

深入理解Spark 2.1 Core（五）：Standalone模式运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/54093333>

概述

前几篇博文都在介绍Spark的调度，这篇博文我们从更加宏观的调度看Spark，讲讲Spark的部署模式。Spark部署模式分以下几种：

- local 模式
- local-cluster 模式
- Standalone 模式
- YARN 模式
- Mesos 模式

我们先来简单介绍下YARN模式，然后深入讲解Standalone模式。

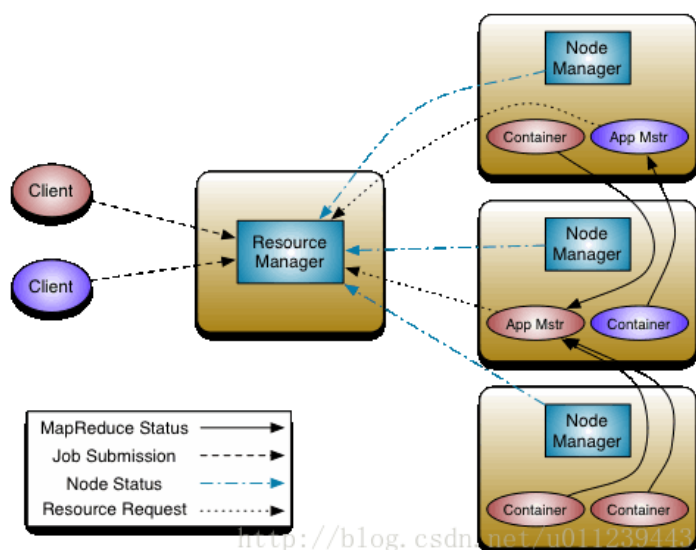
YARN 模式介绍

YARN介绍

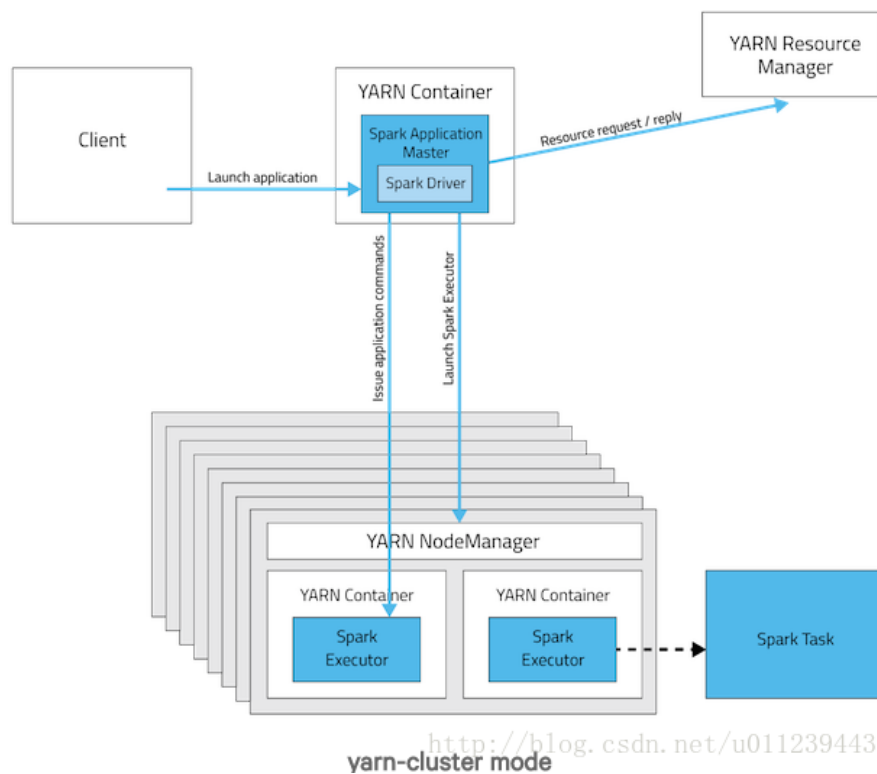
YARN是一个资源管理、任务调度的框架，主要包含三大模块：ResourceManager（RM）、NodeManager（NM）、ApplicationMaster（AM）。

