

深入理解Spark 2.1 Core（七）：任务执行的原理与源码分析

■ 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/54143839>

上篇博文《[深入理解Spark 2.1 Core（六）：资源调度的实现与源码分析](#)》中我们讲解了，AppClient和Executor是如何启动，如何为逻辑上与物理上的资源调度，以及分析了在Spark1.4之前逻辑上资源调度算法的bug。

这篇博文，我们就来讲讲Executor启动后，是如何在Executor上执行Task的，以及其后续处理。

执行Task

我们在《[深入理解Spark 2.1 Core（三）：任务调度器的原理与源码分析](#)》中提到了，任务调度完成后，CoarseGrainedSchedulerBackend.DriverEndpoint会调用launchTasks向CoarseGrainedExecutorBackend发送带着serializedTask的LaunchTask信号。接下来，我们就来讲讲CoarseGrainedExecutorBackend接收到LaunchTask信号后，是如何执行Task的。

调用栈如下：

- CoarseGrainedExecutorBackend.receive
 - Executor.launchTask
 - Executor.TaskRunner.run
 - Executor.updateDependencies
 - Task.run
 - ShuffleMapTask.runTask
 - ResultTask.runTask

CoarseGrainedExecutorBackend.receive

```
1 case LaunchTask(data) =>
2   if (executor == null) {
3     exitExecutor(1, "Received LaunchTask command but executor was null")
4   } else {
5     // 反序列化task描述
6     val taskDesc = ser.deserialize[TaskDescription](data.value)
7     logInfo("Got assigned task " + taskDesc.taskId)
8     // 调用executor.launchTask
9     executor.launchTask(this, taskId = taskDesc.taskId, attemptNumber = taskDesc.attemptNumber,
10      taskDesc.name, taskDesc.serializedTask)
11   }
```

Executor.launchTask

```
1 def launchTask(
2   context: ExecutorBackend,
3   taskId: Long,
4   attemptNumber: Int,
5   taskName: String,
6   serializedTask: ByteBuffer): Unit = {
7   // 创建TaskRunner
8   val tr = new TaskRunner(context, taskId = taskId, attemptNumber = attemptNumber, taskName,
9     serializedTask)
10   // 把taskId 以及 对应的 TaskRunner,
11   // 加入到 ConcurrentHashMap[Long, TaskRunner]
12   runningTasks.put(taskId, tr)
```

```
13 // 线程池 执行 TaskRunner
14 threadPool.execute(tr)
15 }
```

