试

- 启动yarn-session

```
yarn-session.sh -n 3 -s 3 -nm flink-session -d
```


- 通过yarn web ui 找到ApplicationMaster，发现此时的JobManager是在node02启动，现在kill掉JobManager进程 kill -9 进程号

 1586767536855


 1586767559516

 1586767578933

- 再次查看 发现JobManager切换到node03

1586767713253

- 查看node03日志

1586767814649

```
2020-04-08 22:21:36,044 INFO org.apache.flink.yarn.YarnResourceManager
- ResourceManager
akka.tcp://flink@node03:60599/user/resourcemanager was granted leadership
with fencing token 94c94c3d68ed799374303fad7447418b
```

取消job 开始Run a Flink job on YARN模式测试

flink list

flink cancel id

### Run a Flink job on YARN模式测试

- 提交job

```
flink run -m yarn-cluster -yn 3 -ys 3 -ynm flink-job -c
com.msb.stream.WordCount ~/StudyFlink-1.0-SNAPSHOT.jar
```

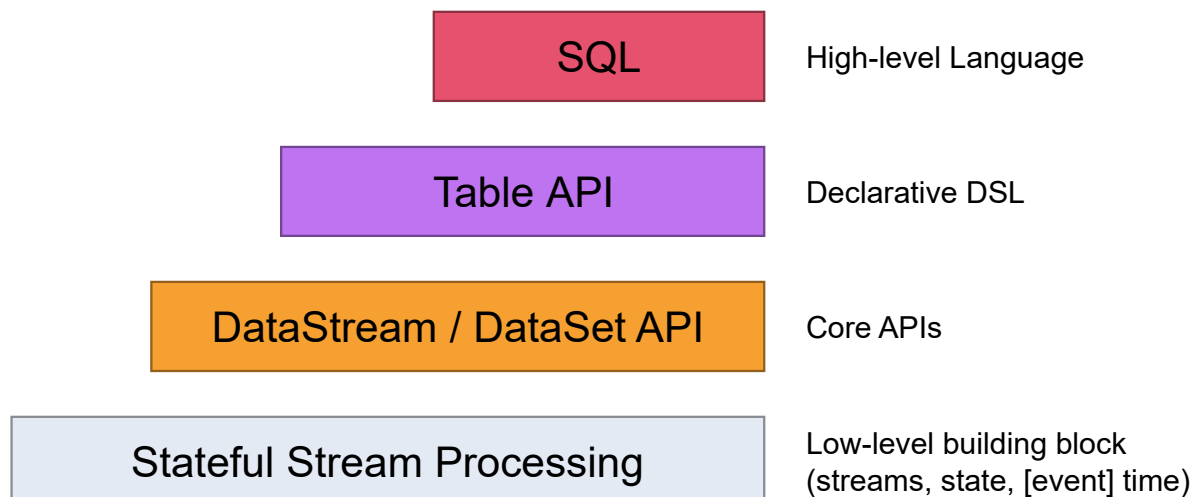
- 停掉JobManager 观察
- 测试完毕，取消job

```
yarn application -kill applicationId
```

## Flink API详解&实操

### Flink API介绍

Flink提供了不同的抽象级别以开发流式或者批处理应用程序



- **Stateful Stream Processing** 最低级的抽象接口是状态化的数据流接口（stateful streaming）。这个接口是通过 ProcessFunction 集成到 DataStream API 中的。该接口允许用户自由的处理来自一个或多个流中的事件，并使用一致的容错状态。另外，用户也可以通过注册 event time 和 processing time 处理回调函数的方法来实现复杂的计算
- **DataStream/DataSet API** DataStream / DataSet API 是 Flink 提供的核心 API，DataSet 处理有界的数据集，DataStream 处理有界或者无界的数据流。用户可以通过各种方法（map / flatmap / window / keyby / sum / max / min / avg / join 等）将数据进行转换 / 计算
- **Table API** Table API 提供了例如 select、project、join、group-by、aggregate 等操作，使用起来却更加简洁,可以在表与 DataStream/DataSet 之间无缝切换，也允许程序将 Table API 与 DataStream 以及 DataSet 混合使用
- **SQL** Flink 提供的最高层级的抽象是 SQL。这一层抽象在语法与表达能力上与 Table API 类似。SQL 抽象与 Table API 交互密切，同时 SQL 查询可以直接在 Table API 定义的表上执行

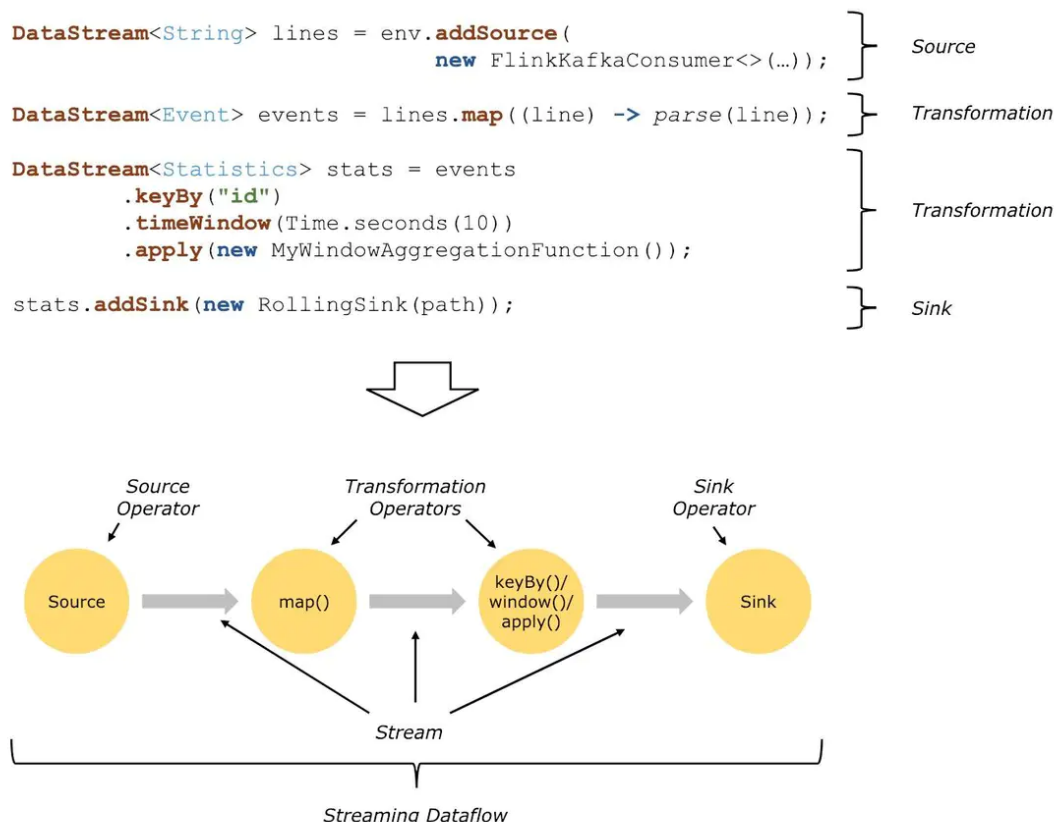
## Dataflows数据流图

在Flink的世界观中，一切都是数据流，所以对于批计算来说，那只是流计算的一个特例而已

Flink Dataflows是由三部分组成，分别是：source、transformation、sink结束

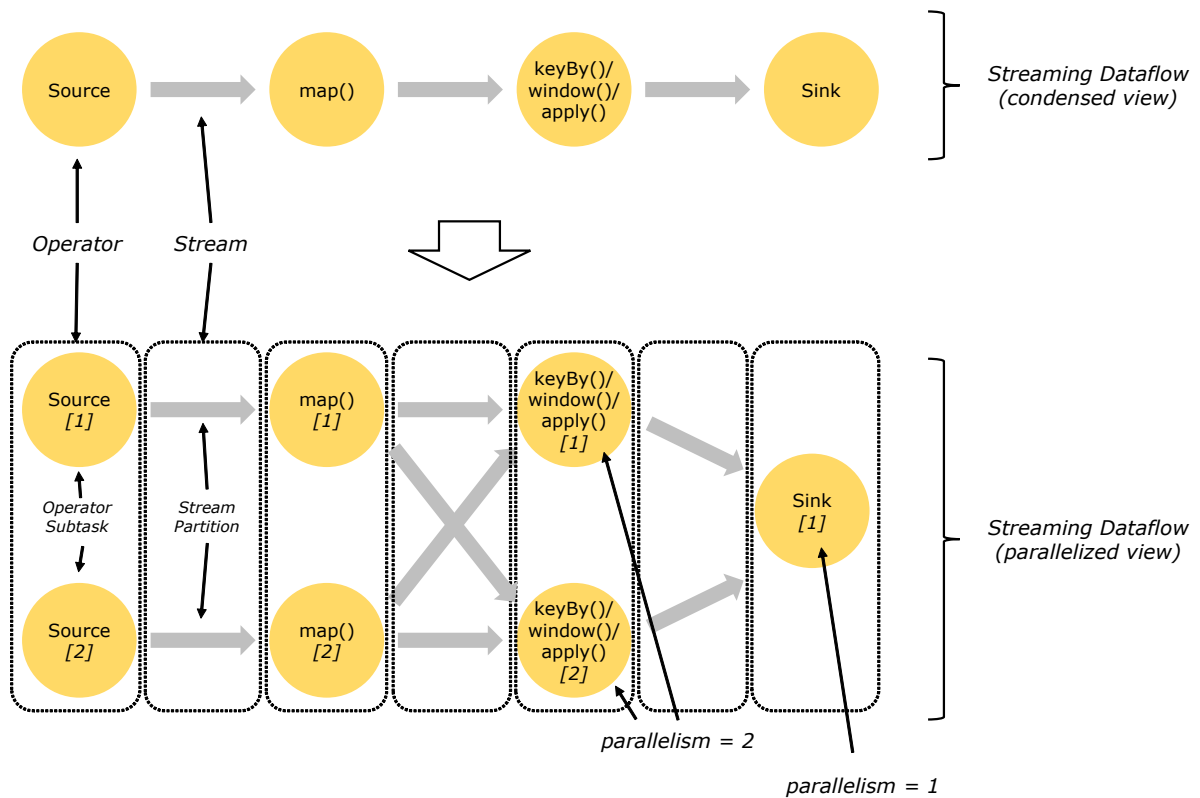
source数据源会源源不断的产生数据，transformation将产生的数据进行各种业务逻辑的数据处理，最终由sink输出到外部（console、kafka、redis、DB.....）

基于Flink开发的程序都能够映射成一个Dataflows



当source数据源的数量比较大或计算逻辑相对比较复杂的情况下，需要提高并行度来处理数据，采用并行数据流

通过设置不同算子的并行度 source并行度设置为2 map也是2.... 代表会启动多个并行的线程来处理数据



## 配置开发环境

每个 Flink 应用都需要依赖一组 Flink 类库。Flink 应用至少需要依赖 Flink APIs。许多应用还会额外依赖连接器类库(比如 Kafka、Cassandra 等)。当用户运行 Flink 应用时(无论是在 IDEA 环境下进行测试，还是部署在分布式环境下)，运行时类库都必须可用

开发工具：IntelliJ IDEA

配置开发Maven依赖：

```
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-scala_2.11</artifactId>
  <version>1.10.0</version>
</dependency>
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-streaming-scala_2.11</artifactId>
  <version>1.10.0</version>
</dependency>
```

注意点：

- 如果要程序打包提交到集群运行，打包的时候不需要包含这些依赖，因为集群环境已经包含了这些依赖，此时依赖的作用域应该设置为 `provided`
- Flink 应用在 IntelliJ IDEA 中运行，这些 Flink 核心依赖的作用域需要设置为 `compile` 而不是 `provided`。否则 IntelliJ 不会添加这些依赖到 classpath，会导致应用运行时抛出 `NoClassDefFoundError` 异常

添加打包插件：

```
<build>
  <plugins>
    <plugin>
```

```

<groupId>org.apache.maven.plugins</groupId>
<artifactId>maven-shade-plugin</artifactId>
<version>3.1.1</version>
<executions>
  <execution>
    <phase>package</phase>
    <goals>
      <goal>shade</goal>
    </goals>
    <configuration>
      <artifactSet>
        <excludes>

<exclude>com.google.code.findbugs:jsr305</exclude>
          <exclude>org.slf4j:*</exclude>
          <exclude>log4j:*</exclude>
        </excludes>
      </artifactSet>
      <filters>
        <filter>
          <!--不要拷贝 META-INF 目录下的签名，
          否则会引起 SecurityExceptions 。 -->
          <artifact>*:*</artifact>
          <excludes>
            <exclude>META-INF/*.SF</exclude>
            <exclude>META-INF/*.DSA</exclude>
            <exclude>META-INF/*.RSA</exclude>
          </excludes>
        </filter>
      </filters>
      <transformers>
        <transformer
implementation="org.apache.maven.plugins.shade.resource.ManifestResourceTransfor
mer">
          <mainClass>my.programs.main.clazz</mainClass>
        </transformer>
      </transformers>
    </configuration>
  </execution>
</executions>
</plugin>
</plugins>
</build>

```

## WordCount流批计算程序

批计算：统计HDFS文件单词出现的次数

读取HDFS数据需要添加Hadoop依赖

```

<dependency>
  <groupId>org.apache.hadoop</groupId>
  <artifactId>hadoop-client</artifactId>
  <version>2.6.5</version>
</dependency>

```

WordCount代码：

```

val env = ExecutionEnvironment.getExecutionEnvironment
val initDS: DataSet[String] =
env.readTextFile("hdfs://node01:9000/flink/data/wc")
val restDS: AggregateDataSet[(String, Int)] = initDS.flatMap(_._split("
")).map(_._1).groupBy(0).sum(1)
restDS.print()

```

流计算：统计数据流中，单词出现的次数

```

//准备环境
/**
 * createLocalEnvironment 创建一个本地执行的环境 local
 * createLocalEnvironmentWithWebUI 创建一个本地执行的环境 同时还开启web UI的查看
端口 8081
 * getExecutionEnvironment 根据你执行的环境创建上下文，比如local cluster
 */
val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setParallelism(1)
/**
 * DataStream: 一组相同类型的元素 组成的数据流
 */
val initStream:DataStream[String] = env.socketTextStream("node01",8888)
val wordStream = initStream.flatMap(_._split(" "))
val pairStream = wordStream.map(_._1)
val keyByStream = pairStream.keyBy(0)
val restStream = keyByStream.sum(1)
restStream.print()

/**
 * 6> (msb,1)
 * 1> (, ,1)
 * 3> (hello,1)
 * 3> (hello,2)
 * 6> (msb,2)
 * 默认就是有状态的计算
 * 6> 代表是哪一个线程处理的
 *
 * 相同的数据一定是由某一个thread处理
 */
//启动Flink 任务
env.execute("first flink job")

```


## WordCount Dataflows 算子链

为了更高效地分布式执行，Flink会尽可能地将operator的subtask链接（chain）在一起形成task。每个task在一个线程中执行。将operators链接成task是非常有效的优化：它能减少线程之间的切换，减少消息的序列化/反序列化，减少数据在缓冲区的交换，减少了延迟的同时提高整体的吞吐量

## Flink任务调度规则

- 不同Task下的subtask分到同一个TaskSlot，提高数据传输效率
- 相同Task下的subtask不会分到同一个TaskSlot，充分利用集群资源



1587035707335

1587035764668

## Flink并行度设置方式

### 1. 在算子上设置

```
val wordStream = initStream.flatMap(_._split(" ")).setParallelism(2)
```

### 2. 在上下文环境中设置

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setParallelism(1)
```

### 3. client提交Job时设置

```
flink run -c com.msb.stream.wordCount -p 3 StudyFlink-1.0-SNAPSHOT.jar
```

### 4. 在flink-conf.yaml配置文件中设置

```
parallelism.default: 1
```

这四种设置并行度的方式，优先级依次递减

## Dataflows DataSource数据源

Flink内嵌支持的数据源非常多，比如HDFS、Socket、Kafka、Collections Flink也提供了addSource方式，可以自定义数据源，本小节将讲解Flink所有内嵌数据源及自定义数据源的原理及API

### File Source

- 通过读取本地、HDFS文件创建一个数据源

如果读取的是HDFS上的文件，那么需要导入Hadoop依赖

```
<dependency>
  <groupId>org.apache.hadoop</groupId>
  <artifactId>hadoop-client</artifactId>
  <version>2.6.5</version>
</dependency>
```

代码：

```
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
//在算子转换的时候，会将数据转换成Flink内置的数据类型，所以需要将隐式转换导入进来，才能自动进行
//类型转换
import org.apache.flink.streaming.api.scala._

object FileSource {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val textStream = env.readTextFile("hdfs://node01:9000/flink/data/wc")
    textStream.flatMap(_.split(" ")).map((_, 1)).keyBy(0).sum(1).print()
    //读完就停止
    env.execute()
  }
}
```

- 每隔10s中读取HDFS指定目录下的新增文件内容，并且进行WordCount

业务场景：在企业中一般都会做实时的ETL，当Flume采集来新的数据，那么基于Flink实时做ETL入仓

```
import org.apache.flink.api.java.io.TextInputFormat
import org.apache.flink.core.fs.Path
import org.apache.flink.streaming.api.functions.source.FileProcessingMode
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
//在算子转换的时候，会将数据转换成Flink内置的数据类型，所以需要将隐式转换导入进来，才能自动进行
//类型转换
import org.apache.flink.streaming.api.scala._

object FileSource {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    //读取hdfs文件
    val filePath = "hdfs://node01:9000/flink/data/"
    val textInputFormat = new TextInputFormat(new Path(filePath))
    //每隔10s中读取 hdfs上新增文件内容
    val textStream =
env.readFile(textInputFormat, filePath, FileProcessingMode.PROCESS_CONTINUOUSLY, 10
)
    //    val textStream = env.readTextFile("hdfs://node01:9000/flink/data/wc")
    textStream.flatMap(_.split(" ")).map((_, 1)).keyBy(0).sum(1).print()
    env.execute()
  }
}
```

**readTextFile底层调用的就是readFile方法，readFile是一个更加底层的方式，使用起来会更加的灵活**

## Collection Source

基于本地集合的数据源，一般用于测试场景，没有太大意义

```
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._

object CollectionSource {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val stream = env.fromCollection(List("hello flink msb", "hello msb msb"))
    stream.flatMap(_._split(" ")).map(_._1).keyBy(0).sum(1).print()
    env.execute()
  }
}
```

## Socket Source

接受Socket Server中的数据，已经讲过

```
val initStream:DataStream[String] = env.socketTextStream("node01",8888)
```

## Kafka Source

Flink接受Kafka中的数据，首先先配置flink与kafka的连接器依赖

官网地址: <https://ci.apache.org/projects/flink/flink-docs-release-1.9/dev/connectors/kafka.html>

maven依赖

```
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-connector-kafka_2.11</artifactId>
  <version>1.9.2</version>
</dependency>
```

代码:

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val prop = new Properties()
prop.setProperty("bootstrap.servers", "node01:9092,node02:9092,node03:9092")
prop.setProperty("group.id", "flink-kafka-id001")
prop.setProperty("key.deserializer", classOf[StringDeserializer].getName)
prop.setProperty("value.deserializer", classOf[StringDeserializer].getName)
/**
 * earliest:从头开始消费，旧数据会频繁消费
 * latest:从最近的数据开始消费，不再消费旧数据
 */
prop.setProperty("auto.offset.reset", "latest")

val kafkaStream = env.addSource(new FlinkKafkaConsumer[(String, String)](
  "flink-kafka", new KafkaDeserializationSchema[(String, String)] {
    override def isEndOfStream(t: (String, String)): Boolean = false

    override def deserialize(consumerRecord: ConsumerRecord[Array[Byte],
Array[Byte]]): (String, String) = {
      val key = new String(consumerRecord.key(), "UTF-8")
      val value = new String(consumerRecord.value(), "UTF-8")
    }
  })
```

```

        (key, value)
    }
    //指定返回数据类型
    override def getProducedType: TypeInformation[(String, String)] =
        createTuple2TypeInformation(createTypeInformation[String],
        createTypeInformation[String])
    }, prop))
    kafkaStream.print()
    env.execute()

```

kafka命令消费key value值

kafka-console-consumer.sh --zookeeper node01:2181 --topic flink-kafka --property print.key=true

默认只是消费value值

KafkaDeserializationSchema: 读取kafka中key、value

SimpleStringSchema: 读取kafka中value

## Custom Source

Sources are where your program reads its input from. You can attach a source to your program by using `StreamExecutionEnvironment.addSource(sourceFunction)`. Flink comes with a number of pre-implemented source functions, but you can always write your own custom sources by implementing the `SourceFunction` for non-parallel sources, or by implementing the `ParallelSourceFunction` interface or extending the `RichParallelSourceFunction` for parallel sources.

