

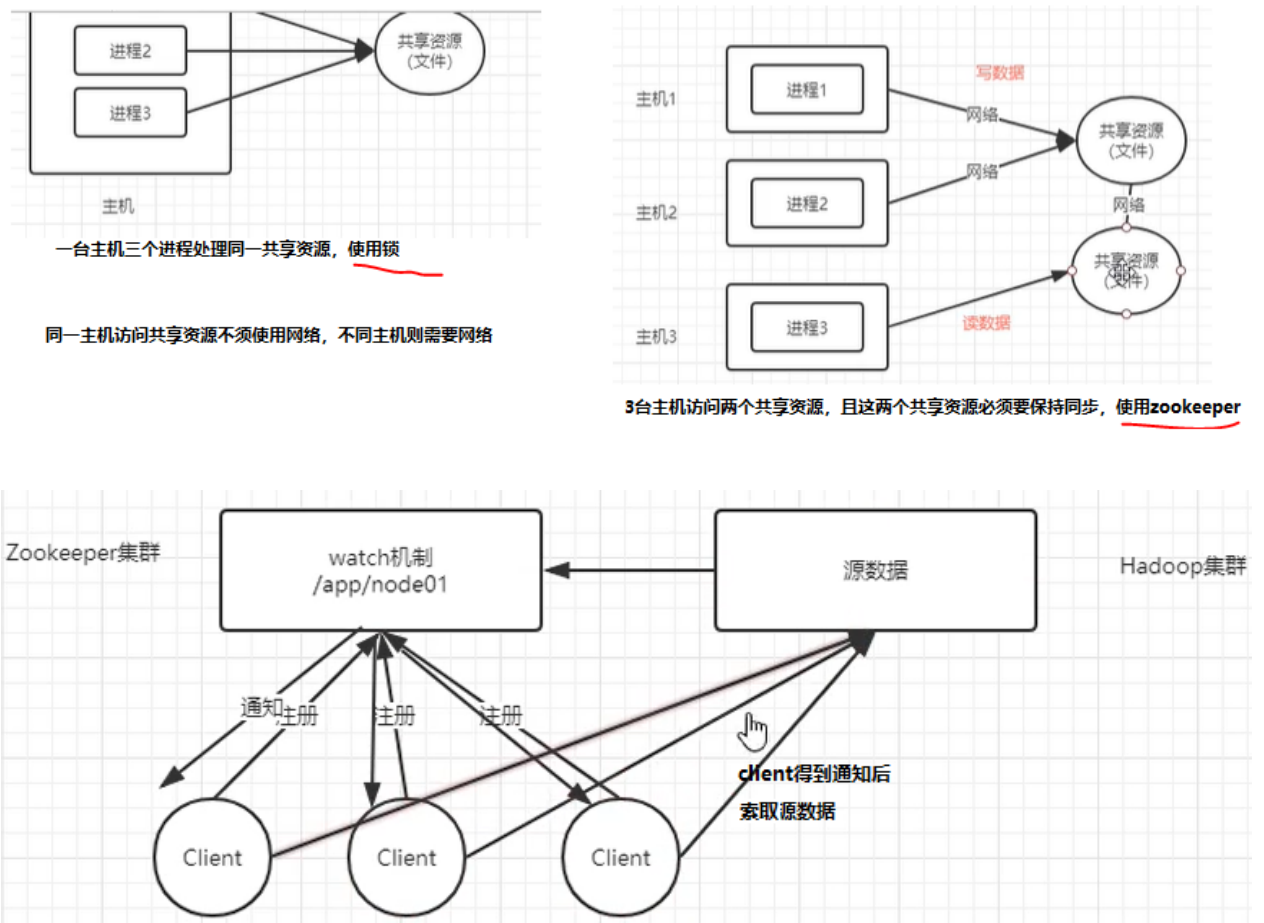
<https://www.bilibili.com/video/BV1ek4y117Yq?p=146> (<https://www.bilibili.com/video/BV1ek4y117Yq?p=146>)

Google三篇论文 1、《分布式文件系统》 2、《分布式计算模型》 3、《BigTable》

hadoop的三大核心：1.分布式存储 2.分布式计算 3.yarn-分布式资源管理调度模块

In []:

