

# Flink DataStream API 中的多面手—— Process Function

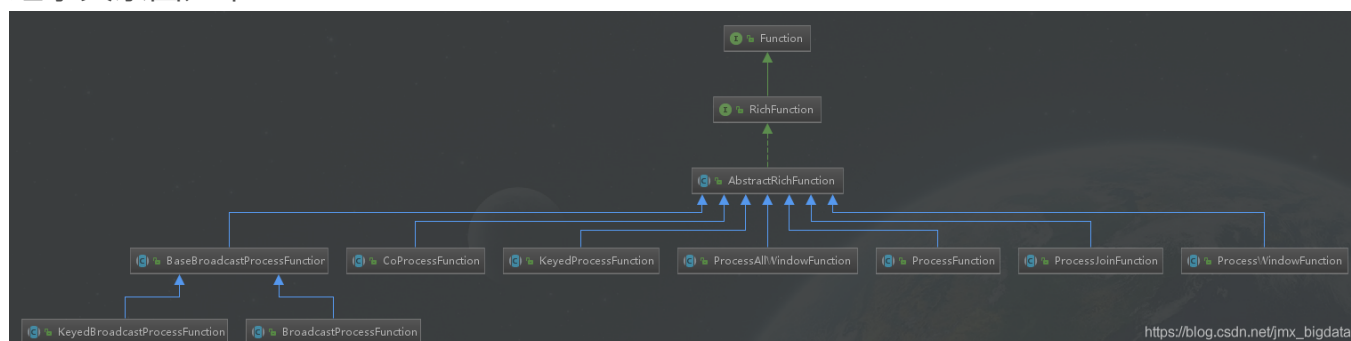
在[Flink的时间与watermarks详解](#)这篇文章中，阐述了Flink的时间与水位线的相关内容。你可能不禁要发问，该如何访问时间戳和水位线呢？首先通过普通的DataStream API是无法访问的，需要借助Flink提供的一个底层的API——Process Function。Process Function不仅能够访问时间戳与水位线，而且还可以注册在将来的某个特定时间触发的计时器(timers)。除此之外，还可以将数据通过Side Outputs发送到多个输出流中。这样以来，可以实现数据分流的功能，同时也是处理迟到数据的一种方式。下面我们将从源码入手，结合具体的使用案例来说明该如何使用Process Function。

## 简介

Flink提供了很多Process Function，每种Process Function都有各自的功能，这些Process Function主要包括：

- ProcessFunction
- KeyedProcessFunction
- CoProcessFunction
- ProcessJoinFunction
- ProcessWindowFunction
- ProcessAllWindowFunction
- BaseBroadcastProcessFunction
  - KeyedBroadcastProcessFunction
- BroadcastProcessFunction

继承关系图如下：



从上面的继承关系中可以看出，都实现了RichFunction接口，所以支持使用 `open()`、`close()`、`getRuntimeContext()` 等方法的调用。从名字上可以看出，这些函数都有不同的适用场景，但是基本的功能是类似的，下面会以KeyedProcessFunction为例来讨论这些函数的通用功能。