其中，ResourceManager负责所有资源的监控、分配和管理；ApplicationMaster负责每一个具体应用程序的调度和协调；NodeManager负责每一个节点的维护。

对于所有的applications，RM拥有绝对的控制权和对资源的分配权。而每个AM则会和RM协商资源，同时和NodeManager通信来执行和监控task。几个模块之间的关系如图所示。



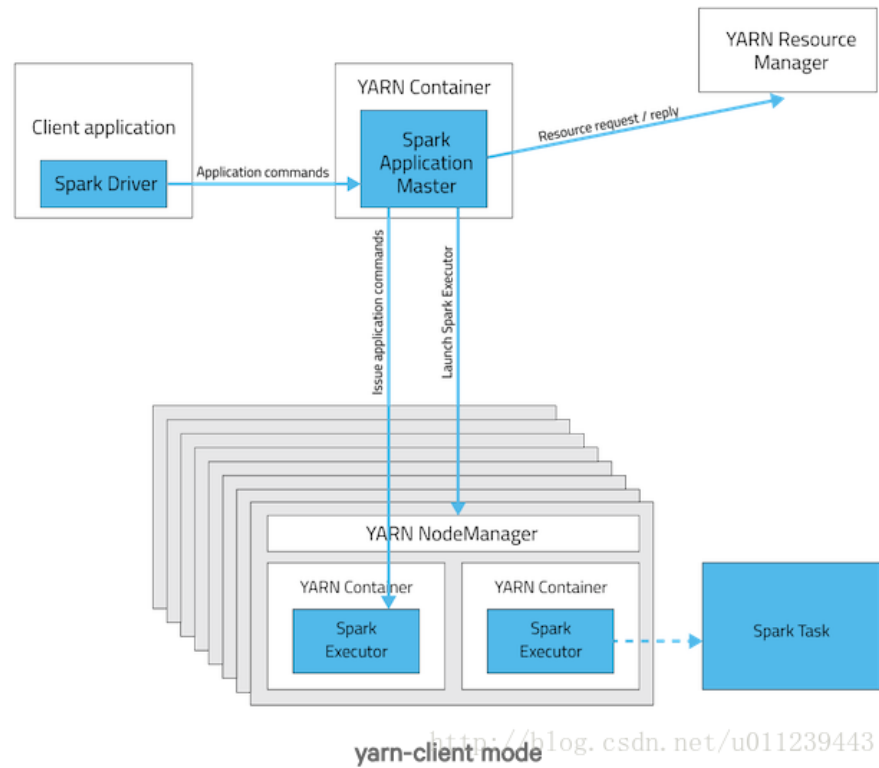
Yarn Cluster 模式



Spark的Yarn Cluster 模式流程如下：

- 本地用YARN Client 提交App 到 Yarn Resource Manager
- Yarn Resource Manager 选个 YARN Node Manager，用它来
 - 创建个ApplicationMaster，SparkContext相当于这个ApplicationMaster管的APP，生成YarnClusterScheduler与YarnClusterSchedulerBackend
 - 选择集群中的容器启动CoarseGrainedExecutorBackend，用来启动spark.executor。
- ApplicationMaster与CoarseGrainedExecutorBackend会有远程调用。