1

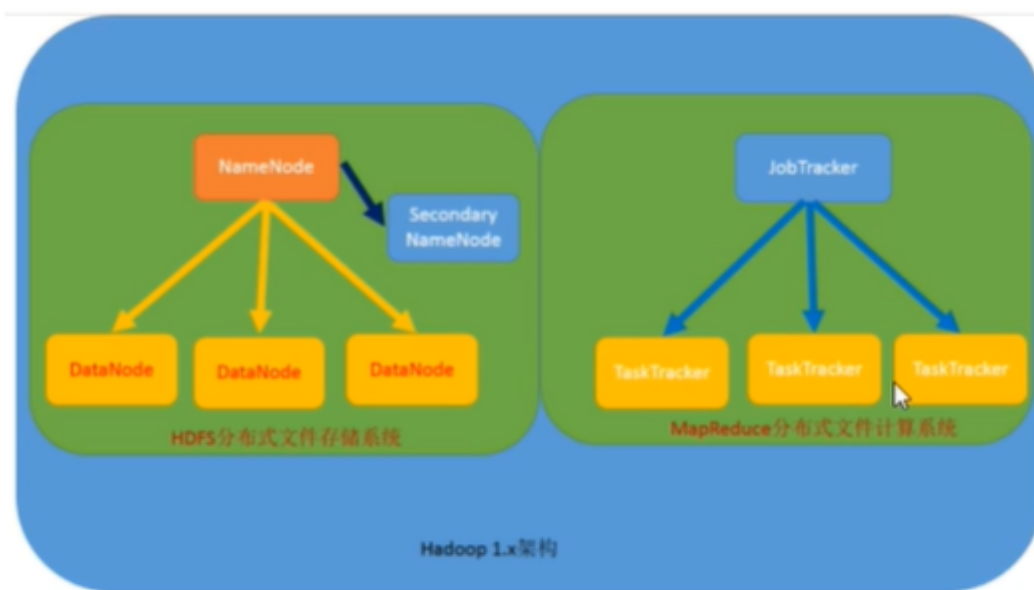
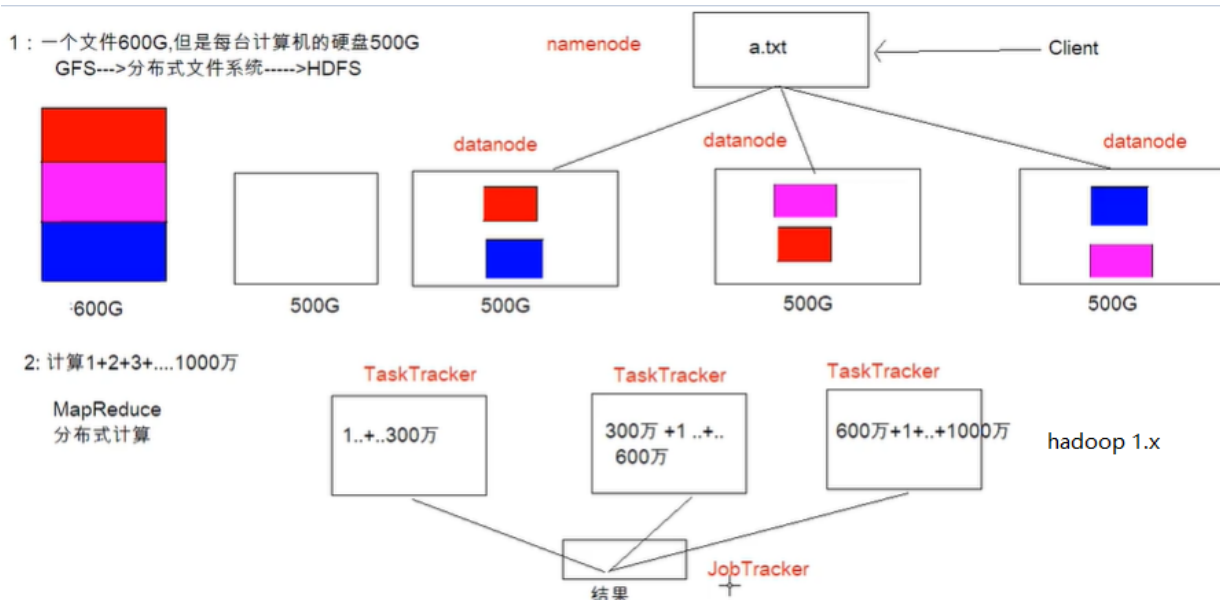


In []:

1

In []:

1



文件系统核心模块:

NameNode: 集群当中的主节点, 管理元数据(文件的大小, 文件的位置, 文件的权限), 主要用于管理集群当中的各种数据

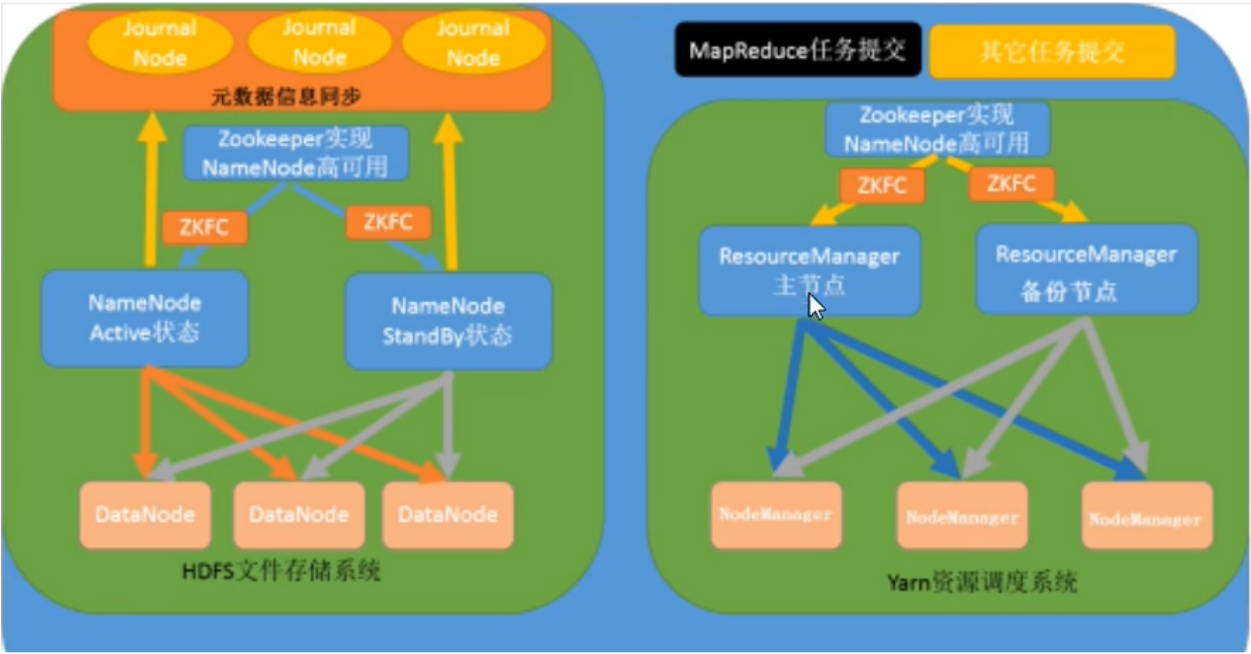
secondaryNameNode: 主要能用于hadoop当中元数据信息的辅助管理

DataNode: 集群当中的从节点, 主要用于存储集群当中的各种数据

数据计算核心模块:

JobTracker: 接收用户的计算请求任务, 并分配任务给从节点

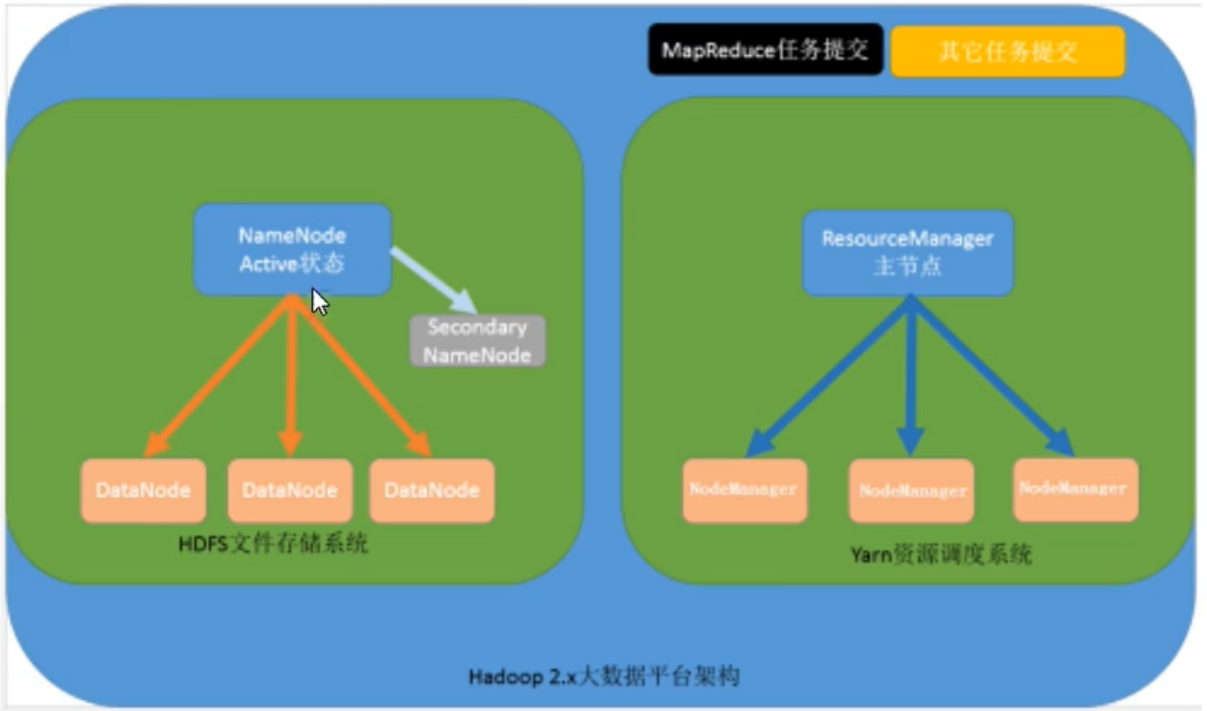
TaskTracker: 负责执行主节点JobTracker分配的任务



hadoop 2.x 生产中使用

