# 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

在上一篇博文《深入理解Spark 2.1 Core （三）： 任务调度器的实现与源码分析 》TaskScheduler在发送任务给executor前的工作就全部完成了。这篇博文，我们来看看当executor计算完任务后，Spark是如何处理获取的计算结果与容错的。

调用栈如下：

- TaskSchedulerImpl.statusUpdate
  - TaskResultGetter.enqueueSuccessfulTask
    - TaskSchedulerImpl.handleSuccessfulTask
      - TaskSetManager.handleSuccessfulTask
        - DAGScheduler.taskEnded
          - DAGSchedulerEventProcessLoop.doOnReceive
            - DAGScheduler.handleTaskCompletion
  - TaskResultGetter.enqueueFailedTask
    - TaskSchedulerImpl.handleFailedTask
      - TaskSetManager.handleFailedTask
        - DAGScheduler.taskEnded
          - DAGSchedulerEventProcessLoop.doOnReceive
            - DAGScheduler.handleTaskCompletion

## TaskSchedulerImpl.statusUpdate

TaskRunner将任务的执行结果发送给DriverEndPoint，DriverEndPoint会转给TaskSchedulerImpl的statusUpdate：

```
 1  def statusUpdate(tid: Long, state: TaskState, serializedData: ByteBuffer) {
 2      var failedExecutor: Option[String] = None
 3      var reason: Option[ExecutorLossReason] = None
 4      synchronized {
 5        try {
 6          taskIdToTaskSetManager.get(tid) match {
 7            case Some(taskSet) =>
 8              //这只针对Mesos调度模式
 9              if (state == TaskState.LOST) {
10                val execId = taskIdToExecutorId.getOrElse(tid, throw new IllegalStateException(
11                  "taskIdToTaskSetManager.contains(tid) <=> taskIdToExecutorId.contains(tid)"))
12                if (executorIdToRunningTaskIds.contains(execId)) {
13                  reason = Some(
14                    SlaveLost(s"Task $tid was lost, so marking the executor as lost as well."))
15                  removeExecutor(execId, reason.get)
16                  failedExecutor = Some(execId)
17                }
18              }
19              //FINISHED KILLED LOST 都属于 isFinished
20              if (TaskState.isFinished(state)) {
21                cleanupTaskState(tid)
```

```
22            taskSet.removeRunningTask(tid)
23            //若FINISHED调用taskResultGetter.enqueueSuccessfulTask,
24            //否则调用taskResultGetter.enqueueFailedTask(taskSet, tid, state, serializedData)
25            if (state == TaskState.FINISHED) {
26              taskResultGetter.enqueueSuccessfulTask(taskSet, tid, serializedData)
27            } else if (Set(TaskState.FAILED, TaskState.KILLED, TaskState.LOST).contains(state)) {
28              taskResultGetter.enqueueFailedTask(taskSet, tid, state, serializedData)
29            }
30          }
31          case None =>
32            logError(
33              ("Ignoring update with state %s for TID %s because its task set is gone (this is " +
34                "likely the result of receiving duplicate task finished status updates) or its " +
35                "executor has been marked as failed.")
36                .format(state, tid))
37        }
38      } catch {
39        case e: Exception => logError("Exception in statusUpdate", e)
40      }
41    }
42    if (failedExecutor.isDefined) {
43      assert(reason.isDefined)
44      dagScheduler.executorLost(failedExecutor.get, reason.get)
45      backend.reviveOffers()
46    }
47  }
```