- 基于SourceFunction接口实现单并行度数据源

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
//source的并行度为1 单并行度source源
val stream = env.addSource(new SourceFunction[String] {
    var flag = true
    override def run(ctx: SourceFunction.SourceContext[String]): Unit = {
        val random = new Random()
        while (flag) {
            ctx.collect("hello" + random.nextInt(1000))
            Thread.sleep(200)
        }
    }
})
//停止产生数据
override def cancel(): Unit = flag = false
})
stream.print()
env.execute()

```

基于ParallelSourceFunction接口实现多并行度数据源

```

public interface ParallelSourceFunction<OUT> extends SourceFunction<OUT> {}

```

```
public abstract class RichParallelSourceFunction<OUT> extends
AbstractRichFunction
    implements ParallelSourceFunction<OUT> {
    private static final long serialVersionUID = 1L;
}
```


实现ParallelSourceFunction接口=继承RichParallelSourceFunction

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val sourceStream = env.addSource(new ParallelSourceFunction[String] {
    var flag = true

    override def run(ctx: SourceFunction.SourceContext[String]): Unit = {
        val random = new Random()
        while (flag) {
            ctx.collect("hello" + random.nextInt(1000))
            Thread.sleep(500)
        }
    }

    override def cancel(): Unit = {
        flag = false
    }
}).setParallelism(2)
```

数据源可以设置为多并行度

1586950908040

## Dataflows Transformations

Transformations算子可以将一个或者多个算子转换成一个新的数据流，使用Transformations算子组合可以进行复杂的业务处理

### Map

DataStream → DataStream

遍历数据流中的每一个元素，产生一个新的元素

### FlatMap

DataStream → DataStream

遍历数据流中的每一个元素，产生N个元素 N=0, 1, 2,.....

### Filter

DataStream → DataStream

过滤算子，根据数据流的元素计算出一个boolean类型的值，true代表保留，false代表过滤掉

### KeyBy

DataStream → KeyedStream

根据数据流中指定的字段来分区，相同指定字段值的数据一定是在同一个分区中，内部分区使用的是HashPartitioner

指定分区字段的方式有三种：

- 1、根据索引号指定
- 2、通过匿名函数来指定
- 3、通过实现KeySelector接口 指定分区字段

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1, 100)
stream
  .map(x => (x % 3, 1))
  //根据索引号来指定分区字段
  //      .keyBy(0)
  //通过传入匿名函数 指定分区字段
  //      .keyBy(x=>x._1)
  //通过实现KeySelector接口 指定分区字段
  .keyBy(new KeySelector[(Long, Int), Long] {
    override def getKey(value: (Long, Int)): Long = value._1
  })
  .sum(1)
  .print()
env.execute()
```

## Reduce

KeyedStream：根据key分组 → DataStream

注意，reduce是基于分区后的流对象进行聚合，也就是说，DataStream类型的对象无法调用reduce方法

```
.reduce((v1,v2) => (v1._1,v1._2 + v2._2))
```

demo01：读取kafka数据，实时统计各个卡口下的车流量

- 实现kafka生产者，读取卡口数据并且往kafka中生产数据

```
val prop = new Properties()
prop.setProperty("bootstrap.servers", "node01:9092,node02:9092,node03:9092")
prop.setProperty("key.serializer", classOf[StringSerializer].getName)
prop.setProperty("value.serializer", classOf[StringSerializer].getName)

val producer = new KafkaProducer[String, String](prop)

val iterator = Source.fromFile("data/carFlow_all_column_test.txt", "UTF-8").getLines()
for (i <- 1 to 100) {
  for (line <- iterator) {
    //将需要的字段值 生产到kafka集群 car_id monitor_id event-time speed
    //车牌号 卡口号 车辆通过时间 通过速度
    val splits = line.split(",")
    val monitorID = splits(0).replace("'", "")
    val car_id = splits(2).replace("'", "")
    val eventTime = splits(4).replace("'", "")
    val speed = splits(6).replace("'", "")
    if (!"00000000".equals(car_id)) {
```

```

        val event = new StringBuilder
        event.append(monitorID + "\t").append(car_id+"\t").append(eventTime +
"\t").append(speed)
        producer.send(new ProducerRecord[String, String]("flink-kafka",
event.toString()))
    }

    Thread.sleep(500)
}
}

```

- 实现代码

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val props = new Properties()
props.setProperty("bootstrap.servers","node01:9092,node02:9092,node03:9092")
props.setProperty("key.deserializer",classOf[StringDeserializer].getName)
props.setProperty("value.deserializer",classOf[StringDeserializer].getName)
props.setProperty("group.id","flink001")
props.setProperty("auto.offset.reset","latest")

val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new
SimpleStringSchema(),props))
stream.map(data => {
    val splits = data.split("\t")
    val carFlow = CarFlow(splits(0),splits(1),splits(2),splits(3).toDouble)
    (carFlow,1)
}).keyBy(_._1.monitorId)
    .sum(1)
    .print()
env.execute()

```

## Aggregations

KeyedStream → DataStream

Aggregations代表的是一类聚合算子，具体算子如下：

```

keyedStream.sum(0)
keyedStream.sum("key")
keyedStream.min(0)
keyedStream.min("key")
keyedStream.max(0)
keyedStream.max("key")
keyedStream.minBy(0)
keyedStream.minBy("key")
keyedStream.maxBy(0)
keyedStream.maxBy("key")

```

demo02: 实时统计各个卡口最先通过的汽车的信息

```

val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new
SimpleStringsSchema(), props))
stream.map(data => {
    val splits = data.split("\t")
    val carFlow = CarFlow(splits(0), splits(1), splits(2), splits(3).toDouble)
    val eventTime = carFlow.eventTime
    val format = new SimpleDateFormat("yyyy-MM-dd HH:mm:ss")
    val date = format.parse(eventTime)
    (carFlow, date.getTime)
}).keyBy(_._1.monitorId)
    .min(1)
    .map(_._1)
    .print()
env.execute()

```

## Union 真合并

**DataStream\*** → DataStream

Union of two or more data streams creating a new stream containing all the elements from all the streams

合并两个或者更多的数据流产生一个新的数据流，这个新的数据流中包含了所合并的数据流的元素

注意：需要保证数据流中元素类型一致

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val ds1 = env.fromCollection(List(("a",1),("b",2),("c",3)))
val ds2 = env.fromCollection(List(("d",4),("e",5),("f",6)))
val ds3 = env.fromCollection(List(("g",7),("h",8)))
// val ds3 = env.fromCollection(List((1,1),(2,2)))
val unionStream = ds1.union(ds2,ds3)
unionStream.print()
env.execute()

```

## Connect 假合并

DataStream,DataStream → ConnectedStreams

合并两个数据流并且保留两个数据流的数据类型，能够共享两个流的状态

```

val ds1 = env.socketTextStream("node01", 8888)
val ds2 = env.socketTextStream("node01", 9999)
val wcStream1 = ds1.flatMap(_ .split(" ")).map(_ , 1)).keyBy(0).sum(1)
val wcStream2 = ds2.flatMap(_ .split(" ")).map(_ , 1)).keyBy(0).sum(1)
val restStream: ConnectedStreams[(String, Int), (String, Int)] =
wcStream2.connect(wcStream1)

```

## CoMap, CoFlatMap

ConnectedStreams → DataStream

CoMap, CoFlatMap并不是具体算子名字，而是一类操作名称

凡是基于ConnectedStreams数据流做map遍历，这类操作叫做CoMap

凡是基于ConnectedStreams数据流做flatMap遍历，这类操作叫做CoFlatMap



### CoMap第一种实现方式:

```
restStream.map(new CoMapFunction[(String,Int),(String,Int),(String,Int)] {  
    //对第一个数据流做计算  
    override def map1(value: (String, Int)): (String, Int) = {  
        (value._1+":first",value._2+100)  
    }  
    //对第二个数据流做计算  
    override def map2(value: (String, Int)): (String, Int) = {  
        (value._1+":second",value._2*100)  
    }  
}).print()
```

### CoMap第二种实现方式:

```
restStream.map(  
    //对第一个数据流做计算  
    x=>{(x._1+":first",x._2+100)}  
    //对第二个数据流做计算  
    ,y=>{(y._1+":second",y._2*100)}  
).print()
```

### CoFlatMap第一种实现方式:

```
ds1.connect(ds2).flatMap((x,c:Collector[String])=>{  
    //对第一个数据流做计算  
    x.split(" ").foreach(w=>{  
        c.collect(w)  
    })  
}  
    //对第二个数据流做计算  
    ,(y,c:Collector[String])=>{  
    y.split(" ").foreach(d=>{  
        c.collect(d)  
    })  
}).print
```

### CoFlatMap第二种实现方式:

```
ds1.connect(ds2).flatMap(  
    //对第一个数据流做计算  
    x=>{  
        x.split(" ")  
    }  
    //对第二个数据流做计算  
    ,y=>{  
        y.split(" ")  
    }  
).print()
```

### CoFlatMap第三种实现方式:

```
ds1.connect(ds2).flatMap(new CoFlatMapFunction[String,String,(String,Int)] {  
    //对第一个数据流做计算
```

```

override def flatMap1(value: String, out: Collector[(String, Int)]): Unit =
{
    val words = value.split(" ")
    words.foreach(x=>{
        out.collect((x,1))
    })
}

//对第二个数据流做计算
override def flatMap2(value: String, out: Collector[(String, Int)]): Unit =
{
    val words = value.split(" ")
    words.foreach(x=>{
        out.collect((x,1))
    })
}
}).print()

```

demo03: 现有一个配置文件存储车牌号与车主的真实姓名，通过数据流中的车牌号实时匹配出对应的车主姓名（注意：配置文件可能实时改变）

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setParallelism(1)
val filePath = "data/carId2Name"
val carId2NameStream = env.readFile(new TextInputFormat(new
Path(filePath)), filePath, FileProcessingMode.PROCESS_CONTINUOUSLY, 10)
val dataStream = env.socketTextStream("node01", 8888)
dataStream.connect(carId2NameStream).map(new CoMapFunction[String, String, String]
{
    private val hashMap = new mutable.HashMap[String, String]()
    override def map1(value: String): String = {
        hashMap.getOrElse(value, "not found name")
    }

    override def map2(value: String): String = {
        val splits = value.split(" ")
        hashMap.put(splits(0), splits(1))
        value + "加载完毕..."
    }
}).print()
env.execute()

```

此demo仅限深度理解connect算子和CoMap操作，后期还需使用广播流优化

## Split

DataStream → SplitStream

根据条件将一个流分成两个或者更多的流

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,100)
val splitStream = stream.split(
  d => {
    d % 2 match {
      case 0 => List("even")
      case 1 => List("odd")
    }
  }
)
splitStream.select("even").print()
env.execute()

```

@deprecated Please use side output instead

## Select

SplitStream → DataStream

从SplitStream中选择一个或者多个数据流

```
splitStream.select("even").print()
```

## side output侧输出流

流计算过程，可能遇到根据不同的条件来分隔数据流。filter分割造成不必要的数据复制

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.socketTextStream("node01",8888)
val gtTag = new OutputTag[String]("gt")
val processStream = stream.process(new ProcessFunction[String, String] {
  override def processElement(value: String, ctx: ProcessFunction[String,
String]#Context, out: Collector[String]): Unit = {
    try {
      val longVar = value.toLong
      if (longVar > 100) {
        out.collect(value)
      } else {
        ctx.output(gtTag, value)
      }
    } catch {
      case e => e.getMessage
      ctx.output(gtTag, value)
    }
  }
})
val sideStream = processStream.getSideOutput(gtTag)
sideStream.print("sideStream")
processStream.print("mainStream")
env.execute()

```

## Iterate (比较重要)

DataStream → IterativeStream → DataStream

Iterate算子提供了对数据流迭代的支持

迭代由两部分组成：迭代体、终止迭代条件


不满足终止迭代条件的数据流会返回到stream流中，进行下一次迭代

满足终止迭代条件的数据流继续往下游发送

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val initStream = env.socketTextStream("node01", 8888)
val stream = initStream.map(_._1.toLong)
stream.iterate {
    iteration => {
        //定义迭代逻辑
        val iterationBody = iteration.map { x => {
            println(x)
            if(x > 0) x - 1
            else x
        } }
        //> 0 大于0的值继续返回到stream流中,当 <= 0 继续往下游发送
        (iterationBody.filter(_ > 0), iterationBody.filter(_ <= 0))
    }
}.print()
env.execute()
```

## 函数类和富函数类

在使用Flink算子的时候，可以通过传入匿名函数和函数类对象 例如：

1587212235441

函数类分为：普通函数类、富函数类（自行划分）

富函数类相比于普通的函数，可以获取运行环境的上下文（Context），拥有一些生命周期方法，管理状态，可以实现更加复杂的功能

普通函数类	富函数类
MapFunction	RichMapFunction
FlatMapFunction	RichFlatMapFunction
FilterFunction	RichFilterFunction
.....	.....

- 使用普通函数类过滤掉车速高于100的车辆信息

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.readTextFile("./data/carFlow_all_column_test.txt")
stream.filter(new FilterFunction[String] {
    override def filter(value: String): Boolean = {
        if (value != null && !"".equals(value)) {
            val speed = value.split(",")(6).replace("'", "").toLong
            if (speed > 100)
                false
            else
                true
        }
    }
}).print()
env.execute()
```

```

        true
    }else
        false
    }
}).print()
env.execute()

```

- 使用富函数类，将车牌号转化成车主真实姓名，映射表存储在Redis中

```

@Public
public abstract class RichMapFunction<IN, OUT> extends AbstractRichFunction
implements MapFunction<IN, OUT> {

    private static final long serialVersionUID = 1L;

    @Override
    public abstract OUT map(IN value) throws Exception;
}

public abstract class AbstractRichFunction implements RichFunction, Serializable
{
    @Override
    public void open(Configuration parameters) throws Exception {}

    @Override
    public void close() throws Exception {}
}

```

## abstract class RichMapFunction实现MapFunction接口

map函数是抽象方法，需要实现

添加redis依赖

wordcount数据写入到redis

```

<dependency>
    <groupId>redis.clients</groupId>
    <artifactId>jedis</artifactId>
    <version>${redis.version}</version>
</dependency>

```

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.socketTextStream("node01", 8888)
stream.map(new RichMapFunction[String, String] {

```

```

    private var jedis: Jedis = _

```

//初始化函数 在每一个thread启动的时候（处理元素的时候，会调用一次）

//在open中可以创建连接redis的连接

```

override def open(parameters: Configuration): Unit = {

```

//getRuntimeContext可以获取flink运行的上下文环境 AbstractRichFunction抽象类提供的

```

    val taskName = getRuntimeContext.getTaskName

```

```

    val subtasks = getRuntimeContext.getTaskNameWithSubtasks
    println("=====open====="+taskName:" + taskName +
"\tsubtasks:"+subtasks)
    jedis = new Jedis("node01", 6379)
    jedis.select(3)
}

//每处理一个元素，就会调用一次
override def map(value: String): String = {
    val name = jedis.get(value)
    if(name == null){
        "not found name"
    }else
        name
}

//元素处理完毕后，会调用close方法
//关闭redis连接
override def close(): Unit = {
    jedis.close()
}
}).setParallelism(2).print()

env.execute()

```

## 底层API(ProcessFunctionAPI)

1587216807922

属于低层次的API，我们前面讲的map、filter、flatMap等算子都是基于这层高层封装出来的

越低层次的API，功能越强大，用户能够获取的信息越多，比如可以拿到元素状态信息、事件时间、设置定时器等

- 监控每辆汽车，车速超过100迈，5s钟后发出超速的警告通知

```

object MonitorOverSpeed02 {
    case class CarInfo(carId:String,speed:Long)
    def main(args: Array[String]): Unit = {
        val env = StreamExecutionEnvironment.getExecutionEnvironment
        val stream = env.socketTextStream("node01",8888)
        stream.map(data => {
            val splits = data.split(" ")
            val carId = splits(0)
            val speed = splits(1).toLong
            CarInfo(carId,speed)
        }).keyBy(_._1.carId)
        //KeyedStream调用process需要传入KeyedProcessFunction
        //DataStream调用process需要传入ProcessFunction
        .process(new KeyedProcessFunction[String,CarInfo,String] {

            override def processElement(value: CarInfo, ctx:
            KeyedProcessFunction[String, CarInfo, String]#Context, out:
            Collector[String]): Unit = {
                val currentTime = ctx.timerService().currentProcessingTime()
                if(value.speed > 100 ){
                    val timerTime = currentTime + 2 * 1000
                    ctx.timerService().registerProcessingTimeTimer(timerTime)
                }
            }
        })
    }
}

```

```

    }
}

    override def onTimer(timestamp: Long, ctx:
    KeyedProcessFunction[String, CarInfo, String]#OnTimerContext, out:
    Collector[String]): Unit = {
        var warnMsg = "warn... time:" + timestamp + "   carID:" +
    ctx.getCurrentKey
        out.collect(warnMsg)
    }
    }).print()

    env.execute()
}
}

```

## 总结

使用Map Filter....算子的适合，可以直接传入一个匿名函数、普通函数类对象(MapFunction FilterFunction)

富函数类对象 (RichMapFunction、RichFilterFunction)

传入的富函数类对象：可以拿到任务执行的上下文，生命周期方法、管理状态.....

如果业务比较复杂，通过Flink提供这些算子无法满足我们的需求，通过process算子直接使用比较底层 API (使用这套API 上下文、生命周期方法、测输出流、时间服务)

KeyedDataStream调用process    KeyedProcessFunction

DataStream调用process    ProcessFunction

具体写代码的适合，看提示就行

## Dataflows分区策略

### shuffle

场景：增大分区、提高并行度，解决数据倾斜

DataStream → DataStream

分区元素随机均匀分发到下游分区，网络开销比较大

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(1)
println(stream.getParallelism)
stream.shuffle.print()
env.execute()

```

console result: 上游数据比较随意的分发到下游

```
2> 1
1> 4
7> 10
4> 6
6> 3
5> 7
8> 2
1> 5
1> 8
1> 9
```

## rebalance

场景：增大分区、提高并行度，解决数据倾斜

DataStream → DataStream

轮询分区元素，均匀的将元素分发到下游分区，下游每个分区的数据比较均匀，在发生数据倾斜时非常有用，网络开销比较大

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setParallelism(3)
val stream = env.generateSequence(1,100)
val shuffleStream = stream.rebalance
shuffleStream.print()
env.execute()
```

console result:上游数据比较均匀的分发到下游

```
8> 6
3> 1
5> 3
7> 5
1> 7
2> 8
6> 4
4> 2
3> 9
4> 10
```

## rescale

场景：减少分区 防止发生大量的网络传输 不会发生全量的重分区

DataStream → DataStream

通过轮询分区元素，将一个元素集合从上游分区发送给下游分区，发送单位是集合，而不是一个个元素

注意：rescale发生的是本地数据传输，而不需要通过网络传输数据，比如taskmanager的槽数。简单来说，上游的数据只会发送给本TaskManager中的下游

1587204455957



```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(2)
stream.writeAsText("./data/stream1").setParallelism(2)
stream.rescale.writeAsText("./data/stream2").setParallelism(4)
env.execute()
```

console result: stream1:1内容 分发给stream2:1和stream2:2

stream1:1

