## 深入理解Spark 2.1 Core （二）：DAG调度器的原理与源码分析

上一篇《深入理解Spark 2.0 （一）： RDD实现及源码分析 》的5.2 Spark任务调度器我们省略过去了，这篇我们就来讲讲Spark的调度器。

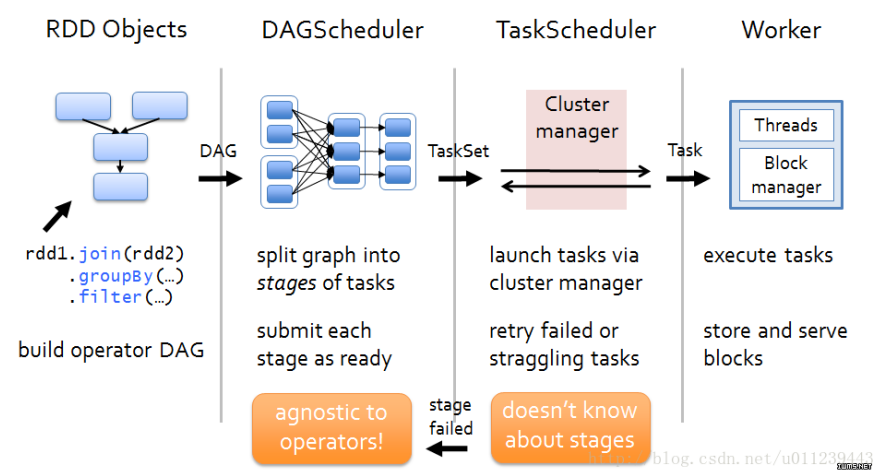# 概述

上一篇《深入理解Spark（一）： RDD实现及源码分析 》提到：

> 定义RDD之后，程序员就可以在动作（注：即action操作）中使用RDD了。动作是向应用程序返回值，或向存储系统导出数据的那些操作，例如，count（返回RDD中的元素个数），collect（返回元素本身），save（将RDD输出到存储系统）。在Spark中，只有在动作第一次使用RDD时，才会计算RDD（即延迟计算）。这样在构建RDD的时候，运行时通过管道的方式传输多个转换。

一次action操作会触发RDD的延迟计算，我们把这样的一次计算称作一个Job。我们还提到了窄依赖和宽依赖的概念：

> 窄依赖指的是：每个parent RDD 的 partition 最多被 child RDD的一个partition使用
> 宽依赖指的是：每个parent RDD 的 partition 被多个 child RDD的partition使用
> 窄依赖每个child RDD 的partition的生成操作都是可以并行的，而宽依赖则需要所有的parent partition shuffle结果得到后再进行。

由于在RDD的一系类转换中，若其中一些连续的转换都是窄依赖，那么它们是可以并行的，而有宽依赖则不行。所有，Spark将宽依赖为划分界限，将Job换分为多个Stage。而一个Stage里面的转换任务，我们可以把它抽象成TaskSet。一个TaskSet中有很多个Task，它们的转换操作都是相同的，不同只是操作的对象是对数据集中的不同子数据集。

接下来，Spark就可以提交这些任务了。但是，如何对这些任务进行调度和资源分配呢？如何通知worker去执行这些任务呢？接下来，我们会一一讲解。



根据以上两个阶段，我们会来详细介绍两个Scheduler，一个是DAGScheduler，另外一个是TaskScheduler。

我们先来看一来在SparkContext中是如何创建它们的：

```
1   val (sched, ts) = SparkContext.createTaskScheduler(this, master, deployMode)
2     _schedulerBackend = sched
3     _taskScheduler = ts
4     _dagScheduler = new DAGScheduler(this)
```

可以看到，我们是先用函数createTaskScheduler创建了taskScheduler，再new了一个DAGScheduler。这个顺序可以改变吗？答案是否定的，我们看下DAGScheduler类就知道了：

```
1   class DAGScheduler(
2       private[scheduler] val sc: SparkContext,
3       private[scheduler] val taskScheduler: TaskScheduler,
4       listenerBus: LiveListenerBus,
5       mapOutputTracker: MapOutputTrackerMaster,
6       blockManagerMaster: BlockManagerMaster,
7       env: SparkEnv,
8       clock: Clock = new SystemClock())
9     extends Logging {
10
11    def this(sc: SparkContext, taskScheduler: TaskScheduler) = {
12      this(
13        sc,
14        taskScheduler,
15        sc.listenerBus,
16        sc.env.mapOutputTracker.asInstanceOf[MapOutputTrackerMaster],
17        sc.env.blockManager.master,
18        sc.env)
19    }
20
21    def this(sc: SparkContext) = this(sc, sc.taskScheduler)
22
23  ***
24
25    }
```

SparkContext中创建的TaskScheduler，会传入DAGScheduler赋值给它的成员变量，再DAG阶段结束后，使用它进行下一步对任务调度等的操作。

