

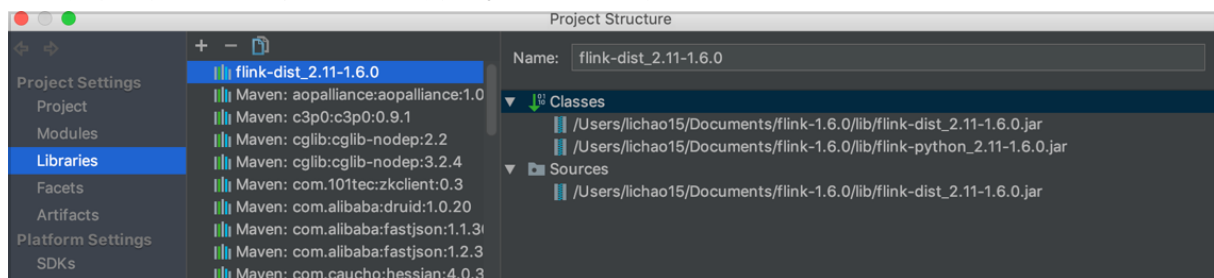
Flink学习

#学习

基础

相关概念

- flink程序能实现在分布式的结合上进行各种转换操作，集合通常来自订阅的来源（文件，kafka,local,in-memory），结果被返回到sinks里（大多数写入分布式文件系统，或者标准输出）
- **flink 的operator** 会尽量按照链式方式分在同一个**task slot**里面，以优化执行
- DataSet and DataStream
 - DataSet和DataStream在flink中都代表一种数据结构，是不可变且包含重复记录的集合。区别在于DataSet是有限的集合，而DataStream是无界的
- flink 配置interlij ideal 在本地运行调试
 - 只需要将flink依赖的包引入项目中即可启动项目



- 讲解Flink怎么序列化objects，怎么分配内存 [Apache Flink: Juggling with Bits and Bytes](#)

DataStream

- [Apache Flink 1.7 Documentation: Flink DataStream API Programming Guide](#)
- datasource（数据源）：
 - File-based: readTextFile, readFile...
 - Socket-based: socketTextStream
 - Collection-based: fromCollection, fromElements
 - custom: addSource, FlinkKafkaConsumer08 or other connectors

DataSet

- [Apache Flink 1.7 Documentation: Flink DataSet API Programming Guide](#)
- 不得不会的Flink Dataset的Delta 迭代操作

savepoint

- [Apache Flink 1.7 Documentation: Savepoints](#)
- Savepoints are created, owned, and deleted by the user.

- 目前savepoint和checkpoint实现和format方式都相同（除了checkpoint选择了rocksdb作为state backend，这样format会有些微不同）
- Operations:
 - Triggering Savepoints: FsStateBackend or RocksDBStateBackend:
 - Trigger a Savepoint
 - Cancel Job with Savepoint
 - `bin/flink cancel -s [:targetDirectory] :jobId`
 - Resuming from Savepoints
 - `$ bin/flink run -s :savepointPath [:runArgs]`
 - Disposing Savepoints
 - `$ bin/flink savepoint -d :savepointPath`

checkpoint

- [Apache Flink 1.7 Documentation: Checkpoints](#)
- 生命周期是由Flink管理，checkpoint的管理，创建以及释放统一通过Flink，而不需要用户干预
- **Checkpoints are usually dropped（随应用退出被清除）** after the job was terminated by the user (except if explicitly configured as retained Checkpoints)

Retained Checkpoints

Checkpoints are by default not retained and are only used to resume a job from failures. They are deleted when a program is cancelled. You can, however, configure periodic checkpoints to be retained. Depending on the configuration these *retained* checkpoints are *not* automatically cleaned up when the job fails or is canceled. This way, you will have a checkpoint around to resume from if your job fails.

```
CheckpointConfig config = env.getCheckpointConfig();
config.enableExternalizedCheckpoints(ExternalizedCheckpointCleanup.RETAIN_ON_CANCELLATION);
```

The ExternalizedCheckpointCleanup mode configures what happens with checkpoints when you cancel the job:

- **ExternalizedCheckpointCleanup.RETAIN_ON_CANCELLATION:** Retain the checkpoint when the job is cancelled. Note that you have to manually clean up the checkpoint state after cancellation in this case.
- **ExternalizedCheckpointCleanup.DELETE_ON_CANCELLATION:** Delete the checkpoint when the job is cancelled. The checkpoint state will only be available if the job fails.

- retained checkpoints 不会随着应用清除，所以需要手动清理，与savepoint有一点区别，不支持rescaling
- resuming from checkpoints

```
bin/flink run -s :checkpointMetaDataPath [:runArgs]
```

- checkpoint 优化 [Apache Flink 1.7 Documentation: Tuning Checkpoints and Large State](#)
 - state 双写：一份在distributed storage(HDFS)；一份在local
 - task-local recovery: 默认是关闭的状态,可以通过 `state.backend.local-recovery` 打开

Barriers

- [Apache Flink 1.8-SNAPSHOT Documentation: Data Streaming Fault Tolerance](#)