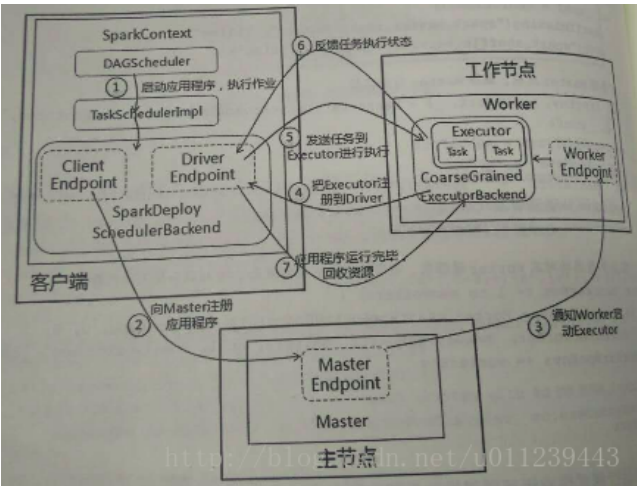
Yarn Client 模式



Spark的Yarn Client 模式流程如下：

- 本地启动SparkContext，生成YarnClientClusterScheduler 和 YarnClientClusterSchedulerBackend
- YarnClientClusterSchedulerBackend启动yarn.Client，用它提交App 到 Yarn Resource Manager
- Yarn Resource Manager 选个 YARN Node Manager，用它来选择集群中的容器启动CoarseCrainedExecutorBackend，用来启动 spark.executor
- YarnClientClusterSchedulerBackend与 CoarseCrainedExecutorBackend会有远程调用。

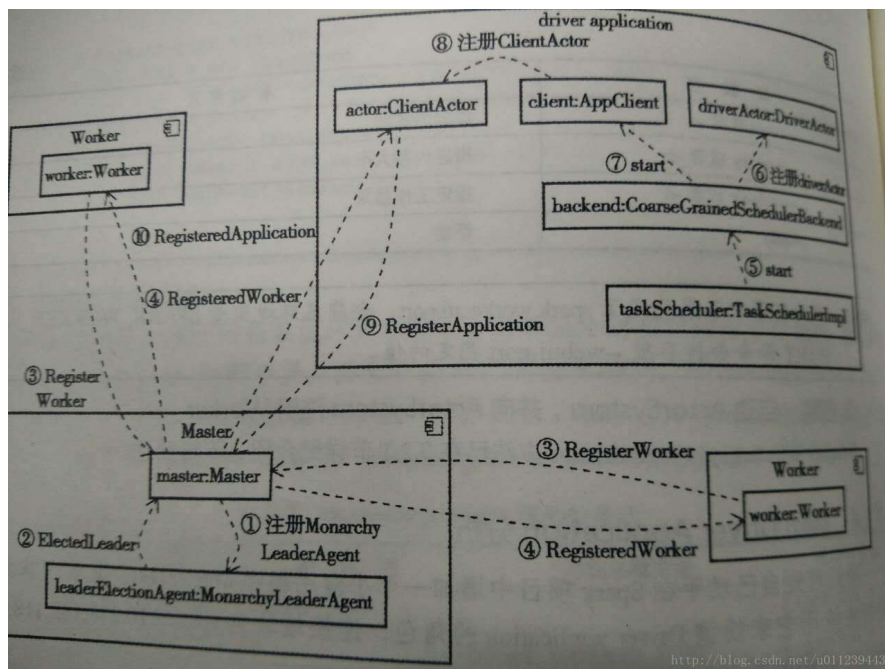
Standalone 模式介绍



1. 启动app，在SparkContxt启动过程中，先初始化DAGScheduler 和 TaskScheduler，并初始化 SparkDeploySchedulerBacken d，并在其内部启动DriverEndpoint和ClientEndpoint。

- ClientEndpoint想Master注册app，Master收到注册信息后把该app加入到等待运行app列表中，等待由Master分配给该app worker。
- app获取到worker后，Master通知Worker的WorkerEndpoint创建CoarseGrainedExecutorBackend进程，在该进程中创建执行容器executor
- executor创建完毕后发送信息给Master和DriverEndpoint，告知Executor创建完毕，在SparkContext注册，后等待DriverEndpoint发送执行任务的消息。
- SparkContext分配TaskSet给CoarseGrainedExecutorBackend，按一定调度策略在executor执行。详见：《深入理解Spark 2.1 Core（二）：DAG调度器的实现与源码分析》与《深入理解Spark 2.1 Core（三）：任务调度器的实现与源码分析》
- CoarseGrainedExecutorBackend在Task处理的过程中，把处理Task的状态发送给DriverEndpoint，Spark根据不同的执行结果来处理。若处理完毕，则继续发送其他TaskSet。详见：《深入理解Spark 2.1 Core（四）：运算结果处理和容错的实现与源码分析》
- app运行完成后，SparkContext会进行资源回收，销毁Worker的CoarseGrainedExecutorBackend进程，然后注销自己。