# 提交Job

调用栈如下：

- rdd.count

  - SparkContext.runJob

    - DAGScheduler.runJob

      - DAGScheduler.submitJob

        - DAGSchedulerEventProcessLoop.doOnReceive

          - DAGScheduler.handleJobSubmitted

接下来，我们来逐个深入：

## rdd.count

RDD的一些action操作都会触发SparkContext的runJob函数，如count()

```
1   def count(): Long = sc.runJob(this, Utils.getIteratorSize _).sum
```

## SparkContext.runJob

SparkContext的runJob会触发 DAGScheduler的runJob：

```
1   def runJob[T, U: ClassTag](
2       rdd: RDD[T],
3       func: (TaskContext, Iterator[T]) => U,
4       partitions: Seq[Int],
5       resultHandler: (Int, U) => Unit): Unit = {
```

```
 6        if (stopped.get()) {
 7          throw new IllegalStateException("SparkContext has been shutdown")
 8        }
 9        val callSite = getCallSite
10        val cleanedFunc = clean(func)
11        logInfo("Starting job: " + callSite.shortForm)
12        if (conf.getBoolean("spark.logLineage", false)) {
13          logInfo("RDD's recursive dependencies:\n" + rdd.toDebugString)
14        }
15        dagScheduler.runJob(rdd, cleanedFunc, partitions, callSite, resultHandler, localProperties.get)
16        progressBar.foreach(_.finishAll())
17        rdd.doCheckpoint()
18      }
```

这里的rdd.doCheckpoint()并不是对自己Checkpoint，而是递归的回溯parent rdd 检查checkpointData是否被定义了，若定义了就将该rdd Checkpoint：

```
 1    private[spark] def doCheckpoint(): Unit = {
 2      RDDOperationScope.withScope(sc, "checkpoint", allowNesting = false, ignoreParent = true) {
 3        if (!doCheckpointCalled) {
 4          doCheckpointCalled = true
 5          if (checkpointData.isDefined) {
 6            if (checkpointAllMarkedAncestors) {
 7              //若想要把checkpointData定义过的RDD的parents也进行checkpoint的话，
 8              //那么我们需要先对parents checkpoint。
 9              //这是因为，如果RDD把自己checkpoint了，
10              //那么它就将lineage中它的parents给切除了。
11              dependencies.foreach(_.rdd.doCheckpoint())
12            }
13            checkpointData.get.checkpoint()
14          } else {
15            dependencies.foreach(_.rdd.doCheckpoint())
16          }
17        }
18      }
19    }
```

具体的checkpoint实现可见上一篇博文。

# DAGScheduler.runJob

DAGScheduler的runJob会触发DAGScheduler的submitJob：

```
 1   /**
 2     * 参数介绍：
 3     * @param rdd:  执行任务的目标TDD
 4     * @param func:  在RDD的分区上所执行的函数
 5     * @param partitions:  需要执行的分区集合;有些job并不会对RDD的所有分区都进行计算的，比如说first()
 6     * @param callSite：用户程序的调用点
 7     * @param resultHandler：回调结果
 8     * @param properties：关于这个job的调度器特征，比如说公平调度的pool名字，这个会在后续讲到
 9     */
10   def runJob[T, U](
11       rdd: RDD[T],
12       func: (TaskContext, Iterator[T]) => U,
13       partitions: Seq[Int],
14       callSite: CallSite,
15       resultHandler: (Int, U) => Unit,
16       properties: Properties): Unit = {
17     val start = System.nanoTime
18     val waiter = submitJob(rdd, func, partitions, callSite, resultHandler, properties)
19
20      ***
```

```
21        waiter.completionFuture.value.get match {
22          case scala.util.Success(_) =>
23            logInfo("Job %d finished: %s, took %f s".format
24              (waiter.jobId, callSite.shortForm, (System.nanoTime - start) / 1e9))
25          case scala.util.Failure(exception) =>
26            logInfo("Job %d failed: %s, took %f s".format
27              (waiter.jobId, callSite.shortForm, (System.nanoTime - start) / 1e9))
28            val callerStackTrace = Thread.currentThread().getStackTrace.tail
29            exception.setStackTrace(exception.getStackTrace ++ callerStackTrace)
30            throw exception
31        }
32      }
```