Executor.TaskRunner.run

```
1  override def run(): Unit = {
2      val threadMXBean = ManagementFactory.getThreadMXBean
3      val taskMemoryManager = new TaskMemoryManager(env.memoryManager, taskId)
4      // 记录开始反序列化的时间
5      val deserializeStartTime = System.currentTimeMillis()
6      // 记录开始反序列化的时的Cpu时间
7      val deserializeStartCpuTime = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
8          threadMXBean.getCurrentThreadCpuTime
9      } else 0L
10     Thread.currentThread().setContextClassLoader(replClassLoader)
11     val ser = env.closureSerializer.newInstance()
12     logInfo(s"Running $taskName (TID $taskId)")
13     execBackend.statusUpdate(taskId, TaskState.RUNNING, EMPTY_BYTE_BUFFER)
14     var taskStart: Long = 0
15     var taskStartCpu: Long = 0
16     // 开始GC的时间
17     startGCtime = computeTotalGcTime()
18
19     try {
20         //反序列化任务信息
21         val (taskFiles, taskJars, taskProps, taskBytes) =
22             Task.deserializeWithDependencies(serializedTask)
23
24         // 根据taskProps设置executor属性
25         Executor.taskDeserializationProps.set(taskProps)
26         // 根据taskFiles和taskJars,
27         // 下载任务所需的File 和 加载所需的Jar包
28         updateDependencies(taskFiles, taskJars)
29         // 根据taskBytes生成task
30         task = ser.deserialize[Task[Any]](taskBytes, Thread.currentThread().getContextClassLoader)
31         //设置task属性
32         task.localProperties = taskProps
33         //设置task内存管理
34         task.setTaskMemoryManager(taskMemoryManager)
35
36         // 若在反序列化之前Task就被kill了,
37         // 抛出异常
38         if (killed) {
39             throw new TaskKilledException
40         }
41
42         logDebug("Task " + taskId + "'s epoch is " + task.epoch)
43         //更新mapOutputTracker Epoch 为task epoch
44         env.mapOutputTracker.updateEpoch(task.epoch)
45
46         // 记录任务开始时间
47         taskStart = System.currentTimeMillis()
48         // 记录任务开始时的cpu时间
49         taskStartCpu = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
50             threadMXBean.getCurrentThreadCpuTime
51         } else 0L
52         var threwException = true
53         val value = try {
54             // 运行Task
55             val res = task.run(
56                 taskAttemptId = taskId,
57                 attemptNumber = attemptNumber,
58                 metricsSystem = env.metricsSystem)
```

```
59     throwException = false
60     res
61 } finally {
62     val releasedLocks = env.blockManager.releaseAllLocksForTask(taskId)
63     val freedMemory = taskMemoryManager.cleanUpAllAllocatedMemory()
64
65     if (freedMemory > 0 && !throwException) {
66         val errMsg = s"Managed memory leak detected; size = $freedMemory bytes, TID = $taskId"
67         if (conf.getBoolean("spark.unsafe.exceptionOnMemoryLeak", false)) {
68             throw new SparkException(errMsg)
69         } else {
70             logWarning(errMsg)
71         }
72     }
73
74     if (releasedLocks.nonEmpty && !throwException) {
75         val errMsg =
76             s"${releasedLocks.size} block locks were not released by TID = $taskId:\n" +
77             releasedLocks.mkString("[", ", ", "]")
78         if (conf.getBoolean("spark.storage.exceptionOnPinLeak", false)) {
79             throw new SparkException(errMsg)
80         } else {
81             logWarning(errMsg)
82         }
83     }
84 }
85 // 记录任务结束时间
86 val taskFinish = System.currentTimeMillis()
87 // 记录任务结束时的cpu时间
88 val taskFinishCpu = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
89     threadMXBean.getCurrentThreadCpuTime
90 } else 0L
91
92 // 若task在运行中被kill了
93 // 则抛出异常
94 if (task.killed) {
95     throw new TaskKilledException
96 }
97
98
99 val resultSer = env.serializer.newInstance()
100 // 结果记录序列化开始的系统时间
101 val beforeSerialization = System.currentTimeMillis()
102 // 序列化结果
103 val valueBytes = resultSer.serialize(value)
104 // 结果记录序列化完成的系统时间
105 val afterSerialization = System.currentTimeMillis()
106
107 // 反序列化发生在两个地方:
108 // 1. 在该函数下反序列化Task信息以及Task实例。
109 // 2. 在任务启动后, Task.run 反序列化 RDD 和 函数
110
111 // 计算task的反序列化费时
112 task.metrics.setExecutorDeserializeTime(
113     (taskStart - deserializeStartTime) + task.executorDeserializeTime)
114 // 计算task的反序列化cpu费时
115 task.metrics.setExecutorDeserializeCpuTime(
116     (taskStartCpu - deserializeStartCpuTime) + task.executorDeserializeCpuTime)