Standalone 启动集群



启动Master

master.Master

我们先看下Master对象的main函数做了什么：

```

1 private[deploy] object Master extends Logging {
2   val SYSTEM_NAME = "sparkMaster"
3   val ENDPOINT_NAME = "Master"
4
5   def main(argStrings: Array[String]) {
6     Utils.initDaemon(log)
7     //创建SparkConf
8     val conf = new SparkConf
9     //解析SparkConf参数
10    val args = new MasterArguments(argStrings, conf)

```

```
11  val (rpcEnv, _, _) = startRpcEnvAndEndpoint(args.host, args.port, args.webUiPort, conf)
12  rpcEnv.awaitTermination()
13  }
14
15  def startRpcEnvAndEndpoint(
16    host: String,
17    port: Int,
18    webUiPort: Int,
19    conf: SparkConf): (RpcEnv, Int, Option[Int]) = {
20    val securityMgr = new SecurityManager(conf)
21    val rpcEnv = RpcEnv.create(SYSTEM_NAME, host, port, conf, securityMgr)
22    //创建Master
23    val masterEndpoint = rpcEnv.setupEndpoint(ENDPOINT_NAME,
24      new Master(rpcEnv, rpcEnv.address, webUiPort, securityMgr, conf))
25    val portsResponse = masterEndpoint.askWithRetry[BoundPortsResponse](BoundPortsRequest)
26    //返回 Master RpcEnv,
27    //web UI 端口,
28    //其他服务的端口
29    (rpcEnv, portsResponse.webUIPort, portsResponse.restPort)
30  }
31 }
```

master.MasterArguments

接下来我们看看master是如何解析参数的：

```
1  private[master] class MasterArguments(args: Array[String], conf: SparkConf) extends Logging {
2    //默认配置
3    var host = Utils.localHostName()
4    var port = 7077
5    var webUiPort = 8080
6    //Spark属性文件
7    //默认为 spark-default.conf
8    var propertiesFile: String = null
9
10   // 检查环境变量
11   if (System.getenv("SPARK_MASTER_IP") != null) {
12     logWarning("SPARK_MASTER_IP is deprecated, please use SPARK_MASTER_HOST")
13     host = System.getenv("SPARK_MASTER_IP")
14   }
15
16   if (System.getenv("SPARK_MASTER_HOST") != null) {
17     host = System.getenv("SPARK_MASTER_HOST")
18   }
19   if (System.getenv("SPARK_MASTER_PORT") != null) {
20     port = System.getenv("SPARK_MASTER_PORT").toInt
21   }
22   if (System.getenv("SPARK_MASTER_WEBUI_PORT") != null) {
23     webUiPort = System.getenv("SPARK_MASTER_WEBUI_PORT").toInt
24   }
25
26   parse(args.toList)
27
28   // 转变SparkConf
29   propertiesFile = Utils.loadDefaultSparkProperties(conf, propertiesFile)
30   //环境变量的SPARK_MASTER_WEBUI_PORT
31   //会被Spark属性spark.master.ui.port所覆盖
32   if (conf.contains("spark.master.ui.port")) {
33     webUiPort = conf.get("spark.master.ui.port").toInt
34   }
35
36   //解析命令行参数
37   //命令行参数会把环境变量和Spark属性都覆盖
38   @tailrec
```

```

39 private def parse(args: List[String]): Unit = args match {
40   case ("--ip" | "-i") :: value :: tail =>
41     Utils.checkHost(value, "ip no longer supported, please use hostname " + value)
42     host = value
43     parse(tail)
44
45   case ("--host" | "-h") :: value :: tail =>
46     Utils.checkHost(value, "Please use hostname " + value)
47     host = value
48     parse(tail)
49
50   case ("--port" | "-p") :: IntParam(value) :: tail =>
51     port = value
52     parse(tail)
53
54   case "--webui-port" :: IntParam(value) :: tail =>
55     webUiPort = value
56     parse(tail)
57
58   case ("--properties-file") :: value :: tail =>
59     propertiesFile = value
60     parse(tail)
61
62   case ("--help") :: tail =>
63     printUsageAndExit(0)
64
65   case Nil =>
66
67   case _ =>
68     printUsageAndExit(1)
69 }
70
71
72 private def printUsageAndExit(exitCode: Int) {
73   System.err.println(
74     "Usage: Master [options]\n" +
75     "\n" +
76     "Options:\n" +
77     "  -i HOST, --ip HOST      Hostname to listen on (deprecated, please use --host or -h) \n" +
78     "  -h HOST, --host HOST    Hostname to listen on\n" +
79     "  -p PORT, --port PORT    Port to listen on (default: 7077)\n" +
80     "  --webui-port PORT      Port for web UI (default: 8080)\n" +
81     "  --properties-file FILE  Path to a custom Spark properties file.\n" +
82     "                          Default is conf/spark-defaults.conf.")
83   System.exit(exitCode)
84 }
85 }