## 源码

### KeyedProcessFunction

```
1  /**
2   * 处理KeyedStream流的低级API函数
3   * 对于输入流中的每个元素都会触发调用processElement方法。该方法会产生0个或多个输出。
4   * 其实现类可以通过Context访问数据的时间戳和计时器(timers)。当计时器(timers)触发时
5   * ，会回调onTimer方法。
6   * onTimer方法会产生0个或者多个输出，并且会注册一个未来的计时器。
7   *
8   * 注意：如果要访问keyed state和计时器(timers)，必须在KeyedStream上使用KeyedPro
9   * cessFunction。
10  * 另外，KeyedProcessFunction的父类AbstractRichFunction实现了RichFunction接
11  * 口，所以，可以使用
12  * open(), close()及getRuntimeContext()方法。
13  *
14  * @param <K> key的类型
15  * @param <I> 输入元素的数据类型
16  * @param <O> 输出元素的数据类型
17  */
18  @PublicEvolving
19  public abstract class KeyedProcessFunction<K, I, O> extends AbstractRich
20  Function {
21
22      private static final long serialVersionUID = 1L;
23      /**
24       * 处理输入流中的每个元素
25       * 该方法会输出0个或者多个输出，类似于FlatMap的功能
26       * 除此之外，该方法还可以更新内部状态或者设置计时器(timer)
27       * @param value 输入元素
28       * @param ctx Context，可以访问输入元素的时间戳，并其可以获取一个时间服
29       * 务器(TimerService)，用于注册计时器(timers)并查询时间
30       * Context只有在processElement被调用期间有效。
31       * @param out 返回的结果值
32       * @throws Exception
33       */
34      public abstract void processElement(I value, Context ctx, Collec
35      tor<O> out) throws Exception;
36
37      /**
38       * 是一个回调函数，当在TimerService中注册的计时器(timers)被触发时，会回调
39
```

```

40 该函数
41      * @param timestamp 触发计时器(timers)的时间戳
42      * @param ctx OnTimerContext, 允许访问时间戳, TimeDomain枚举类提供了
43 两种时间类型:
44      * EVENT_TIME与PROCESSING_TIME
45      * 并其可以获取一个时间服务器(TimerService), 用于注册计时器(timers)并查
46 询时间
47      * OnTimerContext只有在onTimer方法被调用期间有效
48      * @param out 结果输出
49      * @throws Exception
50      */
51      public void onTimer(long timestamp, OnTimerContext ctx, Collecto
52 r<0> out) throws Exception {}
53      /**
54      * 仅仅在processElement()方法或者onTimer方法被调用期间有效
55      */
56      public abstract class Context {
57
58          /**
59          * 当前被处理元素的时间戳, 或者是触发计时器(timers)时的时间戳
60          * 该值可能为null, 比如当程序中设置的时间语义为: TimeCharacteri
61 stic#ProcessingTime
62          * @return
63          */
64          public abstract Long timestamp();
65
66          /**
67          * 访问时间和注册的计时器(timers)
68          * @return
69          */
70          public abstract TimerService timerService();
71
72          /**
73          * 将元素输出到side output (侧输出)
74          * @param outputTag 侧输出的标记
75          * @param value 输出的记录
76          * @param <X>
77          */
78          public abstract <X> void output(OutputTag<X> outputTag,
79 X value);
80          /**
81          * 获取被处理元素的key
82          * @return
83          */
84          public abstract K getCurrentKey();
85      }
86      /**
87      * 当onTimer方法被调用时, 才可以使用OnTimerContext
88      */

```

```

89 |         public abstract class OnTimerContext extends Context {
            /**
            * 触发计时器(timers)的时间类型, 包括两种: EVENT_TIME与PROCESS
            ING_TIME
            * @return
            */
            public abstract TimeDomain timeDomain();
            /**
            * 获取触发计时器(timer)元素的key
            * @return
            */
            @Override
            public abstract K getCurrentKey();
        }
    }

```

上面的源码中, 主要有两个方法, 分析如下:

- processElement(l value, Context ctx, Collector out)

该方法会对流中的每条记录都调用一次, 输出0个或者多个元素, 类似于FlatMap的功能, 通过Collector将结果发出。除此之外, 该函数有一个Context 参数, 用户可以通过Context 访问时间戳、当前记录的key值以及TimerService(关于TimerService, 下面会详细解释)。另外还可以使用output方法将数据发送到side output, 实现分流或者处理迟到数据的功能。

- onTimer(long timestamp, OnTimerContext ctx, Collector out)

该方法是一个回调函数, 当在TimerService中注册的计时器(timers)被触发时, 会回调该函数。其中 `@param timestamp` 参数表示触发计时器(timers)的时间戳, Collector可以将记录发出。细心的你可能会发现, 这两个方法都有一个上下文参数, 上面的方法传递的是Context 参数, onTimer方法传递的是OnTimerContext参数, 这两个参数对象可以实现相似的功能。OnTimerContext还可以返回触发计时器的时间域(EVENT\_TIME与PROCESSING\_TIME)。

## TimerService

在KeyedProcessFunction源码中, 使用TimerService来访问时间和计时器, 下面来看一下源码:

```

1 | @PublicEvolving
2 | public interface TimerService {
3 |     String UNSUPPORTED_REGISTER_TIMER_MSG = "Setting timers is only
4 |     supported on a keyed streams.";
5 |     String UNSUPPORTED_DELETE_TIMER_MSG = "Deleting timers is only s

```

```

6  upported on a keyed streams.";
7      // 返回当前的处理时间
8      long currentProcessingTime();
9      // 返回当前event-time水位线(watermark)
10     long currentWatermark();
11
12     /**
13      * 注册一个计时器(timers), 当processing time的时间等于该计时器时钟时会被
14 调用
15      * @param time
16      */
17     void registerProcessingTimeTimer(long time);
18
19     /**
20      * 注册一个计时器(timers), 当event time的水位线(watermark) 到达该时间时会
21 被触发
22      * @param time
23      */
24     void registerEventTimeTimer(long time);
25
26     /**
27      * 根据给定的触发时间(trigger time)来删除processing-time计时器
28      * 如果这个timer不存在, 那么该方法不会起作用,
29      * 即该计时器(timer)之前已经被注册了, 并且没有过时
30      *
31      * @param time
32      */
33     void deleteProcessingTimeTimer(long time);
34
35     /**
36      * 根据给定的触发时间(trigger time)来删除event-time 计时器
37      * 如果这个timer不存在, 那么该方法不会起作用,
38      * 即该计时器(timer)之前已经被注册了, 并且没有过时
39      *
39      * @param time
40      */
41     void deleteEventTimeTimer(long time);
42 }

