

## 深入理解Spark 2.1 Core（六）：资源调度的原理与源码分析

■ 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/54098376>

<http://blog.csdn.net/u011239443/article/details/54098376>

在上篇博文《深入理解Spark 2.1 Core（五）：Standalone模式运行的实现与源码分析》中，我们讲到了如何启动Master和Worker，还讲到了如何回收资源。但是，我们没有将AppClient是如何启动的，其实它们的启动也涉及到了资源是如何调度的。这篇博文，我们就来讲一下AppClient的启动和逻辑与物理上的资源调度。

### 启动AppClient

调用栈如下：

- StandaloneSchedulerBackend.start
  - StandaloneAppClient.start
    - StandaloneAppClient.ClientEndpoint.onStart
      - StandaloneAppClient.registerWithMaster
        - StandaloneAppClient.tryRegisterAllMasters
- Master.receive
  - Master.createApplication
  - Master.registerApplication
- StandaloneAppClient.ClientEndpoint.receive

### StandaloneSchedulerBackend.start

在Standalone模式下，SparkContext中的backend是StandaloneSchedulerBackend。在StandaloneSchedulerBackend.start中可以看到：

```
1  ***
2  val appUIAddress = sc.ui.map(_.appUIAddress).getOrElse("")
3  val coresPerExecutor = conf.getOption("spark.executor.cores").map(_.toInt)
4
5  val initialExecutorLimit =
6    if (Utils.isDynamicAllocationEnabled(conf)) {
7      Some(0)
8    } else {
9      None
10   }
11  val appDesc = new ApplicationDescription(sc.appName, maxCores, sc.executorMemory, command,
12    appUIAddress, sc.eventLogDir, sc.eventLogCodec, coresPerExecutor, initialExecutorLimit)
13  //创建AppClient
14  client = new StandaloneAppClient(sc.env.rpcEnv, masters, appDesc, this, conf)
15  //启动AppClient
16  client.start()
17  ***
```

### StandaloneAppClient.start

```

1  def start() {
2      //生成了ClientEndpoint, 于是调用其onStart
3      endpoint.set(rpcEnv.setupEndpoint("AppClient", new ClientEndpoint(rpcEnv)))
4  }

```

## StandaloneAppClient.ClientEndpoint.onStart

调用registerWithMaster

```

1  override def onStart(): Unit = {
2      try {
3          registerWithMaster(1)
4      } catch {
5          case e: Exception =>
6              logWarning("Failed to connect to master", e)
7              markDisconnected()
8              stop()
9      }
10 }

```

## StandaloneAppClient.registerWithMaster

```

1
2  private def registerWithMaster(nthRetry: Int) {
3      //向所有的Master注册当前App
4      //一旦成功连接的一个master, 其他将被取消
5      registerMasterFutures.set(tryRegisterAllMasters())
6      registrationRetryTimer.set(registrationRetryThread.schedule(new Runnable {
7          override def run(): Unit = {
8              if (registered.get) {
9
10                 registerMasterFutures.get.foreach(_.cancel(true))
11                 registerMasterThreadPool.shutdownNow()
12             }
13             //若达到最大尝试次数, 则标志死亡, 默认为3
14             else if (nthRetry >= REGISTRATION_RETRIES) {
15                 markDead("All masters are unresponsive! Giving up.")
16             } else {
17                 registerMasterFutures.get.foreach(_.cancel(true))
18                 registerWithMaster(nthRetry + 1)
19             }
20         }
21     }, REGISTRATION_TIMEOUT_SECONDS, TimeUnit.SECONDS))
22 }

```

## StandaloneAppClient.tryRegisterAllMasters

给Master发送RegisterApplication信号:

```

1  private def tryRegisterAllMasters(): Array[JFuture[_]] = {
2      for (masterAddress <- masterRpcAddresses) yield {
3          registerMasterThreadPool.submit(new Runnable {
4              override def run(): Unit = try {
5                  if (registered.get) {
6                      return
7                  }
8                  logInfo("Connecting to master " + masterAddress.toSparkURL + "...")
9                  val masterRef = rpcEnv.setupEndpointRef(masterAddress, Master.ENDPOINT_NAME)
10                 masterRef.send(RegisterApplication(appDescription, self))
11             } catch {
12                 case ie: InterruptedException => // Cancelled
13                 case NonFatal(e) => logWarning(s"Failed to connect to master $masterAddress", e)