# 处理执行成功的结果

我们先来看下处理执行成功的结果的运行机制：

## TaskResultGetter.enqueueSuccessfulTask

```
1    def enqueueSuccessfulTask(
2        taskSetManager: TaskSetManager,
3        tid: Long,
4        serializedData: ByteBuffer): Unit = {
5        //通过线程池来获取结果
6      getTaskResultExecutor.execute(new Runnable {
7        override def run(): Unit = Utils.logUncaughtExceptions {
8          try {
9            val (result, size) = serializer.get().deserialize[TaskResult[_]](serializedData) match {
10            //可以直接获取到的结果
11              case directResult: DirectTaskResult[_] =>
12            //判断大小是否符合要求
13                if (!taskSetManager.canFetchMoreResults(serializedData.limit())) {
14                  return
15                }
16
17          directResult.value(taskResultSerializer.get())
18                (directResult, serializedData.limit())
19                //若不能直接获取到结果
20              case IndirectTaskResult(blockId, size) =>
21                if (!taskSetManager.canFetchMoreResults(size)) {
22                  //  判断大小是否符合要求,
23                  //若不符合则远程的删除计算结果
24                  sparkEnv.blockManager.master.removeBlock(blockId)
25                  return
26                }
27                logDebug("Fetching indirect task result for TID %s".format(tid))
28                scheduler.handleTaskGettingResult(taskSetManager, tid)
29                val serializedTaskResult = sparkEnv.blockManager.getRemoteBytes(blockId)
30                //从远程获取计算结果
```

```
31        if (!serializedTaskResult.isDefined) {
32            //若在任务执行结束后与我们去获取结果之间机器出现故障了
33            //或者block manager 不得不刷新结果了
34            //那么我们将不能够获取到结果
35              scheduler.handleFailedTask(
36                taskSetManager, tid, TaskState.FINISHED, TaskResultLost)
37              return
38          }
39          val deserializedResult = serializer.get().deserialize[DirectTaskResult[_]](
40            serializedTaskResult.get.toByteBuffer)
41          // 反序列化
42          deserializedResult.value(taskResultSerializer.get())
43          sparkEnv.blockManager.master.removeBlock(blockId)
44          (deserializedResult, size)
45        }
46
47
48        result.accumUpdates = result.accumUpdates.map { a =>
49          if (a.name == Some(InternalAccumulator.RESULT_SIZE)) {
50            val acc = a.asInstanceOf[LongAccumulator]
51            assert(acc.sum == 0L, "task result size should not have been set on the executors")
52            acc.setValue(size.toLong)
53            acc
54          } else {
55            a
56          }
57        }
58
59        //处理获取到的计算结果
60        scheduler.handleSuccessfulTask(taskSetManager, tid, result)
61      } catch {
62        case cnf: ClassNotFoundException =>
63          val loader = Thread.currentThread.getContextClassLoader
64          taskSetManager.abort("ClassNotFound with classloader: " + loader)
65        case NonFatal(ex) =>
66          logError("Exception while getting task result", ex)
67          taskSetManager.abort("Exception while getting task result: %s".format(ex))
68      }
69    }
70  })
71  }
```

## TaskSchedulerImpl.handleSuccessfulTask

调用taskSetManager.handleSuccessfulTask

```
1  def handleSuccessfulTask(
2      taskSetManager: TaskSetManager,
3      tid: Long,
4      taskResult: DirectTaskResult[_]): Unit = synchronized {
5    taskSetManager.handleSuccessfulTask(tid, taskResult)
6  }
```

## TaskSetManager.handleSuccessfulTask

```
1  def handleSuccessfulTask(tid: Long, result: DirectTaskResult[_]): Unit = {
2    val info = taskInfos(tid)
3    val index = info.index
4    info.markFinished(TaskState.FINISHED)
5    //从RunningTask中移除该task
6    removeRunningTask(tid)
7   //通知dagScheduler该task完成
8    sched.dagScheduler.taskEnded(tasks(index), Success, result.value(), result.accumUpdates, info)
```

```
9        //杀死所有其他与之相同的task的尝试
10       for (attemptInfo <- taskAttempts(index) if attemptInfo.running) {
11         logInfo(s"Killing attempt ${attemptInfo.attemptNumber} for task ${attemptInfo.id} " +
12           s"in stage ${taskSet.id} (TID ${attemptInfo.taskId}) on ${attemptInfo.host} " +
13           s"as the attempt ${info.attemptNumber} succeeded on ${info.host}")
14         sched.backend.killTask(attemptInfo.taskId, attemptInfo.executorId, true)
15       }
16       if (!successful(index)) {
17         //计数
18         tasksSuccessful += 1
19         logInfo(s"Finished task ${info.id} in stage ${taskSet.id} (TID ${info.taskId}) in" +
20           s" ${info.duration} ms on ${info.host} (executor ${info.executorId})" +
21           s" ($tasksSuccessful/$numTasks)")
22         //若果有所task成功了，
23         //那么标记successful，并且停止
24         successful(index) = true
25         if (tasksSuccessful == numTasks) {
26           isZombie = true
27         }
28       } else {
29         logInfo("Ignoring task-finished event for " + info.id + " in stage " + taskSet.id +
30           " because task " + index + " has already completed successfully")
31       }
32       maybeFinishTaskSet()
33     }
```

