

深入理解Spark 2.1 Core （三）：任务调度器的原理与源码分析

■ 版权声明：本文为博主原创文章，转载请附上原文地址。 <http://blog.csdn.net/u011239443/article/details/53996141>

上一篇博文《深入理解Spark 2.1 Core （二）：DAG调度器的实现与源码分析》讲到了DAGScheduler.submitMissingTasks中最终调用了taskScheduler.submitTasks来提交任务。

这篇我们就从taskScheduler.submitTasks开始讲，深入理解TaskScheduler的运行过程。

提交Task

调用栈如下：

- TaskSchedulerImpl.submitTasks
 - CoarseGrainedSchedulerBackend.reviveOffers
- CoarseGrainedSchedulerBackend.DriverEndpoint.makeOffers
 - TaskSchedulerImpl.resourceOffers
 - TaskSchedulerImpl.resourceOfferSingleTaskSet
- CoarseGrainedSchedulerBackend.DriverEndpoint.launchTasks

TaskSchedulerImpl.submitTasks

TaskSchedulerImpl是TaskScheduler的子类，重写了submitTasks：

```
1  override def submitTasks(taskSet: TaskSet) {
2    val tasks = taskSet.tasks
3    logInfo("Adding task set " + taskSet.id + " with " + tasks.length + " tasks")
4    this.synchronized {
5      //生成TaskSetManager
6      val manager = createTaskSetManager(taskSet, maxTaskFailures)
7      val stage = taskSet.stageId
8      val stageTaskSets =
9        taskSetsByStageIdAndAttempt.getOrElseUpdate(stage, new HashMap[Int, TaskSetManager])
10     stageTaskSets(taskSet.stageAttemptId) = manager
11     val conflictingTaskSet = stageTaskSets.exists { case (_, ts) =>
12       ts.taskSet != taskSet && !ts.isZombie
13     }
14     if (conflictingTaskSet) {
15       throw new IllegalStateException(s"more than one active taskSet for stage $stage:" +
16         s" ${stageTaskSets.toSeq.map{_, ts => ts.taskSet.id}.mkString(",")}")
17     }
18     //将manager等信息放入调度器
19     schedulableBuilder.addTaskSetManager(manager, manager.taskSet.properties)
20
21     if (!isLocal && !hasReceivedTask) {
22       starvationTimer.scheduleAtFixedRate(new TimerTask() {
23         override def run() {
24           if (!hasLaunchedTask) {
25             logWarning("Initial job has not accepted any resources; " +
26               "check your cluster UI to ensure that workers are registered " +
27               "and have sufficient resources")
28           } else {
29             this.cancel()
```

```

30         }
31     }
32     }, STARVATION_TIMEOUT_MS, STARVATION_TIMEOUT_MS)
33 }
34     hasReceivedTask = true
35 }
36 //分配资源
37 backend.reviveOffers()
38 }

```

CoarseGrainedSchedulerBackend.reviveOffers

下面我们来讲讲上一节代码中最后一句：

```
1 backend.reviveOffers()
```

我们先回过头来看TaskScheduler是如何启动的：

```

1  override def start() {
2      backend.start()
3
4      if (!isLocal && conf.getBoolean("spark.speculation", false)) {
5          logInfo("Starting speculative execution thread")
6          speculationScheduler.scheduleAtFixedRate(new Runnable {
7              override def run(): Unit = Utils.tryOrStopSparkContext(sc) {
8                  checkSpeculatableTasks()
9              }
10         }, SPECULATION_INTERVAL_MS, SPECULATION_INTERVAL_MS, TimeUnit.MILLISECONDS)
11     }
12 }

```

我们可以看到TaskScheduler.start会调用backend.start()。

backend 是一个 SchedulerBackend 接口。SchedulerBackend 接口由 CoarseGrainedSchedulerBackend 类实现。我们看下 CoarseGrainedSchedulerBackend 的start：

```

1  override def start() {
2      val properties = new ArrayBuffer[(String, String)]
3      for ((key, value) <- scheduler.sc.conf.getAll) {
4          if (key.startsWith("spark.")) {
5              properties += ((key, value))
6          }
7      }
8
9      driverEndpoint = createDriverEndpointRef(properties)
10 }