```

```
14         }
15     })
16 }
17 }
```

## Master.receive

Master.receive接收并处理RegisterApplication信号

```
1  case RegisterApplication(description, driver) =>
2      // 若之前注册过
3      if (state == RecoveryState.STANDBY) {
4          // 忽略
5      } else {
6          logInfo("Registering app " + description.name)
7          //创建app
8          val app = createApplication(description, driver)
9          //注册app
10         registerApplication(app)
11         logInfo("Registered app " + description.name + " with ID " + app.id)
12         //持久化
13         persistenceEngine.addApplication(app)
14         //回复RegisteredApplication信号
15         driver.send(RegisteredApplication(app.id, self))
16         //资源调度
17         schedule()
18     }
```

让我们深入来看下Master是如何注册app的。

## Master.createApplication

先创建app:

```
1  private def createApplication(desc: ApplicationDescription, driver: RpcEndpointRef):
2      ApplicationInfo = {
3      val now = System.currentTimeMillis()
4      val date = new Date(now)
5      //根据日期生成appId
6      val appId = newApplicationId(date)
7      //传入 时间, appId, 描述信息, 日期, driver, 默认核数,
8      //生成app信息
9      new ApplicationInfo(now, appId, desc, date, driver, defaultCores)
10 }
```

## Master.registerApplication

再注册app:

```
1  private def registerApplication(app: ApplicationInfo): Unit = {
2      //若已有这个app地址,
3      //则返回
4      val appAddress = app.driver.address
5      if (addressToApp.contains(appAddress)) {
6          logInfo("Attempted to re-register application at same address: " + appAddress)
7          return
8      }
9
10     //向 applicationMetricsSystem 注册appSource
11     applicationMetricsSystem.registerSource(app.appSource)
12     //将app加入到 集合
13     //HashSet[ApplicationInfo]
```

```

14     apps += app
15     //更新 id到App
16     //HashMap[String, ApplicationInfo]
17     idToApp(app.id) = app
18     //更新 endpoint到App
19     // HashMap[RpcEndpointRef, ApplicationInfo]
20     endpointToApp(app.driver) = app
21     //更新 address到App
22     // HashMap[RpcAddress, ApplicationInfo]
23     addressToApp(appAddress) = app
24     // 加入到等待数组中
25     //ArrayBuffer[ApplicationInfo]
26     waitingApps += app
27     if (reverseProxy) {
28         webUi.addProxyTargets(app.id, app.desc.appUiUrl)
29     }
30 }

```

## StandaloneAppClient.ClientEndpoint.receive

```

1     case RegisteredApplication(appId_, masterRef) =>
2         //这里的代码有两个缺陷:
3         //1. 一个Master可能接收到多个注册请求,
4         // 并且回复多个RegisteredApplication信号,
5         //这会导致网络不稳定。
6         //2.若master正在变化,
7         //则会接收到多个RegisteredApplication信号
8         //设置appId
9         appId.set(appId_)
10        //编辑已经注册
11        registered.set(true)
12        //创建master信息
13        master = Some(masterRef)
14        //绑定监听
15        listener.connected(appId.get)

```

## 逻辑资源调度

我们可以看到在上一章，Master.receive接收并处理RegisterApplication信号时的最后一行代码：

```

1         //资源调度
2         schedule()

```

下面，我们就来讲讲资源调度。

调用栈如下：

- Master.schedule
  - Master.startExecutorsOnWorkers
    - Master.scheduleExecutorsOnWorkers
    - Master.allocateWorkerResourceToExecutors

## Master.schedule

该方法主要来在等待的app之间调度资源。每次有新的app加入或者可用资源改变的时候，这个方法都会被调用：

```

1     private def schedule(): Unit = {
2         if (state != RecoveryState.ALIVE) {
3             return

```

```

4      }
5      // 得到活的Worker,
6      // 并打乱它们
7      val shuffledAliveWorkers = Random.shuffle(workers.toSeq.filter(_.state == WorkerState.ALIVE))
8      // worker数量
9      val numWorkersAlive = shuffledAliveWorkers.size
10     var curPos = 0
11     //为driver分配资源,
12     //该调度策略为FIFO的策略,
13     //先来的driver会先满足其资源所需的条件
14     for (driver <- waitingDrivers.toList) {
15         var launched = false
16         var numWorkersVisited = 0
17         while (numWorkersVisited < numWorkersAlive && !launched) {
18             val worker = shuffledAliveWorkers(curPos)
19             numWorkersVisited += 1
20             if (worker.memoryFree >= driver.desc.mem && worker.coresFree >= driver.desc.cores) {
21                 launchDriver(worker, driver)
22                 waitingDrivers -= driver
23                 launched = true
24             }
25             curPos = (curPos + 1) % numWorkersAlive
26         }
27     }
28     //启动worker上的executor
29     startExecutorsOnWorkers()
30 }