```

我们可以看到上述参数设置的优先级别为：

系统环境变量 < *spark - default.conf* 中的属性 < 命令行参数 < 应用级代码中的参数设置

启动Worker

worker.Worker

我们先看下Worker对象的主函数做了什么：

```

1 private[deploy] object Worker extends Logging {
2   val SYSTEM_NAME = "sparkWorker"
3   val ENDPOINT_NAME = "Worker"
4
5   def main(argStrings: Array[String]) {
6     Utils.initDaemon(log)

```

```

7      //创建SparkConf
8      val conf = new SparkConf
9      //解析SparkConf参数
10     val args = new WorkerArguments(argStrings, conf)
11     val rpcEnv = startRpcEnvAndEndpoint(args.host, args.port, args.webUiPort, args.cores,
12         args.memory, args.masters, args.workDir, conf = conf)
13     rpcEnv.awaitTermination()
14 }
15
16 def startRpcEnvAndEndpoint(
17     host: String,
18     port: Int,
19     webUiPort: Int,
20     cores: Int,
21     memory: Int,
22     masterUrls: Array[String],
23     workDir: String,
24     workerNumber: Option[Int] = None,
25     conf: SparkConf = new SparkConf): RpcEnv = {
26
27
28     val systemName = SYSTEM_NAME + workerNumber.map(_.toString).getOrElse("")
29     val securityMgr = new SecurityManager(conf)
30     val rpcEnv = RpcEnv.create(systemName, host, port, conf, securityMgr)
31     val masterAddresses = masterUrls.map(RpcAddress.fromSparkURL(_))
32     //创建Worker
33     rpcEnv.setupEndpoint(ENDPOINT_NAME, new Worker(rpcEnv, webUiPort, cores, memory,
34         masterAddresses, ENDPOINT_NAME, workDir, conf, securityMgr))
35     rpcEnv
36 }
37
38 ***