```

我们可以看到CoarseGrainedSchedulerBackend的start会生成driverEndpoint，它是一个rpc的终端，一个RpcEndpoint接口，它由ThreadSafeRpcEndpoint接口实现，而ThreadSafeRpcEndpoint由CoarseGrainedSchedulerBackend的内部类DriverEndpoint实现。

CoarseGrainedSchedulerBackend的reviveOffers就是发送给这个rpc的终端ReviveOffers信号。

```

1  override def reviveOffers() {
2      driverEndpoint.send(ReviveOffers)
3  }

```

CoarseGrainedSchedulerBackend.DriverEndpoint.makeOffers

DriverEndpoint有两种发送信息的函数。一个是send，发送信息后不需要对方回复。一个是ask，发送信息后需要对方回复。对应着，也有两种接收信息的函数。一个是receive，接收后不回复对方：

```

1  override def receive: PartialFunction[Any, Unit] = {
2      case StatusUpdate(executorId, taskId, state, data) =>
3          scheduler.statusUpdate(taskId, state, data.value)
4          if (TaskState.isFinished(state)) {
5              executorDataMap.get(executorId) match {
6                  case Some(executorInfo) =>
7                      executorInfo.freeCores += scheduler.CPUS_PER_TASK
8                      makeOffers(executorId)
9                  case None =>
10
11                      logWarning(s"Ignored task status update ($taskId state $state) " +
12                          s"from unknown executor with ID $executorId")
13              }
14          }
15
16      case ReviveOffers =>
17          makeOffers()
18
19      case KillTask(taskId, executorId, interruptThread) =>
20          executorDataMap.get(executorId) match {
21              case Some(executorInfo) =>
22                  executorInfo.executorEndpoint.send(KillTask(taskId, executorId, interruptThread))
23              case None =>
24
25                  logWarning(s"Attempted to kill task $taskId for unknown executor $executorId.")
26          }
27  }

```

另外一个`receiveAndReply`，接收后回复对方：

```

1  override def receiveAndReply(context: RpcCallContext): PartialFunction[Any, Unit] = {
2
3      case RegisterExecutor(executorId, executorRef, hostname, cores, logUrls) =>
4          if (executorDataMap.contains(executorId)) {
5              executorRef.send(RegisterExecutorFailed("Duplicate executor ID: " + executorId))
6              context.reply(true)
7          } else {
8
9              val executorAddress = if (executorRef.address != null) {
10                  executorRef.address
11              } else {
12                  context.senderAddress
13              }
14              logInfo(s"Registered executor $executorRef ($executorAddress) with ID $executorId")
15              addressToExecutorId(executorAddress) = executorId
16              totalCoreCount.addAndGet(cores)
17              totalRegisteredExecutors.addAndGet(1)
18              val data = new ExecutorData(executorRef, executorRef.address, hostname,
19                  cores, cores, logUrls)
20              CoarseGrainedSchedulerBackend.this.synchronized {
21                  executorDataMap.put(executorId, data)
22                  if (currentExecutorIdCounter < executorId.toInt) {
23                      currentExecutorIdCounter = executorId.toInt
24                  }
25                  if (numPendingExecutors > 0) {
26                      numPendingExecutors -= 1
27                      logDebug(s"Decrement number of pending executors ($numPendingExecutors left)")
28                  }
29              }
30              executorRef.send(RegisteredExecutor)
31              context.reply(true)
32              listenerBus.post(
33                  SparkListenerExecutorAdded(System.currentTimeMillis(), executorId, data))
34              makeOffers()

```

```

35     }
36
37     case StopDriver =>
38         context.reply(true)
39         stop()
40
41     case StopExecutors =>
42         logInfo("Asking each executor to shut down")
43         for ((_, executorData) <- executorDataMap) {
44             executorData.executorEndpoint.send(StopExecutor)
45         }
46         context.reply(true)
47
48     case RemoveExecutor(executorId, reason) =>
49
50         executorDataMap.get(executorId).foreach(_.executorEndpoint.send(StopExecutor))
51         removeExecutor(executorId, reason)
52         context.reply(true)
53
54     case RetrieveSparkAppConfig =>
55         val reply = SparkAppConfig(sparkProperties,
56             SparkEnv.get.securityManager.getIOEncryptionKey())
57         context.reply(reply)
58 }
59
60
61 private def makeOffers() {
62
63     val activeExecutors = executorDataMap.filterKeys(executorIsAlive)
64     val workOffers = activeExecutors.map { case (id, executorData) =>
65         new WorkerOffer(id, executorData.executorHost, executorData.freeCores)
66     }.toIndexedSeq
67     launchTasks(scheduler.resourceOffers(workOffers))
68 }