```

## Master.startExecutorsOnWorkers

接下来我们来看下executor的启动:

```

1 private def startExecutorsOnWorkers(): Unit = {
2     // 这里还是使用的FIFO的调度方式
3     for (app <- waitingApps if app.coresLeft > 0) {
4         val coresPerExecutor: Option[Int] = app.desc.coresPerExecutor
5         // 过滤掉资源不够启动executor的worker
6         // 并按资源逆序排序
7         val usableWorkers = workers.toArray.filter(_.state == WorkerState.ALIVE)
8             .filter(worker => worker.memoryFree >= app.desc.memoryPerExecutorMB &&
9                 worker.coresFree >= coresPerExecutor.getOrElse(1))
10        .sortBy(_.coresFree).reverse
11        //调度worker上的executor,
12        //确定在每个worker上给这个app分配多少核
13        val assignedCores = scheduleExecutorsOnWorkers(app, usableWorkers, spreadOutApps)
14
15        //分配
16        for (pos <- 0 until usableWorkers.length if assignedCores(pos) > 0) {
17            allocateWorkerResourceToExecutors(
18                app, assignedCores(pos), coresPerExecutor, usableWorkers(pos))
19        }
20    }
21 }

```

## Master.scheduleExecutorsOnWorkers

接下来我们就来讲讲核心的worker上的executor资源调度。在将现在的Spark代码之前, 我们看看在Spark1.4之前, 这部分逻辑是如何实现的:

```

1 ***
2     val numUsable = usableWorkers.length
3     // 用来记录每个worker已经分配的核数
4     val assigned = new Array[Int](numUsable)

```

```

5      var toAssign = math.min(app.coresLeft, usableWorkers.map(_._coresFree).sum)
6      var pos = 0
7      while (toAssign > 0) {
8          //遍历worker,
9          //若当前worker还存在资源,
10         //则分配掉1个核。
11         //直到workers的资源全都被分配掉,
12         //或者是app所需要的资源被满足。
13         if (usableWorkers(pos)._coresFree - assigned(pos) > 0) {
14             toAssign -= 1
15             assigned(pos) += 1
16         }
17         pos = (pos + 1) % numUsable
18     }
19     ***

```

在Spark1.4的时候，这段代码被修改了。我们来想一下，以上代码有什么问题？

**问题就在于，core是一个一个的被分配的。**设想，一个集群中有4 worker，每个worker有16个core。用户想启动3个executor，且每个executor拥有16个core。于是，他会这样配置参数：

```

1 spark.cores.max = 48
2 spark.executor.cores = 16

```

显然，我们集群的资源是能满足用户的需求的。但如果一次只能分配一个core，那最终的结果是每个worker上都分配了12个core。由于 $12 < 16$ ，所有没有一个executor能够启动。

下面，我们回过头来看现在的源码中是如何实现这部分逻辑的：

```

1 private def scheduleExecutorsOnWorkers(
2     app: ApplicationInfo,
3     usableWorkers: Array[WorkerInfo],
4     spreadOutApps: Boolean): Array[Int] = {
5     val coresPerExecutor = app.desc.coresPerExecutor
6     val minCoresPerExecutor = coresPerExecutor.getOrElse(1)
7     val oneExecutorPerWorker = coresPerExecutor.isEmpty
8     val memoryPerExecutor = app.desc.memoryPerExecutorMB
9     val numUsable = usableWorkers.length
10    // 用来记录每个worker已经分配的核数
11    val assignedCores = new Array[Int](numUsable)
12    // 用来记录每个worker已经分配的executor数
13    val assignedExecutors = new Array[Int](numUsable)
14    // 剩余总共资源
15    var coresToAssign = math.min(app.coresLeft, usableWorkers.map(_._coresFree).sum)
16
17    //判断是否能启动Executor
18    def canLaunchExecutor(pos: Int): Boolean = {
19        //先省略
20    }
21
22
23    //过滤去能启动executor的Worker
24    var freeWorkers = (0 until numUsable).filter(canLaunchExecutor)
25    //调度资源,
26    //直到worker上的executor被分配完
27    while (freeWorkers.nonEmpty) {
28        freeWorkers.foreach { pos =>
29            var keepScheduling = true
30            while (keepScheduling && canLaunchExecutor(pos)) {
31                // minCoresPerExecutor 是用户设置的 spark.executor.cores
32                coresToAssign -= minCoresPerExecutor
33                assignedCores(pos) += minCoresPerExecutor
34            }
35        }
36    }
37
38    //返回每个worker已经分配的核数
39    assignedCores
40 }