```
1
3
5
7
9
```

stream1:2

```
2
4
6
8
10
```

stream2:1

```
1
5
9
```

stream2:2

```
3
7
```

stream2:3

```
2
6
10
```

stream2:4

```
4
8
```

## broadcast

场景：需要使用映射表、并且映射表会经常发生变动的场景

DataStream → DataStream

上游中每一个元素内容广播到下游每一个分区中

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(2)
stream.writeAsText("./data/stream1").setParallelism(2)
stream.broadcast.writeAsText("./data/stream2").setParallelism(4)
env.execute()
```

console result: stream1:1、2内容广播到了下游每个分区中

stream1:1

```
1
3
5
7
9
```

stream1:2

```
2
4
6
8
10
```

stream2:1

```
1
3
5
7
9
2
4
6
8
10
```

## global

场景：并行度降为1

DataStream → DataStream

上游分区的数据只分发给下游的第一个分区

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(2)
stream.writeAsText("./data/stream1").setParallelism(2)
stream.global.writeAsText("./data/stream2").setParallelism(4)
env.execute()
```

console result: stream1:1、2内容只分发给stream2:1

stream1:1

```
1
3
5
7
9
```

stream1:2

```
2
4
6
8
10
```

stream2:1

```
1
3
5
7
9
2
4
6
8
10
```

## forward

场景：一对一的数据分发，map、flatMap、filter 等都是这种分区策略

DataStream → DataStream

上游分区数据分发到下游对应分区中

partition1->partition1

partition2->partition2

注意：必须保证上下游分区数（并行度）一致，不然会有如下异常：

```
Forward partitioning does not allow change of parallelism
* Upstream operation: Source: Sequence Source-1 parallelism: 2,
* downstream operation: Sink: Unnamed-4 parallelism: 4
* stream.forward.writeAsText("./data/stream2").setParallelism(4)
```

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(2)
stream.writeAsText("./data/stream1").setParallelism(2)
stream.forward.writeAsText("./data/stream2").setParallelism(2)
env.execute()
```

console result: stream1:1->stream2:1、stream1:2->stream2:2

stream1:1

```
1
3
5
7
9
```

stream1:2

```
2
4
6
8
10
```

stream2:1

```
1
3
5
7
9
```

stream2:2

```
2
4
6
8
10
```

## keyBy

场景：与业务场景匹配

DataStream → DataStream

根据上游分区元素的Hash值与下游分区数取模计算出，将当前元素分发到下游哪一个分区

```
MathUtils.murmurHash(keyHash)（每个元素的Hash值） % maxParallelism（下游分区数）
```

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.generateSequence(1,10).setParallelism(2)
stream.writeAsText("./data/stream1").setParallelism(2)
stream.keyBy(0).writeAsText("./data/stream2").setParallelism(2)
env.execute()
```

console result：根据元素Hash值分发到下游分区中

## PartitionCustom

DataStream → DataStream

通过自定义的分区器，来决定元素是如何从上游分区分发到下游分区

```
object ShuffleOperator {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setParallelism(2)
    val stream = env.generateSequence(1,10).map(_ , 1))
    stream.writeAsText("./data/stream1")
    stream.partitionCustom(new customPartitioner(),0)
      .writeAsText("./data/stream2").setParallelism(4)
    env.execute()
  }
  class customPartitioner extends Partitioner[Long]{
    override def partition(key: Long, numPartitions: Int): Int = {
      key.toInt % numPartitions
    }
  }
}
```

## Dataflows Sink

Flink内置了大量sink，可以将Flink处理后的数据输出到HDFS、kafka、Redis、ES、MySQL等等工程场景中，会经常消费kafka中数据，处理结果存储到Redis或者MySQL中

### Redis Sink

Flink处理的数据可以存储到Redis中，以便实时查询

Flink内嵌连接Redis的连接器，只需要导入连接Redis的依赖就可以

```
<dependency>
  <groupId>org.apache.bahir</groupId>
  <artifactId>flink-connector-redis_2.11</artifactId>
  <version>1.0</version>
</dependency>
```

WordCount写入到Redis中，选择的是HSET数据类型

代码如下：

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
val stream = env.socketTextStream("node01",8888)
val result = stream.flatMap(_ .split(" "))
  .map(_ , 1))
  .keyBy(0)
  .sum(1)

//若redis是单机
val config = new
FlinkJedisPoolConfig.Builder().setDatabase(3).setHost("node01").setPort(6379).build()

//如果是 redis集群
/*val addresses = new util.HashSet[InetSocketAddress]()
addresses.add(new InetSocketAddress("node01",6379))
addresses.add(new InetSocketAddress("node01",6379))
val clusterConfig = new
FlinkJedisClusterConfig.Builder().setNodes(addresses).build()*/
```

```

    result.addSink(new RedisSink[(String,Int)](config,new
RedisMapper[(String,Int)] {

    override def getCommandDescription: RedisCommandDescription = {
        new RedisCommandDescription(RedisCommand.HSET,"wc")
    }

    override def getKeyFromData(t: (String, Int)) = {
        t._1
    }

    override def getValueFromData(t: (String, Int)) = {
        t._2 + ""
    }
}))
env.execute()

```

## Kafka Sink

处理结果写入到kafka topic中，Flink也是默认支持，需要添加连接器依赖，跟读取kafka数据用的连接器依赖相同

之前添加过就不需要再次添加了

```

<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-connector-kafka_2.11</artifactId>
  <version>${flink-version}</version>
</dependency>

```

```

import java.lang
import java.util.Properties

import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.connectors.kafka.{FlinkKafkaProducer,
kafkaSerializationSchema}
import org.apache.kafka.clients.producer.ProducerRecord
import org.apache.kafka.common.serialization.StringSerializer

object Kafkasink {
  def main(args: Array[String]): Unit = {

    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val stream = env.socketTextStream("node01",8888)
    val result = stream.flatMap(_.split(" "))
      .map((_, 1))
      .keyBy(0)
      .sum(1)

    val props = new Properties()
    props.setProperty("bootstrap.servers","node01:9092,node02:9092,node03:9092")
    // props.setProperty("key.serializer",classOf[StringSerializer].getName)
    // props.setProperty("value.serializer",classOf[StringSerializer].getName)

```

```

    /**
    public FlinkKafkaProducer(
        FlinkKafkaProducer(defaultTopic: String, serializationSchema:
        KafkaSerializationSchema[IN], producerConfig: Properties, semantic:
        FlinkKafkaProducer.Semantic)
    */
    result.addSink(new FlinkKafkaProducer[(String,Int)]("wc",new
    KafkaSerializationSchema[(String, Int)] {
        override def serialize(element: (String, Int), timestamp: lang.Long):
        ProducerRecord[Array[Byte], Array[Byte]] = {
            new ProducerRecord("wc",element._1.getBytes(),
            (element._2+"").getBytes())
        }
    },props,FlinkKafkaProducer.Semantic.EXACTLY_ONCE))

    env.execute()
}
}

```

## MySQL Sink (幂等性)

Flink处理结果写入到MySQL中，这并不是Flink默认支持的，需要添加MySQL的驱动依赖

```

<dependency>
  <groupId>mysql</groupId>
  <artifactId>mysql-connector-java</artifactId>
  <version>5.1.44</version>
</dependency>

```

因为不是内嵌支持的，所以需要基于RichSinkFunction自定义sink

不要基于SinkFunction自定义sink why? 看源码

消费kafka中数据，统计各个卡口的流量，并且存入到MySQL中

注意点：需要去重，操作MySQL需要幂等性

```

import java.sql.{Connection, DriverManager, PreparedStatement}
import java.util.Properties

import org.apache.flink.api.common.functions.ReduceFunction
import org.apache.flink.api.common.typeinfo.TypeInformation
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.functions.sink.{RichSinkFunction, SinkFunction}
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.connectors.kafka.{FlinkKafkaConsumer, KafkaDeserializationSchema}
import org.apache.kafka.clients.consumer.ConsumerRecord
import org.apache.kafka.common.serialization.StringSerializer

object MySQLSink {

  case class CarInfo(monitorId: String, carId: String, eventTime: String, Speed: Long)

```

```

def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
"node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    //第一个参数 : 消费的topic名
    val stream = env.addSource(new FlinkKafkaConsumer[(String, String)]("flink-kafka", new KafkaDeserializationSchema[(String, String)] {
        //什么时候停止, 停止条件是什么
        override def isEndOfStream(t: (String, String)): Boolean = false

        //要进行序列化的字节流
        override def deserialize(consumerRecord: ConsumerRecord[Array[Byte], Array[Byte]]): (String, String) = {
            val key = new String(consumerRecord.key(), "UTF-8")
            val value = new String(consumerRecord.value(), "UTF-8")
            (key, value)
        }
    })

    //指定一下返回的数据类型 Flink提供的类型
    override def getProducedType: TypeInformation[(String, String)] = {
        createTuple2TypeInformation(createTypeInformation[String],
createTypeInformation[String])
    }
    }, props))

    stream.map(data => {
        val value = data._2
        val splits = value.split("\t")
        val monitorId = splits(0)
        (monitorId, 1)
    }).keyBy(_._1)
        .reduce(new ReduceFunction[(String, Int)] {
            //t1:上次聚合完的结果 t2:当前的数据
            override def reduce(t1: (String, Int), t2: (String, Int)): (String, Int)
= {
                (t1._1, t1._2 + t2._2)
            }
        }).addSink(new MySQLCustomSink)

    env.execute()
}

//幂等性写入外部数据库MySQL
class MySQLCustomSink extends RichSinkFunction[(String, Int)] {
    var conn: Connection = _
    var insertPst: PreparedStatement = _
    var updatePst: PreparedStatement = _

```



```

//每来一个元素都会调用一次
override def invoke(value: (String, Int), context: SinkFunction.Context[_]):
Unit = {
    println(value)
    updatePst.setInt(1, value._2)
    updatePst.setString(2, value._1)
    updatePst.execute()
    println(updatePst.getUpdateCount)
    if(updatePst.getUpdateCount == 0){
        println("insert")
        insertPst.setString(1, value._1)
        insertPst.setInt(2, value._2)
        insertPst.execute()
    }
}

//thread初始化的时候执行一次
override def open(parameters: Configuration): Unit = {
    conn = DriverManager.getConnection("jdbc:mysql://node01:3306/test",
"root", "123123")
    insertPst = conn.prepareStatement("INSERT INTO car_flow(monitorId,count)
VALUES(?,?)")
    updatePst = conn.prepareStatement("UPDATE car_flow SET count = ? WHERE
monitorId = ?")
}

//thread关闭的时候 执行一次
override def close(): Unit = {
    insertPst.close()
    updatePst.close()
    conn.close()
}
}
}

```

## Socket Sink

Flink处理结果发送到套接字 (Socket)

基于RichSinkFunction自定义sink

```

import java.io.PrintStream
import java.net.{InetAddress, Socket}
import java.util.Properties

import org.apache.flink.api.common.functions.ReduceFunction
import org.apache.flink.api.common.typeinfo.TypeInformation
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.functions.sink.{RichSinkFunction,
SinkFunction}
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment,
createTuple2TypeInformation, createTypeInformation}
import org.apache.flink.streaming.connectors.kafka.{FlinkKafkaConsumer,
KafkaDeserializationSchema}
import org.apache.kafka.clients.consumer.ConsumerRecord
import org.apache.kafka.common.serialization.StringSerializer

```

```

//sink 到 套接字 socket
object SocketsSink {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
"node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    //第一个参数 : 消费的topic名
    val stream = env.addSource(new FlinkKafkaConsumer[(String, String)]("flink-kafka", new KafkaDeserializationSchema[(String, String)] {
      //什么时候停止, 停止条件是什么
      override def isEndOfStream(t: (String, String)): Boolean = false

      //要进行序列化的字节流
      override def deserialize(consumerRecord: ConsumerRecord[Array[Byte], Array[Byte]]): (String, String) = {
        val key = new String(consumerRecord.key(), "UTF-8")
        val value = new String(consumerRecord.value(), "UTF-8")
        (key, value)
      }
    })

    //指定一下返回的数据类型 Flink提供的类型
    override def getProducedType: TypeInformation[(String, String)] = {
      createTuple2TypeInformation(createTypeInformation[String],
createTypeInformation[String])
    }
  }, props))

  stream.map(data => {
    val value = data._2
    val splits = value.split("\t")
    val monitorId = splits(0)
    (monitorId, 1)
  }).keyBy(_._1)
  .reduce(new ReduceFunction[(String, Int)] {
    //t1:上次聚合完的结果 t2:当前的数据
    override def reduce(t1: (String, Int), t2: (String, Int)): (String, Int)
= {
      (t1._1, t1._2 + t2._2)
    }
  }).addSink(new SocketCustomSink("node01", 8888))

  env.execute()
}

class SocketCustomSink(host:String,port:Int) extends
RichSinkFunction[(String,Int)]{
  var socket: Socket = _
  var writer:PrintStream = _

```

```

    override def open(parameters: Configuration): Unit = {
        socket = new Socket(InetAddress.getByName(host), port)
        writer = new PrintStream(socket.getOutputStream)
    }

    override def invoke(value: (String, Int), context: SinkFunction.Context[_]):
Unit = {
        writer.println(value._1 + "\t" + value._2)
        writer.flush()
    }

    override def close(): Unit = {
        writer.close()
        socket.close()
    }
}

```

## File Sink

Flink处理的结果保存到文件，这种使用方式不是很常见

支持分桶写入，每一个桶就是一个目录，默认每隔一个小时会产生一个分桶，每个桶下面会存储每一个Thread的处理结果，可以设置一些文件滚动的策略（文件打开、文件大小等），防止出现大量的小文件，代码中详解

Flink默认支持，导入连接文件的连接器依赖

```

<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-connector-fileSystem_2.11</artifactId>
  <version>1.9.2</version>
</dependency>

```

```

import org.apache.flink.api.common.functions.ReduceFunction
import org.apache.flink.api.common.serialization.SimpleStringEncoder
import org.apache.flink.api.common.typeinfo.TypeInformation
import org.apache.flink.core.fs.Path
import org.apache.flink.streaming.api.functions.sink.filesystem.StreamingFileSink
import org.apache.flink.streaming.api.functions.sink.filesystem.rollingpolicies.DefaultRollingPolicy
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, createTuple2TypeInformation, createTypeInformation}
import org.apache.flink.streaming.connectors.kafka.{FlinkKafkaConsumer, KafkaDeserializationSchema}
import org.apache.kafka.clients.consumer.ConsumerRecord
import org.apache.kafka.common.serialization.StringSerializer

object FileSink {
    def main(args: Array[String]): Unit = {
        val env = StreamExecutionEnvironment.getExecutionEnvironment

        //设置连接kafka的配置信息
        val props = new Properties()

```

```

//注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
props.setProperty("bootstrap.servers",
"node01:9092,node02:9092,node03:9092")
props.setProperty("group.id", "flink-kafka-001")
props.setProperty("key.deserializer", classOf[StringSerializer].getName)
props.setProperty("value.deserializer", classOf[StringSerializer].getName)

//第一个参数 : 消费的topic名
val stream = env.addSource(new FlinkKafkaConsumer[(String, String)]("flink-kafka", new KafkaDeserializationSchema[(String, String)] {
    //什么时候停止, 停止条件是什么
    override def isEndOfStream(t: (String, String)): Boolean = false

    //要进行序列化的字节流
    override def deserialize(consumerRecord: ConsumerRecord[Array[Byte], Array[Byte]]): (String, String) = {
        val key = new String(consumerRecord.key(), "UTF-8")
        val value = new String(consumerRecord.value(), "UTF-8")
        (key, value)
    }

    //指定一下返回的数据类型 Flink提供的类型
    override def getProducedType: TypeInformation[(String, String)] = {
        createTuple2TypeInformation(createTypeInformation[String], createTypeInformation[String])
    }
}, props))

val restStream = stream.map(data => {
    val value = data._2
    val splits = value.split("\t")
    val monitorId = splits(0)
    (monitorId, 1)
}).keyBy(_._1)
    .reduce(new ReduceFunction[(String, Int)] {
        //t1:上次聚合完的结果 t2:当前的数据
        override def reduce(t1: (String, Int), t2: (String, Int)): (String, Int)
    = {
        (t1._1, t1._2 + t2._2)
    }
}).map(x=>x._1 + "\t" + x._2)

//设置文件滚动策略
val rolling:DefaultRollingPolicy[String,String] =
DefaultRollingPolicy.create()
    //当文件超过2s没有写入新数据, 则滚动产生一个小文件
    .withInactivityInterval(2000)
    //文件打开时间超过2s 则滚动产生一个小文件 每隔2s产生一个小文件
    .withRolloverInterval(2000)
    //当文件大小超过256 则滚动产生一个小文件
    .withMaxPartSize(256*1024*1024)
    .build()

/**
 * 默认:
 * 每一个小时对应一个桶(文件夹), 每一个thread处理的结果对应桶下面的一个小文件
 * 当小文件大小超过128M或者小文件打开时间超过60s, 滚动产生第二个小文件

```

```

    */
    val sink: StreamingFileSink[String] = StreamingFileSink.forRowFormat(
        new Path("d:/data/rests"),
        new SimpleStringEncoder[String]("UTF-8"))
        .withBucketCheckInterval(1000)
        .withRollingPolicy(rolling)
        .build()

    //    val sink = StreamingFileSink.forBulkFormat(
    //        new Path("./data/rest"),
    //        ParquetAvroWriters.forSpecificRecord(classOf[String]))
    //    ).build()

    restStream.addSink(sink)
    env.execute()
}
}

```

## HBase Sink

计算结果写入sink 两种实现方式:

1. map算子写入 频繁创建hbase连接
2. process写入 适合批量写入hbase

导入HBase依赖包

```

<dependency>
    <groupId>org.apache.hbase</groupId>
    <artifactId>hbase-client</artifactId>
    <version>${hbase.version}</version>
</dependency>
<dependency>
    <groupId>org.apache.hbase</groupId>
    <artifactId>hbase-common</artifactId>
    <version>${hbase.version}</version>
</dependency>
<dependency>
    <groupId>org.apache.hbase</groupId>
    <artifactId>hbase-server</artifactId>
    <version>${hbase.version}</version>
</dependency>

```

读取kafka数据, 统计卡口流量保存至HBase数据库中

1. HBase中创建对应的表

```
create 'car_flow',{NAME => 'count', VERSIONS => 1}
```

2. 实现代码

```

import java.util.{Date, Properties}

import com.msb.stream.util.{DateUtils, HBaseUtil}
import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.configuration.Configuration

```

```

import org.apache.flink.streaming.api.functions.ProcessFunction
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.flink.util.Collector
import org.apache.hadoop.hbase.HBaseConfiguration
import org.apache.hadoop.hbase.client.{HTable, Put}
import org.apache.hadoop.hbase.util.Bytes
import org.apache.kafka.common.serialization.StringSerializer

object HBaseSinkTest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
      "node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new
      SimpleStringSchema(), props))

    stream.map(row => {
      val arr = row.split("\t")
      (arr(0), 1)
    }).keyBy(_._1)
      .reduce((v1: (String, Int), v2: (String, Int)) => {
        (v1._1, v1._2 + v2._2)
      }).process(new ProcessFunction[(String, Int), (String, Int)] {

        var htab: HTable = _

        override def open(parameters: Configuration): Unit = {
          val conf = HBaseConfiguration.create()
          conf.set("hbase.zookeeper.quorum",
            "node01:2181,node02:2181,node03:2181")
          val hbaseName = "car_flow"
          htab = new HTable(conf, hbaseName)
        }

        override def close(): Unit = {
          htab.close()
        }

        override def processElement(value: (String, Int), ctx:
          ProcessFunction[(String, Int), (String, Int)]#Context, out: Collector[(String,
            Int)]): Unit = {
          // rowkey:monitorid 时间戳(分钟) value: 车流量
          val min = DateUtils.getMin(new Date())
          val put = new Put(Bytes.toBytes(value._1))

```

```

        put.addColumn(Bytes.toBytes("count"), Bytes.toBytes(min),
Bytes.toBytes(value._2))
        htab.put(put)
    }
})
env.execute()
}
}

```

## Flink State状态

Flink是一个有状态的流式计算引擎，所以会将中间计算结果(状态)进行保存，默认保存到TaskManager的堆内存中，但是当task挂掉，那么这个task所对应的状态都会被清空，造成了数据丢失，无法保证结果的正确性，哪怕想要得到正确结果，所有数据都要重新计算一遍，效率很低。想要保证At-least-once和Exactly-once，需要把数据状态持久化到更安全的存储介质中，Flink提供了堆内内存、堆外内存、HDFS、RocksDB等存储介质