```

worker.WorkerArguments

worker.WorkerArguments与master.MasterArguments类似：

```

1 private[worker] class WorkerArguments(args: Array[String], conf: SparkConf) {
2     var host = Utils.localHostName()
3     var port = 0
4     var webUiPort = 8081
5     var cores = inferDefaultCores()
6     var memory = inferDefaultMemory()
7     var masters: Array[String] = null
8     var workDir: String = null
9     var propertiesFile: String = null
10
11     // 检查环境变量
12     if (System.getenv("SPARK_WORKER_PORT") != null) {
13         port = System.getenv("SPARK_WORKER_PORT").toInt
14     }
15     if (System.getenv("SPARK_WORKER_CORES") != null) {
16         cores = System.getenv("SPARK_WORKER_CORES").toInt
17     }
18     if (conf.getenv("SPARK_WORKER_MEMORY") != null) {
19         memory = Utils.memoryStringToMb(conf.getenv("SPARK_WORKER_MEMORY"))
20     }
21     if (System.getenv("SPARK_WORKER_WEBUI_PORT") != null) {
22         webUiPort = System.getenv("SPARK_WORKER_WEBUI_PORT").toInt
23     }
24     if (System.getenv("SPARK_WORKER_DIR") != null) {
25         workDir = System.getenv("SPARK_WORKER_DIR")
26     }
27

```

```
28 parse(args.toList)
29
30 // 转变SparkConf
31 propertiesFile = Utils.loadDefaultSparkProperties(conf, propertiesFile)
32
33 if (conf.contains("spark.worker.ui.port")) {
34     webUiPort = conf.get("spark.worker.ui.port").toInt
35 }
36
37 checkWorkerMemory()
38
39 @tailrec
40 private def parse(args: List[String]): Unit = args match {
41     case ("--ip" | "-i") :: value :: tail =>
42         Utils.checkHost(value, "ip no longer supported, please use hostname " + value)
43         host = value
44         parse(tail)
45
46     case ("--host" | "-h") :: value :: tail =>
47         Utils.checkHost(value, "Please use hostname " + value)
48         host = value
49         parse(tail)
50
51     case ("--port" | "-p") :: IntParam(value) :: tail =>
52         port = value
53         parse(tail)
54
55     case ("--cores" | "-c") :: IntParam(value) :: tail =>
56         cores = value
57         parse(tail)
58
59     case ("--memory" | "-m") :: MemoryParam(value) :: tail =>
60         memory = value
61         parse(tail)
62
63     //工作目录
64     case ("--work-dir" | "-d") :: value :: tail =>
65         workDir = value
66         parse(tail)
67
68     case "--webui-port" :: IntParam(value) :: tail =>
69         webUiPort = value
70         parse(tail)
71
72     case ("--properties-file") :: value :: tail =>
73         propertiesFile = value
74         parse(tail)
75
76     case ("--help") :: tail =>
77         printUsageAndExit(0)
78
79     case value :: tail =>
80         if (masters != null) { // Two positional arguments were given
81             printUsageAndExit(1)
82         }
83         masters = Utils.parseStandaloneMasterUrls(value)
84         parse(tail)
85
86     case Nil =>
87         if (masters == null) { // No positional argument was given
88             printUsageAndExit(1)
89         }
90
91     case _ =>
```



```
92     printUsageAndExit(1)
93   }
94
95   ***
```

资源回收

我们在概述中提到了“app运行完成后，SparkContext会进行资源回收，销毁Worker的CoarseGrainedExecutorBackend进程，然后注销自己。”接下来我们就来讲解下Master和Executor是如何感知到Application的退出的。

调用栈如下：

- SparkContext.stop
 - DAGScheduler.stop
 - TaskSchedulerImpl.stop
 - CoarseGrainedSchedulerBackend.stop
 - CoarseGrainedSchedulerBackend.stopExecutors
 - CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply
 - CoarseGrainedExecutorBackend.receive
 - Executor.stop
- CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply

SparkContext.stop

SparkContext.stop会调用DAGScheduler.stop

```
1  ***
2    if (_dagScheduler != null) {
3      Utils.tryLogNonFatalError {
4        _dagScheduler.stop()
5      }
6      _dagScheduler = null
7    }
8
9  ***
```