## DAGScheduler.submitJob

我们接下来看看submitJob里面做了什么：

```
1    def submitJob[T, U](
2        rdd: RDD[T],
3        func: (TaskContext, Iterator[T]) => U,
4        partitions: Seq[Int],
5        callSite: CallSite,
6        resultHandler: (Int, U) => Unit,
7        properties: Properties): JobWaiter[U] = {
8      // 确认没在不存在的partition上执行任务
9      val maxPartitions = rdd.partitions.length
10     partitions.find(p => p >= maxPartitions || p < 0).foreach { p =>
11       throw new IllegalArgumentException(
12         "Attempting to access a non-existent partition: " + p + ". " +
13           "Total number of partitions: " + maxPartitions)
14     }
15     //递增得到jobId
16     val jobId = nextJobId.getAndIncrement()
17     if (partitions.size == 0) {
18       //若Job没对任何一个partition执行任务，
19       //则立即返回
20       return new JobWaiter[U](this, jobId, 0, resultHandler)
21     }
22
23     assert(partitions.size > 0)
24     val func2 = func.asInstanceOf[(TaskContext, Iterator[_]) => _]
25     val waiter = new JobWaiter(this, jobId, partitions.size, resultHandler)
26     eventProcessLoop.post(JobSubmitted(
27       jobId, rdd, func2, partitions.toArray, callSite, waiter,
28       SerializationUtils.clone(properties)))
29     waiter
30   }
```

## DAGSchedulerEventProcessLoop.doOnReceive

eventProcessLoop是一个DAGSchedulerEventProcessLoop类对象，即一个DAG调度事件处理的监听。eventProcessLoop中调用doOnReceive来进行监听

```
1    private def doOnReceive(event: DAGSchedulerEvent): Unit = event match {
2      //当事件为JobSubmitted时，
3      //会调用DAGScheduler.handleJobSubmitted
4      case JobSubmitted(jobId, rdd, func, partitions, callSite, listener, properties) =>
5        dagScheduler.handleJobSubmitted(jobId, rdd, func, partitions, callSite, listener, properties)
6    ***
7    }
```

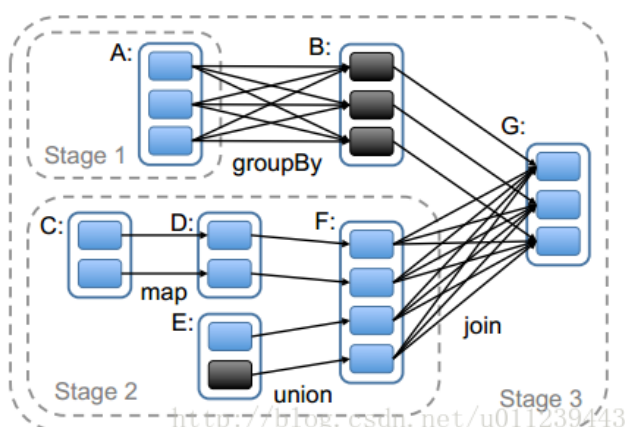## DAGScheduler.handleJobSubmitted

自此Job的提交就完成了：

```scala
private[scheduler] def handleJobSubmitted(jobId: Int,
    finalRDD: RDD[_],
    func: (TaskContext, Iterator[_]) => _,
    partitions: Array[Int],
    callSite: CallSite,
    listener: JobListener,
    properties: Properties) {
  var finalStage: ResultStage = null
  try {
    finalStage = newResultStage(finalRDD, func, partitions, jobId, callSite)
  } catch {
    case e: Exception =>
      logWarning("Creating new stage failed due to exception - job: " + jobId, e)
      listener.jobFailed(e)
      return
  }

  val job = new ActiveJob(jobId, finalStage, callSite, listener, properties)
  clearCacheLocs()
  logInfo("Got job %s (%s) with %d output partitions".format(
    job.jobId, callSite.shortForm, partitions.length))
  logInfo("Final stage: " + finalStage + " (" + finalStage.name + ")")
  logInfo("Parents of final stage: " + finalStage.parents)
  logInfo("Missing parents: " + getMissingParentStages(finalStage))

  val jobSubmissionTime = clock.getTimeMillis()
  jobIdToActiveJob(jobId) = job
  activeJobs += job
  finalStage.setActiveJob(job)
  val stageIds = jobIdToStageIds(jobId).toArray
  val stageInfos = stageIds.flatMap(id => stageIdToStage.get(id).map(_.latestInfo))
  listenerBus.post(
    SparkListenerJobStart(job.jobId, jobSubmissionTime, stageInfos, properties))
  submitStage(finalStage)

  submitWaitingStages()
}
```