先来看下Flink提供的状态有哪些？

Flink中状态分为两种类型

- Keyed State

基于KeyedStream上的状态，这个状态是跟特定的Key绑定，KeyedStream流上的每一个Key都对应一个State，每一个Operator可以启动多个Thread处理，但是相同Key的数据只能由同一个Thread处理，因此一个Keyed状态只能存在于某一个Thread中，一个Thread会有多个Keyed state

- Non-Keyed State (Operator State)

Operator State与Key无关，而是与Operator绑定，整个Operator只对应一个State。比如：Flink中的Kafka Connector就使用了Operator State，它会在每个Connector实例中，保存该实例消费Topic的所有(partition, offset)映射

Flink针对Keyed State提供了以下可以保存State的数据结构

- ValueState:类型为T的单值状态，这个状态与对应的Key绑定，最简单的状态，通过update更新值，通过value获取状态值
- ListState: Key上的状态值为一个列表，这个列表可以通过add方法往列表中添加值，也可以通过get()方法返回一个Iterable来遍历状态值
- ReducingState: 每次调用add()方法添加值的时候，会调用用户传入的reduceFunction，最后合并到一个单一的状态值
- MapState<UK, UV>:状态值为一个Map，用户通过put或putAll方法添加元素，get(key)通过指定的key获取value，使用entries()、keys()、values()检索
- AggregatingState <IN, OUT>:保留一个单值，表示添加到状态的所有值的聚合。和ReducingState相反的是，聚合类型可能与添加到状态的元素类型不同。使用 add(IN) 添加的元素会调用用户指定的 AggregateFunction 进行聚合
- FoldingState<T, ACC>:已过时建议使用AggregatingState 保留一个单值，表示添加到状态的所有值的聚合。与 ReducingState 相反，聚合类型可能与添加到状态的元素类型不同。使用 add(T) 添加的元素会调用用户指定的 FoldFunction 折叠成聚合值

案例1：使用ValueState keyed state检查车辆是否发生了急加速

```

object ValueStateTest {

    case class CarInfo(carId: String, speed: Long)

```

```

def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val stream = env.socketTextStream("node01", 8888)
    stream.map(data => {
        val arr = data.split(" ")
        CarInfo(arr(0), arr(1).toLong)
    }).keyBy(_.carId)
    .map(new RichMapFunction[CarInfo, String]() {

        //保存上一次车速
        private var lastTempState: ValueState[Long] = _

        override def open(parameters: Configuration): Unit = {
            val lastTempStateDesc = new ValueStateDescriptor[Long]
            ("lastTempState", createTypeInfo[Long])
            lastTempState = getRuntimeContext.getState(lastTempStateDesc)
        }

        override def map(value: CarInfo): String = {
            val lastSpeed = lastTempState.value()
            this.lastTempState.update(value.speed)
            if ((value.speed - lastSpeed).abs > 30 && lastSpeed != 0)
                "over speed" + value.toString
            else
                value.carId
        }
    }).print()
    env.execute()
}

```

案例2: 使用MapState 统计单词出现次数 仅供大家理解MapState

```

import org.apache.flink.api.common.functions.RichMapFunction
import org.apache.flink.api.common.state.{MapState, MapStateDescriptor}
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._

//MapState 实现WordCount
object KeyedStateTest {
    def main(args: Array[String]): Unit = {
        val env = StreamExecutionEnvironment.getExecutionEnvironment
        val stream = env.fromCollection(List("I love you", "hello spark", "hello
flink", "hello hadoop"))
        val pairStream = stream.flatMap(_.split(" ")).map((_, 1)).keyBy(_._1)
        pairStream.map(new RichMapFunction[(String, Int), (String, Int)] {

            private var map: MapState[String, Int] = _
            override def open(parameters: Configuration): Unit = {
                //定义map state存储的数据类型
                val desc = new MapStateDescriptor[String, Int]
                ("sum", createTypeInfo[String], createTypeInfo[Int])
                //注册map state
                map = getRuntimeContext.getMapState(desc)
            }
        })
    }
}

```



```

override def map(value: (String, Int)): (String, Int) = {
  val key = value._1
  val v = value._2
  if(map.contains(key)){
    map.put(key,map.get(key) + 1)
  }else{
    map.put(key,1)
  }
  val iterator = map.keys().iterator()
  while (iterator.hasNext){
    val key = iterator.next()
    println("word:" + key + "\t count:" + map.get(key))
  }
  value
}
}).setParallelism(3)
env.execute()
}
}

```

案例3：使用ReducingState统计每辆车的速度总和

```

import com.msb.state.ValueStateTest.CarInfo
import org.apache.flink.api.common.functions.{ReduceFunction, RichMapFunction}
import org.apache.flink.api.common.state.{ReducingState,
ReducingStateDescriptor}
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._

//统计每辆车的速度总和
object ReduceStateTest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val stream = env.socketTextStream("node01", 8888)
    stream.map(data => {
      val arr = data.split(" ")
      CarInfo(arr(0), arr(1).toLong)
    }).keyBy(_._1.carId)
      .map(new RichMapFunction[CarInfo, CarInfo] {
        private var reduceState: ReducingState[Long] = _

        override def map(elem: CarInfo): CarInfo = {
          reduceState.add(elem.speed)
          println("carId:" + elem.carId + " speed count:" + reduceState.get())
          elem
        }

        override def open(parameters: Configuration): Unit = {
          val reduceDesc = new ReducingStateDescriptor[Long]("reduceSpeed", new
ReduceFunction[Long] {
            override def reduce(value1: Long, value2: Long): Long = value1 +
value2
          }, createTypeInfo[Long])
          reduceState = getRuntimeContext.getReducingState(reduceDesc)
        }
      })
  }
}

```

```

    }
  })
  env.execute()
}
}

```

案例4：使用AggregatingState统计每辆车的速度总和

```

import com.msb.state.ValueStateTest.CarInfo
import org.apache.flink.api.common.functions.{AggregateFunction, ReduceFunction,
RichMapFunction}
import org.apache.flink.api.common.state.{AggregatingState,
AggregatingStateDescriptor, ReducingState, ReducingStateDescriptor}
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._

//统计每辆车的速度总和
object ReduceStateTest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    val stream = env.socketTextStream("node01", 8888)
    stream.map(data => {
      val arr = data.split(" ")
      CarInfo(arr(0), arr(1).toLong)
    }).keyBy(_.carId)
      .map(new RichMapFunction[CarInfo, CarInfo] {
        private var aggState: AggregatingState[Long, Long] = _

        override def map(elem: CarInfo): CarInfo = {
          aggState.add(elem.speed)
          println("carId:" + elem.carId + " speed count:" + aggState.get())
          elem
        }

        override def open(parameters: Configuration): Unit = {
          val aggDesc = new AggregatingStateDescriptor[Long, Long, Long]("agg", new
AggregateFunction[Long, Long, Long] {
            //初始化累加器值
            override def createAccumulator(): Long = 0

            //往累加器中累加值
            override def add(value: Long, acc: Long): Long = acc + value

            //返回最终结果
            override def getResult(accumulator: Long): Long = accumulator

            //合并两个累加器值
            override def merge(a: Long, b: Long): Long = a+b
          }, createTypeInfo[Long])

          aggState = getRuntimeContext.getAggregatingState(aggDesc)
        }
      })
    env.execute()
  }
}

```

案例5：统计每辆车的运行轨迹 所谓运行轨迹就是这辆车的信息 按照时间排序，卡口号串联起来

```
import java.text.SimpleDateFormat
import java.util.Properties

import org.apache.flink.api.common.functions.RichMapFunction
import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.api.common.state.{ListState, ListStateDescriptor}
import org.apache.flink.configuration.Configuration
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.kafka.common.serialization.StringSerializer

import scala.collection.JavaConverters._

/**
 *统计每辆车的运行轨迹
 * 所谓运行轨迹就是这辆车的信息 按照时间排序，卡口号串联起来
 */

object ListStateTest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers", "node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new SimpleStringSchema(), props))
    val format = new SimpleDateFormat("yyyy-MM-dd HH:mm:ss")

    stream.map(data => {
      val arr = data.split("\t")
      //卡口、车牌、事件时间、车速
      val time = format.parse(arr(2)).getTime
      (arr(0), arr(1), time, arr(3).toLong)
    }).keyBy(_._2)
      .map(new RichMapFunction[(String, String, Long, Long), (String, String)] {
        //event-time monitor_id
        private var speedInfos: ListState[(Long, String)] = _

        override def map(elem: (String, String, Long, Long)): (String, String) = {
          speedInfos.add(elem._3, elem._1)
          val infos = speedInfos.get().asScala.seq
          val sortList = infos.toList.sortBy(x => x._1).reverse
          val builder = new StringBuilder
          for (elem <- sortList) {
            builder.append(elem._2 + "\t")
          }
        }
      })
```

```

        (elem._2, builder.toString())
    }

    override def open(parameters: Configuration): Unit = {
        val listStateDesc = new ListStateDescriptor[(Long, String)]
        ("speedInfos", createTypeInfoInformation[(Long, String)])
        speedInfos = getRuntimeContext.getListState(listStateDesc)
    }
    }).print()

    env.execute()
}

```

Operator State绑定每一个Operator

实现CheckpointedFunction接口 来操作算子状态

案例6: 自系统启动以来, 总共处理了多少条数据量

```

import org.apache.flink.api.common.functions.{FlatMapFunction,
RichFlatMapFunction}
import org.apache.flink.api.common.state.{ListState, ListStateDescriptor,
ValueState, ValueStateDescriptor}
import org.apache.flink.runtime.state.{FunctionInitializationContext,
FunctionSnapshotContext}
import org.apache.flink.streaming.api.checkpoint.CheckpointedFunction
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.util.Collector

import scala.collection.JavaConverters._

//统计经过flatMap算子的数据量
object FlatMapOperatorStateTest {
    def main(args: Array[String]): Unit = {
        val env = StreamExecutionEnvironment.getExecutionEnvironment
        val stream = env.fromCollection(List("I love you", "hello spark", "hello
flink", "hello hadoop"))
        stream.map(data => {
            (data, 1)
        }).keyBy(_._1)
            .flatMap(new MyFlatMapFunction())

        env.execute()
    }

    class MyFlatMapFunction extends RichFlatMapFunction[(String, Int),
(String, Int, Int)] with CheckpointedFunction{
        private var operatorCount: Long = _
        private var operatorState: ListState[Long] = _

        override def flatMap(value: (String, Int), out: Collector[(String, Int,
Int)]): Unit = {
            operatorCount += 1
            val subtasks = getRuntimeContext.getTaskNameWithSubtasks
            println(subtasks + "==" + operatorState.get())
        }
    }
}

```

```

//进行checkpoint的时候 会被调用，然后持久化到远端
override def snapshotState(context: FunctionSnapshotContext): Unit = {
    operatorState.clear()
    operatorState.add(operatorCount)
}

//初始化方法
override def initializeState(context: FunctionInitializationContext): Unit =
{
    operatorState = context.getOperatorStateStore.getListState(new
ListStateDescriptor[Long]("operateState",createTypeInfo[Long]))
    if(context.isRestored){
        operatorCount = operatorState.get().asScala.sum
    }
}
}
}

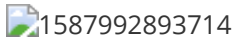
```

## Checkpoint

Flink中基于**异步**轻量级的分布式快照技术提供了Checkpoint容错机制，分布式快照可以将同一时间点Task/Operator的状态数据全局统一快照处理，包括上面提到的用户自定义使用的Keyed State和Operator State，当未来程序出现问题，可以基于保存的快照容错

### Checkpoint原理

Flink会在输入的数据集上间隔性地生成checkpoint barrier，通过栅栏（barrier）将间隔时间段内的数据划分到相应的checkpoint中。当程序出现异常时，Operator就能够从上一次快照中恢复所有算子之前的状态，从而保证数据的一致性。例如在KafkaConsumer算子中维护offset状态，当系统出现问题无法从Kafka中消费数据时，可以将offset记录在状态中，当任务重新恢复时就能够从指定的偏移量开始消费数据。



默认情况Flink不开启检查点，用户需要在程序中通过调用方法配置和开启检查点，另外还可以调整其他相关参数

- Checkpoint开启和时间间隔指定

开启检查点并且指定检查点时间间隔为1000ms，根据实际情况自行选择，如果状态比较大，则建议适当增加该值

```
env.enableCheckpointing(1000)
```

- exactly-ance和at-least-once语义选择

选择exactly-once语义保证整个应用内端到端的数据一致性，这种情况比较适合于数据要求比较高，不允许出现丢数据或者数据重复，与此同时，Flink的性能也相对较弱，而at-least-once语义更适合于时延和吞吐量要求非常高但对数据的一致性要求不高的场景。如下通过setCheckpointingMode()方法来设定语义模式，默认情况下使用的是exactly-once模式

```
env.getCheckpointConfig.setCheckpointingMode(CheckpointingMode.EXACTLY_ONCE)
```

- Checkpoint超时时间

超时时间指定了每次Checkpoint执行过程中的上限时间范围，一旦Checkpoint执行时间超过该阈值，Flink将会中断Checkpoint过程，并按照超时处理。该指标可以通过setCheckpointTimeout方法设定，默认为10分钟

```
env.getCheckpointConfig.setCheckpointTimeout(5 * 60 * 1000)
```

- Checkpoint之间最小时间间隔

该参数主要目的是设定两个Checkpoint之间的最小时间间隔，防止Flink应用密集地触发Checkpoint操作，会占用了大量计算资源而影响到整个应用的性能

```
env.getCheckpointConfig.setMinPauseBetweenCheckpoints(600)
```

- 最大并行执行的Checkpoint数量

在默认情况下只有一个检查点可以运行，根据用户指定的数量可以同时触发多个Checkpoint，进而提升Checkpoint整体的效率

```
env.getCheckpointConfig.setMaxConcurrentCheckpoints(1)
```

- 任务取消后，是否删除Checkpoint中保存的数据

设置为RETAIN\_ON\_CANCELLATION：表示一旦Flink处理程序被cancel后，会保留CheckPoint数据，以便根据实际需要恢复到指定的CheckPoint

设置为DELETE\_ON\_CANCELLATION：表示一旦Flink处理程序被cancel后，会删除CheckPoint数据，只有Job执行失败的时候才会保存CheckPoint

```
env.getCheckpointConfig.enableExternalizedCheckpoints(ExternalizedCheckpointCleanup.RETAIN_ON_CANCELLATION)
```

- 容忍的检查的失败数

设置可以容忍的检查的失败数，超过这个数量则系统自动关闭和停止任务

```
env.getCheckpointConfig.setTolerableCheckpointFailureNumber(1)
```

checkpoint测试：

1. 提交job
2. 取消job
3. 基于checkpoint数据 重启job

```
flink run -c com.msb.state.wordCountCheckpoint -s  
hdfs://node01:9000/flink/sasa/savepoint-917081-0a251a5323b7 ~/StudyFlink-  
1.0-SNAPSHOT.jar
```

如果任务的取消是在第一次checkpoint与第二次checkpoint之间，那么会存在数据的丢失，因为socket是不支持数据回放，如果读取的是kafka 默认支持数据回放

## SavePoint原理

Savepoints 是检查点的一种特殊实现，底层实现其实也是使用Checkpoints的机制。Savepoints是用户以手工命令的方式触发Checkpoint,并将结果持久化到指定的存储路径中，其主要目的是帮助用户在升级和维护集群过程中保存系统中的状态数据，避免因停机运维或者升级应用等正常终止应用的操作而导致系统无法恢复到原有的计算状态的情况，从而无法实现从端到端的 Exactly-Once 语义保证

SavePoint的路径需要在flink-conf.yaml中配置

```
state.savepoints.dir: hdfs://node01:9000/flink/state/savepoint
```

系统的升级顺序

### 1. 先savepoint

```
flink savepoint 91708180bc440568f47ab0ec88087b43
hdfs://node01:9000/flink/sasa
如果在flink-conf.yaml中没有设置SavePoint的路径，可以在进行SavePoint的时候指定路径
```

### 2. cancel job

```
flink cancel 91708180bc440568f47ab0ec88087b43 //job id
```

### 3. 重启job

```
flink run -c com.msb.state.WordCountCheckpoint -s
hdfs://node01:9000/flink/sasa/savepoint-917081-0a251a5323b7 ~/StudyFlink-1.0-SNAPSHOT.jar
```

最佳实战：

为了能够在作业的不同版本之间以及Flink的不同版本之间顺利升级，强烈推荐程序员通过手动给算子赋予ID，这些ID将用于确定每一个算子的状态范围。如果不手动给各算子指定ID，则会由Flink自动给每个算子生成一个ID。而这些自动生成的ID依赖于程序的结构，并且对代码的更改是很敏感的。因此，强烈建议用户手动设置ID

```
stream.flatMap(data => {
    val rest = new ListBuffer[(String, Int)]
    val words = data.split(" ")
    for (word <- words) {
        rest += ((word, 1))
    }
    rest
}).uid("001").keyBy(_._1)
    .reduce((v1: (String, Int), v2: (String, Int)) => {
        (v1._1, v1._2 + v2._2)
    }).uid("002").print()
```

打jar包执行、SavePoint、Cancel job

```
stream.map(_._1).uid("001").keyBy(_._1)
    .reduce((v1: (String, Int), v2: (String, Int)) => {
        (v1._1, v1._2 + v2._2)
    }).uid("002").map(data => {
        println(data + "-savepoint")
    })
```

```
stream.flatMap(_._split(" "))
  .map(_._1, 1))
  .keyBy(_._1)
  .reduce((v1: (String, Int), v2: (String, Int)) => {
    (v1._1, v1._2 + v2._2)
  }).uid("reduce")
  .map(x => {
    println(x + "---savepoint")
    x
  })
  .print()
```

打jar包，提交job（指定SavePoint路径） 根据上次savepoint的各个算子id的状态来恢复

## StateBackend状态后端

在Flink中提供了StateBackend来存储和管理状态数据

Flink一共实现了三种类型的状态管理器：MemoryStateBackend、FsStateBackend、RocksDBStateBackend

### MemoryStateBackend

基于内存的状态管理器将状态数据全部存储在JVM堆内存中。基于内存的状态管理具有非常快速和高效的特点，但也具有非常多的限制，最主要的就是内存的容量限制，一旦存储的状态数据过多就会导致系统内存溢出等问题，从而影响整个应用的正常运行。同时如果机器出现问题，整个主机内存中的状态数据都会丢失，进而无法恢复任务中的状态数据。因此从数据安全的角度建议用户尽可能地避免在生产环境中使用MemoryStateBackend

Flink将MemoryStateBackend作为默认状态后端管理器

```
env.setStateBackend(new MemoryStateBackend(100*1024*1024))
```

注意：聚合类算子的状态会同步到JobManager内存中，因此对于聚合类算子较多的应用会对JobManager的内存造成一定的压力，进而影响集群

### FsStateBackend

和MemoryStateBackend有所不同，FsStateBackend是基于文件系统的一种状态管理器，这里的文件系统可以是本地文件系统，也可以是HDFS分布式文件系统

```
env.setStateBackend(new FsStateBackend("path", true))
```

如果path是本地文件路径，其格式：file:///

如果path是HDFS文件路径，格式为：hdfs://

第二个参数代表是否异步保存状态数据到HDFS，异步方式能够尽可能避免checkpoint的过程中影响流式计算任务。

FsStateBackend更适合任务量比较大的应用，例如：包含了时间范围非常长的窗口计算，或者状态比较大的场景

### RocksDBStateBackend



RocksDBStateBackend是Flink中内置的第三方状态管理器，和前面的状态管理器不同，RocksDBStateBackend需要单独引入相关的依赖包到工程中

```
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-statebackend-rocksdb_2.11</artifactId>
  <version>1.9.2</version>
</dependency>
```

```
env.setStateBackend(new RocksDBStateBackend("hdfs://"))
```

RocksDBStateBackend采用异步的方式进行状态数据的Snapshot，任务中的状态数据首先被写入本地RockDB中，这样在RockDB仅会存储正在进行计算的热数据，而需要进行CheckPoint的时候，会把本地的数据直接复制到远端的FileSystem中。

与FsStateBackend相比，RocksDBStateBackend在性能上要比FsStateBackend高一些，主要是因为借助于RocksDB在本地存储了最新热数据，然后通过异步的方式再同步到文件系统中，但RocksDBStateBackend和MemoryStateBackend相比性能就会较弱一些。RocksDB克服了State受内存限制的缺点，同时又能够持久化到远端文件系统中，推荐在生产中使用