```

```

35 // 若用户没有设置 spark.executor.cores
36 // 则oneExecutorPerWorker就为True
37 // 也就是说, assignedCores中的core都被一个executor使用
38 // 若用户设置了spark.executor.cores,
39 // 则该Worker的assignedExecutors会加1
40 if (oneExecutorPerWorker) {
41     assignedExecutors(pos) = 1
42 } else {
43     assignedExecutors(pos) += 1
44 }
45
46
47 //资源分配算法有两种:
48 // 1. 尽量打散, 将一个app尽可能的分配到不同的节点上,
49 // 这有利于充分利用集群的资源,
50 // 在这种情况下, spreadOutApps设为True,
51 // 于是, 该worker分配好了一个executor之后就退出循环
52 // 轮询到下一个worker
53 // 2. 尽量集中, 将一个app尽可能的分配到同一个的节点上,
54 // 这适合cpu密集型而内存占用比较少的app
55 // 在这种情况下, spreadOutApps设为False,
56 // 于是, 继续下一轮的循环
57 // 在该Worker上分配executor
58 if (spreadOutApps) {
59     keepScheduling = false
60 }
61 }
62 }
63 freeWorkers = freeWorkers.filter(canLaunchExecutor)
64 }
65 assignedCores
66 }

```

接下来看下该函数的内部函数canLaunchExecutor:

```

1 def canLaunchExecutor(pos: Int): Boolean = {
2 // 条件1: 若集群剩余core >= spark.executor.cores
3     val keepScheduling = coresToAssign >= minCoresPerExecutor
4 // 条件2: 若该Worker上的剩余core >= spark.executor.cores
5     val enoughCores = usableWorkers(pos).coresFree - assignedCores(pos) >= minCoresPerExecutor
6
7
8 // 条件3: 若设置了spark.executor.cores
9 // 或者 该Worker还未分配executor
10    val launchingNewExecutor = !oneExecutorPerWorker || assignedExecutors(pos) == 0
11    if (launchingNewExecutor) {
12        val assignedMemory = assignedExecutors(pos) * memoryPerExecutor
13        // 条件4: 若该Worker上的剩余内存 >= spark.executor.memory
14        val enoughMemory = usableWorkers(pos).memoryFree - assignedMemory >= memoryPerExecutor
15        // 条件5: 若分配了该executor后,
16        // 总共分配的core数量 <= spark.cores.max
17        val underLimit = assignedExecutors.sum + app.executors.size < app.executorLimit
18        //若满足 条件3,
19        //且满足 条件1, 条件2, 条件4, 条件5
20        //则返回True
21        keepScheduling && enoughCores && enoughMemory && underLimit
22    } else {
23        //若不满足 条件3,
24        //即一个worker只有一个executor
25        //且满足 条件1, 条件2
26        //也返回True。
27        // 返回后, 不会增加 assignedExecutors
28        keepScheduling && enoughCores
29    }

```

```

30     }
    }

```

通过以上源码，我们可以清楚看到，Spark1.4以后新的逻辑实现其实就是将分配单位从原来的一个core，变为了一个executor（即spark.executor.cores）。而若一个worker上只有一个executor（即没有设置spark.executor.cores），那么就按照原来的逻辑实现。

值得我注意的是：

```

1      //直到worker上的executor被分配完
2      while (freeWorkers.nonEmpty)

```

一个app会尽可能的使用掉集群的所有资源，所以设置spark.cores.max参数是非常有必要的！

## Master.allocateWorkerResourceToExecutors

现在我们回到上述提到的Master.startExecutorsOnWorkers中，深入allocateWorkerResourceToExecutors：

```

1      private def allocateWorkerResourceToExecutors(
2          app: ApplicationInfo,
3          assignedCores: Int,
4          coresPerExecutor: Option[Int],
5          worker: WorkerInfo): Unit = {
6          // 该work上的executor数量
7          // 若没设置 spark.executor.cores
8          // 则为1
9          val numExecutors = coresPerExecutor.map { assignedCores / _ }.getOrElse(1)
10         // 分配给一个executor的core数量
11         // 若没设置 spark.executor.cores
12         // 则为该worker上所分配的所有core是数量
13         val coresToAssign = coresPerExecutor.getOrElse(assignedCores)
14         for (i <- 1 to numExecutors) {
15             //创建该executor信息
16             //并把它加入到app信息中
17             //并返回executor信息
18             val exec = app.addExecutor(worker, coresToAssign)
19             //启动
20             launchExecutor(worker, exec)
21             app.state = ApplicationState.RUNNING
22         }
23     }

