## 官网的说明

```
1  Inside a given Spark application (SparkContext instance), multiple
2  parallel jobs can run simultaneously if they were submitted from
3  separate threads. By "job", in this section, we mean a Spark action
4  (e.g. save, collect) and any tasks that need to run to evaluate that
5  action. Spark's scheduler is fully thread-safe and supports this use
6  case to enable applications that serve multiple requests (e.g. queries
7  for multiple users).
8  By default, Spark's scheduler runs jobs in FIFO fashion. Each job is
9  divided into "stages" (e.g. map and reduce phases), and the first job
10 gets priority on all available resources while its stages have tasks
11 to launch, then the second job gets priority, etc. If the jobs at the
12 head of the queue don't need to use the whole cluster, later jobs can
13 start to run right away, but if the jobs at the head of the queue are
14 large, then later jobs may be delayed significantly.
15 Starting in Spark 0.8, it is also possible to configure fair sharing
16 between jobs. Under fair sharing, Spark assigns tasks between jobs in
17 a "round robin" fashion, so that all jobs get a roughly equal share
18 of cluster resources. This means that short jobs submitted while a
19 long job is running can start receiving resources right away and
20 still get good response times, without waiting for the long job to
21 finish. This mode is best for multi-user settings.
```

在一个spark application(sc 实例)中，多个并行job可以同时运行，如果被一个单独的线程提交的话。

默认情况下，spark的调度器以FIFO运行job，每一个job被切为"stage"(ShuffleMapStage和ReduceStage)，并且当第一个job的stage有task启动时，优先获取所有可用资源，然后第二个job再优先获取，等等。如果第一个job不需要获取整个集群资源，第二个job可以立即运行，但是如果第一个job是很长的，后续的job有非常明显的延迟。

从spark0.8开始，可以在把job配置为FAIR，在FAIR模式下，spark以"轮询"的方式分配job之间的task，以便所有的job获取大致相等的集群资源，这意味着长job开始接收资源的时，提交的短job可以立即执行并获得良好的相应时间，无需等待长job完成。
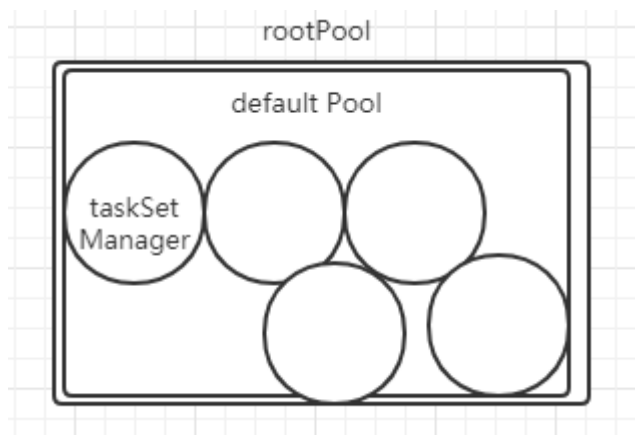
如果启用公平调度的话，只需要将spark.scheduler.mode设为FAIR

```scala
1  val conf = new SparkConf().setMaster(...).setAppName(...)
2  conf.set("spark.scheduler.mode", "FAIR")
```
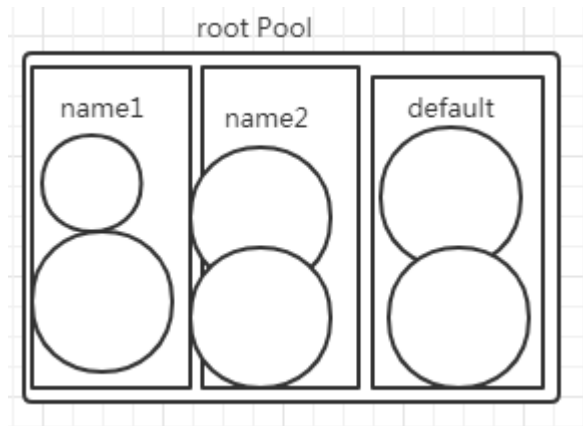
# job提交过程中的 "资源池pool"

　　spark job提交的时候，如果未配置fairScheduler.xml文件(可以去文档末尾查看)，代码未配置spark.scheduler.pool(资源池名)，job会提交的rootPool下的默认资源池，如果配置了，会根据配置参数在rootPool在new Pool，假设我们的application触发了3个job，1个未配置，则构建的调度树为rootPool下有一个默认资源池，两个自定义有name的资源池，当每个job内划分完stage，提交taskSet时，会new taskSetManager管理这组task，并将taskSetManager添加到各自的资源池下，下面用图示来表示该流程。