## 集群级配置StateBackend

全局配置需要需改集群中的配置文件，修改flink-conf.yaml

- 配置FsStateBackend

```
state.backend: filesystem
state.checkpoints.dir: hdfs://namenode-host:port/flink-checkpoints
```

- ☐ FsStateBackend:filesystem
- ☐ MemoryStateBackend:jobmanager
- ☐ RocksDBStateBackend:rocksdb

- 配置MemoryStateBackend

```
state.backend: jobmanager
```

- 配置RocksDBStateBackend

```
state.backend.rocksdb.checkpoint.transfer.thread.num: 1 同时操作RocksDB的线程数
state.backend.rocksdb.localdir: 本地path RocksDB存储状态数据的本地文件路径
```

## Flink Window操作

Flink任务Batch是Streaming的一个特例，因此Flink底层引擎是一个流式引擎，在上面实现了流处理和批处理。而Window就是从Streaming到Batch的桥梁

Window窗口就在一个无界流中设置起始位置和终止位置，让无界流变成有界流，并且在有界流中进行数据处理

Window操作常见的业务场景：统计过去一段时间、最近一些元素的数据指标

## Window窗口分类

Window窗口在无界流中设置起始位置和终止位置的方式可以有两种：

- 根据时间设置
- 根据窗口数据量 (count) 设置

根据窗口的类型划分：

- 滚动窗口
- 滑动窗口

根据数据流类型划分：

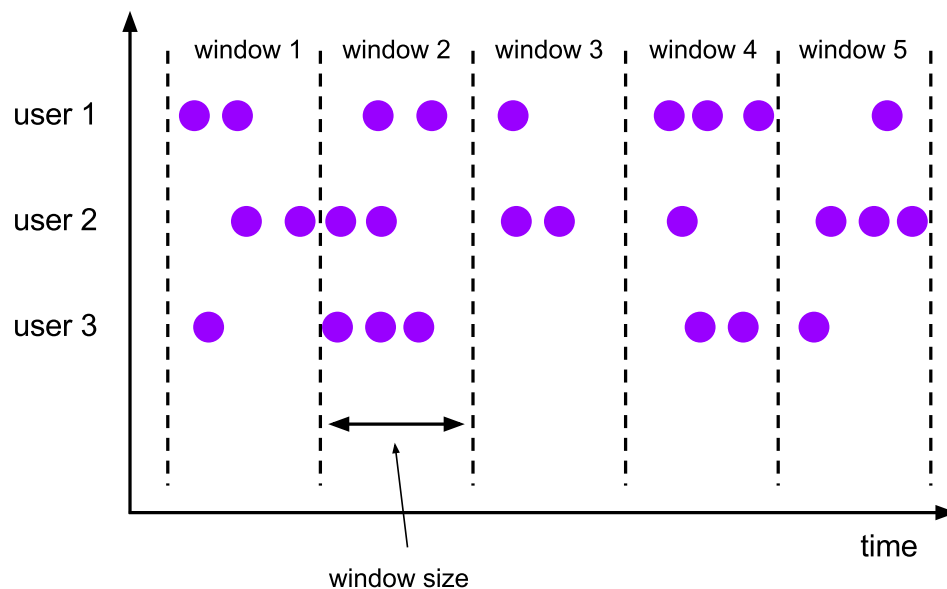
- Keyed Window：基于分组后的数据流之上做窗口操作
- Global Window：基于未分组的数据流之上做窗口操作

根据不同的组合方式，可以组合出来8种窗口类型：

1. 基于分组后的数据流上的时间滚动窗口
2. 基于分组后的数据流上的时间滑动窗口
3. 基于分组后的数据流上的count滚动窗口
4. 基于分组后的数据流上的count滑动窗口
5. 基于未分组的数据流上的时间滚动窗口
6. 基于未分组的数据流上的时间滑动窗口
7. 基于未分组的数据流上的count滚动窗口
8. 基于未分组的数据流上的count滑动窗口

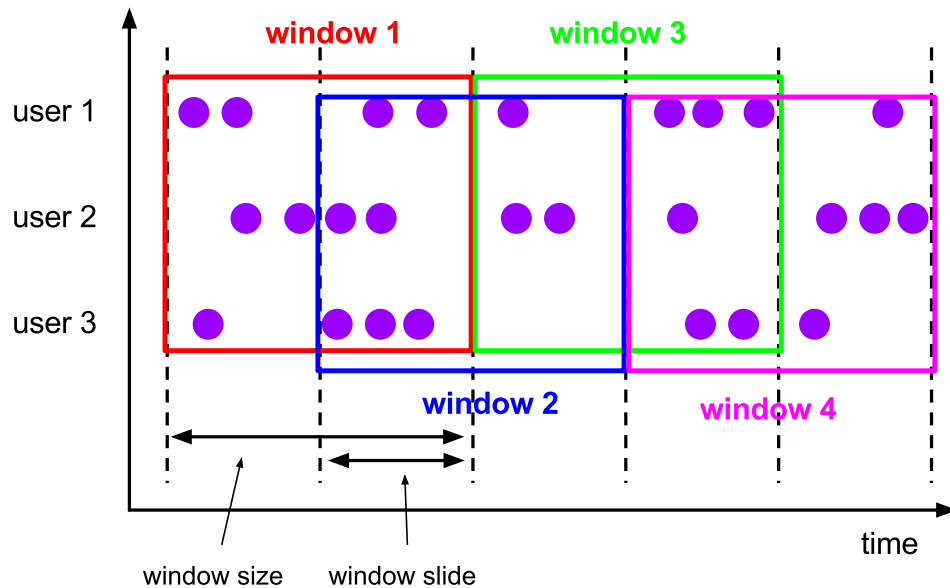
当然我们也可以根据实际业务场景自定义Window，这就是Flink最大的优势：Window种类多，灵活

- Time Window（基于时间的窗口）
  - Tumbling Window：滚动窗口，窗口之间没有数据重叠



- Sliding Window：滑动窗口，窗口内的数据有重叠

在定义滑动窗口的时候，不只是为了定义窗口大小，还要定义窗口的滑动间隔时间（每隔多久滑动一次），如果滑动间隔时间=窗口大小=滚动窗口



## 窗口聚合函数

窗口函数定义了针对窗口内元素的计算逻辑，窗口函数大概分为两类：

1. 增量聚合函数，聚合原理：窗口内保存一个中间聚合结果，随着新元素的加入，不断对该值进行更新

这类函数通常非常节省空间 ReduceFunction、AggregateFunction属于增量聚合函数

2. 全量聚合函数，聚合原理：收集窗口内的所有元素，并且在执行的时候对他们进行遍历，这种聚合函数通常需要占用更多的空间（收集一段时间的数据并且保存），但是它可以支持更复杂的逻辑 ProcessWindowFunction、WindowFunction属于全量窗口函数

注意：这两类函数可以组合搭配使用

## 增量聚合函数

案例1：使用增量聚合函数统计最近20s内，各个卡口的车流量

```
import java.util.Properties

import org.apache.flink.api.common.functions.AggregateFunction
import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.streaming.api.scala.function.ProcessWindowFunction
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment,
createTypeInfoInformation}
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.flink.util.Collector
import org.apache.kafka.common.serialization.StringSerializer

/**
 * 使用增量聚合函数统计最近20s内，各个卡口的车流量
 */
object Demo01StatisCarFlow {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
```

```

val props = new Properties()
//注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
props.setProperty("bootstrap.servers",
"node01:9092,node02:9092,node03:9092")
props.setProperty("group.id", "flink-kafka-001")
props.setProperty("key.deserializer", classOf[StringSerializer].getName)
props.setProperty("value.deserializer", classOf[StringSerializer].getName)

val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new
SimpleStringSchema(), props))

//monitorId + "\t").append(carId + "\t").append(timestamp +
"\t").append(speed)
stream.map(data => {
    val arr = data.split("\t")
    val monitorID = arr(0)
    (monitorID, 1)
}).keyBy(_._1)
    .timewindow(Time.seconds(10))
    //      .reduce(new ReduceFunction[(String, Int)] {
    //          override def reduce(value1: (String, Int), value2: (String,
Int)): (String, Int) = {
    //              (value1._1, value1._2 + value2._2)
    //          }
    //      }).print()
    .aggregate(new AggregateFunction[(String, Int), Int, Int] {
        override def createAccumulator(): Int = 0

        override def add(value: (String, Int), acc: Int): Int = acc + value._2

        override def getResult(acc: Int): Int = acc

        override def merge(a: Int, b: Int): Int = a + b
    },
    //      new WindowFunction[Int, (String, Int), String, Timewindow] {
    //          override def apply(key: String, window: Timewindow, input:
Iterable[Int], out: Collector[(String, Int)]): Unit = {
    //              for (elem <- input) {
    //                  out.collect((key, elem))
    //              }
    //          }
    //      }
    new ProcessWindowFunction[Int, (String, Int), String, Timewindow] {
        override def process(key: String, context: Context, elements:
Iterable[Int], out: Collector[(String, Int)]): Unit = {
            for (elem <- elements) {
                out.collect((key,elem))
            }
        }
    }
    ).print()
env.execute()
}
}

```

ProcessWindowFunction、WindowFunction区别在于ProcessWindowFunction可以获取Flink执行的上下文，可以拿到当前的数据更多信息，比如窗口状态、窗口起始与终止时间、当前水印、时间戳等

案例2：每隔10s统计每辆汽车的平均速度

```
import java.util.Properties

import org.apache.flink.api.common.functions.AggregateFunction
import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.kafka.common.serialization.StringSerializer

object Demo03SpeedAVG {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
      "node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new SimpleStringSchema(), props))

    stream.map(data => {
      val splits = data.split("\t")
      (splits(1), splits(3).toInt)
    }).keyBy(_._1)
      .timewindow(Time.seconds(10))
      .aggregate(new AggregateFunction[(String, Int), (String, Int, Int), (String, Double)] {
        override def createAccumulator(): (String, Int, Int) = ("", 0, 0)

        override def add(value: (String, Int), accumulator: (String, Int, Int)): (String, Int, Int) = {
          (value._1, value._2 + accumulator._2, accumulator._3 + 1)
        }

        override def getResult(accumulator: (String, Int, Int)): (String, Double) = {
          (accumulator._1, accumulator._2.toDouble / accumulator._3)
        }

        override def merge(a: (String, Int, Int), b: (String, Int, Int)): (String, Int, Int) = {
          (a._1, a._2 + b._2, a._3 + b._3)
        }
      }).print()
  }
}
```

```

    env.execute()
  }
}

```

## 全量聚合函数

案例3: 每隔10s对窗口内所有汽车的车速进行排序

```

import java.util.Properties

import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, _}
import org.apache.flink.streaming.api.scala.function.ProcessAllWindowFunction
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.flink.util.Collector
import org.apache.kafka.common.serialization.StringSerializer

object Demo02SortSpeed {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment

    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
      "node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new SimpleStringSchema(), props))

    stream.map(data => {
      val splits = data.split("\t")
      (splits(1), splits(3).toInt)
    }).timewindowAll(Time.seconds(10))
    //注意: 想要全局排序并行度需要设置为1
    .process(new ProcessAllWindowFunction[(String, Int), String, TimeWindow] {
      override def process(context: Context, elements: Iterable[(String, Int)], out: Collector[String]): Unit = {
        val sortList = elements.toList.sortBy(_._2)
        for (elem <- sortList) {
          out.collect(elem._1 + " speed:" + elem._2)
        }
      }
    }).print()
    env.execute()
  }
}

```

案例4: 每隔10s统计出窗口内所有车辆的最大及最小速度

```

import java.util.Properties

```

```

import org.apache.flink.api.common.serialization.SimpleStringSchema
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.connectors.kafka.FlinkKafkaConsumer
import org.apache.kafka.common.serialization.StringSerializer
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.api.scala.function.ProcessAllWindowFunction
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.util.Collector

object Demo04MaxMinSpeed {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    //设置连接kafka的配置信息
    val props = new Properties()
    //注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
    props.setProperty("bootstrap.servers",
      "node01:9092,node02:9092,node03:9092")
    props.setProperty("group.id", "flink-kafka-001")
    props.setProperty("key.deserializer", classOf[StringSerializer].getName)
    props.setProperty("value.deserializer", classOf[StringSerializer].getName)

    val stream = env.addSource(new FlinkKafkaConsumer[String]("flink-kafka", new
      SimpleStringSchema(), props))
    stream.map(data =>{
      val arr = data.split("\t")
      (arr(1),arr(3).toInt)
    }).timewindowAll(Time.seconds(20))
      .process(new ProcessAllWindowFunction[(String,Int),String,TimeWindow] {
        override def process(context: Context, elements: Iterable[(String,
          Int)], out: Collector[String]): Unit = {
          val sortList = elements.toList.sortBy(_._2)
          println(sortList)
          val minSpeedInfo = sortList.head
          val maxSpeedInfo = sortList.last
          val startWindowTime = context.window.getStart
          val endWindowTime = context.window.getEnd
          out.collect(
            "窗口起始时间: "+startWindowTime + "结束时间: "+ endWindowTime + " 最小车辆
            速度车牌号: " + minSpeedInfo._1 + " 车速: "+minSpeedInfo._2 + "\t最大车辆速度车牌号: "
            + maxSpeedInfo._1 + " 车速: " + maxSpeedInfo._2
          )
        }
      }).print()
    env.execute()
  }
}

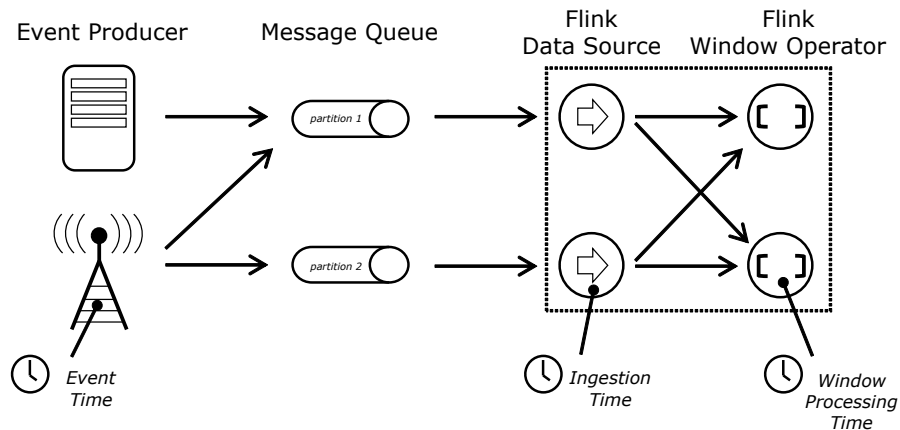
```

## Flink Time时间语义

Flink定义了三类时间

- **处理时间 (Process Time)** 数据进入Flink被处理的系统时间 (Operator处理数据的系统时间)

- **事件时间 (Event Time)** 数据在数据源产生的时间，一般由事件中的时间戳描述，比如用户日志中的Timestamp
- **摄取时间 (Ingestion Time)** 数据进入Flink的时间，记录被Source节点观察到的系统时间



Flink流式计算的时候需要显示定义时间语义，根据不同的时间语义来处理数据，比如指定的时间语义是事件时间，那么我们就切换到事件时间的世界观中，窗口的起始与终止时间都是以事件时间为依据

在Flink中默认使用的是Process Time，如果要使用其他的时间语义，在执行环境中可以设置

```
//设置时间语义为Ingestion Time
env.setStreamTimeCharacteristic(TimeCharacteristic.IngestionTime)
//设置时间语义为Event Time 我们还需要指定一下数据中哪个字段是事件时间（下文会讲）
env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
```

- 基于事件时间的Window操作

```
import org.apache.flink.streaming.api.TimeCharacteristic
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, _}
import org.apache.flink.streaming.api.windowing.assigners.TumblingEventTimeWindows
import org.apache.flink.streaming.api.windowing.time.Time

object EventTimeWindow {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    val stream = env.socketTextStream("node01",
8888).assignAscendingTimestamps(data => {
      val splits = data.split(" ")
      splits(0).toLong
    })

    stream
      .flatMap(x=>x.split(" ").tail)
      .map(_._1)
      .keyBy(_._1)
      // .timeWindow(Time.seconds(10))
      .window(TumblingEventTimeWindows.of(Time.seconds(10)))
      .reduce((v1: (String, Int), v2: (String, Int)) => {
        (v1._1, v1._2 + v2._2)
      })
      .print()

    env.execute()
  }
}
```



```
}  
}
```

## Flink Time Watermark(水印)

Watermark本质就是时间戳

在使用Flink处理数据的时候，数据通常都是按照事件产生的时间（事件时间）的顺序进入到Flink，但是在遇到特殊情况下，比如遇到网络延迟或者使用Kafka（多分区）很难保证数据都是按照事件时间的顺序进入Flink，很有可能是乱序进入。

如果使用的是事件时间这个语义，数据一旦是乱序进入，那么在使用Window处理数据的时候，就会出现延迟数据不会被计算的问题

- 举例：Window窗口长度10s，滚动窗口

001 zs 2020-04-25 10:00:01

001 zs 2020-04-25 10:00:02

001 zs 2020-04-25 10:00:03

001 zs 2020-04-25 10:00:11 窗口触发执行

001 zs 2020-04-25 10:00:05 延迟数据，不会被上一个窗口所计算导致计算结果不正确

Watermark+Window可以很好的解决延迟数据的问题

Flink窗口计算的过程中，如果数据全部到达就会到窗口中的数据做处理，如果过有延迟数据，那么窗口需要等待全部的数据到来之后，再触发窗口执行，需要等待多久？不可能无限期等待，我们用户可以自己来设置延迟时间

这样就可以**尽可能**保证延迟数据被处理

根据用户指定的延迟时间生成水印（Watermak = 最大事件时间-指定延迟时间），当Watermak 大于等于窗口的停止时间，这个窗口就会被触发执行

- 举例：Window窗口长度10s(01-10)，滚动窗口，指定延迟时间3s

001 ls 2020-04-25 10:00:01 wm:2020-04-25 09:59:58

001 ls 2020-04-25 10:00:02 wm:2020-04-25 09:59:59

001 ls 2020-04-25 10:00:03 wm:2020-04-25 10:00:00

001 ls 2020-04-25 10:00:09 wm:2020-04-25 10:00:06

001 ls 2020-04-25 10:00:12 wm:2020-04-25 10:00:09

001 ls 2020-04-25 10:00:08 wm:2020-04-25 10:00:05 延迟数据

001 ls 2020-04-25 10:00:13 wm:2020-04-25 10:00:10 此时wm >= window end time 触发窗口执行 处理的是事件时间01-10的数据，并不是水印时间为01-10的数据 **重点**

讲道理，如果没有Watermark在倒数第三条数据来的时候，就会触发执行，那么倒数第二条的延迟数据就不会被计算，那么有了水印可以处理延迟3s内的数据

**注意：如果数据不会乱序进入Flink，没必要使用Watermark**

- 代码演示

演示数据：

10000 hello msb

14000 hello flink

20000 hello hadoop

```

21000 hello bj
17000 hello sh      迟到数据
23000 hello jjj

import org.apache.flink.streaming.api.TimeCharacteristic
import
org.apache.flink.streaming.api.functions.timestamps.BoundedOutOfOrdernessTimestampExtractor
import org.apache.flink.streaming.api.scala.function.ProcessWindowFunction
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, _}
import
org.apache.flink.streaming.api.windowing.assigners.TumblingEventTimeWindows
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.util.Collector

object EventTimeDelaywindow {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    val stream = env.socketTextStream("node01",
8888).assignTimestampsAndWatermarks(new
BoundedOutOfOrdernessTimestampExtractor[String](Time.seconds(3)) {
      override def extractTimestamp(element: String): Long = {
        element.split(" ")(0).toLong
      }
    })

    stream
      .flatMap(x=>x.split(" ").tail)
      .map(_._1)
      .keyBy(_._1)
      // .timeWindow(Time.seconds(10))
      .window(TumblingEventTimeWindows.of(Time.seconds(10)))
      .process(new ProcessWindowFunction[(String,Int),
(String,Int),String,TimeWindow] {
        override def process(key: String, context: Context, elements:
Iterable[(String, Int)], out: Collector[(String, Int)]): Unit = {
          val start = context.window.getStart
          val end = context.window.getEnd
          var count = 0
          for (elem <- elements) {
            count += elem._2
          }
          println("start:" + start + " end:" + end + " word:" + key + "
count:"+count)
        }
      })
      .print()

    env.execute()
  }
}