接下来我们来看看handleJobSubmitted中的newResultStage，一个非常有趣的划分Stage过程。

# 划分Stage



如我们之前提到的：Spark将宽依赖为划分界限，将Job换分为多个Stage。调用栈为：

- DAGScheduler.newResultStage

  - DAGScheduler.getParentStagesAndId

    - DAGScheduler.getParentStages

      - DAGScheduler.getShuffleMapStage

        - DAGScheduler.getAncestorShuffleDependencies

        - DAGScheduler.newOrUsedShuffleStage

          - DAGScheduler.newShuffleMapStage

接下来，我们来逐个深入：

# DAGScheduler.newResultStage

Spark的Stage调用是从最后一个RDD所在的Stage，ResultStage开始划分的，这里即为G所在的Stage。但是在生成这个Stage之前会生成它的parent Stage，就这样递归的把parent Stage都先生成了。

```
 1    private def newResultStage(
 2        rdd: RDD[_],
 3        func: (TaskContext, Iterator[_]) => _,
 4        partitions: Array[Int],
 5        jobId: Int,
 6        callSite: CallSite): ResultStage = {
 7      val (parentStages: List[Stage], id: Int) = getParentStagesAndId(rdd, jobId)
 8      val stage = new ResultStage(id, rdd, func, partitions, parentStages, jobId, callSite)
 9      stageIdToStage(id) = stage
10      updateJobIdStageIdMaps(jobId, stage)
11      stage
12    }
```

# DAGScheduler.getParentStagesAndId

getParentStagesAndId中得到了ParentStages以及其StageId：

```
 1    private def getParentStagesAndId(rdd: RDD[_], firstJobId: Int): (List[Stage], Int) = {
 2      val parentStages = getParentStages(rdd, firstJobId)
 3      val id = nextStageId.getAndIncrement()
 4      (parentStages, id)
 5    }
```

# DAGScheduler.getParentStages

我们再来深入看看getParentStages做了什么：

```
 1    private def getParentStages(rdd: RDD[_], firstJobId: Int): List[Stage] = {
 2      //将存储ParentStages
 3      val parents = new HashSet[Stage]
 4      //存储已将访问过了的RDD
 5      val visited = new HashSet[RDD[_]]
 6      // 存储需要被处理的RDD
 7      val waitingForVisit = new Stack[RDD[_]]
 8      def visit(r: RDD[_]) {
 9        if (!visited(r)) {
10          //加入访问集合
11          visited += r
12          //遍历该RDD所有的依赖
13          for (dep <- r.dependencies) {
14            dep match {
```

```
15              //若是宽依赖则生成新的Stage
16              case shufDep: ShuffleDependency[_, _, _] =>
17                parents += getShuffleMapStage(shufDep, firstJobId)
18              //若是窄依赖则加入Stack，等待处理
19              case _ =>
20                waitingForVisit.push(dep.rdd)
21            }
22          }
23        }
24      }
25      //在Stack中加入最后一个RDD
26      waitingForVisit.push(rdd)
27      //广度优先遍历
28      while (waitingForVisit.nonEmpty) {
29        visit(waitingForVisit.pop())
30      }
31      //返回ParentStages List
32      parents.toList
33    }
```

其实getParentStages使用的就是广度优先遍历的算法，若知道这点也容易理解了。虽然现在Stage并没有生成，但是我们可以看到划分策略是：广度遍历方式的划分parent RDD 的Stage。

若parent RDD 和 child RDD 为窄依赖，则将parent RDD 纳入 child RDD 所在的Stage中。如图，B被纳入了Stage3中。

若parent RDD 和 child RDD 为宽依赖，则parent RDD将纳入一新的Stage中。如图，F被纳入了Stage2中。

## DAGScheduler.getShuffleMapStage

下面我们来看下getShuffleMapStage是如何生成新的Stage的。
首先shuffleToMapStage中保存了关于Stage的HashMap

```
1  private[scheduler] val shuffleToMapStage = new HashMap[Int, ShuffleMapStage]
```

getShuffleMapStage会先去根据shuffleId去查找shuffleToMapStage

```
1    private def getShuffleMapStage(
2        shuffleDep: ShuffleDependency[_, _, _],
3        firstJobId: Int): ShuffleMapStage = {
4      shuffleToMapStage.get(shuffleDep.shuffleId) match {
5        //若找到则直接返回
6        case Some(stage) => stage
7        case None =>
8          // 检查这个Stage的Parent Stage是否生成
9          // 若没有，则生成它们
10         getAncestorShuffleDependencies(shuffleDep.rdd).foreach { dep =>
11           if (!shuffleToMapStage.contains(dep.shuffleId)) {
12             shuffleToMapStage(dep.shuffleId) = newOrUsedShuffleStage(dep, firstJobId)
13           }
14         }
15         // 生成新的Stage
16         val stage = newOrUsedShuffleStage(shuffleDep, firstJobId)
17         //将新的Stage 加入到 HashMap
18         shuffleToMapStage(shuffleDep.shuffleId) = stage
19         //返回新的Stage
20         stage
21       }
22     }
```