```

我们可以看到之前在CoarseGrainedSchedulerBackend的reviveOffers发送的ReviveOffers信号会在receive中被接收，从而调用makeOffers：

```

1     case ReviveOffers =>
2         makeOffers()

```

makeOffers做的工作为：

```

1     private def makeOffers() {
2         //过滤掉被杀死的Executor
3         val activeExecutors = executorDataMap.filterKeys(executorIsAlive)
4         //根据activeExecutors生成workOffers，
5         //即executor所能提供的资源信息。
6         val workOffers = activeExecutors.map { case (id, executorData) =>
7             new WorkerOffer(id, executorData.executorHost, executorData.freeCores)
8         }.toIndexedSeq
9         //scheduler.resourceOffers分配资源，
10        //并launchTasks发送任务
11        launchTasks(scheduler.resourceOffers(workOffers))
12    }

```

launchTasks主要的实现是向executor发送LaunchTask信号：

```

1     executorData.executorEndpoint.send(LaunchTask(new SerializableBuffer(serializedTask)))

```

TaskSchedulerImpl.resourceOffers

下面我们来深入上节scheduler.resourceOffers分配资源的函数：

```

1  def resourceOffers(offers: IndexedSeq[WorkerOffer]): Seq[Seq[TaskDescription]] = synchronized {
2      //标记每个活的节点并记录它的主机名
3      //并且追踪是否有新的executor加入
4      var newExecAvail = false
5      for (o <- offers) {
6          if (!hostToExecutors.contains(o.host)) {
7              hostToExecutors(o.host) = new HashSet[String]()
8          }
9          if (!executorIdToRunningTaskIds.contains(o.executorId)) {
10             hostToExecutors(o.host) += o.executorId
11             executorAdded(o.executorId, o.host)
12             executorIdToHost(o.executorId) = o.host
13             executorIdToRunningTaskIds(o.executorId) = HashSet[Long]()
14             newExecAvail = true
15         }
16         for (rack <- getRackForHost(o.host)) {
17             hostsByRack.getOrElseUpdate(rack, new HashSet[String]()) += o.host
18         }
19     }
20
21     // 为了避免将Task集中分配到某些机器，随机的打散它们
22     val shuffledOffers = Random.shuffle(offers)
23     // 建立每个worker的TaskDescription数组
24     val tasks = shuffledOffers.map(o => new ArrayBuffer[TaskDescription](o.cores))
25     //记录各个worker的available Cpus
26     val availableCpus = shuffledOffers.map(o => o.cores).toArray
27     //获取按照调度策略排序好的TaskSetManager
28     val sortedTaskSets = rootPool.getSortedTaskSetQueue
29     for (taskSet <- sortedTaskSets) {
30         logDebug("parentName: %s, name: %s, runningTasks: %s".format(
31             taskSet.parent.name, taskSet.name, taskSet.runningTasks))
32         //如果有新的executor加入
33         //则需要从新计算TaskSetManager的就近原则
34         if (newExecAvail) {
35             taskSet.executorAdded()
36         }
37     }
38     // 得到调度序列中的每个TaskSet,
39     // 然后按节点的locality级别增序分配资源
40     // Locality优先序列为: PROCESS_LOCAL, NODE_LOCAL, NO_PREF, RACK_LOCAL, ANY
41     for (taskSet <- sortedTaskSets) {
42         var launchedAnyTask = false
43         var launchedTaskAtCurrentMaxLocality = false
44         //按照就近原则分配
45         for (currentMaxLocality <- taskSet.myLocalityLevels) {
46             do {
47                 //resourceOfferSingleTaskSet为单个TaskSet分配资源,
48                 //若该LocalityLevel的节点下不能再为之分配资源了,
49                 //则返回false
50                 launchedTaskAtCurrentMaxLocality = resourceOfferSingleTaskSet(
51                     taskSet, currentMaxLocality, shuffledOffers, availableCpus, tasks)
52                 launchedAnyTask |= launchedTaskAtCurrentMaxLocality
53             } while (launchedTaskAtCurrentMaxLocality)
54         }
55         if (!launchedAnyTask) {
56             taskSet.abortIfCompletelyBlacklisted(hostToExecutors)
57         }
58     }
59
60     if (tasks.size > 0) {
61         hasLaunchedTask = true
62     }