```

要注意的是

```

1      app.state = ApplicationState.RUNNING

```

这句代码并不是将该app从waitingApp队列中去除。若在该次资源调度中该app并没有启动足够的executor，等到集群资源变化时，会再次资源调度，在waitingApp中遍历到该app，其coresLeft > 0。

```

1      for (app <- waitingApps if app.coresLeft > 0)

```

我们这里做一个实验：

- 我们的实验集群是4\*8核的集群：



State	Cores
ALIVE	8 (0 Used)
ALIVE	8 (0 Used)
ALIVE	8 (0 Used)
ALIVE	8 (0 Used)

- 第1个app，我们申请4个executor，该executor为4个core：

```
1 spark-shell --master spark://cdh03:7077 --total-executor-cores 4 --executor-cores 4
```

可以看到集群资源：

State	Cores
ALIVE	8 (0 Used)
ALIVE	8 (0 Used)
ALIVE	8 (4 Used)
ALIVE	8 (0 Used)

app1的executor：

Executor Summary

ExecutorID	Worker
0	<a href="#">worker-20170102151129</a>

- 第2个app，我们申请4个executor，该executor为6个core：

```
1 spark-shell --master spark://cdh03:7077 --total-executor-cores 24 --executor-cores 6
```

可以看到集群资源：

State	Cores
ALIVE	8 (6 Used)
ALIVE	8 (6 Used)
ALIVE	8 (4 Used)
ALIVE	8 (6 Used)

app2的executor：

Executor Summary

ExecutorID	Worker
2	<a href="#">worker-20170102151232</a>
1	<a href="#">worker-20170102151120</a>
0	<a href="#">worker-20170102151055</a>

我们可以看到，Spark只为app2分配了3个executor。

- 当我们把app1退出

会发现集群资源状态：

State	Cores
ALIVE	8 (6 Used)
ALIVE	8 (6 Used)
ALIVE	8 (6 Used)
ALIVE	8 (6 Used)

app2的executor：

Executor Summary

ExecutorID	Worker
2	<a href="#">worker-20170102151232</a>
1	<a href="#">worker-20170102151120</a>
3	<a href="#">worker-20170102151129</a>
0	<a href="#">worker-20170102151055</a>

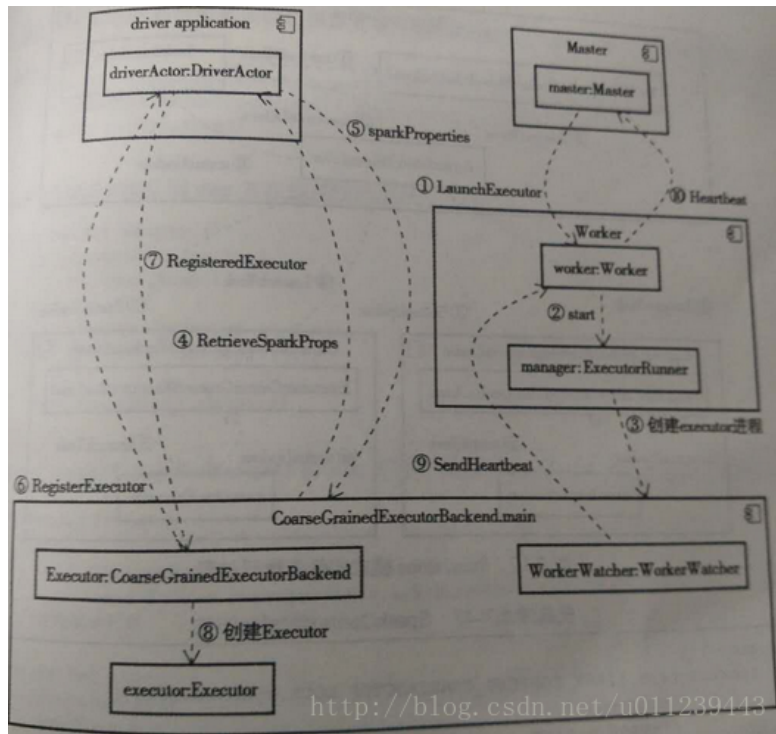
会发现新增加了一个“ worker-20170102151129”的executor。