```

TimerService提供了以下几种方法:

- currentProcessingTime()

返回当前的处理时间

- currentWatermark()

返回当前event-time水位线(watermark)时间戳

- registerProcessingTimeTimer(long time)

针对当前key，注册一个processing time计时器(timers)，当processing time的时间等于该计时器时钟时会被调用

- registerEventTimeTimer(long time)

针对当前key，注册一个event time计时器(timers)，当水位线时间戳大于等于该计时器时钟时会被调用

- deleteProcessingTimeTimer(long time)

针对当前key，删除一个之前注册过的processing time计时器(timers)，如果这个timer不存在，那么该方法不会起作用

- deleteEventTimeTimer(long time)

针对当前key，删除一个之前注册过的event time计时器(timers)，如果这个timer不存在，那么该方法不会起作用

当计时器触发时，会回调onTimer()函数，系统对于ProcessElement()方法和onTimer()方法的调用是同步的

注意:上面的源码中有两个Error 信息,这就说明计时器只能在keyed streams上使用，常见的用途是在某些key值不在使用后清除keyed state，或者实现一些基于时间的自定义窗口逻辑。如果要在一个非KeyedStream上使用计时器，可以使用KeySelector返回一个固定的分区值(比如返回一个常数)，这样所有的数据只会发送到一个分区。

## 使用案例

下面将使用Process Function的side output功能进行分流处理，具体代码如下：

```
1 public class ProcessFunctionExample {
2
3     // 定义side output标签
4     static final OutputTag<UserBehaviors> buyTags = new OutputTag<UserBeh
5 aviors>("buy") {
6     };
7     static final OutputTag<UserBehaviors> cartTags = new OutputTag<UserBe
8 haviors>("cart") {
9     };
10    static final OutputTag<UserBehaviors> favTags = new OutputTag<UserBeh
11 aviors>("fav") {
12    };
```

```

13     static class SplitStreamFunction extends ProcessFunction<UserBehaviors, UserBehaviors> {
14
15
16         @Override
17         public void processElement(UserBehaviors value, Context ctx, Collector<UserBehaviors> out) throws Exception {
18             switch (value.behavior) {
19                 case "buy":
20                     ctx.output(buyTags, value);
21                     break;
22                 case "cart":
23                     ctx.output(cartTags, value);
24                     break;
25                 case "fav":
26                     ctx.output(favTags, value);
27                     break;
28                 default:
29                     out.collect(value);
30             }
31         }
32     }
33
34     public static void main(String[] args) throws Exception {
35         StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment().setParallelism(1);
36
37         // 模拟数据源[userId,behavior,product]
38         SingleOutputStreamOperator<UserBehaviors> splitStream = env.fromElements(
39             new UserBehaviors(1L, "buy", "iphone"),
40             new UserBehaviors(1L, "cart", "huawei"),
41             new UserBehaviors(1L, "buy", "logi"),
42             new UserBehaviors(1L, "fav", "oppo"),
43             new UserBehaviors(2L, "buy", "huawei"),
44             new UserBehaviors(2L, "buy", "onemore"),
45             new UserBehaviors(2L, "fav", "iphone")).process(new SplitStreamFunction());
46
47         // 获取分流之后购买行为的数据
48         splitStream.getSideOutput(buyTags).print("data_buy");
49         // 获取分流之后加购行为的数据
50         splitStream.getSideOutput(cartTags).print("data_cart");
51         // 获取分流之后收藏行为的数据
52         splitStream.getSideOutput(favTags).print("data_fav");
53
54         env.execute("ProcessFunctionExample");
55     }
56 }

```

## 总结

本文首先介绍了Flink提供的几种底层Process Function API，这些API可以访问时间戳和水位线，同时支持注册一个计时器，进行调用回调函数onTimer()。接着从源码的角度解读了这些API的共同部分，详细解释了每个方法的具体含义和使用方式。最后，给出了一个Process Function常见使用场景案例，使用其实现分流处理。除此之外，用户还可以使用这些函数，通过注册计时器，在回调函数中定义处理逻辑，使用非常的灵活。