Window, waterMark, Trigger

- [Window, waterMark, Trigger介绍- 简书](#)
- window
 - **非常重要还未研究** [Windows](#)
 - 滚动窗口：分配器将每个元素分配到一个指定窗口大小的窗口中，并且不会重叠；
`TumblingEventTimeWindows.of(Time.seconds(5))`
 - 滑动窗口：滑动窗口分配器将元素分配到固定长度的窗口中，与滚动窗口类似，窗口的大小由窗口大小参数来配置，另一个窗口滑动参数控制滑动窗口开始的频率；因此可能出现窗口重叠，如果滑动参数小于滚动参数的话；
`SlidingEventTimeWindows.of(Time.seconds(10), Time.seconds(5))`
 - 会话窗口：通过session活动来对元素进行分组，跟滚动窗口和滑动窗口相比，不会有重叠和固定的开始时间和结束时间的情况。当他在一个固定的时间周期内不再收到元素，即非活动间隔产生，那么窗口就会关闭；
 - 一个session窗口通过一个session间隔来配置，这个session间隔定义了非活跃周期的长度。当这个非活跃周期产生，那么当前的session将关闭并且后续的元素将被分配到新的session窗口中去。如：
`EventTimeSessionWindows.withGap(Time.minutes(10))`
- 触发器(Triggers)
 - 触发器定义了一个窗口何时被窗口函数处理
 - `EventTimeTrigger`
 - `ProcessingTimeTrigger`
 - `CountTrigger`
 - `PurgingTrigger`
- 驱逐器(Evictors)

任务提交与停止姿势

- 任务提交
 - **启动命令详解** :[Apache Flink 1.7 Documentation: YARN Setup](#)
 - 参数

Usage:

Required

`-n, --container <arg>` Number of YARN container to allocate (=Number of

Task Managers)

Optional

`-D <arg>` Dynamic properties

`-d, --detached` Start detached

<code>-jm,--jobManagerMemory <arg></code>	Memory for JobManager Container with optional unit (default: MB)
<code>-nm,--name</code>	Set a custom name for the application on YARN
<code>-q,--query</code>	Display available YARN resources (memory, cores)
<code>-qu,--queue <arg></code>	Specify YARN queue.
<code>-s,--slots <arg></code>	Number of slots per TaskManager
<code>-tm,--taskManagerMemory <arg></code>	Memory per TaskManager Container with optional unit (default: MB)
<code>-z,--zookeeperNamespace <arg></code>	Namespace to create the Zookeeper sub-paths for HA mode

- 提交到yarn-cluster上需要以 `y` 或者 `yarn` 作为前缀；如: `ynm=nm`

```
flink run -c com.jacobs.jobs.realtime.wordcount.WindowWordCount target/real-time-jobs-1.0.0-SNAPSHOT.jar
```

```
flink run -m yarn-cluster -ynm SinkToKafkaStream -yn 4 -yjm 1024m -ytm 4096m -ys 4 -yqu feed.prod -c com.weibo.api.feed.dm.stream.TestFlinkStream /data1/dm-flink/feed-dm-flink-1.0.4-SNAPSHOT.jar
```

```
flink run -m yarn-cluster -ynm SinkToKafkaStream -yn 2 -yjm 1024m -ytm 4096m -ys 2 -yqu feed.prod -c com.weibo.api.feed.dm.stream.SinkToKafkaStream /data1/dm-flink/feed-dm-flink-1.0.4-SNAPSHOT.jar
```

• 停止任务

- **关闭或重启flink程序不能直接kill掉**，这样会导致flink来不及制作checkpoint，而应该调用flink提供的cancel语意

//重启正确姿势，with savepoint

1. 调用cancel, cancel之前先触发savepoint

```
bin/flink cancel -s [:targetDirectory] :jobId -yid: yarnAppId
```

```
例子: flink cancel -s hdfs://vcp-yz-nameservice1/user/hcp/hcpsys/feed/flink-checkpoints/test-user-logs_97b4e67859af4bfb1b597355f1c846f3 -yid application_1542801635735_2121
```

2. 从savepoint中恢复flink程序

```
bin/flink run -s :savepointPath [:runArgs]
```

```
例子: flink run -s hdfs://vcp-yz-nameservice1/user/hcp/hcpsys/feed/flink-  
checkpoints/test-user-logs/savepoint-97b4e6-22dd5890dd0c -m yarn-cluster -ynm  
TestSinkUserLogStream -yn 4 -yjm 1024m -ytm 4096m -ys 4 -yqu feed.prod -c  
com.weibo.api.feed.dm.stream.TestFlinkStream /data1/dm-flink/feed-dm-  
flink-1.0.4-SNAPSHOT.jar
```

3. 查看运行中的任务

```
flink list -m yarn-cluster -yid application_1548467718478_0002 -r
```

