

# Spark 初始

讲师：陈博

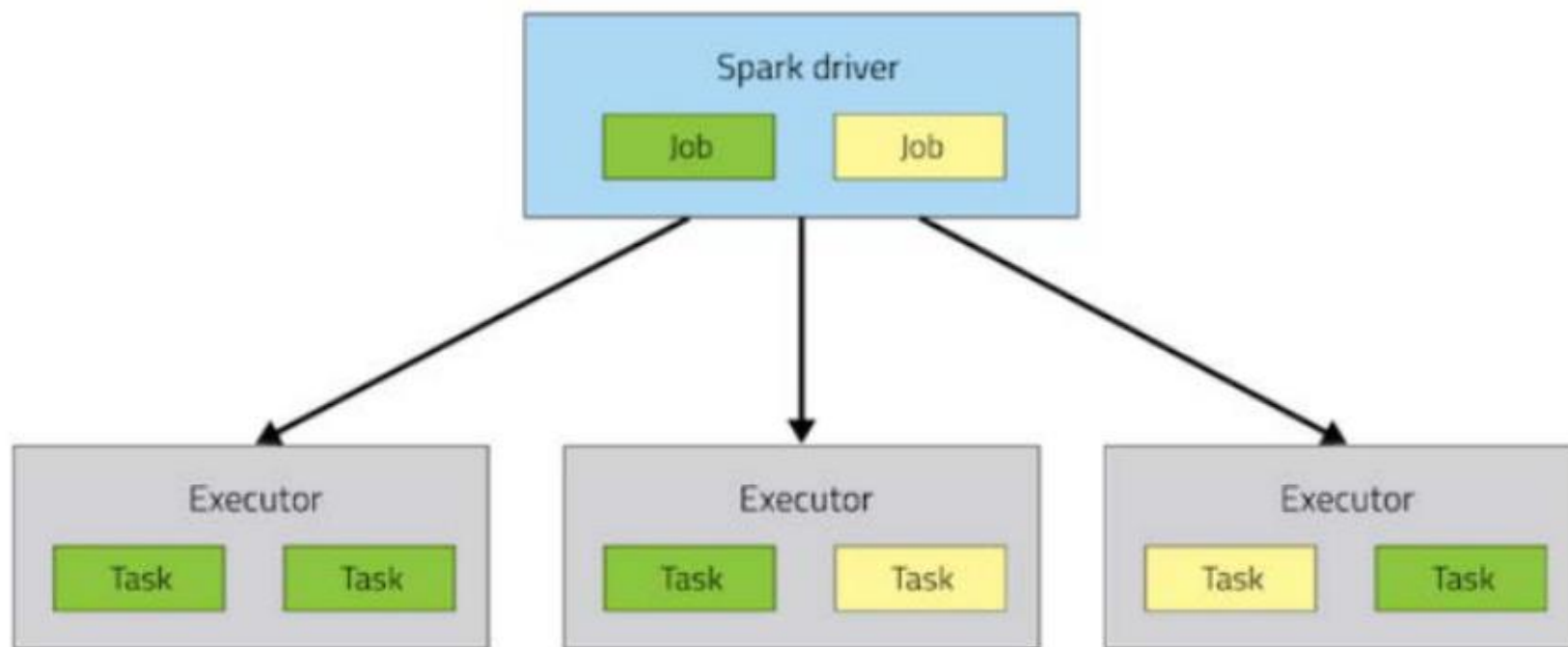


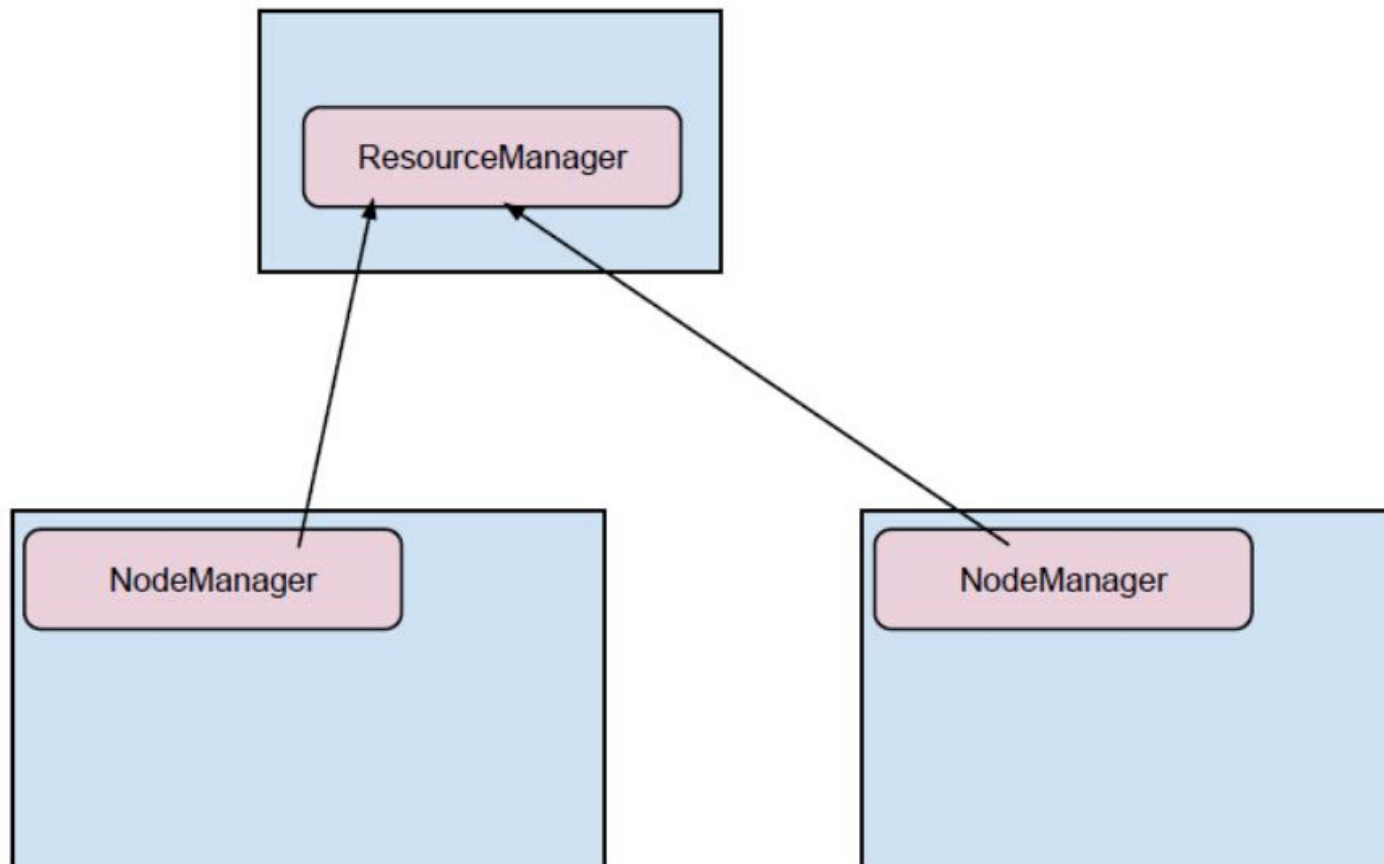
- Why YARN?
- Hadoop MR
- Spark
- MPI
- Storm
- 共享数据，更方便的管理集群