　　假设job都未自定义，则会在root Pool下new 默认Pool，所有job的所有taskSetManager都

会添加到默认Pool来调度，如下图



　　假设application触发3个job，2个自定义了name，如下图，未自定义job的taskSetManager添加到默认Pool下，自定义的添加到自己的Pool下



```
1  //直接从taskScheduler.submitTasks()开始，之前的在别的小模块细说
2  taskScheduler.submitTasks(new TaskSet(
3    ..., properties))
4
5  taskSchedulerImpl.submitTasks(taskSet)
6
7  manager = createTaskSetManager(taskSet, maxTaskFailures)
8
```

```
 9  schedulableBuilder.addTaskSetManager(manager,
10   manager.taskSet.properties)
11
12  case FIFO
13   rootPool.addSchedulable(manager)
14
15  case FAIR
16   会先判断用户是否自定义了资源池，如果未自定义，则new默认资源池，否则根据
17  配置new资源池
18
```

　　文字描述：将task打包成taskSet提交时，会给当前taskSet创建一个taskSetManager(持续跟踪每个task，在失败次数内，如果task失败则重新提交，并且通过延迟调度尽量保证数据本地性),然后schedulableBuilder将此taskSetManager添加到调度树里，如果是FIFO，则直接添加该taskSetManager到根池下。

　　FAIR是重点，下面是详细的match到FAIR的源码

```
 1  override def addTaskSetManager(manager: Schedulable,
 2   properties: Properties) {
 3   val poolName = if (properties != null) {//判断用户是否自定义资源池
 4   //获取配置的资源池名，如果获取不到，取默认值"default"
 5   properties.getProperty(FAIR_SCHEDULER_PROPERTIES,
 6  DEFAULT_POOL_NAME)
 7   } else {
 8  DEFAULT_POOL_NAME
 9   }
10   //根据资源池名获取资源池
11   var parentPool = rootPool.getSchedulableByName(poolName)
12   if (parentPool == null) {
13   //如果为空，则new 资源池
14   parentPool = new Pool(poolName, DEFAULT_SCHEDULING_MODE,
15  DEFAULT_MINIMUM_SHARE, DEFAULT_WEIGHT)
16   //添加到根池下
17   rootPool.addSchedulable(parentPool)
18   }
19   //不为空，添加taskSetManager到该资源池
20   parentPool.addSchedulable(manager)
21  }
```

## 资源池涉及的源码

#trait Schedulable

可调度的实体，有资源池pool和taskSetManager

```scala
1  var parent: Pool//父资源池
2
3  // child queues 资源池下面添加的子队列
4  //ConcurrentLinkedQueue并发编程中，保证线程安全的
5  def schedulableQueue: ConcurrentLinkedQueue[Schedulable]
6
7  def schedulingMode: SchedulingMode//调度模式
8  def weight: Int//权重
9  def minShare: Int//需要最小的core数
10 def runningTasks: Int//需要运行的task数
11 def priority: Int//job ID
12 def stageId: Int//stage ID
13 def name: String//资源池的名字
14
15 //添加可调度的实体
16 def addSchedulable(schedulable: Schedulable): Unit
17 //移除可调度的实体
18 def removeSchedulable(schedulable: Schedulable): Unit
19 //通过资源池name获取资源池
20 def getSchedulableByName(name: String): Schedulable
```

#class FIFOSchedulingAlgorithm

//FIFO调度策略

```scala
1  override def comparator(s1: Schedulable, s2: Schedulable): Boolean = {
2   val priority1 = s1.priority//获取jobID，越早start的job jobID越小
3   val priority2 = s2.priority
4
5   //signum 大于0 返回1 小于0 返回-1 等于0 返回 0
6   var res = math.signum(priority1 - priority2)
7   if (res == 0) {//如果jobID相同，同一job，比较stageID
8   val stageId1 = s1.stageId//获取stageID
9   val stageId2 = s2.stageId
10  res = math.signum(stageId1 - stageId2)
11  }
12  res < 0 //再比较res是否小于0
13 }
```