## DAGScheduler.taskEnded

我们再深入看下是如何通知dagScheduler该task完成的：

```
1    def taskEnded(
2        task: Task[_],
3        reason: TaskEndReason,
4        result: Any,
5        accumUpdates: Seq[AccumulatorV2[_, _]],
6        taskInfo: TaskInfo): Unit = {
7      //发送CompletionEvent信号
8      eventProcessLoop.post(
9        CompletionEvent(task, reason, result, accumUpdates, taskInfo))
10   }
```

## DAGSchedulerEventProcessLoop.doOnReceive

上一篇博文讲过，DAGSchedulerEventProcessLoop的doOnReceive会对信号进行监听：

```
1      case completion: CompletionEvent =>
2        dagScheduler.handleTaskCompletion(completion)
```

## DAGScheduler.handleTaskCompletion

我们来看下DAGScheduler.handleTaskCompletion部分核心代码：

```
1    ***
2      //根据stageId 得到stage
3      val stage = stageIdToStage(task.stageId)
4      //这里的event就是completion
5      event.reason match {
6      //这里只看成功的流程
7        case Success =>
8        //将这个task 从stage等待处理分区中删去
9          stage.pendingPartitions -= task.partitionId
10         task match {
```

```
11      //若是最后一个Stage的task
12        case rt: ResultTask[_, _] =>
13      //将stage 转为 ResultStage
14          val resultStage = stage.asInstanceOf[ResultStage]
15          resultStage.activeJob match {
16          //获取这Stage的job
17            case Some(job) =>
18              if (!job.finished(rt.outputId)) {
19                updateAccumulators(event)
20                //标记状态
21                job.finished(rt.outputId) = true
22                //计数
23                job.numFinished += 1
24                // 若Job的所有partition都完成了，
25                //移除这个Job
26                if (job.numFinished == job.numPartitions) {
27                  markStageAsFinished(resultStage)
28                  cleanupStateForJobAndIndependentStages(job)
29                  listenerBus.post(
30                    SparkListenerJobEnd(job.jobId, clock.getTimeMillis(), JobSucceeded))
31                }
32                //通知 JobWaiter 有任务成功
33                //但 taskSucceeded 会运行用户自定义的代码
34                //因此可能抛出异常
35                try {
36                  job.listener.taskSucceeded(rt.outputId, event.result)
37                } catch {
38                  case e: Exception =>
39                    // 标记为失败
40                    job.listener.jobFailed(new SparkDriverExecutionException(e))
41                }
42              }
43            case None =>
44              logInfo("Ignoring result from " + rt + " because its job has finished")
45          }
46      //若不是最后一个Stage的Task
47        case smt: ShuffleMapTask =>
48          val shuffleStage = stage.asInstanceOf[ShuffleMapStage]
49          updateAccumulators(event)
50          val status = event.result.asInstanceOf[MapStatus]
51          val execId = status.location.executorId
52          logDebug("ShuffleMapTask finished on " + execId)
53          if (failedEpoch.contains(execId) && smt.epoch <= failedEpoch(execId)) {
54            logInfo(s"Ignoring possibly bogus $smt completion from executor $execId")
55          } else {
56          //将Task的partitionId和status
57          //追加到OutputLoc
58            shuffleStage.addOutputLoc(smt.partitionId, status)
59          }
60
61          if (runningStages.contains(shuffleStage) && shuffleStage.pendingPartitions.isEmpty) {
62            markStageAsFinished(shuffleStage)
63            logInfo("looking for newly runnable stages")
64            logInfo("running: " + runningStages)
65            logInfo("waiting: " + waitingStages)
66            logInfo("failed: " + failedStages)
67
68      //将outputLoc信息注册到mapOutputTracker
69      //上篇博文中有提到：
70      //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mapOutputTra
71      // 所以 DAGScheduler.newOrUsedShuffleStage需要先判断Stage是否已经被计算过
72      ///若计算过，DAGScheduler.newOrUsedShuffleStage则把结果复制到新创建的stage
73      //如果没计算过，DAGScheduler.newOrUsedShuffleStage就向注册mapOutputTracker Stage，为存储元数据占位
74            mapOutputTracker.registerMapOutputs(
```

```
75              shuffleStage.shuffleDep.shuffleId,
76              shuffleStage.outputLocInMapOutputTrackerFormat(),
77              changeEpoch = true)
78
79          clearCacheLocs()
80
81          if (!shuffleStage.isAvailable) {
82            //若Stage不可用（一些任务失败），则从新提交Stage
83            logInfo("Resubmitting " + shuffleStage + " (" + shuffleStage.name +
84              ") because some of its tasks had failed: " +
85              shuffleStage.findMissingPartitions().mkString(", "))
86            submitStage(shuffleStage)
87          } else {
88            //  若该Stage的所有分区都完成了
89            if (shuffleStage.mapStageJobs.nonEmpty) {
90              val stats = mapOutputTracker.getStatistics(shuffleStage.shuffleDep)
91              //将各个Task的标记为Finished
92              for (job <- shuffleStage.mapStageJobs) {
93                markMapStageJobAsFinished(job, stats)
94              }
95            }
96            //提交该Stage的正在等在的Child Stages
97            submitWaitingChildStages(shuffleStage)
98          }
99        }
100       }
101   ***
```