运行模式

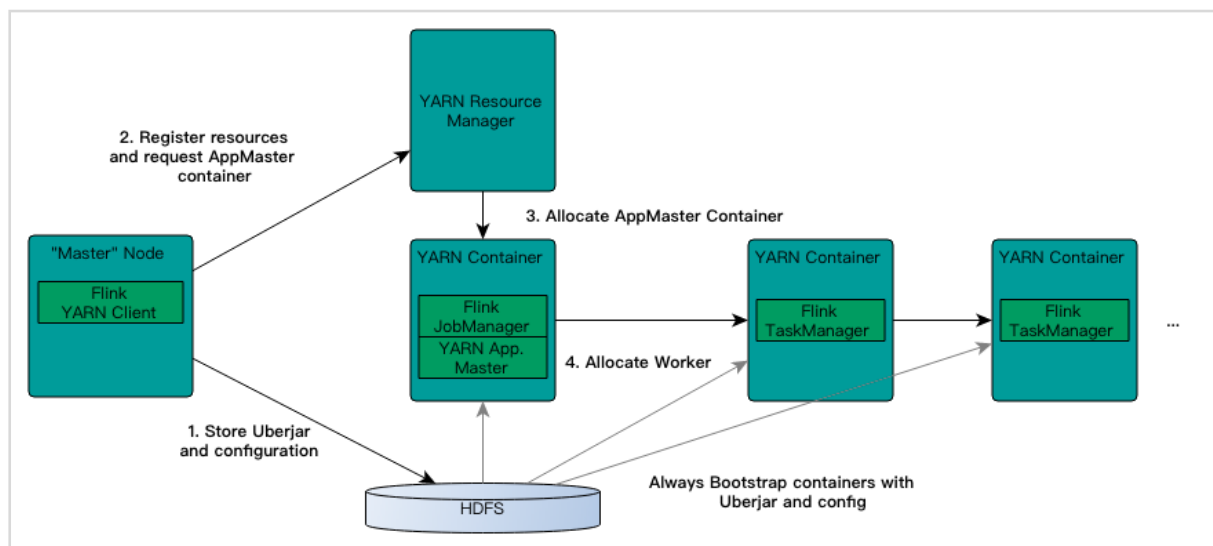
Standalone

- standalone 启动cluster

```
/usr/local/flink-1.6.0/bin;./start-cluster.sh
```

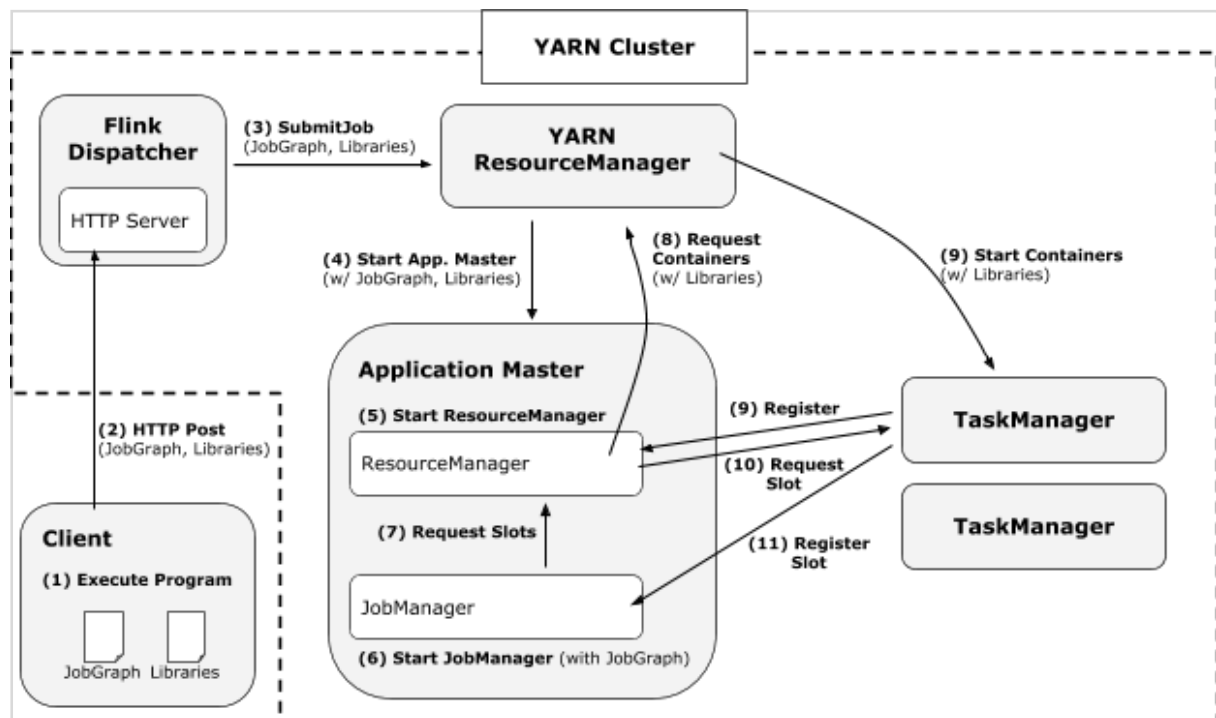
On Yarn Cluster

- Apache Flink 1.7 Documentation: YARN
- 参考文章[Flink1.6 - flink on yarn分布式部署架构 - 深山含笑](#)
- 架构图