#class FairSchedulingAlgorithm

//公平调度策略

```scala
1  override def comparator(s1: Schedulable, s2: Schedulable): Boolean = {
```

```scala
2   val minShare1 = s1.minShare//调度池要求最少的cpu core数
3   val minShare2 = s2.minShare
4   val runningTasks1 = s1.runningTasks//需要运行的task数
5   val runningTasks2 = s2.runningTasks
6
7   //是否需要运行的task数小于要求cpu core数
8   val s1Needy = runningTasks1 < minShare1
9   val s2Needy = runningTasks2 < minShare2
10
11  //求task数/cpu数，如果需要task数越少，cpu越大，值越小
12  val minShareRatio1 = runningTasks1.toDouble
13  / math.max(minShare1, 1.0)
14  val minShareRatio2 = runningTasks2.toDouble
15  / math.max(minShare2, 1.0)
16
17  //task数/权重 权重越高，值越小
18  val taskToWeightRatio1 = runningTasks1.toDouble
19  / s1.weight.toDouble
20  val taskToWeightRatio2 = runningTasks2.toDouble
21  / s2.weight.toDouble
22
23  var compare = 0
24  if (s1Needy && !s2Needy) {//如果s1 task数小于cpu数，先运行s1
25  return true
26  } else if (!s1Needy && s2Needy) {
27  return false
28  } else if (s1Needy && s2Needy) {//如果都小于
29  //比较比值，小的先运行
30  compare = minShareRatio1.compareTo(minShareRatio2)
31  } else {
32  //还相等，比较给的权重
33  compare = taskToWeightRatio1.compareTo(taskToWeightRatio2)
34  }
35  if (compare < 0) {
36  true
37  } else if (compare > 0) {
38  false
39  } else {
40  s1.name < s2.name//上述如果都相等，比较给定的资源池name
41  }
```

# #class Pool

## //可调度的资源池

```
1  //调度队列
2  val schedulableQueue = new ConcurrentLinkedQueue[Schedulable]
3  #该类的主要方法有：
4  1.根据传入的schedulingMode来匹配是FAIR或FIFO
5  2.添加/移除可调度的实体，如果是根池，可以添加资源池，
6    如果不是根池，只能添加tasksetmanager
7  3.根据传入的task数，增加或减少要运行的task数
```

# #trait SchedulableBuilder

## //构建调度树的接口

```
1  def rootPool: Pool//返回根资源池
2  def buildPools(): Unit//构建树节点
3  def addTaskSetManager(manager: Schedulable, properties: Properties)
4    : Unit//构建叶节点
```

# #class FIFOSchedulableBuilder

## //构建FIFO调度树

```
1  override def addTaskSetManager(manager: Schedulable,
2   properties: Properties) {
3    rootPool.addSchedulable(manager)//在根池子里添加tasksetmanager
4  }
```

# #class FairSchedulableBuilder

## //构建FAIR调度树

```
1  #三个主要的方法
2  def buildPools()
3  如果用户创建了fairScheduler.xml，则取读取内容，根据配置创建资源池，如果配置
4  有误，则创建默认资源池
5
6  def buildDefaultPool()
7  先判断是否有name为default的资源池(大池子/根池子)，没有则创建默认资源池，
8  有的话将tasksetmanager添加到下面
9
10  def buildFairSchedulerPool()
11  构建公平调度的池子，加载fairScheduler.xml配置文件，遍历结构，用
12  获取到的参数new 资源池，
13
14  rootPool.addSchedulable(new Pool(poolName, schedulingMode,
15   minShare, weight))
```

fairScheduler.xml的配置如下(官网案例)

可以配置资源池名，调度模式，权重，需要的cpu core数，可以在job触发前设置spark.scheduler.pool，spark会根据设置的name去xml文件找对应配置参数

```xml
<?xml version="1.0"?>
<allocations>
<pool name="production">
<schedulingMode>FAIR</schedulingMode>
<weight>1</weight>
<minShare>2</minShare>
</pool>
<pool name="test">
<schedulingMode>FIFO</schedulingMode>
<weight>2</weight>
<minShare>3</minShare>
</pool>
</allocations>
```