```

```

63     return tasks
64 }

```

这里涉及到两个排序，首先调度器会对TaskSet进行排序：

```

1 val sortedTaskSets = rootPool.getSortedTaskSetQueue

```

取出每个TaskSet后，我们又会根据从近到远的Locality Level 的来对各个Task进行资源的分配。

TaskSchedulerImpl.resourceOfferSingleTaskSet

接下来我们来看下为单个TaskSet分配资源的具体实现：

```

1 private def resourceOfferSingleTaskSet(
2     taskSet: TaskSetManager,
3     maxLocality: TaskLocality,
4     shuffledOffers: Seq[WorkerOffer],
5     availableCpus: Array[Int],
6     tasks: IndexedSeq[ArrayBuffer[TaskDescription]]) : Boolean = {
7     var launchedTask = false
8     //遍历各个executor
9     for (i <- 0 until shuffledOffers.size) {
10         val execId = shuffledOffers(i).executorId
11         val host = shuffledOffers(i).host
12         if (availableCpus(i) >= CPUS_PER_TASK) {
13             try {
14                 //获取taskSet中，相对于该execId，host所能接收的最大距离maxLocality的task
15                 //maxLocality的值在TaskSchedulerImpl.resourceOffers中从近到远的遍历
16                 for (task <- taskSet.resourceOffer(execId, host, maxLocality)) {
17                     tasks(i) += task
18                     val tid = task.taskId
19                     taskIdToTaskSetManager(tid) = taskSet
20                     taskIdToExecutorId(tid) = execId
21                     executorIdToRunningTaskIds(execId).add(tid)
22                     availableCpus(i) -= CPUS_PER_TASK
23                     assert(availableCpus(i) >= 0)
24                     launchedTask = true
25                 }
26             } catch {
27                 case e: TaskNotSerializableException =>
28                     logError(s"Resource offer failed, task set ${taskSet.name} was not serializable")
29                     return launchedTask
30             }
31         }
32     }
33     return launchedTask
34 }

```

CoarseGrainedSchedulerBackend.DriverEndpoint.launchTasks

我们回到CoarseGrainedSchedulerBackend.DriverEndpoint.makeOffers，看最后一步，发送任务的函数launchTasks：

```

1 private def launchTasks(tasks: Seq[Seq[TaskDescription]]) {
2     for (task <- tasks.flatten) {
3         val serializedTask = ser.serialize(task)
4         //若序列化Task大小达到Rpc限制，
5         //则停止
6         if (serializedTask.limit >= maxRpcMessageSize) {
7             scheduler.taskIdToTaskSetManager.get(task.taskId).foreach { taskSetMgr =>
8                 try {
9                     var msg = "Serialized task %s:%d was %d bytes, which exceeds max allowed: " +
10                         "spark.rpc.message.maxSize (%d bytes). Consider increasing " +
11                         "spark.rpc.message.maxSize or using broadcast variables for large values."

```

```

12         msg = msg.format(task.taskId, task.index, serializedTask.limit, maxRpcMessageSize)
13         taskSetMgr.abort(msg)
14     } catch {
15         case e: Exception => logError("Exception in error callback", e)
16     }
17 }
18 }
19 else {
20     // 减少改task所对应的executor信息的core数量
21     val executorData = executorDataMap(task.executorId)
22     executorData.freeCores -= scheduler.CPU_PER_TASK
23
24     logDebug(s"Launching task ${task.taskId} on executor id: ${task.executorId} hostname: " +
25             s"${executorData.executorHost}.")
26
27     //向executorEndpoint 发送LaunchTask 信号
28     executorData.executorEndpoint.send(LaunchTask(new SerializableBuffer(serializedTask)))
29 }
30 }
31 }