```

DataStream API提供了自定义水印生成器和内置水印生成器

生成水印策略:

- 周期性水印 (Periodic Watermark) 根据事件或者处理时间周期性的触发水印生成器(Assigner), 默认100ms, 每隔100毫秒自动向流里注入一个Watermark

周期性水印API 1:

```
val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
env.getConfig.setAutoWatermarkInterval(100)
val stream = env.socketTextStream("node01",
8888).assignTimestampsAndWatermarks(new
BoundedOutOfOrdernessTimestampExtractor[String](Time.seconds(3)) {
    override def extractTimestamp(element: String): Long = {
        element.split(" ")(0).toLong
    }
})
```

周期性水印API 2:

```
import org.apache.flink.streaming.api.TimeCharacteristic
import
org.apache.flink.streaming.api.functions.AssignerWithPeriodicWatermarks
import org.apache.flink.streaming.api.scala.function.ProcessWindowFunction
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, _}
import org.apache.flink.streaming.api.watermark.Watermark
import
org.apache.flink.streaming.api.windowing.assigners.TumblingEventTimeWindows
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.util.Collector

object EventTimeDelayWindow {

    class MyTimestampAndWatermarks(delayTime: Long) extends
AssignerWithPeriodicWatermarks[String] {

        var maxCurrentWatermark: Long = _

        //水印=最大事件时间-延迟时间 后被调用 水印是递增, 小于上一个水印不会被发射出去
        override def getCurrentWatermark: Watermark = {
            //产生水印
            new Watermark(maxCurrentWatermark - delayTime)
        }

        //获取当前的时间戳 先被调用
        override def extractTimestamp(element: String, previousElementTimestamp:
Long): Long = {
            val currentTimeStamp = element.split(" ")(0).toLong
            maxCurrentWatermark = math.max(currentTimeStamp, maxCurrentWatermark)
            currentTimeStamp
        }
    }

    def main(args: Array[String]): Unit = {
        val env = StreamExecutionEnvironment.getExecutionEnvironment
        env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
        env.getConfig.setAutoWatermarkInterval(100)
    }
}
```

```

    val stream = env.socketTextStream("node01",
8888).assignTimestampsAndWatermarks(new MyTimestampAndWatermarks(3000L))

    stream
      .flatMap(x => x.split(" ").tail)
      .map(_._1)
      .keyBy(_._1)
      //      .timewindow(Time.seconds(10))
      .window(TumblingEventTimeWindows.of(Time.seconds(10)))
      .process(new ProcessWindowFunction[(String, Int), (String, Int),
String, TimeWindow] {
        override def process(key: String, context: Context, elements:
Iterable[(String, Int)], out: Collector[(String, Int)]): Unit = {
          val start = context.window.getStart
          val end = context.window.getEnd
          var count = 0
          for (elem <- elements) {
            count += elem._2
          }
          println("start:" + start + " end:" + end + " word:" + key + "
count:" + count)
        }
      })
      .print()

    env.execute()
  }
}

```

- 间歇性水印生成器

间歇性水印 (Punctuated Watermark) 在观察到事件后，会依据用户指定的条件来决定是否发射水印

比如，在车流量的数据中，001卡口通信经常异常，传回到服务器的数据会有延迟问题，其他的卡口都是正常的，那么这个卡口的数据需要打上水印

```

package com.msb.stream.windowt

import org.apache.flink.streaming.api.TimeCharacteristic
import
org.apache.flink.streaming.api.functions.AssignerWithPunctuatedWatermarks
import org.apache.flink.streaming.api.scala.{StreamExecutionEnvironment, _}
import org.apache.flink.streaming.api.watermark.Watermark
import org.apache.flink.streaming.api.windowing.time.Time

object PunctuatedWatermarkTest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setParallelism(1)
    env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    //卡口号、时间戳
    env
      .socketTextStream("node01", 8888)
      .map(data => {
        val splits = data.split(" ")
        (splits(0), splits(1).toLong)
      })
  }
}

```

```

        .assignTimestampsAndWatermarks(new myWatermark(3000))
        .keyBy(_._1)
        .timewindow(Time.seconds(5))
        .reduce((v1: (String, Long), v2: (String, Long)) => {
            (v1._1 + "," + v2._1, v1._2 + v2._2)
        }).print()

    env.execute()
}

class myWatermark(delay: Long) extends
    AssignerWithPunctuatedWatermarks[(String, Long)] {
    var maxTimeStamp: Long = _

    override def checkAndGetNextWatermark(elem: (String, Long),
        extractedTimestamp: Long): Watermark = {
        maxTimeStamp = extractedTimestamp.max(maxTimeStamp)
        if ("001".equals(elem._1)) {
            new Watermark(maxTimeStamp - delay)
        } else {
            new Watermark(maxTimeStamp)
        }
    }

    override def extractTimestamp(element: (String, Long),
        previousElementTimestamp: Long): Long = {
        element._2
    }
}

```

## AllowedLateness

基于Event-Time的窗口处理流式数据，虽然提供了Watermark机制，却只能在一定程度上解决了数据乱序的问题。但在某些情况下数据可能延时会非常严重，即使通过Watermark机制也无法等到数据全部进入窗口再进行处理。Flink中默认会将这些迟到的数据做丢弃处理，但是有些时候用户希望即使数据延迟并不是很严重的情况下，也能继续窗口计算，不希望对于数据延迟比较严重的数据混入正常的计算流程中，此时就需要使用Allowed Lateness机制来对迟到的数据进行额外的处理。

举例：

例如用户大屏数据展示系统，即使正常的窗口中没有将迟到的数据进行统计，但为了保证页面数据显示的连续性，后来接入到系统中迟到比较严重的数据所统计出来的结果不希望显示在屏幕上，而是将延时数据和结果存储到数据库中，便于后期对延时数据进行分析。对于这种情况需要借助Side Output来处理，通过使用sideOutputLateData (OutputTag) 来标记迟到数据计算的结果，然后使用getSideOutput (lateOutputTag) 从窗口结果中获取lateOutputTag标签对应的数据，之后转成独立的DataStream数据集进行处理，创建late-data的OutputTag，再通过该标签从窗口结果中将迟到数据筛选出来

Flink默认当窗口计算完毕后，窗口元素数据及状态会被清空，但是使用AllowedLateness，可以延迟清除窗口元素数据及状态，以便于当延迟数据到来的时候，能够重新计算当前窗口

Watermark 2s   AllowedLateness 3s

10000 hello  
 11000 spark  
 14000 flink  
 15000 hadoop 此时窗口并不会计算，因为watermark设为2s 此时的watermark是13000 窗口范围10000-15000  
 17000 sqoop 此时窗口会被计算 默认：窗口计算完毕，窗口数据全部会被清空  
 12000 flume 此时窗口重新计算（10000-15000），因为开启了AllowedLateness 3s，当watermark>=window end+ AllowedLateness 3s 窗口数据及状态才会被清除掉，此时的watermark是15000  
 20000 scala 此时上一个窗口（10000-15000）的数据及状态会被清空  
 12000 hdfs 此时窗口不会重新计算，因为现在watermark是18000>=15000+3000,12000数据是迟到非常严重的数据，会被放入到侧输出流中

本来10000-15000的窗口，在15000的时候会计算，但是由于watermark 的原因，等待了2s 17000的时候才会计算，又因为AllowedLateness 3s的原因，10000-15000的窗口会被保存3s（注意这是eventtime时间语义），直到20000出现，才会被删除，所以在20000没有出现之前，凡是事件时间在10000-15000的数据都会重新进行窗口计算

超过5s的数据，称之为迟到非常严重的数据，放入到侧输出流  
 5s以内的数据，称之为迟到不严重的数据，窗口重新计算

```
import org.apache.flink.streaming.api.TimeCharacteristic
import org.apache.flink.streaming.api.functions.ProcessFunction
import org.apache.flink.streaming.api.functions.timestamps.BoundedOutOfOrdernessTimestampExtractor
import org.apache.flink.streaming.api.scala.StreamExecutionEnvironment
import org.apache.flink.streaming.api.scala._
import org.apache.flink.streaming.api.scala.function.ProcessAllWindowFunction
import org.apache.flink.streaming.api.windowing.time.Time
import org.apache.flink.streaming.api.windowing.windows.TimeWindow
import org.apache.flink.util.Collector

object Allowlatest {
  def main(args: Array[String]): Unit = {
    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    env.setParallelism(1)
    val stream = env.socketTextStream("node01", 8888)
    var lateTag = new OutputTag[(Long, String)]("late")
    val value = stream.map(x => {
      val strings = x.split(" ")
      (strings(0).toLong, strings(1))
    }).assignTimestampsAndWatermarks(new
    BoundedOutOfOrdernessTimestampExtractor[(Long, String)](Time.seconds(2)) {
      override def extractTimestamp(element: (Long, String)): Long = element._1
    }).timewindowAll(Time.seconds(5))
      .allowedLateness(Time.seconds(3))
      .sideOutputLateData(lateTag)
      .process(new ProcessAllWindowFunction[(Long, String), (Long, String),
    Timewindow] {
        override def process(context: Context, elements: Iterable[(Long,
    String)], out: Collector[(Long, String)]): Unit = {
          println(context.window.getStart + "---" + context.window.getEnd)
          for (elem <- elements) {
            out.collect(elem)
          }
        }
      })
  }
}
```

```

    }
  }
})
value.print("main")
value.getSideOutput(lateTag).print("late")
env.execute()
}
}

```

问题1：使用AllowedLateness 方法是不是会降低flink计算的吞吐量？ 是的

问题2：直接watermark设置为5 不是也可以代替这一通操作嘛？ 不能代替， watermark设置为5的话，允许延迟5s，每次处理过去5s的窗口数据，延迟比较高，如果使用这通操作，每次处理过去2s的数据，实时性比较高，当有新的延迟数据，即时计算，对于计算实时性比较高的场景还得使用这一通操作

问题3： watermark (5s) +滑动窗口（滑动间隔2s）能够实现这通计算？ 不行

案例：每隔5s统计各个卡口最近5s的车流量（滑动窗口），计算实时性小于2（ps：当10s的数据来了，8s之前的数据必须处理完），允许数据延迟5s，数据延迟超过5s的数据放入到侧输出流中

## Flink关联维表实战

在Flink实际开发过程中，可能会遇到source 进来的数据，需要连接数据库里面的字段，再做后面的处理

比如，想要通过id获取对应的地区名字，这时候需要通过id查询地区维度表，获取具体的地区名

对于不同的应用场景，关联维度表的方式不同

- 场景1：维度表信息基本不发生改变，或者发生改变的频率很低

实现方案：采用Flink提供的CachedFile

Flink提供了一个分布式缓存（CachedFile），类似于hadoop，可以使用户在并行函数中很方便的读取本地文件，并把它放在TaskManager节点中，防止task重复拉取。此缓存的工作机制如下：程序注册一个文件或者目录(本地或者远程文件系统，例如hdfs或者s3)，通过ExecutionEnvironment注册缓存文件并为它起一个名称。当程序执行，Flink自动将文件或者目录复制到所有TaskManager节点的本地文件系统，**仅会执行一次**。用户可以通过这个指定的名称查找文件或者目录，然后从TaskManager节点的本地文件系统访问它

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
env.registerCachedFile("/root/id2city","id2city")

val socketStream = env.socketTextStream("node01",8888)
val stream = socketStream.map(_._1.toInt)
stream.map(new RichMapFunction[Int,String] {

    private val id2CityMap = new mutable.HashMap[Int,String]()
    override def open(parameters: Configuration): Unit = {
        val file =
getRuntimeContext().getDistributedCache().getFile("id2city")
        val str = FileUtils.readFileUtf8(file)
        val strings = str.split("\r\n")
        for(str <- strings){
            val splits = str.split(" ")
            val id = splits(0).toInt
            val city = splits(1)
            id2CityMap.put(id,city)
        }
    }
}


```

```

    }
    override def map(value: Int): String = {
        id2CityMap.getOrElse(value, "not found city")
    }
  }).print()
env.execute()

```

在集群中查看对应TaskManager的log日志，发现注册的file会被拉取到各个TaskManager的工作目录区

 1590302497460

- 场景2：对于维度表更新频率比较高并且对于查询维度表的实时性要求比较高  
实现方案：使用定时器，定时加载外部配置文件或者数据库

```

val env = StreamExecutionEnvironment.getExecutionEnvironment
env.setParallelism(1)
val stream = env.socketTextStream("node01", 8888)

stream.map(new RichMapFunction[String, String] {

    private val map = new mutable.HashMap[String, String]()

    override def open(parameters: Configuration): Unit = {
        println("init data ...")
        query()
        val timer = new Timer(true)
        timer.schedule(new TimerTask {
            override def run(): Unit = {
                query()
            }
        }, 1000, 2000)
    }

    def query() = {
        val source =
Source.fromFile("D:\\code\\StudyFlink\\data\\id2city", "UTF-8")
        val iterator = source.getLines()
        for (elem <- iterator) {
            val vs = elem.split(" ")
            map.put(vs(0), vs(1))
        }
    }

    override def map(key: String): String = {
        map.getOrElse(key, "not found city")
    }
  }).print()

env.execute()


```

如果维度信息在配置文件中存储，那么还有一个解决方案，就是使用readFile读取文件，因为这个方法可以检测内容是否发生改变，之前在讲readFile的时候讲过，不再赘述.....

- 场景3：对于维度表更新频率高并且对于查询维度表的实时性要求高



实现方案：管理员在修改配置文件的时候，需要将更改的信息同步值Kafka配置Topic中，然后将kafka的配置流信息变成广播流，广播到业务流的各个线程中

1590305917005

```
val env = StreamExecutionEnvironment.getExecutionEnvironment

//设置连接kafka的配置信息
val props = new Properties()
//注意 sparkstreaming + kafka (0.10之前版本) receiver模式 zookeeper url (元数据)
props.setProperty("bootstrap.servers", "node01:9092,node02:9092,node03:9092")
props.setProperty("group.id", "flink-kafka-001")
props.setProperty("key.deserializer", classOf[StringSerializer].getName)
props.setProperty("value.deserializer", classOf[StringSerializer].getName)
val consumer = new FlinkKafkaConsumer[String]("configure", new
SimpleStringSchema(), props)
//从topic最开始的数据读取
// consumer.setStartFromEarliest()
//从最新的数据开始读取
consumer.setStartFromLatest()

//动态配置信息流
val configureStream = env.addSource(consumer)
//业务流
val busStream = env.socketTextStream("node01", 8888)

val descriptor = new MapStateDescriptor[String, String]("dynamicConfig",
    BasicTypeInfo.STRING_TYPE_INFO,
    BasicTypeInfo.STRING_TYPE_INFO)
//设置广播流的数据描述信息
val broadcastStream = configureStream.broadcast(descriptor)

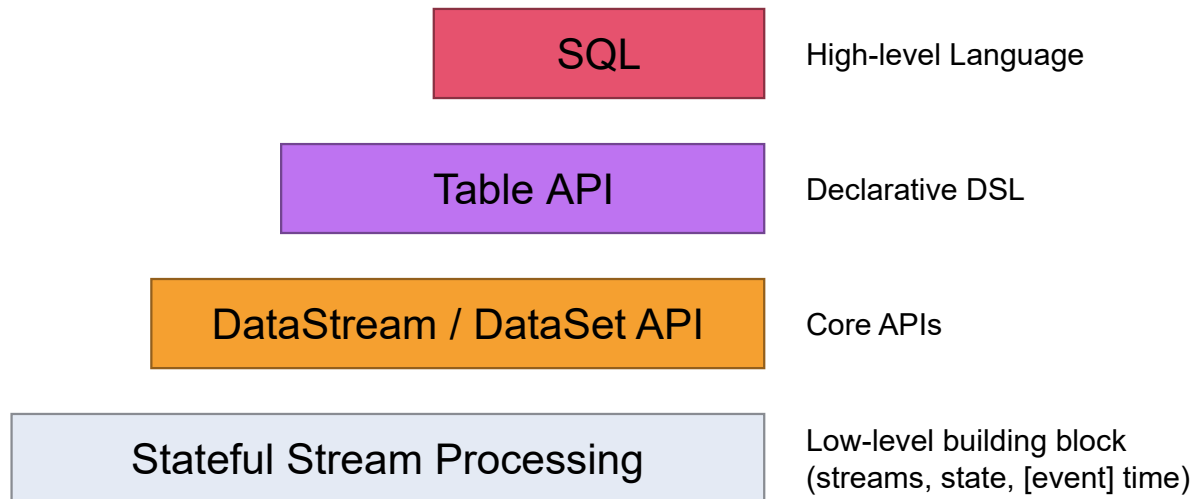
//connect关联业务流与配置信息流，broadcastStream流中的数据会广播到下游的各个线程中
busStream.connect(broadcastStream)
    .process(new BroadcastProcessFunction[String, String, String] {
        override def processElement(line: String, ctx:
BroadcastProcessFunction[String, String, String]#ReadOnlyContext, out:
Collector[String]): Unit = {
            val broadcast = ctx.getBroadcastState(descriptor)
            val city = broadcast.get(line)
            if(city == null){
                out.collect("not found city")
            }else{
                out.collect(city)
            }
        }
    })

//kafka中配置流信息，写入到广播流中
    override def processBroadcastElement(line: String, ctx:
BroadcastProcessFunction[String, String, String]#Context, out:
Collector[String]): Unit = {
        val broadcast = ctx.getBroadcastState(descriptor)
        //kafka中的数据
        val elems = line.split(" ")
        broadcast.put(elems(0), elems(1))
    }
}).print()
```


```
env.execute()
```

## TableAPI和Flink SQL

在Spark中有DataFrame这样的关系型编程接口，因其强大且灵活的表达能力，能够让用户通过非常丰富的接口对数据进行处理，有效降低了用户的使用成本。Flink也提供了关系型编程接口Table API以及基于Table API的SQL API，让用户能够通过使用结构化编程接口高效地构建Flink应用。同时Table API以及SQL能够统一处理批量和实时计算业务，无须切换修改任何应用代码就能够基于同一套API编写流式应用和批量应用，从而达到真正意义的批流统一



在 Flink 1.8 架构里，如果用户需要同时流计算、批处理的场景下，用户需要维护两套业务代码，开发人员也要维护两套技术栈，非常不方便。Flink 社区很早就设想过将批数据看作一个有界流数据，将批处理看作流计算的一个特例，从而实现流批统一，阿里巴巴的 Blink 团队在这方面做了大量的工作，已经实现了 Table API & SQL 层的流批统一。阿里巴巴已经将 Blink 开源回馈给 Flink 社区

1590473668405

## 开发环境构建

在 Flink 1.9 中，Table 模块迎来了核心架构的升级，引入了阿里巴巴Blink团队贡献的诸多功能，取名叫做：Blink Planner。在使用Table API和SQL开发Flink应用之前，通过添加Maven的依赖配置到项目中，在本地工程中引入相应的依赖库，库中包含了Table API和SQL接口

```
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-table-planner_2.11</artifactId>
  <version>1.9.1</version>
</dependency>
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-table-api-scala-bridge_2.11</artifactId>
  <version>1.9.1</version>
</dependency>
```

## TableEnvironment

和DataStream API一样，Table API和SQL中具有相同的基本编程模型。首先需要构建对应的TableEnvironment创建关系型编程环境，才能够在程序中使用Table API和SQL来编写应用程序，另外Table API和SQL接口可以在应用中同时使用，Flink SQL基于Apache Calcite框架实现了SQL标准协议，是构建在Table API之上的更高级接口。

首先需要在环境中创建TableEnvironment对象，TableEnvironment中提供了注册内部表、执行Flink SQL语句、注册自定义函数等功能。根据应用类型的不同，TableEnvironment创建方式也有所不同，但是都是通过调用create()方法创建

流计算环境下创建TableEnvironment：

