

Flink整体执行流程

以Flink源码中自带的WordCount为例，执行的入口从用户程序的execute()函数入手，execute()的源码如下：

```
1 public JobExecutionResult execute(String jobName) throws Exception {
2     StreamGraph streamGraph = getStreamGraph();
3     streamGraph.setJobName(jobName);
4     JobGraph jobGraph = streamGraph.getJobGraph();
5     . . . . .
6     LocalFlinkMiniCluster exec = new LocalFlinkMiniCluster(configuration, true);
7     try {
8         exec.start();
9         return exec.submitJobAndWait(jobGraph, getConfig().isSysoutLoggingEnabled())
10    }
11    finally {
12        transformations.clear();
13        exec.stop();
14    }
15 }
```

函数内部主要有 getStreamGraph()、getJobGraph()、exec.start()、exec.submitJobAndWait() 等。getStreamGraph()的作用是生成StreamGraph图，getJobGraph()的作用是生成JobGraph的图，exec.start()的作用是建立Client、JobManager、TaskManager三者之间通信初始化，exec.submitJobAndWait()的作用提交job并且等待job执行后的结果，该函数提供了任务执行调度执行的入口，进入Client类中，首先执行createUserCodeClassLoader()函数，创建用户代码的加载器，然后执行jobClient.SubmitJobAndWait()，进入JobClient类，在函数内部会执行submit函数，从该函数开始进入AKKA通信阶段，首先会进入JobClientActor，会创建一个jobclientActor来对JobManager和client进行通信，当通信对象创建之后，会执行akka机制的ask函数，该函数的作用是发出一个消息，然后要求收到方给予回复。当消息发出之后，OnReceive()函数会收到actor发出的消息请求，然后调用handleMessage()方法来处理消息请求，该函数内部有connectToJobManager()方法，此方法内部的tryToSubmitJob()函数是正式提交任务的操作，主要做的工作就是uploadUserJars()上传用户程序的jar文件，接着会jobManager.tell()向JobManager发出一个submit消息请求。

当JobManager收到Client发送的消息之后，会执行JobManager内部的submitJob方法，

```
1 case SubmitJob(jobGraph, listeningBehaviour) =>
2     val client = sender()
3
4     val jobInfo = new JobInfo(client, listeningBehaviour, System.currentTimeMillis(),
5                               jobGraph.getSessionTimeout)
6     log.info("liuzf---开始执行JobManager的submitJob()")
7     submitJob(jobGraph, jobInfo)
```

首先会把由client收到的job信息封装在jobinfo中，然后把jobinfo以及job的任务图jobGraph一起发送给submit()去执行，在JobManager的submit函数中处理的函数逻辑比较复杂，比较重要的函数执行过程如下：

```

1 private def submitJob(jobGraph: JobGraph, jobInfo: JobInfo, isRecovery: Boolean = false) {
2     try {
3         libraryCacheManager.registerJob(jobGraph.getJobID, jobGraph.getUserJarBlobKey,
4                                         jobGraph.getClasspaths)
5     }
6     val userCodeLoader = libraryCacheManager.getClassLoader(jobGraph.getJobID)
7
8 }
9 executionGraph = ExecutionGraphBuilder.buildGraph()
10    try {
11        submittedJobGraphs.putJobGraph(new SubmittedJobGraph(jobGraph, jobInfo))
12        jobInfo.notifyClients(
13            decorateMessage(JobSubmitSuccess(jobGraph.getJobID)))
14        log.info(s"开始调度 job $jobId ($jobName).")
15        executionGraph.scheduleForExecution()