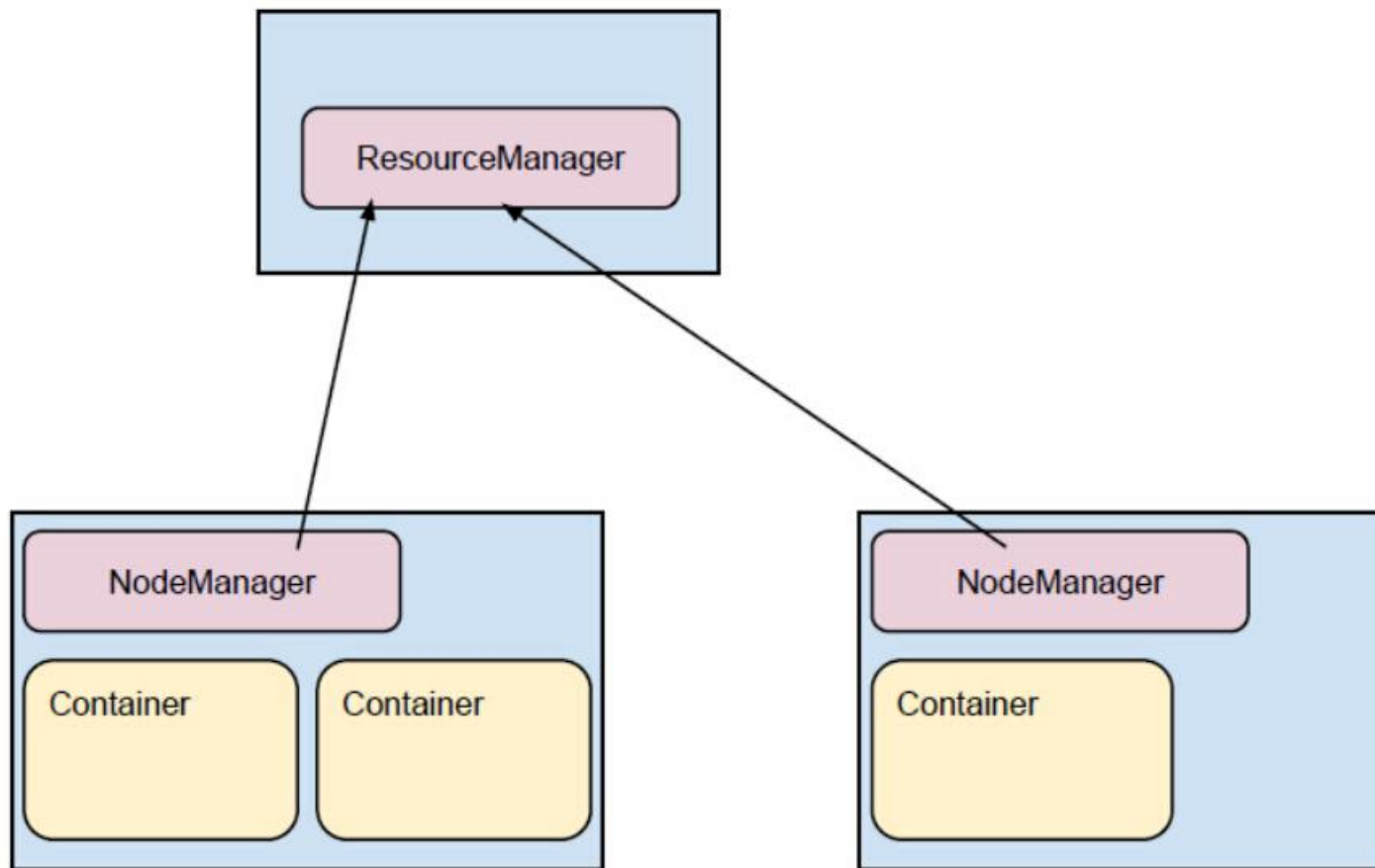


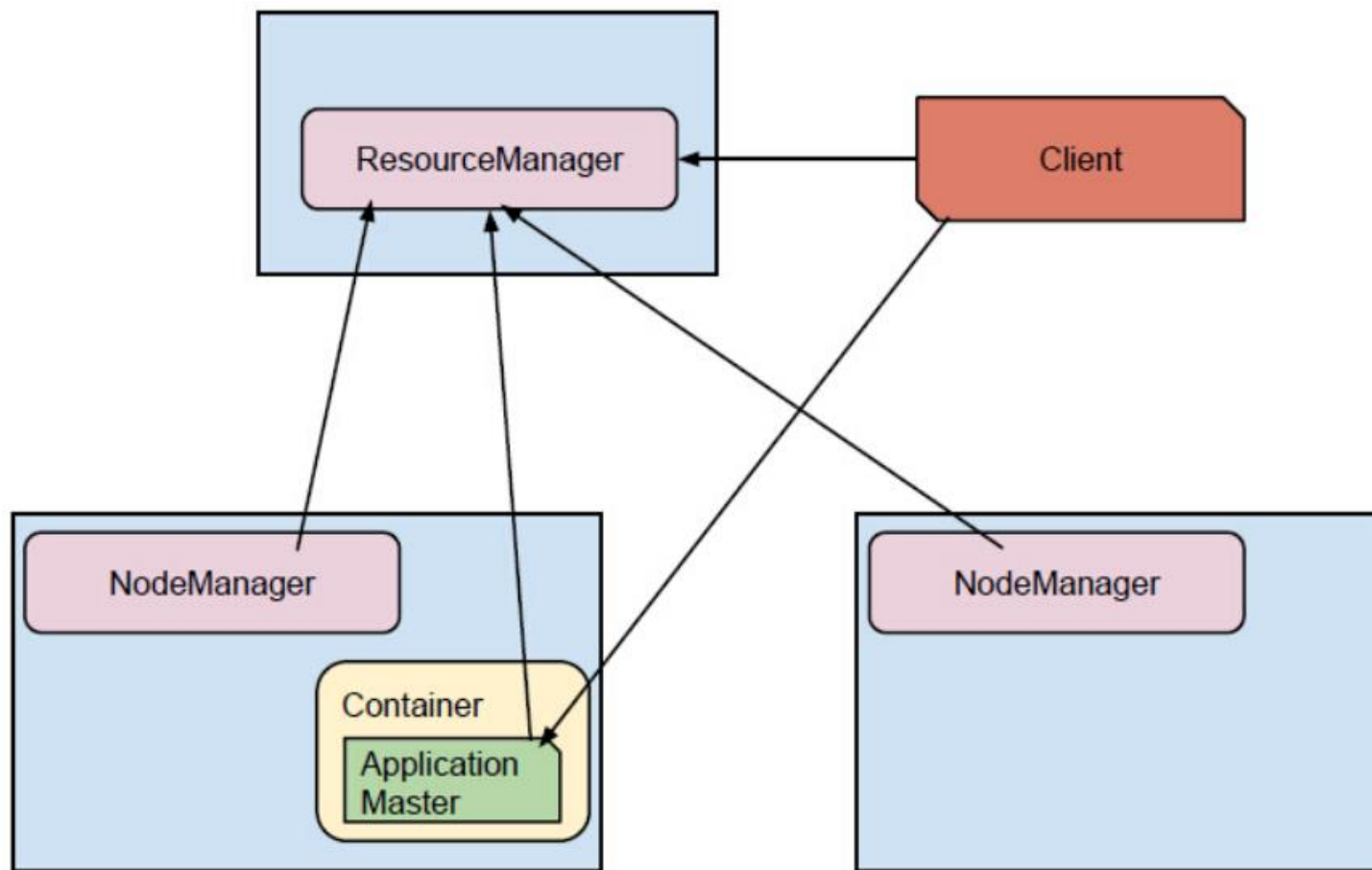
- ResourceManager
- ApplicationMaster
- NodeManager
- Container