其实，只要集群中的app没结束，它们都会在waitingApps中，当该app结束时，才会将这个app从waitingApps中移除

```
1  def removeApplication(app: ApplicationInfo, state: ApplicationState.Value) {
2  ***
3      waitingApps -= app
4  ***
5  }
```

物理资源调度与启动Executor

接下来，我们就来讲逻辑上资源调度完后，该如何物理上资源调度，即启动Executor。



调用栈如下：

- Master.launchExecutor
- Worker.receive
  - ExecutorRunner.start
    - ExecutorRunner.fetchAndRunExecutor
- CoarseGrainedExecutorBackend.main
  - CoarseGrainedExecutorBackend.run
    - CoarseGrainedExecutorBackend.onStart
- CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply
- CoarseGrainedExecutorBackend.receive

## Master.launchExecutor

```

1 private def launchExecutor(worker: WorkerInfo, exec: ExecutorDesc): Unit = {
2   logInfo("Launching executor " + exec.fullId + " on worker " + worker.id)
3   //在worker信息中加入executor信息
4   worker.addExecutor(exec)
5   //给worker发送LaunchExecutor信号
6   worker.endpoint.send(LaunchExecutor(masterUrl,
7     exec.application.id, exec.id, exec.application.desc, exec.cores, exec.memory))
8   //给driver发送ExecutorAdded信号
9   exec.application.driver.send(
10     ExecutorAdded(exec.id, worker.id, worker.hostPort, exec.cores, exec.memory))
11 }

```

## Worker.receive

worker接收到LaunchExecutor信号后处理：

```

1  case LaunchExecutor(masterUrl, appId, execId, appDesc, cores_, memory_) =>
2  if (masterUrl != activeMasterUrl) {
3      logWarning("Invalid Master (" + masterUrl + ") attempted to launch executor.")
4  } else {
5      try {
6          logInfo("Asked to launch executor %s/%d for %s".format(appId, execId, appDesc.name))
7
8          // 创建executor的工作目录
9          // shuffle持久化结果会存在这个目录下
10         // 节点应每块磁盘大小尽可能相同
11         // 并在配置中在每块磁盘上都设置SPARK_WORKER_DIR,
12         // 以增加IO性能
13         val executorDir = new File(workDir, appId + "/" + execId)
14         if (!executorDir.mkdirs()) {
15             throw new IOException("Failed to create directory " + executorDir)
16         }
17
18         // 为app创建本地dir
19         // app完成后, 此目录会被删除
20         val appLocalDirs = appDirectories.getOrElse(appId,
21             Utils.getOrCreateLocalRootDirs(conf).map { dir =>
22                 val appDir = Utils.createDirectory(dir, namePrefix = "executor")
23                 Utils.chmod700(appDir)
24                 appDir.getAbsolutePath()
25             }.toSeq)
26         appDirectories(appId) = appLocalDirs
27         //创建ExecutorRunner
28         val manager = new ExecutorRunner(
29             appId,
30             execId,
31             appDesc.copy(command = Worker.maybeUpdateSSLSettings(appDesc.command, conf)),
32             cores_,
33             memory_,
34             self,
35             workerId,
36             host,
37             webUi.boundPort,
38             publicAddress,
39             sparkHome,
40             executorDir,
41             workerUri,
42             conf,
43             appLocalDirs, ExecutorState.RUNNING)
44         executors(appId + "/" + execId) = manager
45         //启动ExecutorRunner
46         manager.start()
47         coresUsed += cores_
48         memoryUsed += memory_
49         // 向Master发送ExecutorStateChanged信号
50         sendToMaster(ExecutorStateChanged(appId, execId, manager.state, None, None))
51     } catch {
52         case e: Exception =>
53             logError(s"Failed to launch executor $appId/$execId for ${appDesc.name}.", e)
54             if (executors.contains(appId + "/" + execId)) {
55                 executors(appId + "/" + execId).kill()
56                 executors -= appId + "/" + execId
57             }
58             sendToMaster(ExecutorStateChanged(appId, execId, ExecutorState.FAILED,
59                 Some(e.toString), None))
60     }
61 }

```

## ExecutorRunner.start

接下来我们深入看下ExecutorRunner

```
1 private[worker] def start() {
2   //创建worker线程
3   workerThread = new Thread("ExecutorRunner for " + fullId) {
4     override def run() { fetchAndRunExecutor() }
5   }
6   //启动worker线程
7   workerThread.start()
8   // 创建Shutdownhook线程
9   // 用于worker关闭时, 杀掉executor
10  shutdownHook = ShutdownHookManager.addShutdownHook { () =>
11    if (state == ExecutorState.RUNNING) {
12      state = ExecutorState.FAILED
13    }
14    killProcess(Some("Worker shutting down")) }
15 }
```