```

首先执行 libraryCacheManager.registerJob(), 向 CacheManager 进行注册, 请求缓存, 然后执行 getClassLoader() 来加载用户的代码加载器, 接下来会调用 ExecutionGraph 中的 buildGraph() 构造 ExecutionGraph 的并行化版本的执行图, 当逻辑执行图构造完毕之后, 这时候可以通知 Client 任务已经成功提交, 并且提交过程结束。接下来会调用 scheduleForExecution() 来会整体的资源进行调度分配, 主要是每个 TaskManager 中的 slot 的分配, 并且当 slot 分配完成之后, 所有的 task 的任务状态发生改变, 由 CREATED→SCHEDULED。接下分配完之后, 接下来执行 depolyToSlot() 函数, 就要进入部署状态, 同样会执行 transitionState() 函数, 将 SCHEDULED 状态变为 DEPLOYING 状态, 接着的重要函数是 submitTask() 函数, 该函数会通过 AKKA 机制, 向 TaskManager 发出一个 submitTask 的消息请求, TaskManager 收到消息请求后, 会执行 submitTask() 方法, 该函数的重要执行过程如下:

```

1 public submitTask(){
2     val task = new Task(. . . .)
3     log.info(s"Received task ${task.getTaskInfo.getTaskNameWithSubtasks()}")
4     val execId = tdd.getExecutionAttemptId
5     val prevTask = runningTasks.put(execId, task)
6     if (prevTask != null) {
7         runningTasks.put(execId, prevTask)
8         throw new IllegalStateException("TaskManager already contains a task with this id")
9     }
10    task.startTaskThread()
11    sender ! decorateMessage(Acknowledge.get())
12 }

```

首先执行 Task 的构造函数, 生成具体物理执行的相关组件, 比如 ResultPartition 等, 最后创建执行 Task 的线程, 然后调用 startTaskThread() 来启动具体的执行线程, Task 线程内部的 run() 方法承载了被执行的核心逻辑, 该方法具体的内容为:

```

1 public void run() {
2     while (true) {
3         ExecutionState current = this.executionState;
4         if (current == ExecutionState.CREATED) {
5             if (transitionState(ExecutionState.CREATED, ExecutionState.DEPLOYING)) {
6                 break;

```

```

7         }
8     }
9     invocable = loadAndInstantiateInvokable(userCodeClassLoader, nameOfInvokable);
10    network.registerTask(this);
11    Environment env = new RuntimeEnvironment(. . . . );
12    invocable.setEnvironment(env);
13    // -----
14    //  actual task core work
15    if (!transitionState(ExecutionState.DEPLOYING, ExecutionState.RUNNING)) {
16    }
17    // notify everyone that we switched to running
18    notifyObservers(ExecutionState.RUNNING, null);
19    executingThread.setContextClassLoader(userCodeClassLoader);
20    // run the invocable
21    invocable.invoke();
22
23    if (transitionState(ExecutionState.RUNNING, ExecutionState.FINISHED)) {
24        notifyObservers(ExecutionState.FINISHED, null);
25    }
26    Finally{
27        // free the network resources
28        network.unregisterTask(this);
29        // free memory resources
30        if (invocable != null) {
31            memoryManager.releaseAll(invocable);
32        }
33        libraryCache.unregisterTask(jobId, executionId);
34        removeCachedFiles(distributedCacheEntries, fileCache);

```

首先执行 transitionState() 函数将 TaskManager 的状态由 CREATED 转变为 DEPOLYING 状态，然后调用 loadAndTrantiateInvokable() 对用户代码打包成 jar 包，并且生成用户代码加载器，然后执行 network.registerTask()，执行该函数之前，会执行 NetworkEnvironment 的构造函数，该类是 TaskManager 通信的主对象，主要用于跟踪中间结果并负责所有的数据交换，在该类中会创建协助通信的关键部件，比如网络缓冲池，连接管理器，结果分区管理器，结果分区可消费通知器等。当网络对象准备完成后，创建一个运行环境，然后执行 invoke.setEnvironment(env)，将各种配置打包到运行环境中。

当运行环境准备之后，接下来到了具体分析任务执行的时候，首先会调用 transitionState() 函数将任务状态由 DEPOLYING 改为 RUNNING 状态，然后会调用 notifyObservers() 通知所有的 task 观察者也改变状态，然后执行 setContextClassLoader() 将执行的类加载器设置为用户执行的加载器，然后执行 invocable.invoke()，该函数是分界点，执行前用户逻辑没有被触发，执行之后说明用户逻辑已完成。当执行完成之后，调用 transitionState() 函数执行的 RUNNING 状态改成 FINISHED 状态。同样调用 notifyObservers() 来通知其他观察者改变状态，最后，释放资源。

总体的函数执行图如下：因图片太大——>>>>

链接: <https://pan.baidu.com/s/149pxDTNtDX5kAG3ocJsNGg> 提取码: bps8