# 处理执行失败的结果

## TaskResultGetter.enqueueFailedTask

下面，我们回归头来看如何处理失败的结果。

```
1   def enqueueFailedTask(taskSetManager: TaskSetManager, tid: Long, taskState: TaskState,
2     serializedData: ByteBuffer) {
3     var reason : TaskFailedReason = UnknownReason
4     try {
5     //通过线程池来处理结果
6       getTaskResultExecutor.execute(new Runnable {
7         override def run(): Unit = Utils.logUncaughtExceptions {
8           val loader = Utils.getContextOrSparkClassLoader
9           try {
10          //若序列化数据，即TaskFailedReason，存在且长度大于0
11          //则反序列化获取它
12            if (serializedData != null && serializedData.limit() > 0) {
13              reason = serializer.get().deserialize[TaskFailedReason](
14                serializedData, loader)
15            }
16          } catch {
17          //若是ClassNotFoundException,
18          //打印log
19            case cnd: ClassNotFoundException =>
20              logError(
21                "Could not deserialize TaskEndReason: ClassNotFound with classloader " + loader)
22              //若其他异常,
23              //不进行操作
24            case ex: Exception =>
25          }
26          //处理失败的任务
27          scheduler.handleFailedTask(taskSetManager, tid, taskState, reason)
28        }
29      })
```

```
30      } catch {
31        case e: RejectedExecutionException if sparkEnv.isStopped =>
32      }
33    }
```

## TaskSchedulerImpl.handleFailedTask

```
1    def handleFailedTask(
2        taskSetManager: TaskSetManager,
3        tid: Long,
4        taskState: TaskState,
5        reason: TaskFailedReason): Unit = synchronized {
6        //处理失败任务
7      taskSetManager.handleFailedTask(tid, taskState, reason)
8      if (!taskSetManager.isZombie && taskState != TaskState.KILLED) {
9      //handleFailedTask会将失败任务放入待运行的队列等待下一次调度
10     //所以这里开始新的一轮调度
11       backend.reviveOffers()
12     }
13   }
```

## TaskSetManager.handleFailedTask

我们来看下handleFailedTask核心代码：

```
1    ***
2        //调用dagScheduler处理失败任务
3        sched.dagScheduler.taskEnded(tasks(index), reason, null, accumUpdates, info)
4
5        if (successful(index)) {
6          logInfo(
7            s"Task ${info.id} in stage ${taskSet.id} (TID $tid) failed, " +
8            "but another instance of the task has already succeeded, " +
9            "so not re-queuing the task to be re-executed.")
10       } else {
11       //将这个任务重新加入到等待队列中
12         addPendingTask(index)
13       }
14
15       if (!isZombie && reason.countTowardsTaskFailures) {
16         taskSetBlacklistHelperOpt.foreach(_.updateBlacklistForFailedTask(
17           info.host, info.executorId, index))
18         assert (null != failureReason)
19       //计数 这个任务的重试次数
20         numFailures(index) += 1
21       //若大于等于最大重试次数，默认为4，
22       //则取消这个任务
23         if (numFailures(index) >= maxTaskFailures) {
24           logError("Task %d in stage %s failed %d times; aborting job".format(
25             index, taskSet.id, maxTaskFailures))
26           abort("Task %d in stage %s failed %d times, most recent failure: %s\nDriver stacktrace:"
27             .format(index, taskSet.id, maxTaskFailures, failureReason), failureException)
28           return
29         }
30       }
31       maybeFinishTaskSet()
32     }
```