117 // 计算task运行费时
118 task.metrics.setExecutorRunTime((taskFinish - taskStart) - task.executorDeserializeTime)
119 // 计算task运行cpu费时
120 task.metrics.setExecutorCpuTime(
121     (taskFinishCpu - taskStartCpu) - task.executorDeserializeCpuTime)
122 // 计算GC时间
```

```
task.metrics.setJvmGcTime(computeTotalGcTime() - startGcTime)

//计算结果序列化时间
task.metrics.setResultSerializationTime(afterSerialization - beforeSerialization)

val accumUpdates = task.collectAccumulatorUpdates()
// 这里代码存在缺陷:
// value相当于被序列化了两次
val directResult = new DirectTaskResult(valueBytes, accumUpdates)
val serializedDirectResult = ser.serialize(directResult)
// 得到结果的大小
val resultSize = serializedDirectResult.limit

// 对于计算结果, 会根据结果的大小有不同的策略:
// 1.生成结果在(正无穷,1GB):
// 超过1GB的部分结果直接丢弃,
// 可以通过spark.driver.maxResultSize实现
// 默认为1G
// 2.生成结果大小在[1GB,128MB - 200KB]
// 会把该结果以taskId为编号存入BlockManager中,
// 然后把该编号通过Netty发送给Driver,
// 该阈值是Netty框架传输的最大值
// spark.akka.frameSize (默认为128MB) 和Netty的预留空间reservedSizeBytes (200KB) 的差值
// 3.生成结果大小在 (128MB - 200KB,0) :
// 直接通过Netty发送到Driver
val serializedResult: ByteBuffer = {
  if (maxResultSize > 0 && resultSize > maxResultSize) {
    logWarning(s"Finished $taskName (TID $taskId). Result is larger than maxResultSize " +
      s"(${Utils.bytesToString(resultSize)} > ${Utils.bytesToString(maxResultSize)}), " +
      s"dropping it.")
    ser.serialize(new IndirectTaskResult[Any](TaskResultBlockId(taskId), resultSize))
  } else if (resultSize > maxDirectResultSize) {
    val blockId = TaskResultBlockId(taskId)
    env.blockManager.putBytes(
      blockId,
      new ChunkedByteBuffer(serializedDirectResult.duplicate()),
      StorageLevel.MEMORY_AND_DISK_SER)
    logInfo(
      s"Finished $taskName (TID $taskId). $resultSize bytes result sent via BlockManager")
    ser.serialize(new IndirectTaskResult[Any](blockId, resultSize))
  } else {
    logInfo(s"Finished $taskName (TID $taskId). $resultSize bytes result sent to driver")
    serializedDirectResult
  }
}
// 更新execBackend 状态
execBackend.statusUpdate(taskId, TaskState.FINISHED, serializedResult)

} catch {
  case ffe: FetchFailedException =>
    val reason = ffe.toTaskFailedReason
    setTaskFinishedAndClearInterruptStatus()
    execBackend.statusUpdate(taskId, TaskState.FAILED, ser.serialize(reason))

  case _: TaskKilledException =>
    logInfo(s"Executor killed $taskName (TID $taskId)")
    setTaskFinishedAndClearInterruptStatus()
    execBackend.statusUpdate(taskId, TaskState.KILLED, ser.serialize(TaskKilled))

  case _: InterruptedException if task.killed =>
    logInfo(s"Executor interrupted and killed $taskName (TID $taskId)")
    setTaskFinishedAndClearInterruptStatus()
    execBackend.statusUpdate(taskId, TaskState.KILLED, ser.serialize(TaskKilled))
}
```

```

187
188     case CausedBy(cDE: CommitDeniedException) =>
189         val reason = cDE.toTaskFailedReason
190         setTaskFinishedAndClearInterruptStatus()
191         execBackend.statusUpdate(taskId, TaskState.FAILED, ser.serialize(reason))
192
193     case t: Throwable =>
194         logError(s"Exception in $taskName (TID $taskId)", t)
195
196         val accums: Seq[AccumulatorV2[_, _]] =
197             if (task != null) {
198                 task.metrics.setExecutorRunTime(System.currentTimeMillis() - taskStart)
199                 task.metrics.setJvmGCTime(computeTotalGcTime() - startGCtime)
200                 task.collectAccumulatorUpdates(taskFailed = true)
201             } else {
202                 Seq.empty
203             }
204
205         val accUpdates = accums.map(acc => acc.toInfo(Some(acc.value), None))
206
207         val serializedTaskEndReason = {
208             try {
209                 ser.serialize(new ExceptionFailure(t, accUpdates).withAccums(accums))
210             } catch {
211                 case _: NotSerializableException =>
212                     ser.serialize(new ExceptionFailure(t, accUpdates, false).withAccums(accums))
213             }
214         }
215         setTaskFinishedAndClearInterruptStatus()
216         execBackend.statusUpdate(taskId, TaskState.FAILED, serializedTaskEndReason)
217
218         if (Utils.isFatalError(t)) {
219             SparkUncaughtExceptionHandler.uncaughtException(t)
220         }
221
222     } finally {
223         // 任务结束后移除
224         runningTasks.remove(taskId)
225     }
226 }