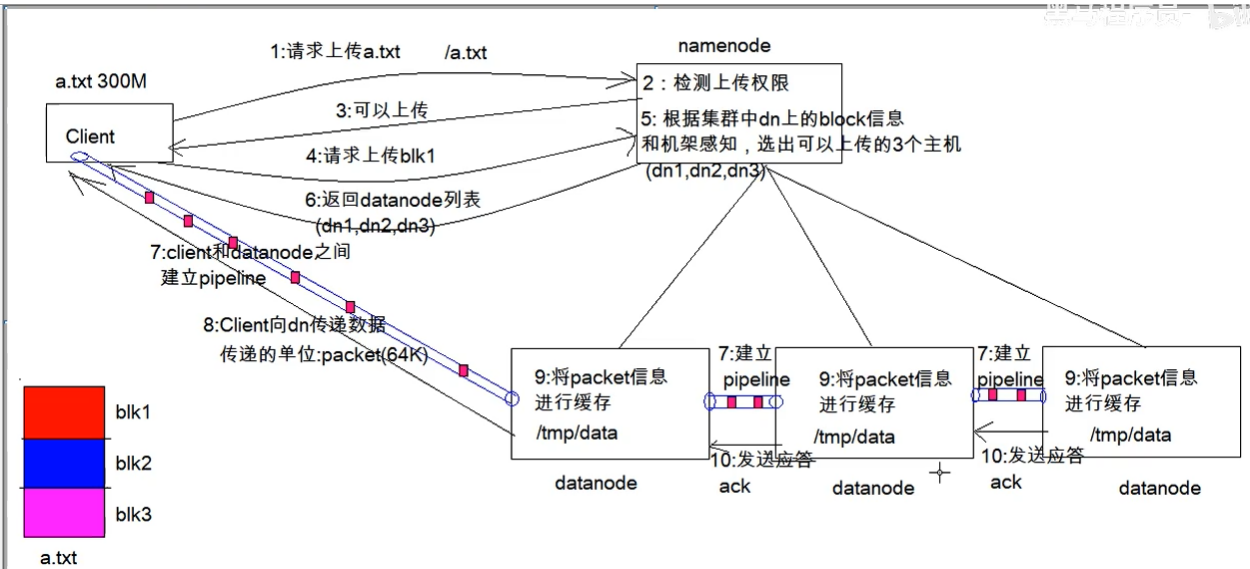
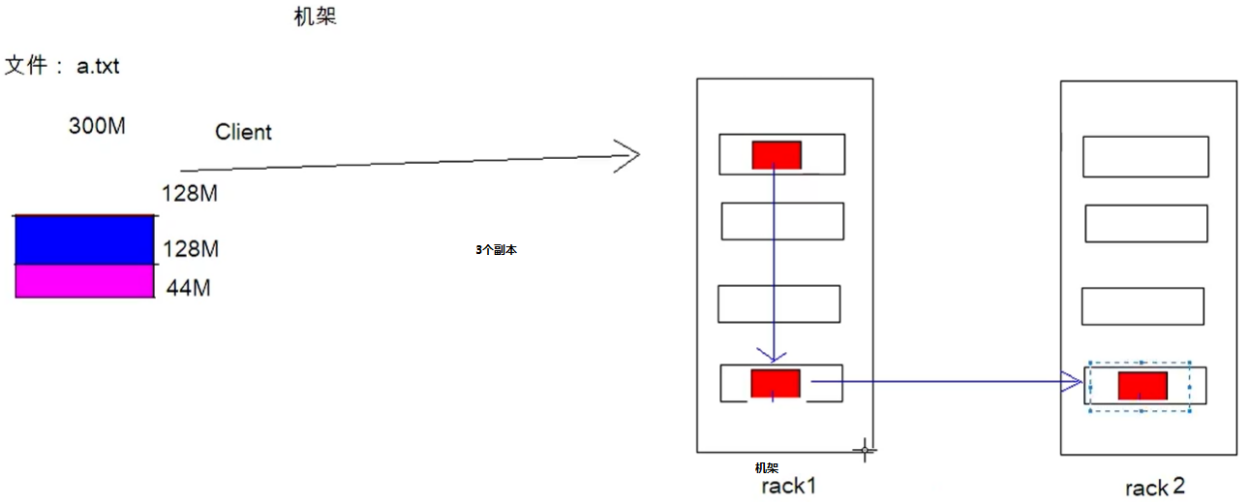
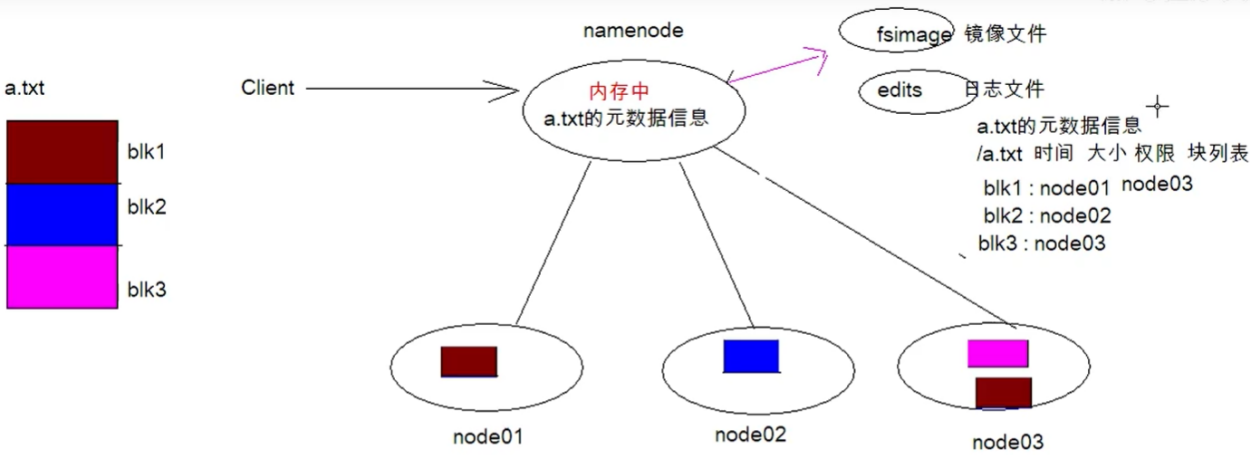
学习使用下:

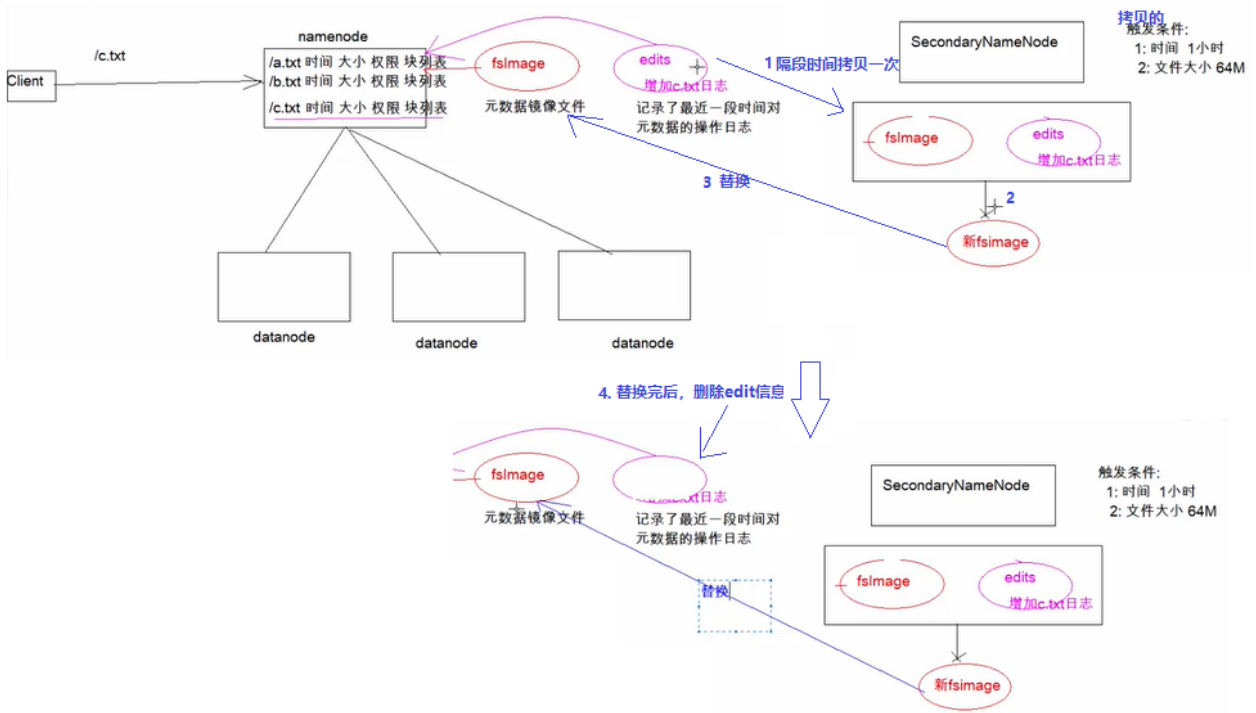
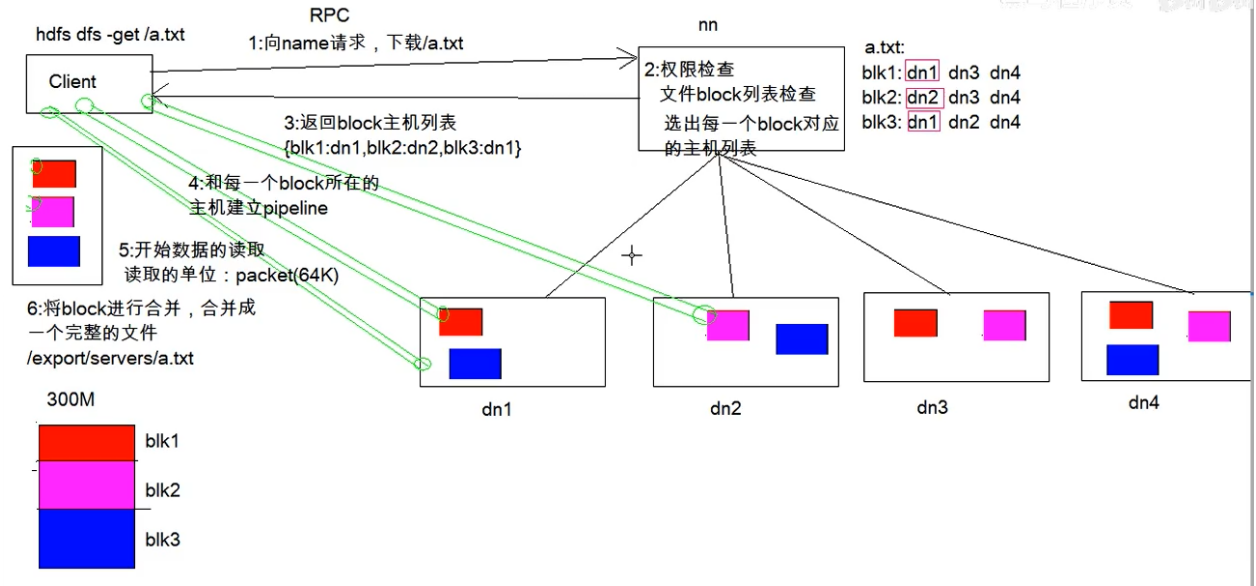
NameNode与ResourceManager单节点架构模型