```
//创建流式计算的上下文环境
val env = StreamExecutionEnvironment.getExecutionEnvironment
//创建Table API的上下文环境
val tableEnv = StreamTableEnvironment.create(env)
```

## Table API

在Flink中创建一张表有两种方法：

- 从一个文件中导入表结构（Structure）（常用于批计算）（静态）
- 从DataStream或者DataSet转换成Table（动态）

### 1)创建Table

Table API中已经提供了TableSource从外部系统获取数据，例如常见的数据库、文件系统和Kafka消息队列等外部系统。

#### 1. 从文件中创建Table（静态表）

Flink允许用户从本地或者分布式文件系统中读取和写入数据，在Table API中可以通过CsvTableSource类来创建，只需指定相应的参数即可。但是文件格式必须是CSV格式的。其他文件格式也支持（在Flink还有Connector的支持其他格式或者自定义TableSource）

```
//创建流式计算的上下文环境
val env = StreamExecutionEnvironment.getExecutionEnvironment
//创建Table API的上下文环境
val tableEnv = StreamTableEnvironment.create(env)

val source = new
CsvTableSource("D:\\code\\StudyFlink\\data\\tableexamples"
    , Array[String]("id", "name", "score")
    , Array(Types.INT, Types.STRING, Types.DOUBLE)
)
//将source注册成一张表 别名: exampleTab
tableEnv.registerTableSource("exampleTab", source)
tableEnv.scan("exampleTab").printSchema()
```

代码最后不需要env.execute()，这并不是一个流式计算任务

#### 2. 从DataStream中创建Table（动态表）

前面已经知道Table API是构建在DataStream API和DataSet API之上的一层更高级的抽象，因此用户可以灵活地使用Table API将Table转换成DataStream或DataSet数据集，也可以将DataStream或DataSet数据集转换成Table，这和Spark中的DataFrame和RDD的关系类似

## 2) 修改Table中字段名

Flink支持把自定义POJOs类的所有case类的属性名字变成字段名，也可以通过基于字段偏移位置和字段名称两种方式重新修改：

```
//导入table库中的隐式转换
import org.apache.flink.table.api.scala._
// 基于位置重新指定字段名称为"field1", "field2", "field3"
val table = tStreamEnv.fromDataStream(stream, 'field1, 'field2, 'field3)
// 将DataStream转换成Table,并且将字段名称重新成别名
val table: Table = tStreamEnv.fromDataStream(stream, 'rowtime as 'newTime,
'id as 'newId,'variable as 'newVariable)
```

**注意：**要导入隐式转换。如果使用as 修改字段，必须修改表中所有的字段。

## 3) 查询和过滤

在Table对象上使用select操作符查询需要获取的指定字段，也可以使用filter或where方法过滤字段和检索条件，将需要的数据检索出来。

```
object TableAPITest {

  def main(args: Array[String]): Unit = {
    val streamEnv: StreamExecutionEnvironment =
      StreamExecutionEnvironment.getExecutionEnvironment
    streamEnv.setParallelism(1)
    //初始化Table API的上下文环境
    val tableEnv = StreamTableEnvironment.create(streamEnv)
    //导入隐式转换，建议写在这里，可以防止IDEA代码提示出错的问题
    import org.apache.flink.streaming.api.scala._
    import org.apache.flink.table.api.scala._
    val data = streamEnv.socketTextStream("hadoop101", 8888)
      .map(line=>{
        var arr = line.split(",")
        new
        StationLog(arr(0).trim, arr(1).trim, arr(2).trim, arr(3).trim, arr(4).trim, arr(5).trim, arr(6).trim)
      })

    val table: Table = tableEnv.fromDataStream(data)
    //查询
    tableEnv.toAppendStream[Row](
      table.select('sid, 'callType as 'type, 'callTime, 'callOut))
      .print()
    //过滤查询
    tableEnv.toAppendStream[Row](
      table.filter('callType === "success") //filter
      .where('callType === "success")) //where
      .print()
    tableEnv.execute("sql")
  }
}
```

其中toAppendStream函数是把Table对象转换成DataStream对象。

## 4) 分组聚合

举例：我们统计每个基站的日志数量。

```
val table: Table = tableEnv.fromDataStream(data)
tableEnv.toRetractStream[Row](
  table.groupBy('sid).select('sid, 'sid.count as 'logCount))
  .filter(_._1==true) //返回的如果是true才是Insert的数据
  .print()
```

在代码中可以看出，使用toAppendStream和toRetractStream方法将Table转换为DataStream[T]数据集，T可以是Flink自定义的数据格式类型Row，也可以是用户指定的数据格式类型。在使用toRetractStream方法时，返回的数据类型结果为DataStream[(Boolean,T)]，Boolean类型代表数据更新类型，True对应INSERT操作更新的数据，False对应DELETE操作更新的数据。

## 5) UDF自定义的函数

用户可以在Table API中自定义函数类，常见的抽象类和接口是：

- ScalarFunction
- TableFunction
- AggregateFunction
- TableAggregateFunction

案例：使用Table完成基于流的WordCount

```
object TableAPITest2 {

  def main(args: Array[String]): Unit = {
    val streamEnv: StreamExecutionEnvironment =
      StreamExecutionEnvironment.getExecutionEnvironment
    streamEnv.setParallelism(1)
    //初始化Table API的上下文环境
    val tableEnv = StreamTableEnvironment.create(streamEnv)
    //导入隐式转换，建议写在这里，可以防止IDEA代码提示出错的问题
    import org.apache.flink.streaming.api.scala._
    import org.apache.flink.table.api.scala._

    val stream: DataStream[String] =
      streamEnv.socketTextStream("hadoop101", 8888)
    val table: Table = tableEnv.fromDataStream(stream, 'words)
    var my_func = new MyFlatMapFunction() //自定义UDF
    val result: Table = table.flatMap(my_func('words)).as('word, 'count)
      .groupBy('word) //分组
      .select('word, 'count.sum as 'c) //聚合
    tableEnv.toRetractStream[Row](result)
      .filter(_._1==true)
      .print()

    tableEnv.execute("table_api")
  }
  //自定义UDF
  class MyFlatMapFunction extends TableFunction[Row]{
    //定义类型
    override def getResultType: TypeInformation[Row] = {
      Types.ROW(Types.STRING, Types.INT)
    }
    //函数主体
  }
```

```

def eval(str:String):Unit ={
  str.trim.split(" ")
    .foreach({word=>{
      var row =new Row(2)
      row.setField(0,word)
      row.setField(1,1)
      collect(row)
    }})
}
}
}
}

```

## 6) Window

Flink支持ProcessTime、EventTime和IngestionTime三种时间概念，针对每种时间概念，Flink Table API中使用Schema中单独的字段来表示时间属性，当时间字段被指定后，就可以在基于时间的操作算子中使用相应的时间属性。

在Table API中通过使用.rowtime来定义EventTime字段，在ProcessTime时间字段名后使用.proctime后缀来指定ProcessTime时间属性

案例：统计最近5秒钟，每个基站的呼叫数量

```

object TableAPITest {

  def main(args: Array[String]): Unit = {
    val streamEnv: StreamExecutionEnvironment =
      StreamExecutionEnvironment.getExecutionEnvironment
    //指定EventTime为时间语义
    streamEnv.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    streamEnv.setParallelism(1)
    //初始化Table API的上下文环境
    val tableEnv =StreamTableEnvironment.create(streamEnv)
    //导入隐式转换，建议写在这里，可以防止IDEA代码提示出错的问题
    import org.apache.flink.streaming.api.scala._
    import org.apache.flink.table.api.scala._

    val data = streamEnv.socketTextStream("hadoop101",8888)
      .map(line=>{
        var arr =line.split(",")
        new
        StationLog(arr(0).trim,arr(1).trim,arr(2).trim,arr(3).trim,arr(4).trim.toLong,ar
        r(5).trim.toLong)
      })
      .assignTimestampsAndWatermarks( //引入watermark
        new BoundedOutOfOrdernessTimestampExtractor[StationLog](Time.seconds(2))
      ){//延迟2秒
        override def extractTimestamp(element: StationLog) = {
          element.callTime
        }
      })

    //设置时间属性
    val table: Table =
      tableEnv.fromDataStream(data,'sid','callOut','callIn','callType','callTime.rowtime)
    //滚动window ,第一种写法

```

```

    val result: Table = table.window(Tumble over 5.second on 'callTime as
'window')
    //第二种写法
    val result: Table =
table.window(Tumble.over("5.second").on("callTime").as("window"))
    .groupBy('window, 'sid)
    .select('sid, 'window.start, 'window.end, 'window.rowtime, 'sid.count)
    //打印结果
    tableEnv.toRetractStream[Row](result)
    .filter(_._1==true)
    .print()

    tableEnv.execute("sql")
}
}

```

上面的案例是滚动窗口，如果是滑动窗口也是一样，代码如下：

```

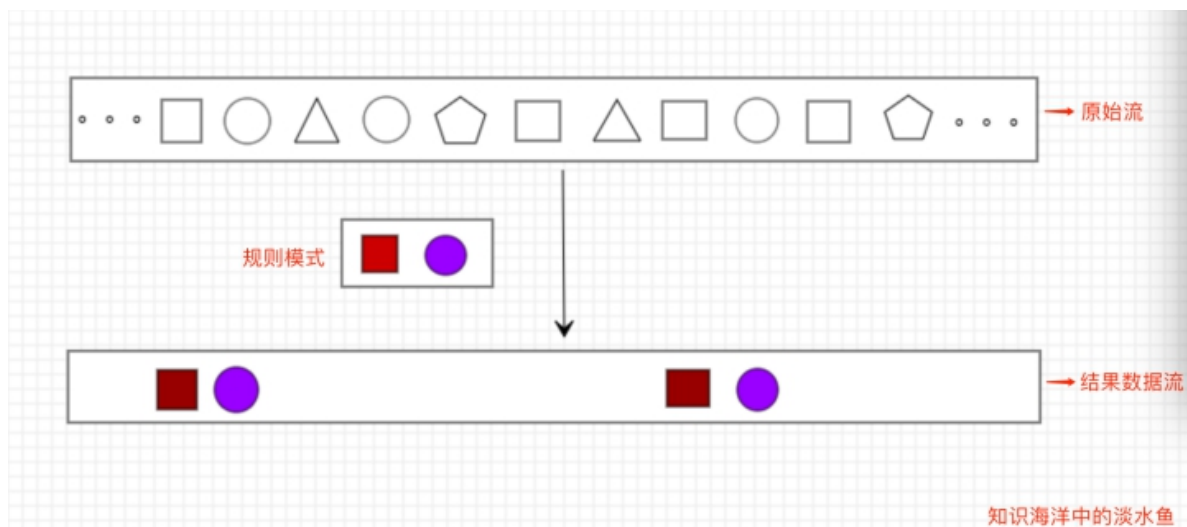
//滑动窗口，窗口大小为：10秒，滑动步长为5秒：第一种写法
table.window(Slide over 10.second every 5.second on 'callTime as 'window)
//滑动窗口第二种写法
table.window(Slide.over("10.second").every("5.second").on("callTime").as("window
"))

```

## Flink的复杂事件处理CEP

复杂事件处理（CEP）是一种基于流处理的技术，将系统数据看作不同类型的事件，通过分析事件之间的关系，建立不同的事件关系序列库，并利用过滤、关联、聚合等技术，最终由简单事件产生高级事件，并通过模式规则的方式对重要信息进行跟踪和分析，从实时数据中发掘有价值的信息。复杂事件处理主要应用于防范网络欺诈、设备故障检测、风险规避和智能营销等领域。Flink基于DataStrem API提供了FlinkCEP组件栈，专门用于对复杂事件的处理，帮助用户从流式数据中发掘有价值的信息。

CEP(Complex Event Processing)就是在无界事件流中检测事件模式，让我们掌握数据中重要的部分。flink CEP是在flink中实现的复杂事件处理库。



### 1. CEP相关概念

#### 1) 配置依赖

在使用FlinkCEP组件之前，需要将FlinkCEP的依赖库引入项目工程中。

```
<dependency>
  <groupId>org.apache.flink</groupId>
  <artifactId>flink-cep-scala_2.11</artifactId>
  <version>1.9.1</version>
</dependency>
```

org.apache.flink flink-cep-scala\_2.11 1.9.1

## 2) 事件定义

简单事件：简单事件存在于现实场景中，主要的特点为处理单一事件，事件的定义可以直接观察出来，处理过程中无须关注多个事件之间的关系，能够通过简单的数据处理手段将结果计算出来。

复杂事件：相对于简单事件，复杂事件处理的不仅是单一的事件，也处理由多个事件组成的复合事件。复杂事件处理监测分析事件流(Event Streaming)，当特定事件发生时来触发某些动作。

复杂事件中事件与事件之间包含多种类型关系，常见的有时序关系、聚合关系、层次关系、依赖关系及因果关系等。

## 2. Pattern API

FlinkCEP中提供了Pattern API用于对输入流数据的复杂事件规则定义，并从事件流中抽取事件结果。包含四个步骤：

- 1、输入事件流的创建
- 2、Pattern的定义
- 3、Pattern应用在事件流上检测
- 4、选取结果

### 1) 模式定义

定义Pattern可以是单次执行模式，也可以是循环执行模式。单词执行模式一次只接受一个事件，循环执行模式可以接收一个或者多个事件。通常情况下，可以通过指定循环次数将单次执行模式变为循环执行模式。每种模式能够将多个条件组合应用到同一事件之上，条件组合可以通过where方法进行叠加。每个Pattern都是通过begin方法定义的

```
val start = Pattern.begin[Event]("start_pattern")
```

```
val start = Pattern.beginEvent
```

下一步通过Pattern.where()方法在Pattern上指定Condition，只有当Condition满足之后，当前的Pattern才会接受事件。

```
start.where(_.getCallType == "success")
```

```
start.where(_.getCallType == "success")
```

#### 1. 设置循环次数

对于已经创建好的Pattern，可以指定循环次数，形成循环执行的Pattern。



- ☐ times: 可以通过times指定固定的循环执行次数。

```
//指定循环触发4次
start.times(4);
//可以执行触发次数范围,让循环执行次数在该范围之内
start.times(2, 4);
```

- ☐ optional: 也可以通过optional关键字指定要么不触发要么触发指定的次数。

```
start.times(4).optional();
start.times(2, 4).optional();
```

- ☐ greedy: 可以通过greedy将Pattern标记为贪婪模式,在Pattern匹配成功的前提下,会尽可能多地触发。

```
//触发2、3、4次,尽可能重复执行
start.times(2, 4).greedy();
//触发0、2、3、4次,尽可能重复执行
start.times(2, 4).optional().greedy();
```

- ☐ oneOrMore: 可以通过oneOrMore方法指定触发一次或多次。

```
// 触发一次或者多次
start.oneOrMore();
//触发一次或者多次,尽可能重复执行
start.oneOrMore().greedy();
// 触发0次或者多次
start.oneOrMore().optional();
// 触发0次或者多次,尽可能重复执行
start.oneOrMore().optional().greedy();
```

- ☐ timesOrMore: 通过timesOrMore方法可以指定触发固定次数以上,例如执行两次以上。

```
// 触发两次或者多次
start.timesOrMore(2);
// 触发两次或者多次,尽可能重复执行
start.timesOrMore(2).greedy();
// 不触发或者触发两次以上,尽可能重复执行
start.timesOrMore(2).optional().greedy();
```

```
// 触发两次或者多次start.timesOrMore(2);// 触发两次或者多次,尽可能重复执行
start.timesOrMore(2).greedy();// 不触发或者触发两次以上,尽可能重复执行
start.timesOrMore(2).optional().greedy();
```

## 2. 定义条件

每个模式都需要指定触发条件,作为事件进入到该模式是否接受的判断依据,当事件中的数值满足了条件时,便进行下一步操作。在FlinkCFP中通过pattern.where()、pattern.or()及pattern.until()方法来为Pattern指定条件,且Pattern条件有Simple Conditions及Combining Conditions等类型。

- 简单条件: Simple Condition继承于Iterative Condition类,其主要根据事件中的字段信息进行判断,决定是否接受该事件。

```
// 把通话成功的事件挑选出来
start.where(_.getCallType == "success")
```

- 组合条件：组合条件是将简单条件进行合并，通常情况下也可以使用where方法进行条件的组合，默认每个条件通过AND逻辑相连。如果需要使用OR逻辑，直接使用or方法连接条件即可。

```
// 把通话成功，或者通话时长大于10秒的事件挑选出来
val start = Pattern.begin[StationLog]("start_pattern")
  .where(_.callType=="success")
  .or(_.duration>10)
```

- 终止条件：如果程序中使用了oneOrMore或者oneOrMore().optional()方法，则必须指定终止条件，否则模式中的规则会一直循环下去，如下终止条件通过until()方法指定。

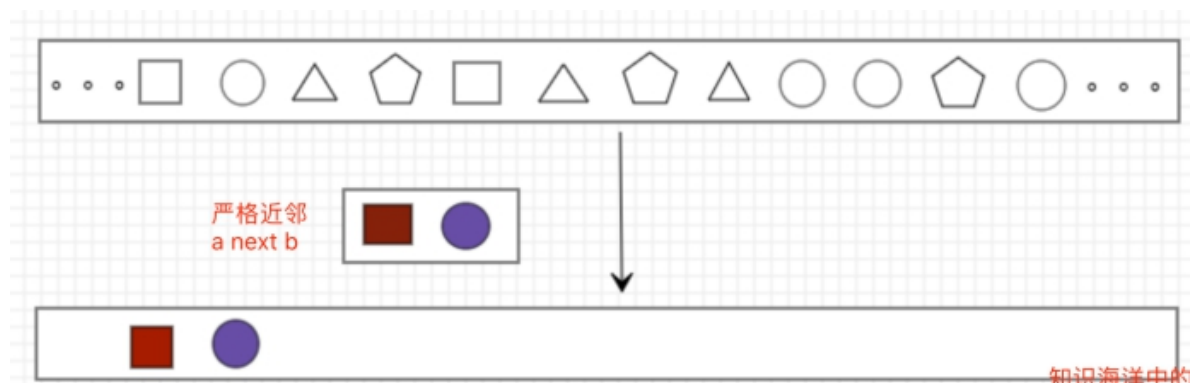
```
pattern.oneOrMore.until(_.callOut.startsWith("186"))
```

pattern.oneOrMore.until(\_.callOut.startsWith("186"))

### 3. 模式序列

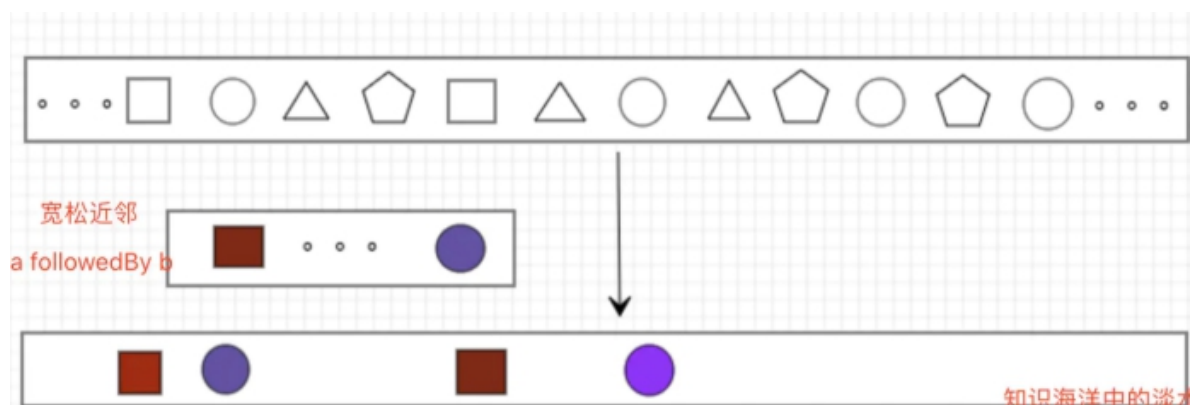
将相互独立的模式进行组合然后形成模式序列。模式序列基本的编写方式和独立模式一致，各个模式之间通过邻近条件进行连接即可，其中有严格邻近、宽松邻近、非确定宽松邻近三种邻近连接条件。

- ☐ 严格邻近：严格邻近条件中，需要所有的事件都按照顺序满足模式条件，不允许忽略任意不满足的模式。



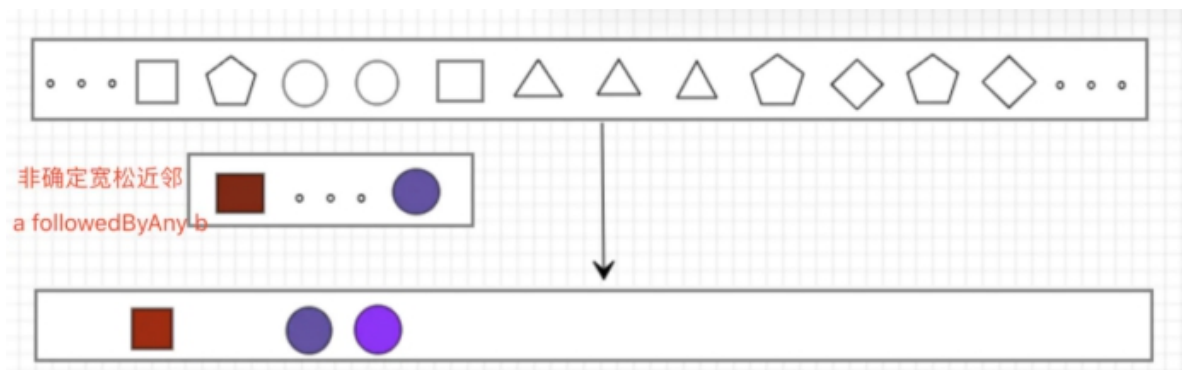
```
val strict: Pattern[Event] = start.next("middle").where(...)
```

- ☐ 宽松邻近：在宽松邻近条件下，会忽略没有成功匹配模式条件，并不会像严格邻近要求得那么高，可以简单理解为OR的逻辑关系。