```

Executor.updateDependencies

接下来, 我们来看看更新executor的依赖, 即下载任务所需的File和加载所需的Jar包:

```

1  private def updateDependencies(newFiles: HashMap[String, Long], newJars: HashMap[String, Long]) {
2      lazy val hadoopConf = SparkHadoopUtil.get.newConfiguration(conf)
3      synchronized {
4          // 下载任务所需的File
5          for ((name, timestamp) <- newFiles if currentFiles.getOrElse(name, -1L) < timestamp) {
6              logInfo("Fetching " + name + " with timestamp " + timestamp)
7              Utils.fetchFile(name, new File(SparkFiles.getRootDirectory()), conf,
8                  env.securityManager, hadoopConf, timestamp, useCache = !isLocal)
9              currentFiles(name) = timestamp
10         }
11         // 加载所需的Jar包
12         for ((name, timestamp) <- newJars) {
13             val localName = name.split("/").last
14             val currentTimeStamp = currentJars.get(name)
15                 .orElse(currentJars.get(localName))
16                 .getOrElse(-1L)
17             if (currentTimeStamp < timestamp) {
18                 logInfo("Fetching " + name + " with timestamp " + timestamp)
19                 Utils.fetchFile(name, new File(SparkFiles.getRootDirectory()), conf,

```

```

20         env.securityManager, hadoopConf, timestamp, useCache = !isLocal)
21         currentJars(name) = timestamp
22         // 把它加入到 class loader
23         val url = new File(SparkFiles.getRootDirectory(), localName).toURI.toURL
24         if (!urlClassLoader.getURLs().contains(url)) {
25             logInfo("Adding " + url + " to class loader")
26             urlClassLoader.addURL(url)
27         }
28     }
29 }
30 }
31 }

```

Task.run

接下来, 我们来看看这篇博文最核心的部分——task运行:

```

1  final def run(
2      taskAttemptId: Long,
3      attemptNumber: Int,
4      metricsSystem: MetricsSystem): T = {
5      SparkEnv.get.blockManager.registerTask(taskAttemptId)
6      //创建TaskContextImpl
7      context = new TaskContextImpl(
8          stageId,
9          partitionId,
10         taskAttemptId,
11         attemptNumber,
12         taskMemoryManager,
13         localProperties,
14         metricsSystem,
15         metrics)
16         //在TaskContext中设置TaskContextImpl
17         TaskContext.setTaskContext(context)
18         taskThread = Thread.currentThread()
19
20         if (!_killed) {
21             kill(interruptThread = false)
22         }
23
24         new CallerContext("TASK", appId, appAttemptId, jobId, Option(stageId), Option(stageAttemptId),
25             Option(taskAttemptId), Option(attemptNumber)).setCurrentContext()
26
27         try {
28             // 调用runTask
29             runTask(context)
30         } catch {
31             case e: Throwable =>
32                 try {
33                     context.markTaskFailed(e)
34                 } catch {
35                     case t: Throwable =>
36                         e.addSuppressed(t)
37                 }
38             throw e
39         } finally {
40             // 标记Task完成
41             context.markTaskCompleted()
42             try {
43                 Utils.tryLogNonFatalError {
44                     // 释放内存
45                     SparkEnv.get.blockManager.memoryStore.releaseUnrollMemoryForThisTask(MemoryMode.ON_HEAP)
46                     SparkEnv.get.blockManager.memoryStore.releaseUnrollMemoryForThisTask(MemoryMode.OFF_HEAP)
47                     val memoryManager = SparkEnv.get.memoryManager

```

```

48         memoryManager.synchronized { memoryManager.notifyAll() }
49     }
50 } finally {
51     //取消TaskContext设置
52     TaskContext.unset()
53 }
54 }
55 }