可以发现这部分的代码和上述的newResultStage部分很像，所以可以看成一种递归的方法。

## DAGScheduler.getAncestorShuffleDependencies

我们再来看下getAncestorShuffleDependencies，可想而知，它应该会和newResultStage中的getParentStages会非常类似：

```scala
private def getAncestorShuffleDependencies(rdd: RDD[_]): Stack[ShuffleDependency[_, _, _]] = {
  val parents = new Stack[ShuffleDependency[_, _, _]]
  val visited = new HashSet[RDD[_]]
  val waitingForVisit = new Stack[RDD[_]]
  def visit(r: RDD[_]) {
    if (!visited(r)) {
      visited += r
      for (dep <- r.dependencies) {
        dep match {
          case shufDep: ShuffleDependency[_, _, _] =>
            if (!shuffleToMapStage.contains(shufDep.shuffleId)) {
              parents.push(shufDep)
            }
          case _ =>
        }
        waitingForVisit.push(dep.rdd)
      }
    }
  }

  waitingForVisit.push(rdd)
  while (waitingForVisit.nonEmpty) {
    visit(waitingForVisit.pop())
  }
  parents
}
```

可以看到的确和newResultStage中的getParentStages会非常类似，不同的是这里会先判断shuffleToMapStage是否存在这个Stage，不存在的话会push到parents这个Stack，最会返回给上述的getShuffleMapStage，调用newOrUsedShuffleStage生成新的Stage。

## DAGScheduler.newOrUsedShuffleStage

那现在就来看newOrUsedShuffleStage是如何生成新的Stage的。
首先ShuffleMapTask的计算结果（其实是计算结果数据所在的位置、大小等元数据信息）都会传给Driver的mapOutputTracker。所以需要先判断Stage是否已经被计算过：

```scala
private def newOrUsedShuffleStage(
    shuffleDep: ShuffleDependency[_, _, _],
    firstJobId: Int): ShuffleMapStage = {
  val rdd = shuffleDep.rdd
  val numTasks = rdd.partitions.length
  //生成新的Stage
  val stage = newShuffleMapStage(rdd, numTasks, shuffleDep, firstJobId, rdd.creationSite)
  //判断Stage是否已经被计算过
  //若计算过，则把结果复制到新的stage
  if (mapOutputTracker.containsShuffle(shuffleDep.shuffleId)) {
    val serLocs = mapOutputTracker.getSerializedMapOutputStatuses(shuffleDep.shuffleId)
    val locs = MapOutputTracker.deserializeMapStatuses(serLocs)
    (0 until locs.length).foreach { i =>
      if (locs(i) ne null) {
        stage.addOutputLoc(i, locs(i))
      }
    }
  } else {
    logInfo("Registering RDD " + rdd.id + " (" + rdd.getCreationSite + ")")
    //如果没计算过，就在注册mapOutputTracker Stage
    //为存储元数据占位
    mapOutputTracker.registerShuffle(shuffleDep.shuffleId, rdd.partitions.length)
  }
  stage
}
```

## DAGScheduler.newShuffleMapStage

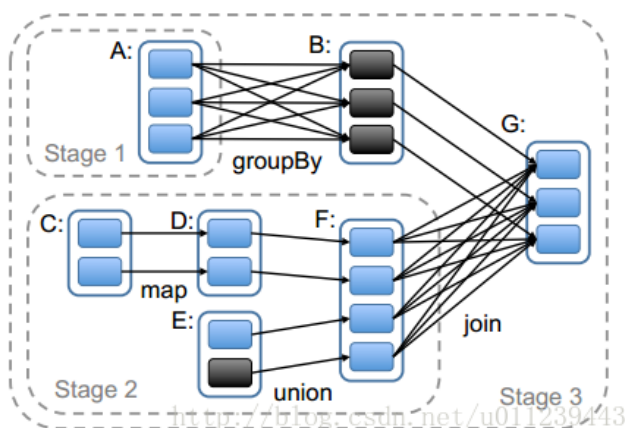递归就发生在newShuffleMapStage，它的实现和最一开始的newResultStage类似，也是先getParentStagesAndId，然后生成一个ShuffleMapStage：

```scala
 1    private def newShuffleMapStage(
 2        rdd: RDD[_],
 3        numTasks: Int,
 4        shuffleDep: ShuffleDependency[_, _, _],
 5        firstJobId: Int,
 6        callSite: CallSite): ShuffleMapStage = {
 7      val (parentStages: List[Stage], id: Int) = getParentStagesAndId(rdd, firstJobId)
 8      val stage: ShuffleMapStage = new ShuffleMapStage(id, rdd, numTasks, parentStages,
 9        firstJobId, callSite, shuffleDep)
10
11      stageIdToStage(id) = stage
12      updateJobIdStageIdMaps(firstJobId, stage)
13      stage
14    }
```

## 回顾

到此，Stage划分过程就结束了。我们在根据一开始的图，举例回顾下：



- 首先，我们想 `newResultStage RDD_G`所在的 `Stage3`

- 但在`new Stage`之前会调用`getParentStagesAndId`

- `getParentStagesAndId`中又会调用`getParentStages`，来广度优先的遍历`RDD_G`所依赖的`RDD`。如果是窄依赖，就纳入G所在的`Stage3`，如`RDD_B`就纳入了`Stage3`