```

executorEndpoint接收到LaunchTask信号（包含SerializableBuffer(serializedTask)）后，会开始执行任务。

调度任务

Pool.getSortedTaskSetQueue

上一章我们讲到TaskSchedulerImpl.resourceOffers中会调用：

```
1 val sortedTaskSets = rootPool.getSortedTaskSetQueue
```

获取按照调度策略排序好的TaskSetManager。接下来我们深入讲解这行代码。

rootPool是一个Pool对象。Pool定义为：一个可调度的实体，代表着Pool的集合或者TaskSet的集合，即Schedulable为一个接口，由Pool类和TaskSetManager类实现

getSortedTaskSetQueue：

```

1 override def getSortedTaskSetQueue: ArrayBuffer[TaskSetManager] = {
2     //生成TaskSetManager数组
3     var sortedTaskSetQueue = new ArrayBuffer[TaskSetManager]
4     //对调度实体进行排序
5     val sortedSchedulableQueue =
6         schedulableQueue.asScala.toSeq.sortWith(taskSetSchedulingAlgorithm.comparator)
7     for (schedulable <- sortedSchedulableQueue) {
8         //从调度实体中取得TaskSetManager数组
9         sortedTaskSetQueue ++= schedulable.getSortedTaskSetQueue
10    }
11    sortedTaskSetQueue
12 }

```

其中调度算法taskSetSchedulingAlgorithm，会在Pool被生成时候根据SchedulingMode被设定为FairSchedulingAlgorithm或者FIFOSchedulingAlgorithm

```

1 var taskSetSchedulingAlgorithm: SchedulingAlgorithm = {
2     schedulingMode match {
3         case SchedulingMode.FAIR =>
4             new FairSchedulingAlgorithm()
5         case SchedulingMode.FIFO =>
6             new FIFOSchedulingAlgorithm()
7         case _ =>
8             val msg = "Unsupported scheduling mode: $schedulingMode. Use FAIR or FIFO instead."

```

```

9         throw new IllegalArgumentException(msg)
10    }
11 }

```

TaskSchedulerImpl.initialize

Pool被生成是什么时候被生成的呢？我们来看下TaskSchedulerImpl的初始化就能发现：

```

1  def initialize(backend: SchedulerBackend) {
2      this.backend = backend
3      // 创建一个名字为空的rootPool
4
5      rootPool = new Pool("", schedulingMode, 0, 0)
6      schedulableBuilder = {
7          schedulingMode match {
8              //TaskSchedulerImpl在初始化时，
9              //根据SchedulingMode来创建不同的schedulableBuilder
10             case SchedulingMode.FIFO =>
11                 new FIFOSchedulableBuilder(rootPool)
12             case SchedulingMode.FAIR =>
13                 new FairSchedulableBuilder(rootPool, conf)
14             case _ =>
15                 throw new IllegalArgumentException(s"Unsupported spark.scheduler.mode: $schedulingMode")
16         }
17     }
18     schedulableBuilder.buildPools()
19 }

```

FIFO 调度

FIFOSchedulableBuilder.addTaskSetManager

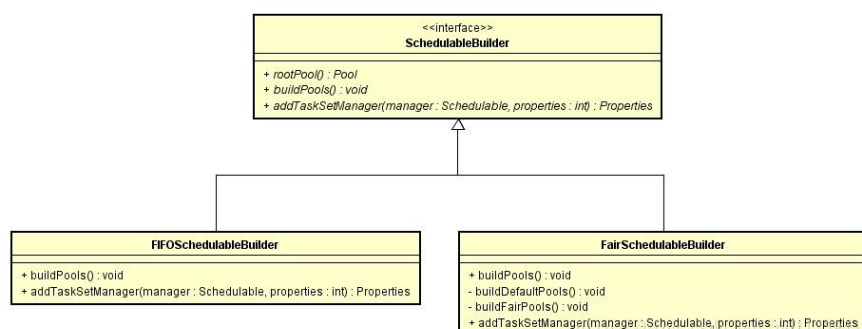
接下来，我们回过头看TaskSchedulerImpl.submitTasks中的schedulableBuilder.addTaskSetManager。

```

1  schedulableBuilder.addTaskSetManager(manager, manager.taskSet.properties)