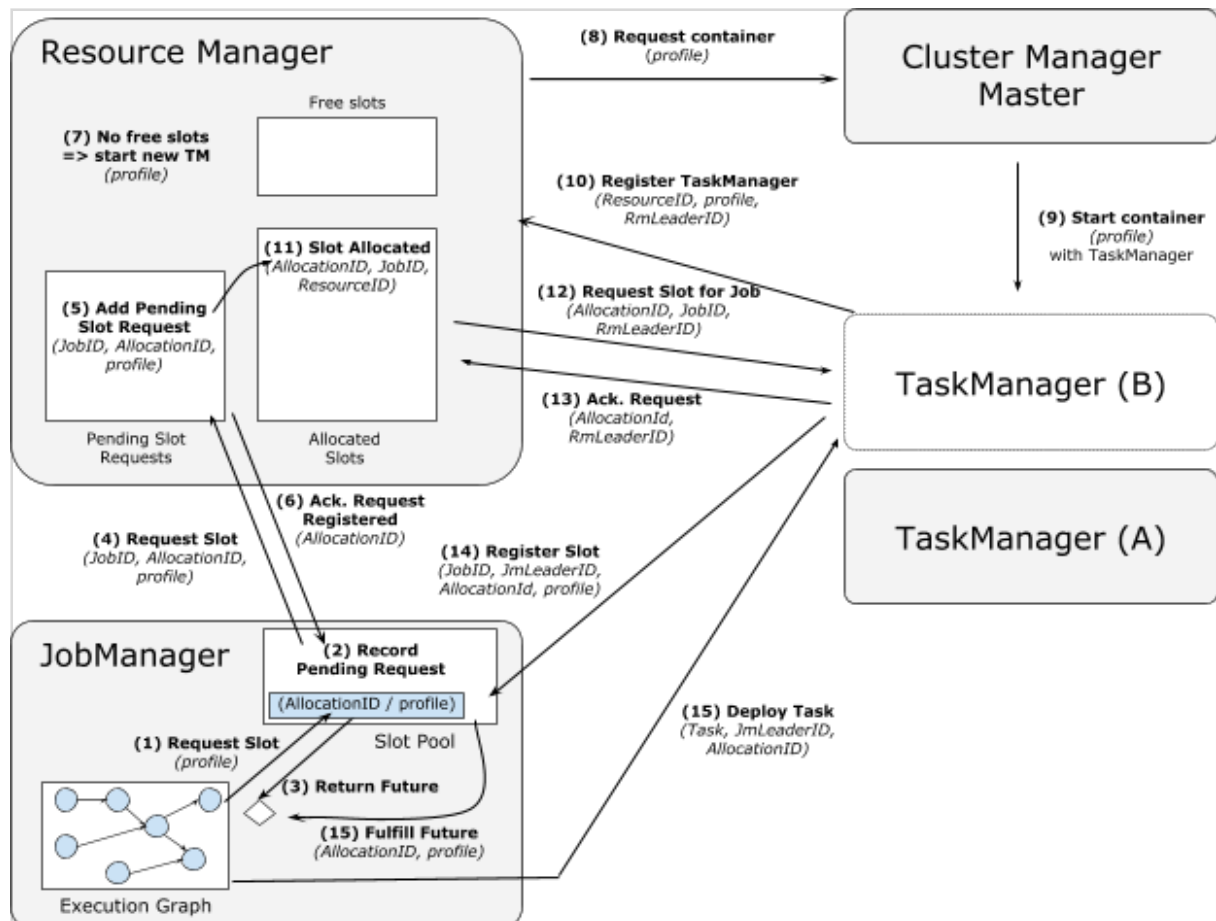


- JobManager 和 ApplicationMaster 运行在同一个JVM里
- on yarn 两种模式
 - session模式：允许运行多个作业在同一个Flink集群中，代价是作业之间没有资源隔离（同一个TM中可能跑多个不同作业的task）

- per-job模式（生产环境）：per-job模式是指在yarn上为每一个Flink作业都分配一个单独的Flink集群，这样就解决了不同作业之间的资源隔离问题
- 摘录参考文章 相比旧的Flink-on-YARN架构（Flink 1.5之前），新的yarn架构带来了以下的优势：
 - client可以直接在yarn上面启动一个作业，不在像以前需要先启动一个固定大小的Flink集群然后把作业提交到这个Flink集群上
 - 按需申请容器（指被同一个作业的不同算子所使用的容器可以有不同的CPU/Memory配置），没有被使用的容器将会被释放



- slot资源申请/分配流程分析
- 请求新TaskExecutor的slot分配



- ResourceManager挂掉：不会挂掉task,不断尝试重新注册ResourceManager详细见[参考文章](#)
- TaskExecutor挂掉
- JobMaster挂掉

资源分配相关？

- 在operator中对并行度的设置将决定任务分配到几个task slot里面去
- 申请资源的时候由-ys 决定向每个nm申请几个cores，而一个job真正能使用到的slot是由

Flink程序运行流程分解

- 基本步骤
 - 1. Obtain an execution environment

```
getExecutionEnvironment()
createLocalEnvironment()
createRemoteEnvironment(host: String, port: Int, jarFiles: String*)
```

- 2. Load/create the initial data

```
val text: DataStream[String] = env.readTextFile("file:///path/to/file")
```

- 3. Specify transformations on this data

```
//create a new DataStream by converting every String in the original collection  
to an integer  
val mapped = input.map { x => x.toInt }
```

- 4. Specify where to put the results of your computations

```
writeAsText(path: String)  
print()
```

