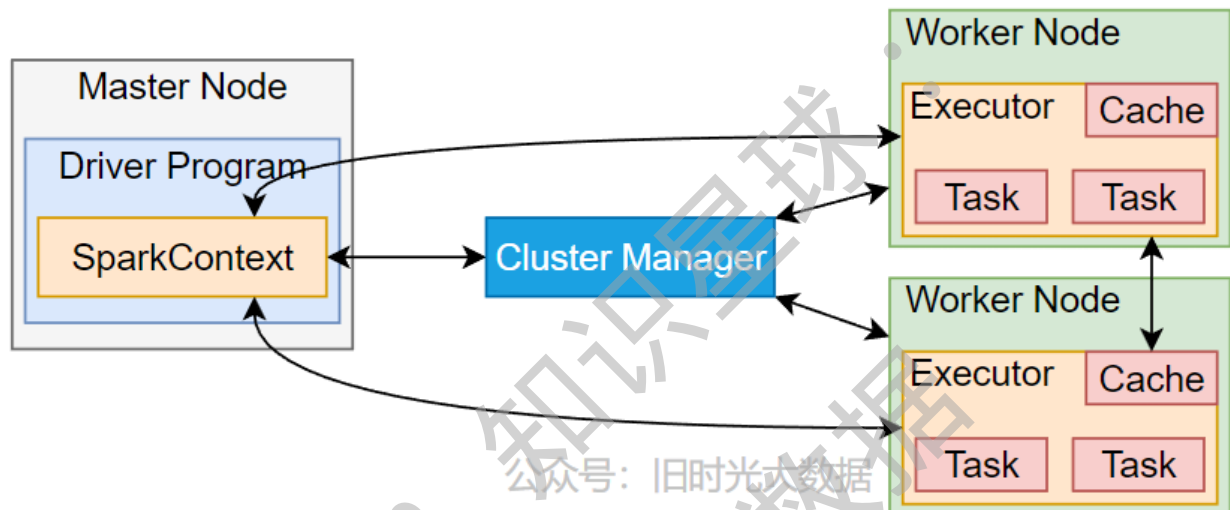


Spark的架构

问过的一些公司：美的杭州研究院(2022.08)，小红书(2022.04)，欢聚(2021.10)，美团(2021.08)，深信服社招(2020.11)，小米(2020.09)，一点资讯(2020.08)，阿里(2018.09)

参考答案：

主要包括五个组件：Driver、Master、Worker、Executor和Task。



1、Driver

Driver是一个进程，我们编写的Spark程序运行在Driver上，由Driver进程执行，Driver是作业的主进程，具有main函数，是程序的入口点，Driver进程启动后，向Master发送请求，进行注册，申请资源，在后面的Executor启动后，会向Driver进行反注册，Driver注册了Executor后，正式执行Spark程序，读取数据源，创建RDD或Dataframe，生成Stage，提交Task到Executor。

2、Master

常驻Master进程，该进程负责管理所有的Worker节点。

分配任务、收集运行信息、监控worker的存活状态。

3、Worker

常驻Worker进程，该进程与Master节点通信，还管理Spark任务的执行。

启动Executor，监控任务运行状态。

4、Executor

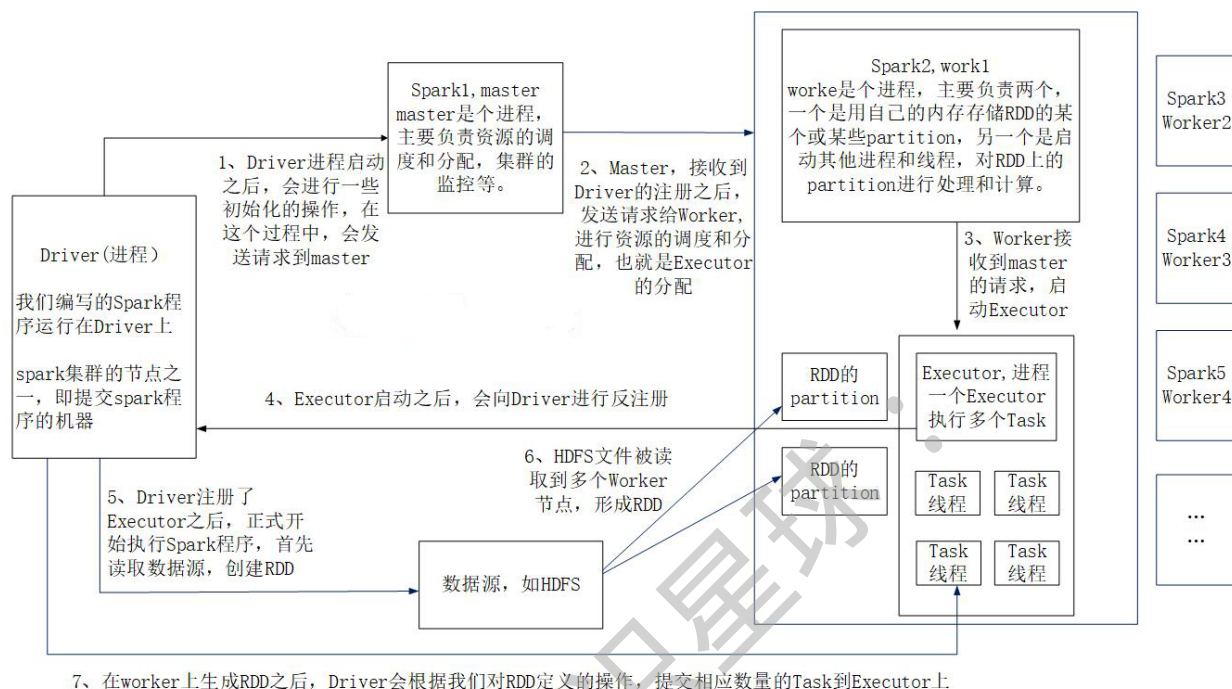
Executor是个进程，一个Executor执行多个Task，多个Executor可以并行执行，可以通过-num-executors来指定Executor的数量，建议Executor最大为集群可用的CPU核数减1。

5、Task

Task是个线程，具体的Spark任务是在Task上运行的，某些并行的算子，有多少个分区就有多少个Task，但是有些算子像Take这样的只有一个Task。

6、详细流程

流程图：



1. Driver进程启动之后，会进行一些初始化的操作，在这个过程中，会发送请求到Master
2. Master接收到Driver的注册之后，发送请求给Worker，进行资源的调度和分配，也就是Executor的分配
3. Worker接收到master的请求，启动Executor
4. Executor启动之后，会向Driver进行反注册
5. Driver注册了Executor之后，正式开始执行Spark程序，首先读取数据源，创建RDD
6. HDFS文件被读取到多个Worker节点，形成RDD
7. 在worker上生成RDD之后，Driver会根据我们对RDD定义的操作，提交相应数量的Task到Executor上

欢迎加入知识星球，获取《大数据面试题 V4.0》以及更多大数据开发学习资料



知识星球

长按扫码领取优惠