```

我们深入讲一下addTaskSetManager：



schedulableBuilder 是一个 SchedulerBuilder 接口，SchedulerBuilder 接口由两个类 FIFOSchedulableBuilder 和 FairSchedulableBuilder实现。

我们这里先讲解FIFOSchedulableBuilder，FIFOSchedulableBuilder的addTaskSetManager：

```

1  override def addTaskSetManager(manager: Schedulable, properties: Properties) {
2      rootPool.addSchedulable(manager)
3  }

```

再看addSchedulable：


```

1  override def addSchedulable(schedulable: Schedulable) {
2      require(schedulable != null)
3      schedulableQueue.add(schedulable)
4      schedulableNameToSchedulable.put(schedulable.name, schedulable)
5      schedulable.parent = this
6  }

```

实际上是将 manager 加入到 schedulableQueue（这里是 FIFO 的 queue），将 manager 的 name 加入到一个名为 schedulableNameToSchedulable 的 ConcurrentHashMap[String, Schedulable] 中，并将 manager 的 parent 设置为 rootPool。

FIFOSchedulableBuilder.buildPools()

上述后一行代码：

```

1  schedulableBuilder.buildPools()

```

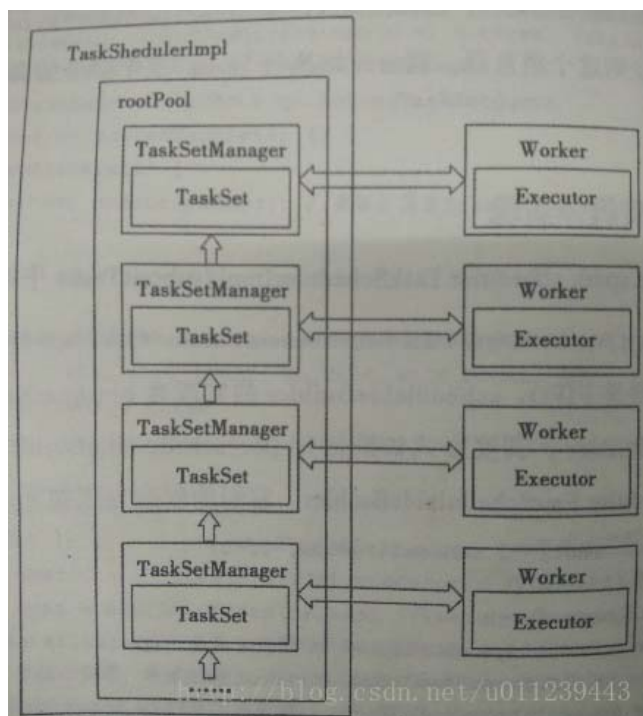
buildPools 会因不同的调度器而异。如果是 FIFOSchedulableBuilder，那么就为空：

```

1  override def buildPools() {
2      // nothing
3  }

```

这是因为 rootPool 里面不包含其他的 Pool，而是像上述所讲的直接将 manager 的 parent 设置为 rootPool。实际上，这是一种 2 层的树形结构，第 0 层为 rootPool，第二层叶子节点为各个 manager：



FIFOSchedulingAlgorithm

一切就绪后，我们可以来看 FIFO 的核心调度算法了：

```

1  private[spark] class FIFOSchedulingAlgorithm extends SchedulingAlgorithm {
2      override def comparator(s1: Schedulable, s2: Schedulable): Boolean = {
3          //priority 实际上是 Job ID
4          val priority1 = s1.priority
5          val priority2 = s2.priority
6          //先比较 Job ID
7          var res = math.signum(priority1 - priority2)
8          if (res == 0) {

```

```

9      //若Job ID相同,
10     //则比较 Stage ID
11     val stageId1 = s1.stageId
12     val stageId2 = s2.stageId
13     res = math.signum(stageId1 - stageId2)
14 }
15 res < 0
16 }
17 }