- 5. Trigger the program execution

Flink watermark机制

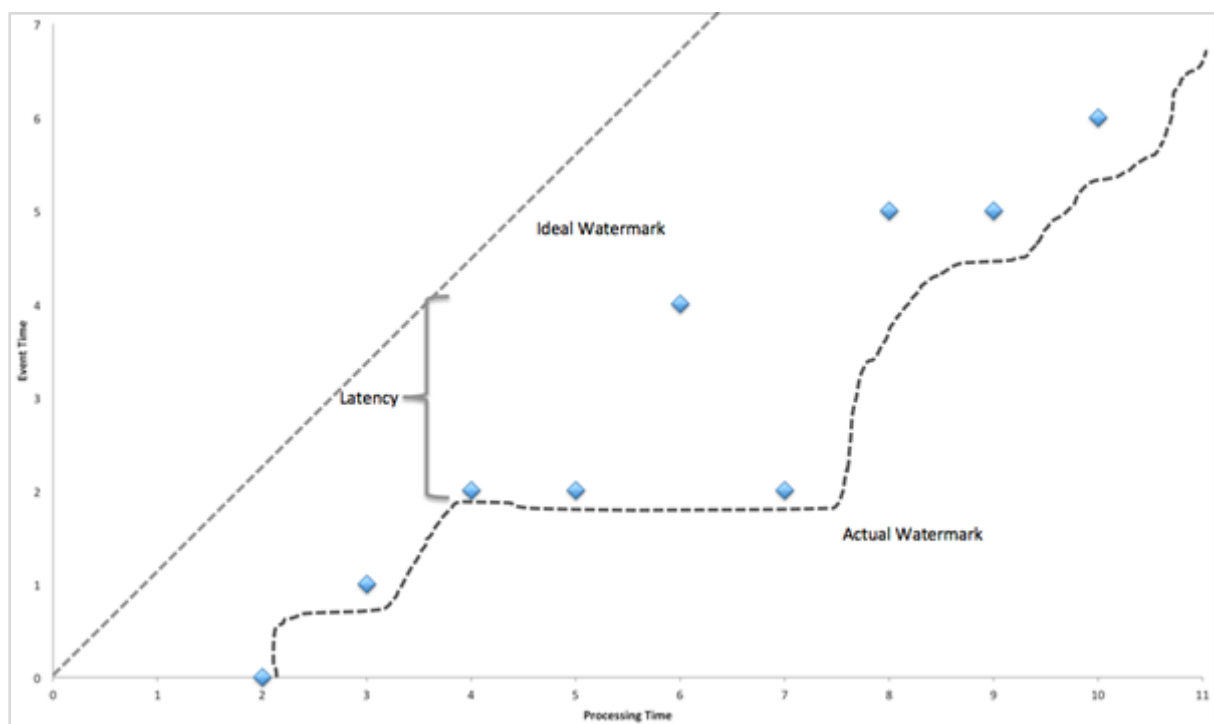
- **【重要】详细讲解watermark:** Flink流计算编程–watermark（水位线）

- window 触发的两个条件

```
1、 watermark时间 >= window_end_time  
2、 在>window_start_time>window_end_time)中有数据存在
```

- 摘录：深入理解Flink核心技术

- 纵坐标为event_time，横坐标为processingTime，理想情况自然是两者一致，但实际情况肯定不可能



- 摘录：使用EventTime与WaterMark进行流数据处理

```
// 这块结合上图理解watermark的值
@Override
public final Watermark getCurrentWatermark() {
    long potentialWM = currentMaxTimestamp - maxOutOfOrderness; // 当前最大
    事件时间戳，减去允许最大延迟到达时间
    if (potentialWM >= lastEmittedWatermark) { // 检查上一次emit的WaterMark时
    间戳，如果比lastEmittedWatermark大则更新其值
        lastEmittedWatermark = potentialWM;
    }
    return new Watermark(lastEmittedWatermark);
}
```

- Windowing, WaterMark, Trigger 三者依赖关系

1、Windowing：就是负责该如何生成Window，比如Fixed Window、Slide Window，当你配置好生成Window的策略时，Window就会根据时间动态生成，最终得到一个一个个的Window，包含一个时间范围：[起始时间，结束时间)，它们是一个一个受限于该时间范围的事件记录的容器，每个Window会收集一堆记录，满足指定条件会触发Window内事件记录集合的计算处理。

2、WaterMark：它其实不太好理解，可以将它定义为一个函数 $E=f(P)$ ，当前处理系统的处理时间 P ，根据一定的策略 f 会映射到一个事件时间 E ，可见 E 在坐标系中的表现形式是一条曲线，根据 f 的不同曲线形状也不同。假设，处理时间12:00:00，我希望映射到事件时间11:59:30，这时对于延迟30秒以内（事件时间范围11:59:30~12:00:00）的事件记录到达处理系统，都指派到时间范围包含处理时间12:00:00这个Window中。事件时间超过12:00:00的就会由Trigger去做补偿了。

3、Trigger：为了满足实际不同的业务需求，对上述事件记录指派给Window未能达到实际效果，而做出的一种补偿，比如事件记录在WaterMark时间戳之后到达事件处理系统，因为已经在对应的Window时间范围之后，我有很多选择：选择丢弃，选择是满足延迟3秒后还是指派给该Window，选择只接受对应的Window时间范围之后的5个事件记录，等等，这都是满足业务需要而制定的触发Window重新计算的策略，所以非常灵活。

Sink Connectors

- Kafka
- Elasticsearch
- RabbitMQ
- Rolling File Sink (HDFS)
 - [Apache Flink 1.7 Documentation: HDFS Connector](#)