Hadoop 2.x大数据平台架构

In []: 1





```
In [ ]: 1
```

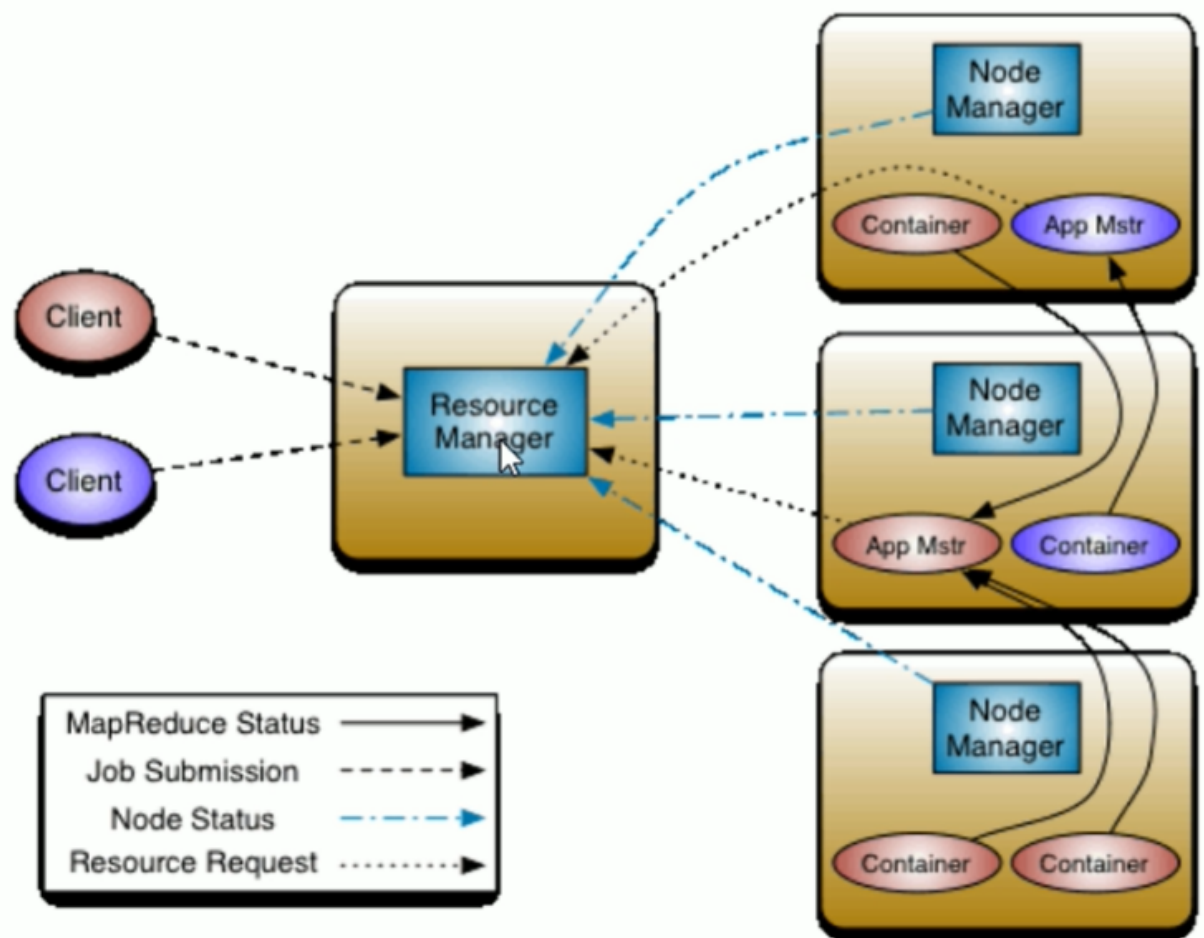
```
In [ ]: 1
```

yarn核心出发点是为了分离资源管理与作业监控，实现分离的做法是拥有一个全局的资源管理（ResourceManager，RM），以及每个应用程序对应一个的应用管理器（ApplicationMaster，AM）

总结一句话就是说：yarn主要就是为了调度资源，管理任务等

其调度分为两个层级来说：

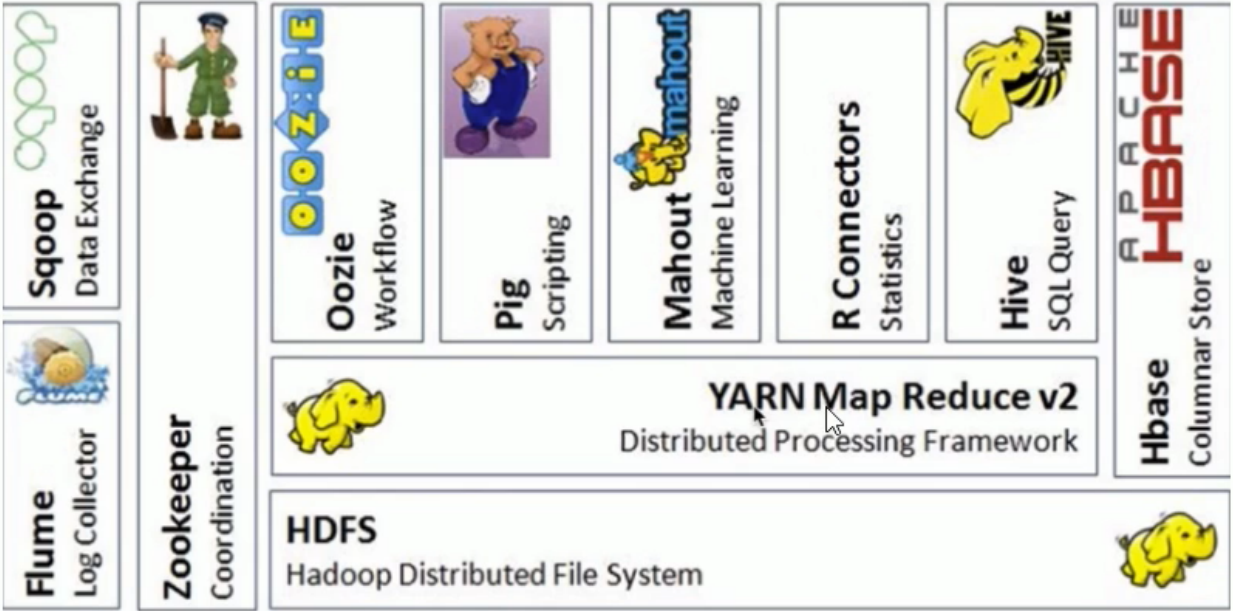
- 一级调度管理：
计算资源管理(CPU,内存，网络IO，磁盘)
- 二级调度管理：
任务内部的计算模型管理 (AppMaster的任务精细化管理)