```

FAIR 调度

FairSchedulableBuilder.addTaskSetManager

FairSchedulableBuilder的addTaskSetManager会比FIFOSchedulableBuilder的复杂:

```

1  override def addTaskSetManager(manager: Schedulable, properties: Properties) {
2      //先生成一个默认的parentPool
3      var poolName = DEFAULT_POOL_NAME
4      var parentPool = rootPool.getSchedulableByName(poolName)
5      //若有配置信息,
6      //则根据配置信息得到poolName
7      if (properties != null) {
8          //FAIR_SCHEDULER_PROPERTIES = "spark.scheduler.pool"
9          poolName = properties.getProperty(FAIR_SCHEDULER_PROPERTIES, DEFAULT_POOL_NAME)
10         parentPool = rootPool.getSchedulableByName(poolName)
11         //若rootPool中没有这个pool
12         if (parentPool == null) {
13             //我们会根据用户在app上的配置生成新的pool,
14             //而不是根据xml 文件
15             parentPool = new Pool(poolName, DEFAULT_SCHEDULING_MODE,
16                 DEFAULT_MINIMUM_SHARE, DEFAULT_WEIGHT)
17             rootPool.addSchedulable(parentPool)
18             logInfo("Created pool %s, schedulingMode: %s, minShare: %d, weight: %d".format(
19                 poolName, DEFAULT_SCHEDULING_MODE, DEFAULT_MINIMUM_SHARE, DEFAULT_WEIGHT))
20         }
21     }
22     //将这个manager加入到这个pool
23     parentPool.addSchedulable(manager)
24     logInfo("Added task set " + manager.name + " tasks to pool " + poolName)
25 }
26 }

```

FairSchedulableBuilder.buildPools()

FairSchedulableBuilder.buildPools需要根据\$SPARK_HOME/conf/fairscheduler.xml文件来构建调度树。配置文件大致如下:

```

1  <allocations>
2      <pool name="production">
3          <schedulingMode>FAIR</schedulingMode>
4          <weight>1</weight>
5          <minShare>2</minShare>
6      </pool>
7      <pool name="test">
8          <schedulingMode>FIFO</schedulingMode>
9          <weight>2</weight>
10         <minShare>3</minShare>
11     </pool>
12 </allocations>

```

buildFairSchedulerPool:

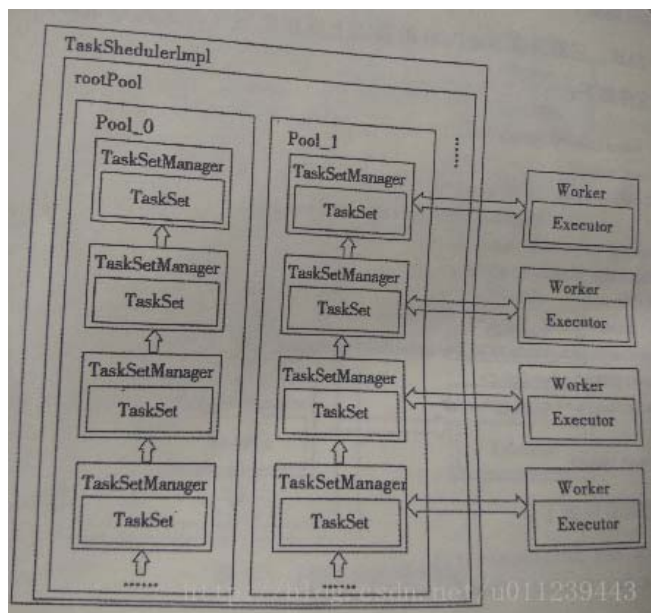
```

1
2 private def buildFairSchedulerPool(is: InputStream) {
3     //加载xml 文件
4     val xml = XML.load(is)
5     //遍历
6     for (poolNode <- (xml \ POOLS_PROPERTY)) {
7
8         val poolName = (poolNode \ POOL_NAME_PROPERTY).text
9         val schedulingMode = DEFAULT_SCHEDULING_MODE
10        val minShare = DEFAULT_MINIMUM_SHARE
11        val weight = DEFAULT_WEIGHT
12
13        val xmlSchedulingMode = (poolNode \ SCHEDULING_MODE_PROPERTY).text
14        if (xmlSchedulingMode != "") {
15            try {
16                schedulingMode = SchedulingMode.withName(xmlSchedulingMode)
17            } catch {
18                case e: NoSuchElementException =>
19                    logWarning(s"Unsupported schedulingMode: $xmlSchedulingMode, " +
20                        s"using the default schedulingMode: $schedulingMode")
21            }
22        }
23
24        val xmlMinShare = (poolNode \ MINIMUM_SHARES_PROPERTY).text
25        if (xmlMinShare != "") {
26            minShare = xmlMinShare.toInt
27        }
28
29        val xmlWeight = (poolNode \ WEIGHT_PROPERTY).text
30        if (xmlWeight != "") {
31            weight = xmlWeight.toInt
32        }
33        //根据xml的配置，
34        //最终生成一个新的Pool
35        val pool = new Pool(poolName, schedulingMode, minShare, weight)
36        //将这个Pool加入到rootPool中
37        rootPool.addSchedulable(pool)
38        logInfo("Created pool %s, schedulingMode: %s, minShare: %d, weight: %d".format(
39            poolName, schedulingMode, minShare, weight))
40    }
41 }