flink-connector-kafka-0.8_2.11	1.0.0	FlinkKafkaConsumer08 FlinkKafkaProducer08	0.8.x	Uses the SimpleConsumer API of Kafka internally. Offsets are committed to ZK by Flink.
flink-connector-kafka-0.9_2.11	1.0.0	FlinkKafkaConsumer09 FlinkKafkaProducer09	0.9.x	Uses the new Consumer API Kafka.
flink-connector-kafka-0.10_2.11	1.2.0	FlinkKafkaConsumer010 FlinkKafkaProducer010	0.10.x	This connector supports Kafka messages with timestamps both for producing and consuming.
flink-connector-kafka-0.11_2.11	1.4.0	FlinkKafkaConsumer011 FlinkKafkaProducer011	0.11.x	Since 0.11.x Kafka does not support scala 2.10. This connector supports Kafka transactional messaging to provide exactly once semantic for the producer.
flink-connector-kafka_2.11	1.7.0	FlinkKafkaConsumer FlinkKafkaProducer	>= 1.0.0	<p>This universal Kafka connector attempts to track the latest version of the Kafka client. The version of the client it uses may change between Flink releases. Modern Kafka clients are backwards compatible with broker versions 0.10.0 or later. However for Kafka 0.11.x and 0.10.x versions, we recommend using dedicated <code>flink-connector-kafka-0.11_2.11</code> and <code>link-connector-kafka-0.10_2.11</code> respectively.</p> <p>Attention: as of Flink 1.7 the universal Kafka connector is considered to be in a BETA status and might not be as stable as the 0.11 connector. In case of problems with the universal connector, you can try to use <code>flink-connector-kafka-0.11_2.11</code> which should be compatible with all of the Kafka versions starting from 0.11.</p>

开启自动发现机制

- 分区自动发现
 - By default, partition discovery is *disabled*. To enable it, set a non-negative value for `flink.partition-discovery.interval-millis` in the provided properties config, representing the discovery interval in *milliseconds*.
- topic 自动发现
 - based on pattern matching on the topic names using regular expressions
 - 同样需要设置`flink.partition-discovery.interval-millis`

```
val env = StreamExecutionEnvironment.getExecutionEnvironment()

val properties = new Properties()
properties.setProperty("bootstrap.servers", "localhost:9092")
properties.setProperty("group.id", "test")

val myConsumer = new FlinkKafkaConsumer08[String](
  java.util.regex.Pattern.compile("test-topic-[0-9]"),
  new SimpleStringSchema,
  properties)

val stream = env.addSource(myConsumer)
```

produce message to kafka with flink

- [docApache Flink 1.7 Documentation: Apache Kafka Connector](#)

```
val stream: DataStream[String] = ...

val myProducer = new FlinkKafkaProducer011[String](
    "localhost:9092", // broker list
    "my-topic",       // target topic
    new SimpleStringSchema) // serialization schema

// versions 0.10+ allow attaching the records' event timestamp when writing them to Kafka;
// this method is not available for earlier Kafka versions
myProducer.setWriteTimestampToKafka(true)

stream.addSink(myProducer)
```

- 如果不指定partitioner的话，默认是FlinkFixedPartitioner
 - maps each Flink Kafka Producer parallel subtask to a single Kafka partition (i.e., all records received by a sink subtask will end up in the same Kafka partition)
- Kafka 0.9 and 0.10
 - setLogFailureOnly=false, setFlushOnCheckpoint=true provide *at-least-once* guarantees
 - 默认retries 为0，是为了防止重试导致消息重复
- Kafka 0.11 and newer
 - provide *exactly-once delivery guarantees*
 - use FlinkKafkaProducer011 with semantic settings
 - Semantic.NONE
 - Semantic.ATLEASTONCE
 - Semantic.EXACTLY_ONCE
 - 同样得设置事务隔离级别，readcommitted, readuncommitted 默认是 read_uncommitted

StreamingFileSink与Kafka 结合如何做到exactly once?

- [An Overview of End-to-End Exactly-Once Processing in Apache Flink \(with Apache Kafka, too!\)](#)
 - 二阶段提交
- partfile 有三种状态：in-progress, pending, finished；part file先被写成in-progress，一旦被关闭写入，会变成pending，当检查点成功之后，pending状态的文件将变成finished；
- 如果失败，将从上一个检查点开始重新store，期间回滚in-progress中的文件，以确保不会重复保存上一个检查点之后的数据