## ExecutorRunner.fetchAndRunExecutor

workerThread执行主要是调用fetchAndRunExecutor，下面我们来看下这个方法：

```
1 private def fetchAndRunExecutor() {
2   try {
3     // 创建进程builder
4     val builder = CommandUtils.buildProcessBuilder(appDesc.command, new SecurityManager(conf),
5       memory, sparkHome.getAbsolutePath, substituteVariables)
6     val command = builder.command()
7     val formattedCommand = command.asScala.mkString("\n", "\n\n", "\n")
8     logInfo(s"Launch command: $formattedCommand")
9     //创建进程builder执行目录
10    builder.directory(executorDir)
11    //为进程builder设置环境变量
12    builder.environment.put("SPARK_EXECUTOR_DIRS", appLocalDirs.mkString(File.pathSeparator))
13    builder.environment.put("SPARK_LAUNCH_WITH_SCALA", "0")
14
15    val baseUrl =
16      if (conf.getBoolean("spark.ui.reverseProxy", false)) {
17        s"/proxy/$workerId/logPage/?appId=$appId&executorId=$execId&logType="
18      } else {
19        s"http://$publicAddress:$webUiPort/logPage/?appId=$appId&executorId=$execId&logType="
20      }
21    builder.environment.put("SPARK_LOG_URL_STDERR", s"${baseUrl}stderr")
22    builder.environment.put("SPARK_LOG_URL_STDOUT", s"${baseUrl}stdout")
23
24    //启动进程builder, 创建进程
25    process = builder.start()
26    val header = "Spark Executor Command: %s\n%s\n\n".format(
27      formattedCommand, "=" * 40)
28
29    // 重定向它的stdout和stderr到文件中
30    val stdout = new File(executorDir, "stdout")
31    stdoutAppender = FileAppender(process.getInputStream, stdout, conf)
32
33    val stderr = new File(executorDir, "stderr")
34    Files.write(header, stderr, StandardCharsets.UTF_8)
35    stderrAppender = FileAppender(process.getErrorStream, stderr, conf)
36
37    // 等待进程退出。
38    // 当driver通知该进程退出
39    // executor会退出并返回0或者非0的exitCode
```

```

40     val exitCode = process.waitFor()
41     state = ExecutorState.EXITED
42     val message = "Command exited with code " + exitCode
43     // 给Worker发送ExecutorStateChanged信号
44     worker.send(ExecutorStateChanged(appId, execId, state, Some(message), Some(exitCode)))
45 } catch {
46     case interrupted: InterruptedException =>
47         logInfo("Runner thread for executor " + fullId + " interrupted")
48         state = ExecutorState.KILLED
49         killProcess(None)
50     case e: Exception =>
51         logError("Error running executor", e)
52         state = ExecutorState.FAILED
53         killProcess(Some(e.toString))
54 }
55 }
56 }

```

## CoarseGrainedExecutorBackend.main

builder start的是CoarseGrainedExecutorBackend实例进程，我们看下它的主函数：

```

1  def main(args: Array[String]) {
2      var driverUrl: String = null
3      var executorId: String = null
4      var hostname: String = null
5      var cores: Int = 0
6      var appId: String = null
7      var workerUrl: Option[String] = None
8      val userClassPath = new mutable.ListBuffer[URL]()
9      // 设置参数
10     var argv = args.toList
11     while (!argv.isEmpty) {
12         argv match {
13             case ("--driver-url") :: value :: tail =>
14                 driverUrl = value
15                 argv = tail
16             case ("--executor-id") :: value :: tail =>
17                 executorId = value
18                 argv = tail
19             case ("--hostname") :: value :: tail =>
20                 hostname = value
21                 argv = tail
22             case ("--cores") :: value :: tail =>
23                 cores = value.toInt
24                 argv = tail
25             case ("--app-id") :: value :: tail =>
26                 appId = value
27                 argv = tail
28             case ("--worker-url") :: value :: tail =>
29                 workerUrl = Some(value)
30                 argv = tail
31             case ("--user-class-path") :: value :: tail =>
32                 userClassPath += new URL(value)
33                 argv = tail
34             case Nil =>
35             case tail =>
36                 System.err.println(s"Unrecognized options: ${tail.mkString(" ")}")
37                 printUsageAndExit()
38         }
39     }
40
41     if (driverUrl == null || executorId == null || hostname == null || cores <= 0 ||
42         appId == null) {