hadoop生态系统

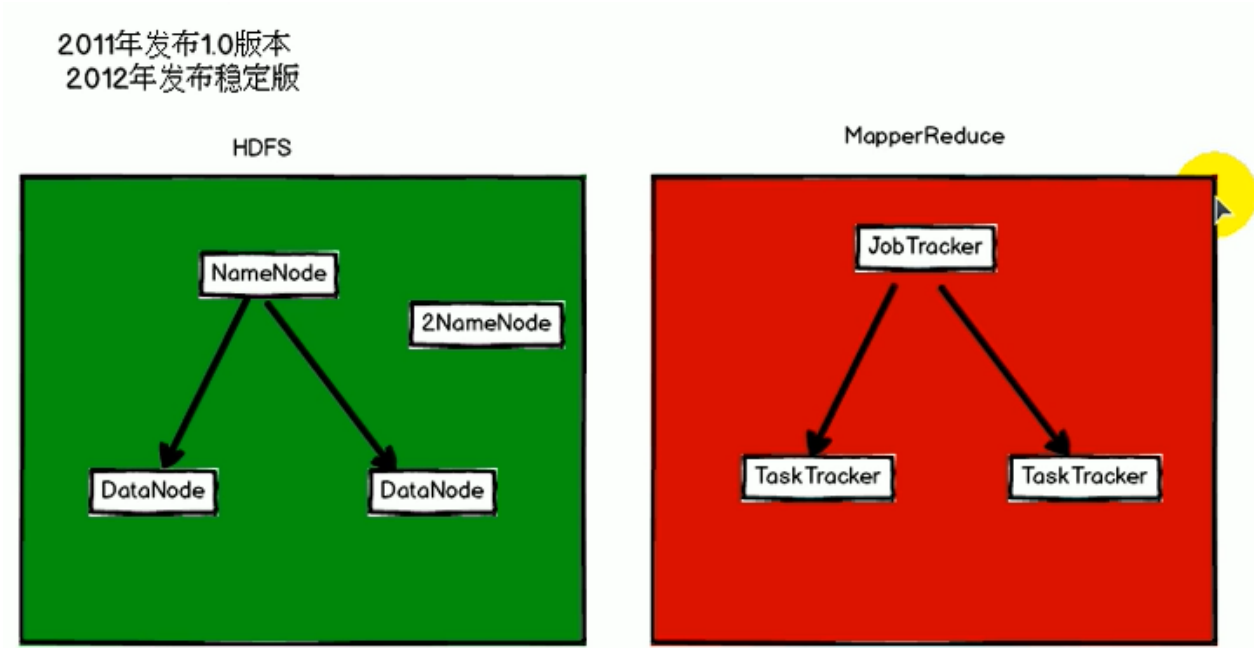
- 广义的Hadoop：指的是Hadoop生态系统，Hadoop生态系统是一个很庞大的概念，hadoop是其中最重要最基础的一个部分，生态系统中每一子系统只解决某一个特定的问题域（甚至可能更窄），不搞统一型的全能系统，而是小而精的多个小系统；



In []:

1

<https://www.bilibili.com/video/BV174411X7Pk?from=search&seid=16140933637673040363>
(<https://www.bilibili.com/video/BV174411X7Pk?from=search&seid=16140933637673040363>)

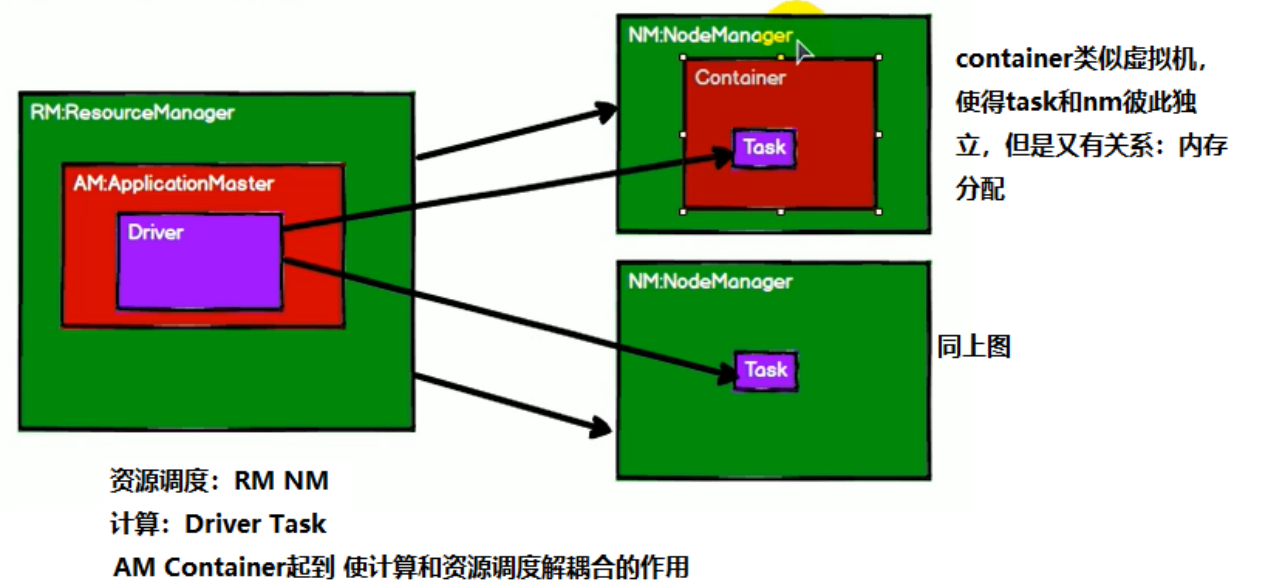


hadoop 1.0版本

MR的缺点：

- mr基于数据集的计算，所以面向数据
- 1. 基本运算规则从存储介质中获取（采集）数据，然后进行计算，最后将结果存储到介质中，所以主要应用于一次性计算，不适合于数据挖掘和机器学习这样的迭代计算和图形挖掘计算。
 - 2. MR基于文件存储介质的操作，所以性能非常的慢
 - 3. MR和hadoop紧密耦合在一起，无法动态替换

2013年发布2.X版本 (Yarn)



In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1