## DAGScheduler.handleTaskCompletion

与处理成功结果的过程相同，接下来也会调用DAGScheduler.taskEnded。DAGSchedulerEventProcessLoop的doOnReceive接收CompletionEvent信号，调用dagScheduler.handleTaskCompletion(completion)

我们来看下DAGScheduler.handleTaskCompletion 处理失败任务部分的核心代码：

```scala
    //重新提交任务
    case Resubmitted =>
      logInfo("Resubmitted " + task + ", so marking it as still running")
      //把任务加入的等待队列
      stage.pendingPartitions += task.partitionId


    //获取结果失败
    case FetchFailed(bmAddress, shuffleId, mapId, reduceId, failureMessage) =>
      val failedStage = stageIdToStage(task.stageId)
      val mapStage = shuffleIdToMapStage(shuffleId)
      //若失败的尝试ID 不是 stage尝试ID,
      //则忽略这个失败
      if (failedStage.latestInfo.attemptId != task.stageAttemptId) {
        logInfo(s"Ignoring fetch failure from $task as it's from $failedStage attempt" +
          s" ${task.stageAttemptId} and there is a more recent attempt for that stage " +
          s"(attempt ID ${failedStage.latestInfo.attemptId}) running")
      } else {
        //若失败的Stage还在运行队列,
        //标记这个Stage完成
        if (runningStages.contains(failedStage)) {
          logInfo(s"Marking $failedStage (${failedStage.name}) as failed " +
            s"due to a fetch failure from $mapStage (${mapStage.name})")
          markStageAsFinished(failedStage, Some(failureMessage))
        } else {
          logDebug(s"Received fetch failure from $task, but its from $failedStage which is no " +
            s"longer running")
        }
        //若不允许重试,
        //则停止这个Stage
        if (disallowStageRetryForTest) {
          abortStage(failedStage, "Fetch failure will not retry stage due to testing config",
            None)
        }
        //若达到最大重试次数,
        //则停止这个Stage
        else if (failedStage.failedOnFetchAndShouldAbort(task.stageAttemptId)) {
          abortStage(failedStage, s"$failedStage (${failedStage.name}) " +
            s"has failed the maximum allowable number of " +
            s"times: ${Stage.MAX_CONSECUTIVE_FETCH_FAILURES}. " +
            s"Most recent failure reason: ${failureMessage}", None)
        } else {
          if (failedStages.isEmpty) {
          //若失败的Stage中，没有个task完成了，
          //则重新提交Stage。
          //若果有完成的task的话，我们不能重新提交Stage,
          //因为有些task已经被调度过了。
          //task级别的重新提交是在TaskSetManager.handleFailedTask进行的
            logInfo(s"Resubmitting $mapStage (${mapStage.name}) and " +
              s"$failedStage (${failedStage.name}) due to fetch failure")
            messageScheduler.schedule(new Runnable {
              override def run(): Unit = eventProcessLoop.post(ResubmitFailedStages)
            }, DAGScheduler.RESUBMIT_TIMEOUT, TimeUnit.MILLISECONDS)
          }
          failedStages += failedStage
          failedStages += mapStage
        }
        // 移除OutputLoc中的数据
        // 取消注册mapOutputTracker
        if (mapId != -1) {
          mapStage.removeOutputLoc(mapId, bmAddress)
          mapOutputTracker.unregisterMapOutput(shuffleId, mapId, bmAddress)
        }
```

```
63
64          //当有executor上发生多次获取结果失败,
65          //则标记这个executor丢失
66          if (bmAddress != null) {
67            handleExecutorLost(bmAddress.executorId, filesLost = true, Some(task.epoch))
68          }
69        }
70
71      //拒绝处理
72      case commitDenied: TaskCommitDenied =>
73        // 不做任何事,
74        //让 TaskScheduler 来决定如何处理
75
76      //异常
77      case exceptionFailure: ExceptionFailure =>
78        // 更新accumulator
79        updateAccumulators(event)
80
81      //task结果丢失
82      case TaskResultLost =>
83      // 不做任何事,
84      // 让 TaskScheduler 处理这些错误和重新提交任务
85
86    // executor 丢失
87    // 任务被杀死
88    // 未知错误
89     case _: ExecutorLostFailure | TaskKilled | UnknownReason =>
90        // 不做任何事,
91        // 若这task不断的错误,
92        // TaskScheduler 会停止 job
```