```

可想而知，FAIR 调度并不是简单的公平调度。我们会先根据xml配置文件生成很多pool加入rootPool中，而每个app会根据配置“spark.scheduler.pool”的poolName，将TaskSetManager加入到某个pool中。其实，rootPool还会对Pool也进程一次调度。

所以，在FAIR调度策略中包含了两层调度。第一层的rootPool内的多个Pool，第二层是Pool内的多个TaskSetManager。fairscheduler.xml文件中，weight（任务权重）和minShare（最小任务数）是用来设置第一层调度的，该调度使用的是FAIR算法。而第二层调度由schedulingMode设置。



但对于Standalone模式下的单个app，FAIR调度的多个Pool显得鸡肋，因为app只能选择一个Pool。但是我们可以代码级别硬编码的去分配：

```
1 sc.setLocalProperty("spark.scheduler.pool", "Pool_1")
```

FAIRSchedulingAlgorithm

接下来，我们就来讲解FAIR算法：

```
1 private[spark] class FairSchedulingAlgorithm extends SchedulingAlgorithm {
2   override def comparator(s1: Schedulable, s2: Schedulable): Boolean = {
3     val minShare1 = s1.minShare
4     val minShare2 = s2.minShare
5     val runningTasks1 = s1.runningTasks
6     val runningTasks2 = s2.runningTasks
7     //若s1运行的任务数小于s1的最小任务数
8     val s1Needy = runningTasks1 < minShare1
9     //若s2运行的任务数小于s2的最小任务数
10    val s2Needy = runningTasks2 < minShare2
11    //minShareRatio = 运行的任务数/最小任务数
12    //代表着负载程度，越小，负载越小
13    val minShareRatio1 = runningTasks1.toDouble / math.max(minShare1, 1.0)
14    val minShareRatio2 = runningTasks2.toDouble / math.max(minShare2, 1.0)
15    //taskToWeightRatio = 运行的任务数/权重
16    //权重越大，越优先
17    //即taskToWeightRatio 越小 越优先
18    val taskToWeightRatio1 = runningTasks1.toDouble / s1.weight.toDouble
19    val taskToWeightRatio2 = runningTasks2.toDouble / s2.weight.toDouble
20
21    var compare = 0
22    //若s1运行的任务小于s1的最小任务数，而s2不然
23    //则s1优先
24    if (s1Needy && !s2Needy) {
25      return true
26    }
27    //若s2运行的任务小于s2的最小任务数，而s1不然
28    //则s2优先
29    else if (!s1Needy && s2Needy) {
30      return false
31    }
32    //若s1 s2 运行的任务都小于自己的最小任务数
```

```
33 //比较minShareRatio, 哪个小, 哪个优先
34 else if (s1Needy && s2Needy) {
35     compare = minShareRatio1.compareTo(minShareRatio2)
36 }
37 //若s1 s2 运行的任务都不小于自己的最小任务数
38 //比较taskToWeightRatio, 哪个小, 哪个优先
39 else {
40     compare = taskToWeightRatio1.compareTo(taskToWeightRatio2)
41 }
42 if (compare < 0) {
43     true
44 } else if (compare > 0) {
45     false
46 } else {
47     s1.name < s2.name
48 }
49 }
50 }
```

至此, TaskScheduler在发送任务给executor前的工作就全部完成了。