```
val relaxed: Pattern[Event, _] = start.followedBy("middle").where(...)
```

- ❑ 非确定宽松邻近：和宽松邻近条件相比，非确定宽松邻近条件指在模式匹配过程中可以忽略已经匹配的条件。



```
val nonDetermin: Pattern[Event, _] = start.followedByAny("middle").where(...)
```

- ❑ 除以上模式序列外，还可以定义“不希望出现某种近邻关系”：

.notNext() —— 不想让某个事件严格紧邻前一个事件发生

.notFollowedBy() —— 不想让某个事件在两个事件之间发生

注意：

- 1、所有模式序列必须以 .begin() 开始
- 2、模式序列不能以 .notFollowedBy() 结束
- 3、“not”类型的模式不能被 optional 所修饰
- 4、此外，还可以为模式指定时间约束，用来要求在多长时间内匹配有效

//指定模式在10秒内有效pattern.within(Time.seconds(10)); //窗口

## 2) 模式检测

调用 CEP.pattern(), 给定输入流和模式，就能得到一个 PatternStream

```
//cep 做模式检测
val patternStream = CEP.pattern[EventLog](dataStream.keyBy(_.id), pattern)
```

## 3) 选择结果

得到PatternStream类型的数据集后，接下来数据获取都基于PatternStream进行。该数据集中包含了所有的匹配事件。目前在FlinkCEP中提供select和flatSelect两种方法从PatternStream提取事件结果事件。

1. 通过Select Funciton抽取正常事件

可以通过在PatternStream的Select方法中传入自定义Select Funciton完成对匹配事件的转换与输出。其中Select Funciton的输入参数为Map[String, Iterable[IN]]，Map中的key为模式序列中的Pattern名称，Value为对应Pattern所接受的事件集合，格式为输入事件的数据类型。

```
def selectFunction(pattern : Map[String, Iterable[IN]]): OUT = {
  //获取pattern中的startEvent
  val startEvent = pattern.get("start_pattern").get.next
  //获取Pattern中middleEvent
  val middleEvent = pattern.get("middle").get.next
  //返回结果
  OUT(startEvent, middleEvent)}

```

## 2. 通过Flat Select Function抽取正常事件

Flat Select Function和Select Function相似，不过Flat Select Function在每次调用可以返回任意数量的结果。因为Flat Select Function使用Collector作为返回结果的容器，可以将需要输出的事件都放置在Collector中返回。

```
def flatSelectFn(pattern : Map[String, Iterable[IN]], collector :
Collector[OUT]) = {    //获取pattern中startEvent
  val startEvent = pattern.get("start_pattern").get.next
  //获取Pattern中middleEvent
  val middleEvent = pattern.get("middle").get.next
  //并根据startEvent的value数量进行返回
  for (i <- 0 to startEvent.getValue) {
    collector.collect(OUT(startEvent, middleEvent))
  }
}

```

## 3. 通过Select Function抽取超时事件

如果模式中有within(time)，那么就很有可能有过时的数据存在，通过PatternStream.select方法分别获取超时事件和正常事件。首先需要创建OutputTag来标记超时事件，然后在PatternStream.select方法中使用OutputTag，就可以将超时事件从PatternStream中抽取出来。

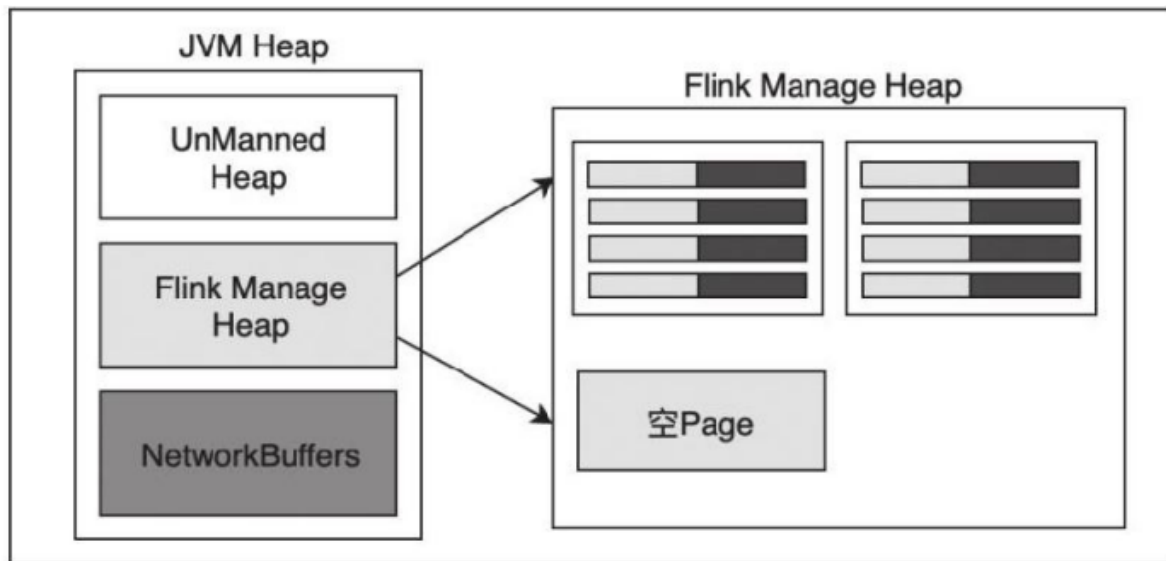
```
// 通过CEP.pattern方法创建
PatternStream val patternStream: PatternStream[Event] = CEP.pattern(input,
pattern) //创建OutputTag,并命名为timeout-output
val timeoutTag = OutputTag[String]("timeout-output")
//调用PatternStream select()并指定timeoutTag val result:
SingleOutputStreamOperator[NormalEvent] = patternStream.select(timeoutTag){
//超时事件获取
(pattern: Map[String, Iterable[Event]], timestamp: Long) =>
TimeoutEvent()//返回异常事件
} {
//正常事件获取
pattern: Map[String, Iterable[Event]] =>
NormalEvent()
//返回正常事件
}
//调用getSideOutput方法,并指定timeoutTag将超时事件输出val timeoutResult:
DataStream[TimeoutEvent] = result.getSideOutput(timeoutTag)

```

# Flink内存优化

在大数据领域，大多数开源框架（Hadoop、Spark、Flink）都是基于JVM运行，但是JVM的内存管理机制往往存在着诸多类似OutOfMemoryError的问题，主要是因为创建过多的对象实例而超过JVM的最大堆内存限制，却没有被有效回收掉，这在很大程度上影响了系统的稳定性，尤其对于大数据应用，面对大量的数据对象产生，仅仅靠JVM所提供的各种垃圾回收机制很难解决内存溢出的问题。在开源框架中有很多框架都实现了自己的内存管理，例如Apache Spark的Tungsten项目，在一定程度上减轻了框架对JVM垃圾回收机制的依赖，从而更好地使用JVM来处理大规模数据集。

Flink也基于JVM实现了自己的内存管理，将JVM根据内存区分为Unmanned Heap、Flink Managed Heap、Network Buffers三个区域。在Flink内部对Flink Managed Heap进行管理，在启动集群的过程中直接将堆内存初始化成Memory Pages Pool，也就是将内存全部以二进制数组的方式占用，形成虚拟内存使用空间。新创建的对象都是以序列化成二进制数据的方式存储在内存页面池中，当完成计算后数据对象Flink就会将Page置空，而不是通过JVM进行垃圾回收，保证数据对象的创建永远不会超过JVM堆内存大小，也有效地避免了因为频繁GC导致的系统稳定性问题。



## 1) JobManager配置

JobManager在Flink系统中主要承担管理集群资源、接收任务、调度Task、收集任务状态以及管理TaskManager的功能，JobManager本身并不直接参与数据的计算过程中，因此JobManager的内存配置项不是特别多，只要指定JobManager堆内存大小即可。

jobmanager.heap.size：设定JobManager堆内存大小，默认为1024MB。

## 2) TaskManager配置

TaskManager作为Flink集群中的工作节点，所有任务的计算逻辑均执行在TaskManager之上，因此对TaskManager内存配置显得尤为重要，可以通过以下参数配置对TaskManager进行优化和调整。

- taskmanager.heap.size：设定TaskManager堆内存大小，默认值为1024M，如果在Yarn的集群中，TaskManager取决于Yarn分配给TaskManager Container的内存大小，且Yarn环境下一般会减掉一部分内存用于Container的容错。
- taskmanager.jvm-exit-on-oom：设定TaskManager是否因为JVM发生内存溢出而停止，默认为false，当TaskManager发生内存溢出时，也不会导致TaskManager停止。
- taskmanager.memory.size：设定TaskManager内存大小，默认为0，如果不设定该值将会使用taskmanager.memory.fraction作为内存分配依据。
- taskmanager.memory.fraction：设定TaskManager堆中去除Network Buffers内存后的内存分配比例。该内存主要用于TaskManager任务排序、缓存中间结果等操作。例如，如果设定为0.8，则代表TaskManager保留80%内存用于中间结果数据的缓存，剩下20%的内存用于创建用户定义函数中的数据对象存储。注意，该参数只有在taskmanager.memory.size不设定的情况下才生效。

- taskmanager.memory.off-heap: 设置是否开启堆外内存供Managed Memory或者Network Buffers使用。
- taskmanager.memory.preallocate: 设置是否在启动TaskManager过程中直接分配TaskManager管理内存。
- taskmanager.numberOfTaskSlots: 每个TaskManager分配的slot数量。

### 3. Flink的网络缓存优化

Flink将JVM堆内存切分为三个部分，其中一部分为Network Buffers内存。Network Buffers内存是Flink数据交互层的关键内存资源，主要目的是缓存分布式数据处理过程中的输入数据。。通常情况下，比较大的Network Buffers意味着更高的吞吐量。如果系统出现“Insufficient number of network buffers”的错误，一般是因为Network Buffers配置过低导致，因此，在这种情况下需要适当调整TaskManager上Network Buffers的内存大小，以使得系统能够达到相对较高的吞吐量。

目前Flink能够调整Network Buffer内存大小的方式有两种：一种是通过直接指定Network Buffers内存数量的方式，另外一种是通过配置内存比例的方式。

#### 1) 设定Network Buffer内存数量\*\*（过时了）\*\*

直接设定Network Buffer数量需要通过如下公式计算得出：

$$\text{NetworkBuffersNum} = \text{total-degree-of-parallelism} * \text{intra-node-parallelism} * n$$

其中total-degree-of-parallelism表示每个TaskManager的总并发数量，intra-node-parallelism表示每个TaskManager输入数据源的并发数量，n表示在预估计算过程中Repartitioning或Broadcasting操作并行的数量。intra-node-parallelism通常情况下与Task-Manager的所占有的CPU数一致，且Repartitioning和Broadcasting一般下不会超过4个并发。可以将计算公式转化如下：

$$\text{NetworkBuffersNum} = 8^2 * \text{TM} * 4$$

其中slots-per-TM是每个TaskManager上分配的slots数量，TMs是TaskManager的总数量。对于一个含有20个TaskManager，每个TaskManager含有8个Slot的集群来说，总共需要的Network Buffer数量为 $8^2 * 20 = 1280$ 个，因此集群中配置Network Buffer内存的大小约为160M较为合适。

计算完Network Buffer数量后，可以通过添加如下两个参数对Network Buffer内存进行配置。其中segment-size为每个Network Buffer的内存大小，默认为32KB，一般不需要修改，通过设定numberOfBuffers参数以达到计算出的内存大小要求。

- taskmanager.network.numberOfBuffers: 指定Network堆栈Buffer内存块的数量。
- taskmanager.memory.segment-size.: 内存管理器和Network栈使用的内存Buffer大小，默认为32KB。

#### 2) 设定Network内存比例\*\*（推荐）\*\*

从1.3版本开始，Flink就提供了通过指定内存比例的方式设置Network Buffer内存大小。

- taskmanager.network.memory.fraction: JVM中用于Network Buffers的内存比例。
- taskmanager.network.memory.min: 最小的Network Buffers内存大小，默认为64MB。
- taskmanager.network.memory.max: 最大的Network Buffers内存大小，默认1GB。
- taskmanager.memory.segment-size: 内存管理器和Network栈使用的Buffer大小，默认为32KB。

## 常见面试问题汇总

## 1. 面试题一：应用架构

---

问题：公司怎么提交的实时任务，有多少 Job Manager？

解答：1. 我们使用 yarn -cluster 模式提交任务。每次提交都会创建一个新的 Flink 集群，为每一个 job 提供一个 yarn-session，Job 之间互相独立，互不影响，方便管理。任务执行完成之后创建的集群也会消失。线上命令脚本如下：

```
bin/yarn-session.sh -n 7 -s 8 -jm 3072 -tm 32768 -qu root.. -nm - -d
```

```
bin/flink run -c . . . . jar 路径
```

其中申请 7 个 taskManager，每个 8 核，每个 taskmanager 有 32768M 内存。

2. 集群默认只有一个 Job Manager。但为了防止单点故障，我们配置了高可用。我们公司一般配置一个主 Job Manager，两个备用 Job Manager，然后结合 ZooKeeper 的使用，来达到高可用。

## 2. 面试题二：压测和监控

---

问题：怎么做压力测试和监控？

解答：我们一般碰到的压力来自以下几个方面：

一，产生数据流的速度如果过快，而下游的算子消费不过来的话，会产生背压。背压的监控可以使用 Flink Web UI (localhost:8081) 来可视化监控，一旦报警就能知道。一般情况下背压问题的产生可能是由于 sink 这个操作符没有优化好，做一下优化就可以了。比如如果是写入 Elasticsearch，那么可以改成批量写入，可以调大 Elasticsearch 队列的大小等等策略。

二，设置 watermark 的最大延迟时间这个参数，如果设置的过大，可能会造成内存的压力。可以设置最大延迟时间小一些，然后把迟到元素发送到侧输出流中去。晚一点更新结果。或者使用类似于 RocksDB 这样的状态后端，RocksDB 会开辟堆外存储空间，但 IO 速度会变慢，需要权衡。

三，还有就是滑动窗口的长度如果过长，而滑动距离很短的话，Flink 的性能会下降的很厉害。我们主要通过时间分片的方法，将每个元素只存入一个“重叠窗口”，这样就可以减少窗口处理中状态的写入。参见链接：[https://www.infoq.cn/article/slhs\\_qY6HCpMQNbITl9M](https://www.infoq.cn/article/slhs_qY6HCpMQNbITl9M)

四，状态后端使用 RocksDB，还没有碰到被撑爆的问题。

## 3. 面试题三：为什么用 Flink

---

问题：为什么使用 Flink 替代 Spark？

解答：主要考虑的是 flink 的低延迟、高吞吐量和对流式数据应用场景更好的支持；另外，flink 可以很好地处理乱序数据，而且可以保证 exactly-once 的状态一致性。详见上面文档，有 Flink 和 Spark 的详细对比。

## 4. 面试题四：checkpoint 的存储

---

问题：Flink 的 checkpoint 存在哪里？

解答：可以是内存，文件系统，或者 RocksDB。详见上面文档。

## 5. 面试题五：exactly-once 的保证

---

问题：如果下级存储不支持事务，Flink 怎么保证 exactly-once？



解答：整体端到端的exactly-once对sink要求比较高，具体实现主要有幂等写入和事务性写入两种方式。幂等写入的场景依赖于业务逻辑，更常见的是用事务性写入。而事务性写入又有预写日志（WAL）和两阶段提交（2PC）两种方式。

如果外部系统不支持事务，那么可以用预写日志的方式，把结果数据先当成状态保存，然后在收到checkpoint 完成的通知时，一次性写入 sink 系统。

## 6. 面试题六：状态机制

---

问题：说一下 Flink 状态机制？

解答：Flink内置的很多算子，包括源source，数据存储sink都是有状态的。在Flink中，状态始终与特定算子相关联。Flink会以checkpoint的形式对各个任务的状态进行快照，用于保证故障恢复时的状态一致性。Flink通过状态后端来管理状态和checkpoint的存储，状态后端可以有不同的配置选择。详见上面文档。

## 7. 面试题七：海量key去重

---

问题：怎么去重？考虑一个实时场景：双十一场景，滑动窗口长度为 1 小时，滑动距离为 10 秒钟，亿级用户，怎样计算 UV？

解答：使用类似于 scala 的 set 数据结构或者 redis 的 set 显然是不行的，因为可能有上亿个 Key，内存放不下。所以可以考虑使用布隆过滤器（Bloom Filter）来去重。

## 8. 面试题八：checkpoint与spark比较

---

问题：Flink 的 checkpoint 机制对比 spark 有什么不同和优势？

解答：spark streaming 的 checkpoint 仅仅是针对 driver 的故障恢复做了数据和元数据的 checkpoint。而 flink 的 checkpoint 机制 要复杂了很多，它采用的是轻量级的分布式快照，实现了每个算子的快照，及流动中的数据的快照。详见上面文档及文章链接：<https://cloud.tencent.com/developer/article/1189624>

## 9. 面试题九：watermark机制

---

问题：请详细解释一下Flink的Watermark机制。

解答：Watermark本质是Flink中衡量EventTime进展的一个机制，主要用来处理乱序数据。详见上面文档。

## 10. 面试题十：exactly-once如何实现

---

问题：Flink中exactly-once语义是如何实现的，状态是如何存储的？

解答：Flink依靠checkpoint机制来实现exactly-once语义，如果要想实现端到端的exactly-once，还需要外部source和sink满足一定的条件。状态的存储通过状态后端来管理，Flink中可以配置不同的状态后端。

## 11. 面试题十一：CEP

---

问题：Flink CEP 编程中当状态没有到达的时候会将数据保存在哪里？



解答：在流式处理中，CEP 当然是要支持 EventTime 的，那么相对应的也要支持数据的迟到现象，也就是watermark的处理逻辑。CEP对未匹配成功的事件序列的处理，和迟到数据是类似的。在 Flink CEP的处理逻辑中，状态没有满足的和迟到的数据，都会存储在一个Map数据结构中，也就是说，如果我们限定判断事件序列的时长为5分钟，那么内存中就会存储5分钟的数据，这在我看来，也是对内存的极大损伤之一。

## 12. 面试题十二：三种时间语义

---

问题：Flink 三种时间语义是什么，分别说出应用场景？

解答：

1. Event Time：这是实际应用最常见的时间语义，具体见文档第七章。
2. Processing Time：没有事件时间的情况下，或者对实时性要求超高的情况下。
3. Ingestion Time：存在多个Source Operator的情况下，每个Source Operator可以使用自己本地系统时钟指派 Ingestion Time。后续基于时间相关的各种操作，都会使用数据记录中的 Ingestion Time。

## 13. 面试题十三：数据高峰的处理

---

问题：Flink程序在面对数据高峰期时如何处理？

解答：使用大容量的 Kafka 把数据先放到消息队列里面作为数据源，再使用 Flink 进行消费，不过这样会影响到一点实时性。