```

Task有两个子类，一个是最最后的Stage的Task，ShuffleMapTask；一个是最最后的Stage的Task，ResultTask。它们都覆盖了Task的runTask方法，接下来我们就分别来讲下它们的runTask方法。

ShuffleMapTask.runTask

根据每个Stage的partition数量来生成ShuffleMapTask，ShuffleMapTask会根据下游的Partition数量和Shuffle的策略来生成一系列文件。

```

1  override def runTask(context: TaskContext): MapStatus = {
2
3      val threadMXBean = ManagementFactory.getThreadMXBean
4      // 记录反序列化开始时间
5      val deserializeStartTime = System.currentTimeMillis()
6      // 记录反序列化开始时的Cpu时间
7      val deserializeStartCpuTime = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
8          threadMXBean.getCurrentThreadCpuTime
9      } else 0L
10     val ser = SparkEnv.get.closureSerializer.newInstance()
11     // 反序列化rdd 及其 依赖
12     val (rdd, dep) = ser.deserialize[(RDD[_], ShuffleDependency[_], _, _)](
13         ByteBuffer.wrap(taskBinary.value), Thread.currentThread.getContextClassLoader)
14     // 计算 反序列化费时
15     _executorDeserializeTime = System.currentTimeMillis() - deserializeStartTime
16     // 计算 反序列化Cpu费时
17     _executorDeserializeCpuTime = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
18         threadMXBean.getCurrentThreadCpuTime - deserializeStartCpuTime
19     } else 0L
20
21     var writer: ShuffleWriter[Any, Any] = null
22     try {
23         //获取shuffleManager
24         val manager = SparkEnv.get.shuffleManager
25         // writer
26         writer = manager.getWriter[Any, Any](dep.shuffleHandle, partitionId, context)
27         // 调用writer.write 开始计算RDD,
28         // 这部分 我们会在后续博文讲解
29         writer.write(rdd.iterator(partition, context).asInstanceOf[Iterator[_ <: Product2[Any, Any]]])
30         // 停止计算, 并返回结果
31         writer.stop(success = true).get
32     } catch {
33         case e: Exception =>
34             try {
35                 if (writer != null) {
36                     writer.stop(success = false)
37                 }
38             } catch {
39                 case e: Exception =>
40                     log.debug("Could not stop writer", e)
41             }
42             throw e
43     }
44 }

```

ResultTask.runTask

```
1  override def runTask(context: TaskContext): U = {
2      val threadMXBean = ManagementFactory.getThreadMXBean
3      // 记录反序列化开始时间
4      val deserializeStartTime = System.currentTimeMillis()
5      // 记录反序列化开始时的Cpu时间
6      val deserializeStartCpuTime = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
7          threadMXBean.getCurrentThreadCpuTime
8      } else 0L
9      val ser = SparkEnv.get.closureSerializer.newInstance()
10     // 反序列化rdd 及其 作用于RDD的结果函数
11     val (rdd, func) = ser.deserialize[(RDD[T], (TaskContext, Iterator[T]) => U)](
12         ByteBuffer.wrap(taskBinary.value), Thread.currentThread.getContextClassLoader)
13     // 计算 反序列化费时
14     _executorDeserializeTime = System.currentTimeMillis() - deserializeStartTime
15     // 计算 反序列化Cpu费时
16     _executorDeserializeCpuTime = if (threadMXBean.isCurrentThreadCpuTimeSupported) {
17         threadMXBean.getCurrentThreadCpuTime - deserializeStartCpuTime
18     } else 0L
19     // 这部分 我们会在后续博文讲解
20     func(context, rdd.iterator(partition, context))
21 }
```

后续处理

计量统计

对各个费时的统计，上章已经讲解。

回收内存

这在上章Task.run也已经讲解。

处理执行结果

Executor.TaskRunner.run的execBackend.statusUpdate，在《[深入理解Spark 2.1 Core （四）：运算结果处理和容错的原理与源码分析](#)》中我们已经讲解过。