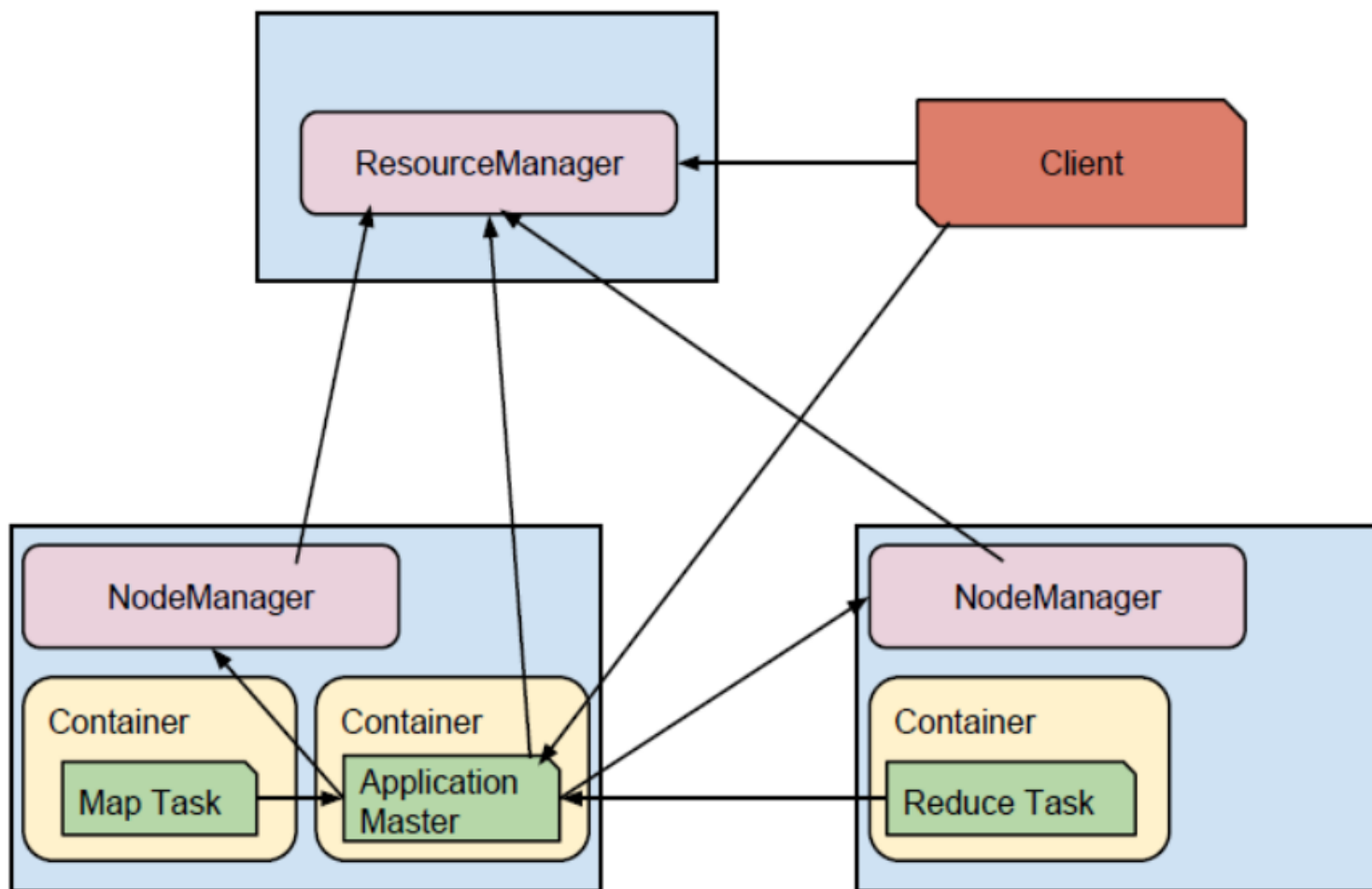






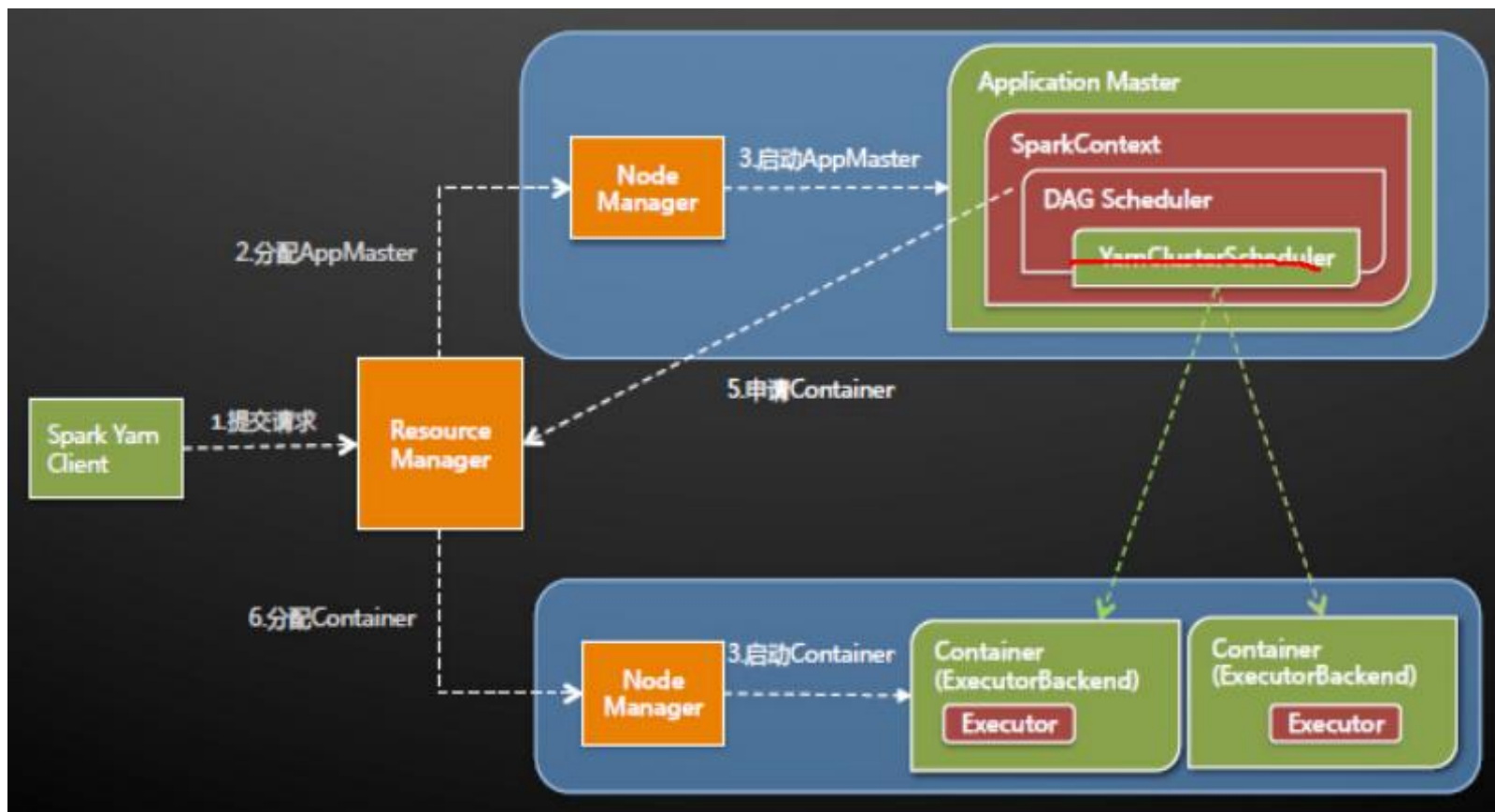


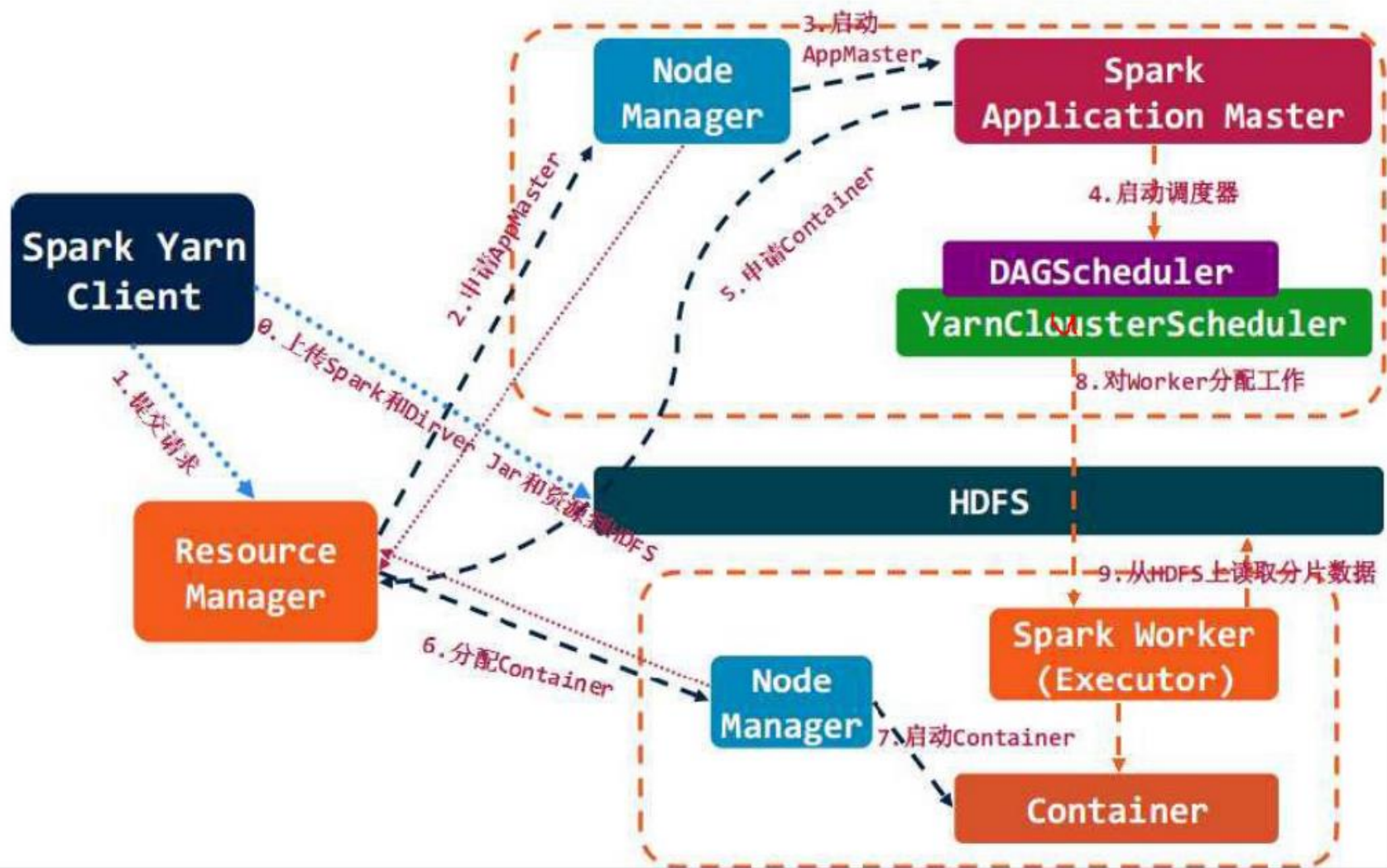
# YARN架构图



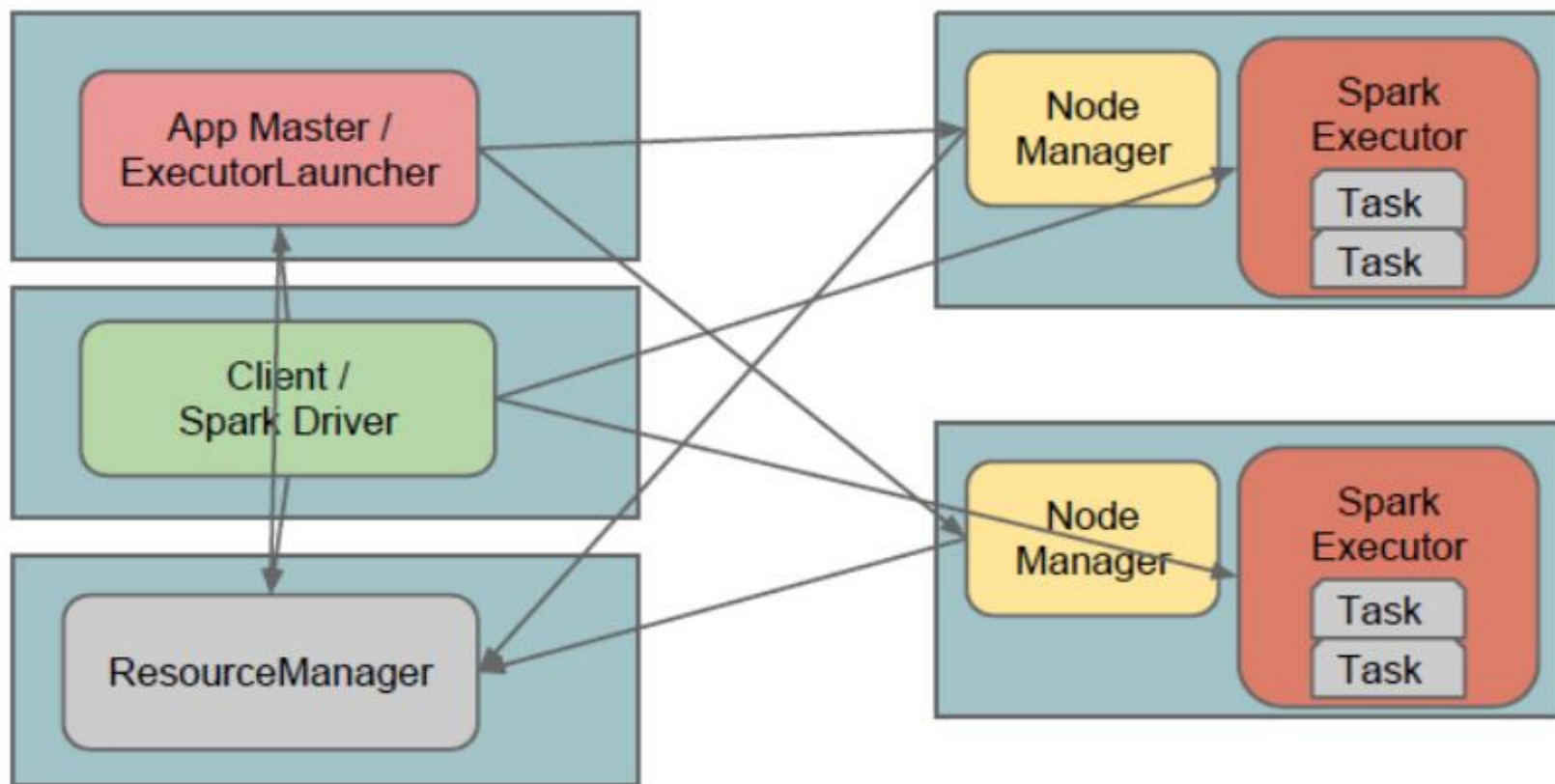


- 每个SparkContext对应一个ApplicationMaster
- 每个Executor对应一个Container

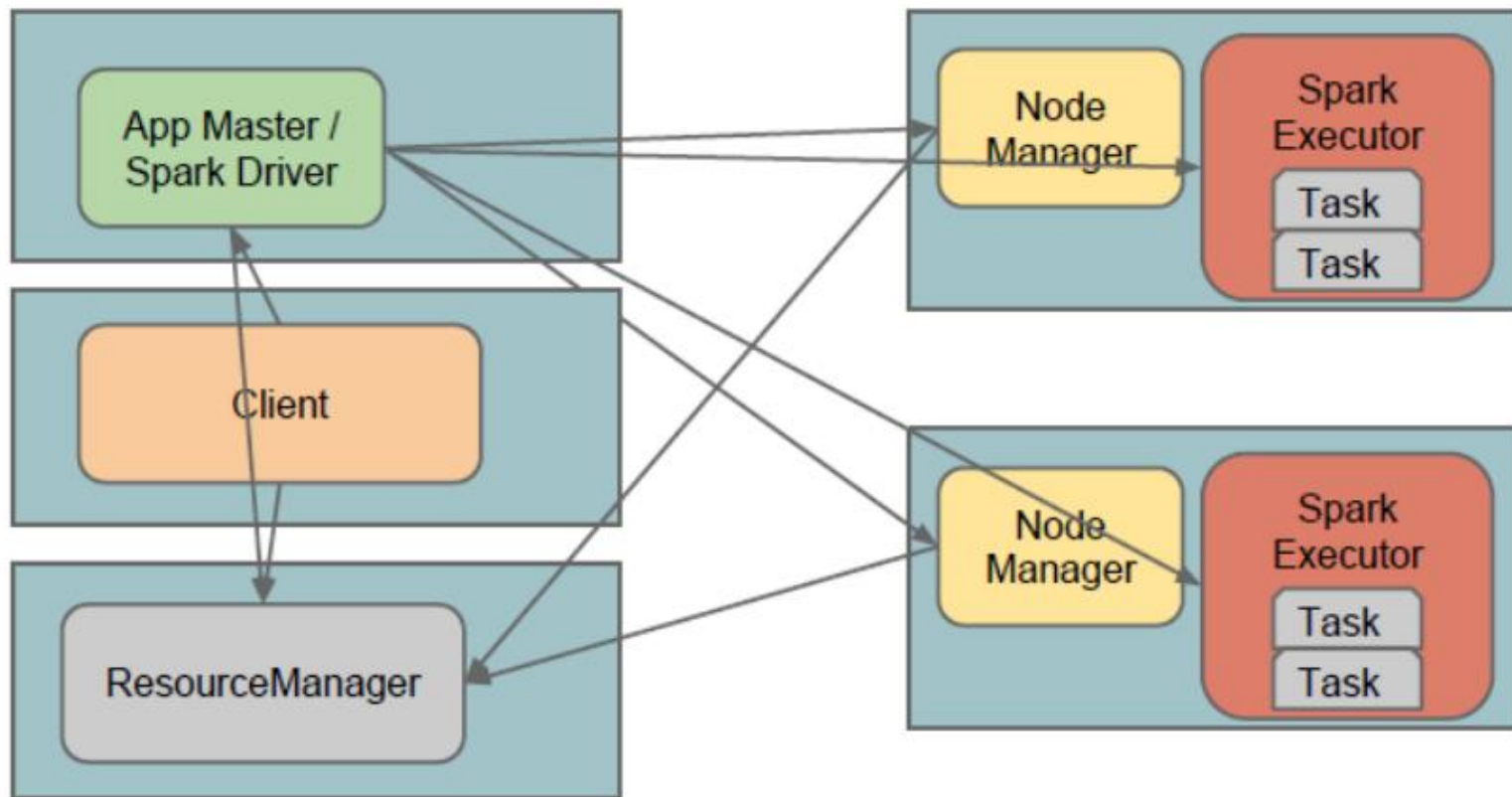