```

```

43     printUsageAndExit()
44 }
45 //调用run方法
46 run(driverUrl, executorId, hostname, cores, appId, workerUrl, userClassPath)
47 System.exit(0)
48 }

```

## CoarseGrainedExecutorBackend.run

```

1  private def run(
2      driverUrl: String,
3      executorId: String,
4      hostname: String,
5      cores: Int,
6      appId: String,
7      workerUrl: Option[String],
8      userClassPath: Seq[URL]) {
9
10     Utils.initDaemon(log)
11
12     SparkHadoopUtil.get.runAsSparkUser { () =>
13         Utils.checkHost(hostname)
14
15         val executorConf = new SparkConf
16         val port = executorConf.getInt("spark.executor.port", 0)
17         val fetcher = RpcEnv.create(
18             "driverPropsFetcher",
19             hostname,
20             port,
21             executorConf,
22             new SecurityManager(executorConf),
23             clientMode = true)
24         val driver = fetcher.setupEndpointRefByURI(driverUrl)
25         // 给driver发送RetrieveSparkAppConfig信号,
26         // 并根据返回的信息创建属性
27         val cfg = driver.askWithRetry[SparkAppConfig](RetrieveSparkAppConfig)
28         val props = cfg.sparkProperties ++ Seq[(String, String)](("spark.app.id", appId))
29         fetcher.shutdown()
30
31         // 根据这些属性来创建SparkEnv
32         val driverConf = new SparkConf()
33         for ((key, value) <- props) {
34             if (SparkConf.isExecutorStartupConf(key)) {
35                 driverConf.setIfMissing(key, value)
36             } else {
37                 driverConf.set(key, value)
38             }
39         }
40         if (driverConf.contains("spark.yarn.credentials.file")) {
41             logInfo("Will periodically update credentials from: " +
42                 driverConf.get("spark.yarn.credentials.file"))
43             SparkHadoopUtil.get.startCredentialUpdater(driverConf)
44         }
45
46         val env = SparkEnv.createExecutorEnv(
47             driverConf, executorId, hostname, port, cores, cfg.ioEncryptionKey, isLocal = false)
48
49         // 创建CoarseGrainedExecutorBackend Endpoint
50         env.rpcEnv.setupEndpoint("Executor", new CoarseGrainedExecutorBackend(
51             env.rpcEnv, driverUrl, executorId, hostname, cores, userClassPath, env))
52         // 创建WorkerWatcher Endpoint
53         // 用来给worker发送心跳,
54         // 告诉worker 这个进程还活着
55         workerUrl.foreach { url =>

```

```
56     env.rpcEnv.setupEndpoint("WorkerWatcher", new WorkerWatcher(env.rpcEnv, url))
57   }
58   env.rpcEnv.awaitTermination()
59   SparkHadoopUtil.get.stopCredentialUpdater()
60 }
61 }
```

## CoarseGrainedExecutorBackend.onStart

new CoarseGrainedExecutorBackend 会调用CoarseGrainedExecutorBackend.onStart:

```
1  override def onStart() {
2    logInfo("Connecting to driver: " + driverUrl)
3    rpcEnv.asyncSetupEndpointRefByURI(driverUrl).flatMap { ref =>
4      //向driver端发送RegisterExecutor信号
5      driver = Some(ref)
6      ref.ask[Boolean](RegisterExecutor(executorId, self, hostname, cores, extractLogUrls))
7    }(ThreadUtils.sameThread).onComplete {
8      case Success(msg) =>
9      case Failure(e) =>
10       exitExecutor(1, s"Cannot register with driver: $driverUrl", e, notifyDriver = false)
11    }(ThreadUtils.sameThread)
12  }
```

## CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply

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



makeOffers()所做的逻辑，在《[深入理解Spark 2.1 Core （三）：任务调度器的原理与源码分析](#)》里已经讲解过。主要是调度任务，并向executor发送任务。

## CoarseGrainedExecutorBackend.receive

CoarseGrainedExecutorBackend接收到来自driver的RegisteredExecutor信号后：

```
1      case RegisteredExecutor =>
2          logInfo("Successfully registered with driver")
3          try {
4              //创建executor
5              executor = new Executor(executorId, hostname, env, userClassPath, isLocal = false)
6          } catch {
7              case NonFatal(e) =>
8                  exitExecutor(1, "Unable to create executor due to " + e.getMessage, e)
9          }
```

至此，Executor就成功的启动了！