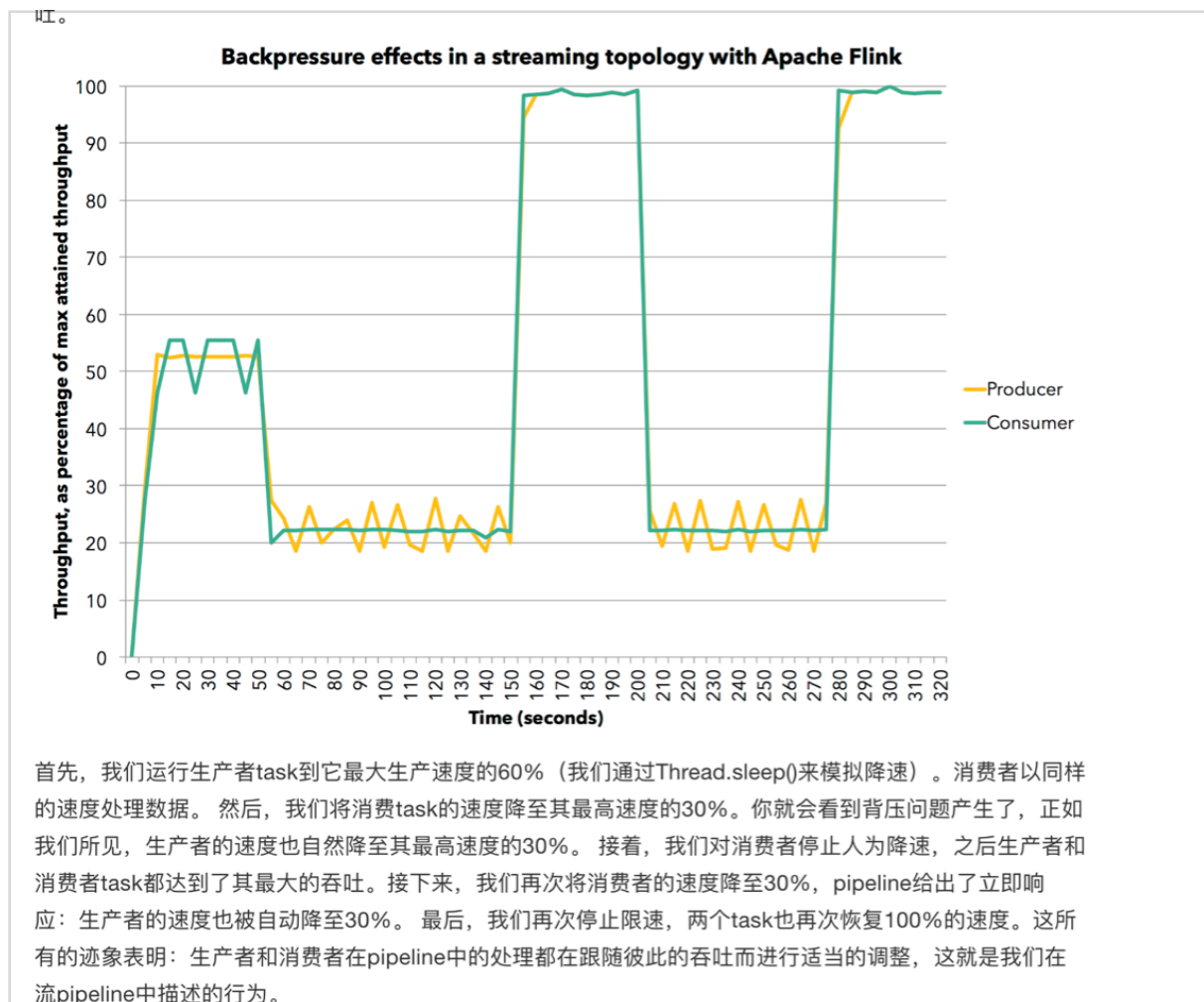
flink如何控制kafka offset提交与checkpoint&&savepoint相结合

- [FlinkKafkaConsumer使用详解](#)
- [Flink 小贴士 \(2\): Flink 如何管理 Kafka 消费位点 | Jark's Blog](#)
- 关闭checkpoint(**Checkpointingdisabled**):
 - 此时，Flink Kafka Consumer依赖于它使用的具体的Kafka client的自动定期提交offset的行为，相应的设置是 Kafka properties中的 enable.auto.commit (或者 auto.commit.enable 对于Kafka 0.8) 以及 auto.commit.interval.ms

- 开启checkpoint(**Checkpointingenabled**):
 - 在这种情况下，Flink Kafka Consumer会将offset存到checkpoint中
 - **制作完checkpoint 一并提交offsets** 当checkpoint 处于completed的状态时（**整个job的所有operator都收到了这个checkpoint的barrier**）。将offset记录起来并提交，从而保证exactly-once
- exactly once的两个风险点：可结合savepoint来做
 - 1. 异常退出的情况，没法来得及做checkpoint, 而checkpoint间隔太长会导致丢失大量数据；可以通过airflow周期性手动触发savepoint恢复；封装hflink脚本
 - 解决思路是结合savepoint来做，通过**airflow定时的触发savepoint**操作，防止因checkpoint未及时做数据丢失
 - 规定一分钟savepoint一次，这样即使分钟级别的数据丢失也是可以容忍
 - 2. 第一点利用savepoint来做也有风险：在做savepoint的时候，如果异常退出，parfile未及时关闭导致数据丢失
 - **暂时可以认为问题较小？**

flink如何控制背压

- 如何做到挂很久之后重新启动时限制拉取的消息量？（类似spark max.partition）
 - 背压通过task slot 的stackTrace判断
 - 可以在kafka source那层控制一次性消费量，类似于spark
 - 其实flink天然就是背压的状态
 - 不用担心一次性拉取全部的量，有个默认的max.partition.fetch.bytes=1048576；可以背压从sink到source 的整个pipeline,同时对source进行限流来适配整个pipeline中最慢组件的速度，从而获得系统的稳定状态。
 - 也不用担心一次拉取的量不够，通过背压可以调整



Flink 如何管理 Kafka 消费位点

- [Flink 小贴士 \(2\): Flink 如何管理 Kafka 消费位点 | Jark's Blog](#)

Flink 高性能部署

- [Apache Flink 1.8-SNAPSHOT Documentation: JobManager High Availability \(HA\)](#)

Flink 内存管理

- [001待整理：flink的背压机制·flink简明实战教程](#)
- [Flink 原理与实现：内存管理 - 简书](#)
- [Flink序列化框架分析_慕课手记](#)
- [Flink 原理与实现：内存管理 | Jark's Blog](#)
 - 有效避免OOM所有的运行时数据结构和算法只能通过内存池申请内存，保证了其使用的内存大小是固定的，不会因为运行时数据结构和算法而发生OOM。在内存吃紧的情况下，算法（sort/join等）会高效地将一大批内存块写到磁盘，之后再读回来。因此，OutOfMemoryErrors可以有效地被避免
 - 堆外内存存在写磁盘或网络传输时是 zero-copy，而堆内存的话，至少需要 copy 一次