DAGScheduler.stop

DAGScheduler.stop会调用TaskSchedulerImpl.stop

```
1  def stop() {
2    //停止消息调度
3    messageScheduler.shutdownNow()
4    //停止事件处理循环
5    eventProcessLoop.stop()
6    //调用TaskSchedulerImpl.stop
7    taskScheduler.stop()
8  }
```

TaskSchedulerImpl.stop

TaskSchedulerImpl.stop会调用CoarseGrainedSchedulerBackend.stop

```

1  override def stop() {
2      //停止推断
3      speculationScheduler.shutdown()
4      //调用CoarseGrainedSchedulerBackend.stop
5      if (backend != null) {
6          backend.stop()
7      }
8      //停止结果获取
9      if (taskResultGetter != null) {
10         taskResultGetter.stop()
11     }
12     starvationTimer.cancel()
13 }

```

CoarseGrainedSchedulerBackend.stop

```

1  override def stop() {
2      //调用stopExecutors()
3      stopExecutors()
4      try {
5          if (driverEndpoint != null) {
6              //发送StopDriver信号
7              driverEndpoint.askWithRetry[Boolean](StopDriver)
8          }
9      } catch {
10         case e: Exception =>
11             throw new SparkException("Error stopping standalone scheduler's driver endpoint", e)
12     }
13 }

```

CoarseGrainedSchedulerBackend.stopExecutors

我们先看下CoarseGrainedSchedulerBackend.stopExecutors

```

1  def stopExecutors() {
2      try {
3          if (driverEndpoint != null) {
4              logInfo("Shutting down all executors")
5              //发送StopExecutors信号
6              driverEndpoint.askWithRetry[Boolean](StopExecutors)
7          }
8      } catch {
9          case e: Exception =>
10             throw new SparkException("Error asking standalone scheduler to shut down executors", e)
11     }
12 }

```

CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply

DriverEndpoint接收并回应该信号：

```

1      case StopExecutors =>
2          logInfo("Asking each executor to shut down")
3          for ((_, executorData) <- executorDataMap) {
4              //给CoarseGrainedExecutorBackend发送StopExecutor信号
5              executorData.executorEndpoint.send(StopExecutor)
6          }
7          context.reply(true)

```

CoarseGrainedExecutorBackend.receive

CoarseGrainedExecutorBackend接收该信号：

```
1 case StopExecutor =>
2   stopping.set(true)
3   logInfo("Driver commanded a shutdown")
4   //这里并没有直接关闭Executor,
5   //因为Executor必须先返回确认帧给CoarseGrainedSchedulerBackend
6   //所以, 这的策略是给自己再发一个Shutdown信号, 然后处理
7   self.send(Shutdown)
8
9 case Shutdown =>
10  stopping.set(true)
11  new Thread("CoarseGrainedExecutorBackend-stop-executor") {
12    override def run(): Unit = {
13      // executor.stop() 会调用 `SparkEnv.stop()`
14      // 直到 RpcEnv 彻底结束
15      // 但是, 如果 `executor.stop()` 运行在和RpcEnv相同的线程里面,
16      // RpcEnv 会等到`executor.stop()`结束后才能结束,
17      // 这就产生了死锁
18      // 因此, 我们需要新建一个线程
19      executor.stop()
20    }
  }
```

Executor.stop

```
1 def stop(): Unit = {
2   env.metricsSystem.report()
3   //关闭心跳
4   heartbeater.shutdown()
5   heartbeater.awaitTermination(10, TimeUnit.SECONDS)
6   //关闭线程池
7   threadPool.shutdown()
8   if (!isLocal) {
9     //停止SparkEnv
10    env.stop()
11  }
12 }
```

CoarseGrainedSchedulerBackend.DriverEndpoint.receiveAndReply

我们回过头来看CoarseGrainedSchedulerBackend.stop, 调用stopExecutors()结束后, 会给 driverEndpoint发送StopDriver信号。CoarseGrainedSchedulerBackend.DriverEndpoint接收信号并回复：

```
1 case StopDriver =>
2   context.reply(true)
3   //停止driverEndpoint
4   stop()
```