- 若过是宽依赖，我们这里以`RDD_F`为例（与`RDD_A`处理过程相同）。我们就会调用`getShuffleMapStage`，来判断`RDD_F`所在的`Stage2`是否已经生成了，如果生成了就直接返回。

- 若还没生成，我们先调用`getAncestorShuffleDependencies`。`getAncestorShuffleDependencies`类似于`getParentStages`，也是用广度优先的遍历`RDD_F`所依赖的`RDD`。如果是窄依赖，如`RDD_C`、`RDD_D`和`RDD_E`，都被纳入了F所在的`Stage`2。但是假设`RDD_E`有个`parent RDD ``RDD_H`，`RDD_H`和`RDD_E`之间是宽依赖，那么该怎么办呢？我们会先判断RDD_H所在的Stage是否已经生成。若还没生成，我们把它put到一个parents Stack 中，最后返回。

- 对于那些返回的还没生成的Stage我们会调用`newOrUsedShuffleStage`

- `newOrUsedShuffleStage`会调用`newShuffleMapStage`，来生成新的Stage。而`newShuffleMapStage`的实现类似于`newResultStage`。这样我们就可以递归下去，使得每个Stage所依赖的Stage都已经生成了，再来生成这个的Stage。如这里，会将RDD_H所在的St

age生成了，然后在再生成Stage2。

- `newOrUsedShuffleStage`生成新的Stage后，会判断Stage是否被计算过。若已经被计算过，就从`mapOutPutTracker`中复制计算结果。若没计算过，则向`mapOutPutTracker`注册占位。

- 最后，回到`newResultStage`中，`new ResultStage`，这里即生成了`Stage3`。至此，`Stage`划分过程就结束了。

# 生成任务

调用栈如下：

- DAGScheduler.handleJobSubmitted

  - DAGScheduler.submitStage

    - DAGScheduler.getMissingParentStages

    - DAGScheduler.submitMissingTasks

## DAGScheduler.handleJobSubmitted

我们再回过头来看**"提交Job"**的最后一步handleJobSubmitted：

```
1    private[scheduler] def handleJobSubmitted(jobId: Int,
2        finalRDD: RDD[_],
3        func: (TaskContext, Iterator[_]) => _,
4        partitions: Array[Int],
5        callSite: CallSite,
6        listener: JobListener,
7        properties: Properties) {
8      var finalStage: ResultStage = null
9      try {
10       finalStage = newResultStage(finalRDD, func, partitions, jobId, callSite)
11     } catch {
12       case e: Exception =>
13         logWarning("Creating new stage failed due to exception - job: " + jobId, e)
14         listener.jobFailed(e)
15         return
16     }
17     ***
18   }
```

在**"划分Stage"**中我们已经深入的讲解了finalStage的生成：

```
1  finalStage = newResultStage(finalRDD, func, partitions, jobId, callSite)
```

接下来，我们继续往下看handleJobSubmitted的代码：

```
1      //生成新的job
2      val job = new ActiveJob(jobId, finalStage, callSite, listener, properties)
3      clearCacheLocs()
4      logInfo("Got job %s (%s) with %d output partitions".format(
5        job.jobId, callSite.shortForm, partitions.length))
6      logInfo("Final stage: " + finalStage + " (" + finalStage.name + ")")
7      logInfo("Parents of final stage: " + finalStage.parents)
8      logInfo("Missing parents: " + getMissingParentStages(finalStage))
9      //得到job提交的时间
10     val jobSubmissionTime = clock.getTimeMillis()
11     //得到job id
12     jobIdToActiveJob(jobId) = job
13     //添加到activeJobs HashSet
```

```
14      activeJobs += job
15      //将finalStage甚至ActiveJob为该job
16      finalStage.setActiveJob(job)
17      //得到stage 的id 信息
18      val stageIds = jobIdToStageIds(jobId).toArray
19      val stageInfos = stageIds.flatMap(id => stageIdToStage.get(id).map(_.latestInfo))
20      //监听
21      listenerBus.post(
22        SparkListenerJobStart(job.jobId, jobSubmissionTime, stageInfos, properties))
23      //提交
24      submitStage(finalStage)
25      //等待
26      submitWaitingStages()
```