Flink Broadcast State

- [Flink 小贴士 \(6\): 使用 Broadcast State 的 4 个注意事项 | Jark's Blog](#)
- [A Practical Guide to Broadcast State in Apache Flink](#)
- [Flink Broadcast 广播变量应用案例实战-Flink牛刀小试 - 掘金](#)

相关问题

- 不允许**null field**或**instance**在**operator**间传递，即使代码写了filter；这一点跟spark处理机制不同[Re: Cannot pass objects with null-valued fields to the next operator](#)

metric监控rest api

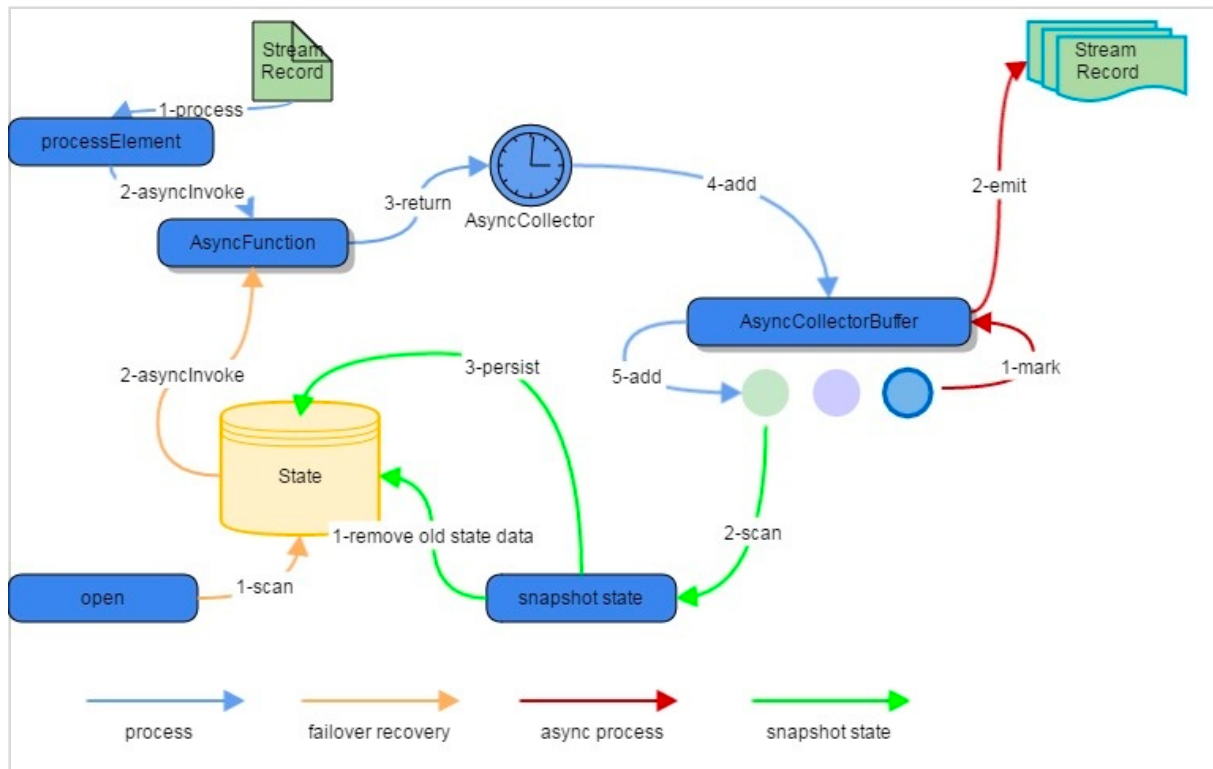
- [Apache Flink 1.8-SNAPSHOT Documentation: Monitoring REST API](#)
- Flink 监控上游Kafka的情况
 - 从 Kafka version 0.9开始，kafka consumer 把所有的metric暴露了出来，所有已经暴露的指标见文档 [The Kafka documentation lists all exported metrics](#)
 - fetch-size-avg: The average number of bytes fetched per request for a topic
 -

Restart Strategies

- **doc** [Apache Flink 1.7 Documentation: Restart Strategies](#)
- Fixed Delay Restart Strategy
- Failure Rate Restart Strategy
- No Restart Strategy
- Fallback Restart Strategy

Flink with async IO

- [Apache Flink 1.7 Documentation: Asynchronous I/O for External Data Access](#)
- 设计思想：[FLIP-12: Asynchronous I/O Design and Implementation - Apache Flink - Apache Software Foundation](#)



Flink with table api

- [Apache Flink 1.7 Documentation: Concepts & Common API](#)
- [Table Api操作符 sum, max...]([Apache Flink 1.7 Documentation: Built-In Functions](#))

Flink for Cloudera

- flink 默认不支持cdh版本Hadoop，需要重新编译一下flink源码
 - <https://ci.apache.org/projects/flink/flink-docs-release-1.7/flinkDev/building.html>

```
mvn clean install -DskipTests -Pvendor-repos -Dhadoop.version=3.0.0-cdh6.0.1
```

- 配置maven仓库 [解决maven仓库默认不支持cdh - 这里的天空比较蓝 - CSDN博客](#)

```
<mirror>
  <id>nexus-aliyun</id>
  <mirrorOf>*,!cloudera</mirrorOf>
  <name>Nexus aliyun</name>
  <url>
    http://maven.aliyun.com/nexus/content/groups/public
  </url>
</mirror>
```