- YarnClientClusterScheduler
- Client和Driver运行行在一起，AM只用用来获取资源



- YarnClusterScheduler
- Driver和AM运行行在一一起，Client可以跑远





- 确保HADOOP\_CONF\_DIR或YARN\_CONF\_DIR指向hadoop集群上含有客户端配置文档的目录。
- 这里有两种发布模式可以被用于在YARN上发布Spark应用，在yarn-cluster模式中，Spark驱动程序跑在一个集群上由YARN管理的Application Master线程中，并且客户端会消失当初始化应用后。在yarn-client模式中，驱动程序跑在客户端线程中，并且Application Master仅被用于请求YARN上面的资源。
- 不像Spark的standalone和Mesos模式，他们需要指定master的地址通过具体的“master”参数。在YARN模式中，资源管理ResourceManager的地址是自动由hadoop配置读取出来的，因此，主要的参数就是简单的yarn-client或者yarn-cluster。
- 去发布一个Spark应用通过yarn-cluster模式：
- `./bin/spark-submit --class path.to.your.Class --master yarn-cluster [options] <app jar> [app options]`



- ./bin/spark-submit --class  
org.apache.spark.examples.SparkPi \
- --master yarn-cluster \
- --num-executors 3 \
- --driver-memory 4g \
- --executor-memory 2g \
- --executor-cores 1 \
- Lib/spark-examples\*.jar \
- 10



- 上面的代码将开始一个YARN客户端程序有一个默认的Application Master.然后SparkPi将会被作为一个Application Master的子线程运行。客户端将周期地将更新的状态告诉Application Master并将他们显示出来。客户端将推出一旦你的应用跑完了。
- 去发布一个Spark应用在yarn-client模式，一样的，只是用yarn-client取代yarn-cluster。去运行spark-shell：
- `./bin/spark-shell -master yarn-client`