## DAGScheduler.submitStage

接下来我们来看Stage是如何提交的。我们需要找到哪些parent Stage缺失，然后我们先运行生成这些Stage。这是一个深度优先遍历的过程：

```
1    private def submitStage(stage: Stage) {
2      val jobId = activeJobForStage(stage)
3      if (jobId.isDefined) {
4        logDebug("submitStage(" + stage + ")")
5        if (!waitingStages(stage) && !runningStages(stage) && !failedStages(stage)) {
6          //得到缺失的Parent Stage
7          val missing = getMissingParentStages(stage).sortBy(_.id)
8          logDebug("missing: " + missing)
9          if (missing.isEmpty) {
10            logInfo("Submitting " + stage + " (" + stage.rdd + "), which has no missing parents")
11            //如果没有缺失的Parent Stage，
12            //那么代表着该Stage可以运行了
13            //submitMissingTasks会完成DAGScheduler最后的工作，
14            //向TaskScheduler 提交 Task
15            submitMissingTasks(stage, jobId.get)
16          } else {
17          //深度优先遍历
18            for (parent <- missing) {
19              submitStage(parent)
20            }
21            waitingStages += stage
22          }
23        }
24      } else {
25        abortStage(stage, "No active job for stage " + stage.id, None)
26      }
27    }
```

## DAGScheduler.getMissingParentStages

getMissingParentStages类似于getParentStages，也是使用广度优先遍历：

```
1    private def getMissingParentStages(stage: Stage): List[Stage] = {
2      val missing = new HashSet[Stage]
3      val visited = new HashSet[RDD[_]]
4      val waitingForVisit = new Stack[RDD[_]]
5      def visit(rdd: RDD[_]) {
6        if (!visited(rdd)) {
7          visited += rdd
8          val rddHasUncachedPartitions = getCacheLocs(rdd).contains(Nil)
9          if (rddHasUncachedPartitions) {
10            for (dep <- rdd.dependencies) {
11              dep match {
12                //若是宽依赖 并且 不可用 ，
```

```
13          //则加入 missing HashSet
14            case shufDep: ShuffleDependency[_, _, _] =>
15              val mapStage = getShuffleMapStage(shufDep, stage.firstJobId)
16              if (!mapStage.isAvailable) {
17                missing += mapStage
18              }
19            //若是窄依赖
20            //则加入等待访问的 HashSet
21            case narrowDep: NarrowDependency[_] =>
22              waitingForVisit.push(narrowDep.rdd)
23          }
24        }
25      }
26    }
27    }
28    waitingForVisit.push(stage.rdd)
29    while (waitingForVisit.nonEmpty) {
30      visit(waitingForVisit.pop())
31    }
32    missing.toList
33  }
```

## DAGScheduler.submitMissingTasks

最后，我们来看下DAGScheduler最后的工作，提交Task：

```
1  private def submitMissingTasks(stage: Stage, jobId: Int) {
2    logDebug("submitMissingTasks(" + stage + ")")
3    // pendingPartitions 是 HashSet[Int]
4    //存储待处理的Task
5    stage.pendingPartitions.clear()
6
7    // 找出还未就算的Partition
8    val partitionsToCompute: Seq[Int] = stage.findMissingPartitions()
9
10   //从一个ActiveJob中得到关于这个Stage的
11   //调度池，job组描述等信息
12   val properties = jobIdToActiveJob(jobId).properties
13   // runningStages 是 HashSet[Stage]
14   //将当前Stage加入到运行中Stage集合
15   runningStages += stage
16
17   stage match {
18     case s: ShuffleMapStage =>
19       outputCommitCoordinator.stageStart(stage = s.id, maxPartitionId = s.numPartitions - 1)
20     case s: ResultStage =>
21       outputCommitCoordinator.stageStart(
22         stage = s.id, maxPartitionId = s.rdd.partitions.length - 1)
23   }
24   val taskIdToLocations: Map[Int, Seq[TaskLocation]] = try {
25     stage match {
26       case s: ShuffleMapStage =>
27         partitionsToCompute.map { id => (id, getPreferredLocs(stage.rdd, id))}.toMap
28       case s: ResultStage =>
29         partitionsToCompute.map { id =>
30           val p = s.partitions(id)
31           (id, getPreferredLocs(stage.rdd, p))
32         }.toMap
33     }
34   } catch {
35     case NonFatal(e) =>
36       stage.makeNewStageAttempt(partitionsToCompute.size)
37       listenerBus.post(SparkListenerStageSubmitted(stage.latestInfo, properties))
38       abortStage(stage, s"Task creation failed: $e\n${Utils.exceptionString(e)}", Some(e))
```
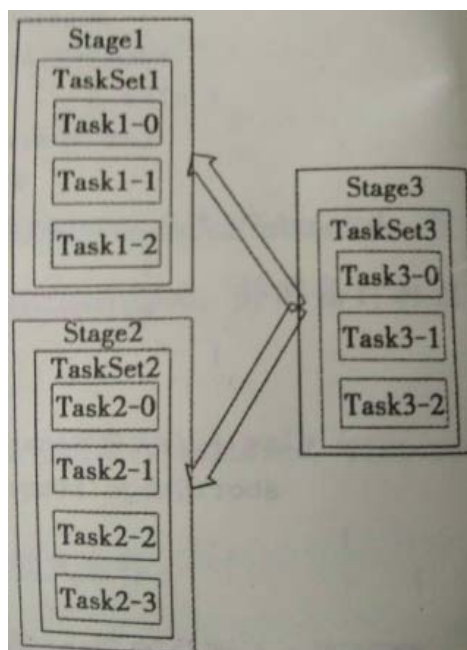
```scala
39        runningStages -= stage
40        return
41    }
42
43    stage.makeNewStageAttempt(partitionsToCompute.size, taskIdToLocations.values.toSeq)
44  //向listenerBus发送SparkListenerStageSubmitted事件
45  listenerBus.post(SparkListenerStageSubmitted(stage.latestInfo, properties))
46
47    var taskBinary: Broadcast[Array[Byte]] = null
48    try {
49  //对于最后一个Stage的Task，
50  //序列化并广播(rdd, func)。
51  //若是其他的Stage的Task，
52  //序列化并广播(rdd, shuffleDep)
53      val taskBinaryBytes: Array[Byte] = stage match {
54        case stage: ShuffleMapStage =>
55          JavaUtils.bufferToArray(
56            closureSerializer.serialize((stage.rdd, stage.shuffleDep): AnyRef))
57        case stage: ResultStage =>
58          JavaUtils.bufferToArray(closureSerializer.serialize((stage.rdd, stage.func): AnyRef))
59      }
60
61      taskBinary = sc.broadcast(taskBinaryBytes)
62    } catch {
63  //若序列化失败，停止这个stage
64      case e: NotSerializableException =>
65        abortStage(stage, "Task not serializable: " + e.toString, Some(e))
66        runningStages -= stage
67
68        // 停止执行
69        return
70      case NonFatal(e) =>
71        abortStage(stage, s"Task serialization failed: $e\n${Utils.exceptionString(e)}", Some(e))
72        runningStages -= stage
73        return
74    }
75
76    val tasks: Seq[Task[_]] = try {
77  //对于最后一个Stage的Task，
78  //则创建ResultTask。
79  //若是其他的Stage的Task，
80  //则创建ShuffleMapTask。
81      stage match {
82        case stage: ShuffleMapStage =>
83          partitionsToCompute.map { id =>
84            val locs = taskIdToLocations(id)
85            val part = stage.rdd.partitions(id)
86            new ShuffleMapTask(stage.id, stage.latestInfo.attemptId,
87              taskBinary, part, locs, stage.latestInfo.taskMetrics, properties, Option(jobId),
88              Option(sc.applicationId), sc.applicationAttemptId)
89          }
90
91        case stage: ResultStage =>
92          partitionsToCompute.map { id =>
93            val p: Int = stage.partitions(id)
94            val part = stage.rdd.partitions(p)
95            val locs = taskIdToLocations(id)
96            new ResultTask(stage.id, stage.latestInfo.attemptId,
97              taskBinary, part, locs, id, properties, stage.latestInfo.taskMetrics,
98              Option(jobId), Option(sc.applicationId), sc.applicationAttemptId)
99          }
100       }
101     } catch {
102       case NonFatal(e) =>
```

```
103          abortStage(stage, s"Task creation failed: $e\n${Utils.exceptionString(e)}", Some(e))
104          runningStages -= stage
105          return
106       }
107
108       if (tasks.size > 0) {
109          logInfo("Submitting " + tasks.size + " missing tasks from " + stage + " (" + stage.rdd + ")")
110          stage.pendingPartitions ++= tasks.map(_.partitionId)
111          logDebug("New pending partitions: " + stage.pendingPartitions)
112          //创建TaskSet并提交
113          taskScheduler.submitTasks(new TaskSet(
114             tasks.toArray, stage.id, stage.latestInfo.attemptId, jobId, properties))
115          stage.latestInfo.submissionTime = Some(clock.getTimeMillis())
116       } else {
117          markStageAsFinished(stage, None)
118
119          val debugString = stage match {
120             case stage: ShuffleMapStage =>
121                s"Stage ${stage} is actually done; " +
122                   s"(available: ${stage.isAvailable}," +
123                   s"available outputs: ${stage.numAvailableOutputs}," +
124                   s"partitions: ${stage.numPartitions})"
125             case stage : ResultStage =>
126                s"Stage ${stage} is actually done; (partitions: ${stage.numPartitions})"
127          }
128          logDebug(debugString)
129
130          submitWaitingChildStages(stage)
131       }
132    }
```



TaskSet保存了Stage包含的一组完全相同的Task，每个Task的处理逻辑完全相同，不同的是处理的数据，每个Task负责一个Partition。

至此，DAGScheduler就完成了它的任务了。接下来一篇博文，我们会从上述代码中的：

```
1    taskScheduler.submitTasks(new TaskSet(
2       tasks.toArray, stage.id, stage.latestInfo.attemptId, jobId, properties))
```

开始讲起，深入理解TaskScheduler的工作过程。