

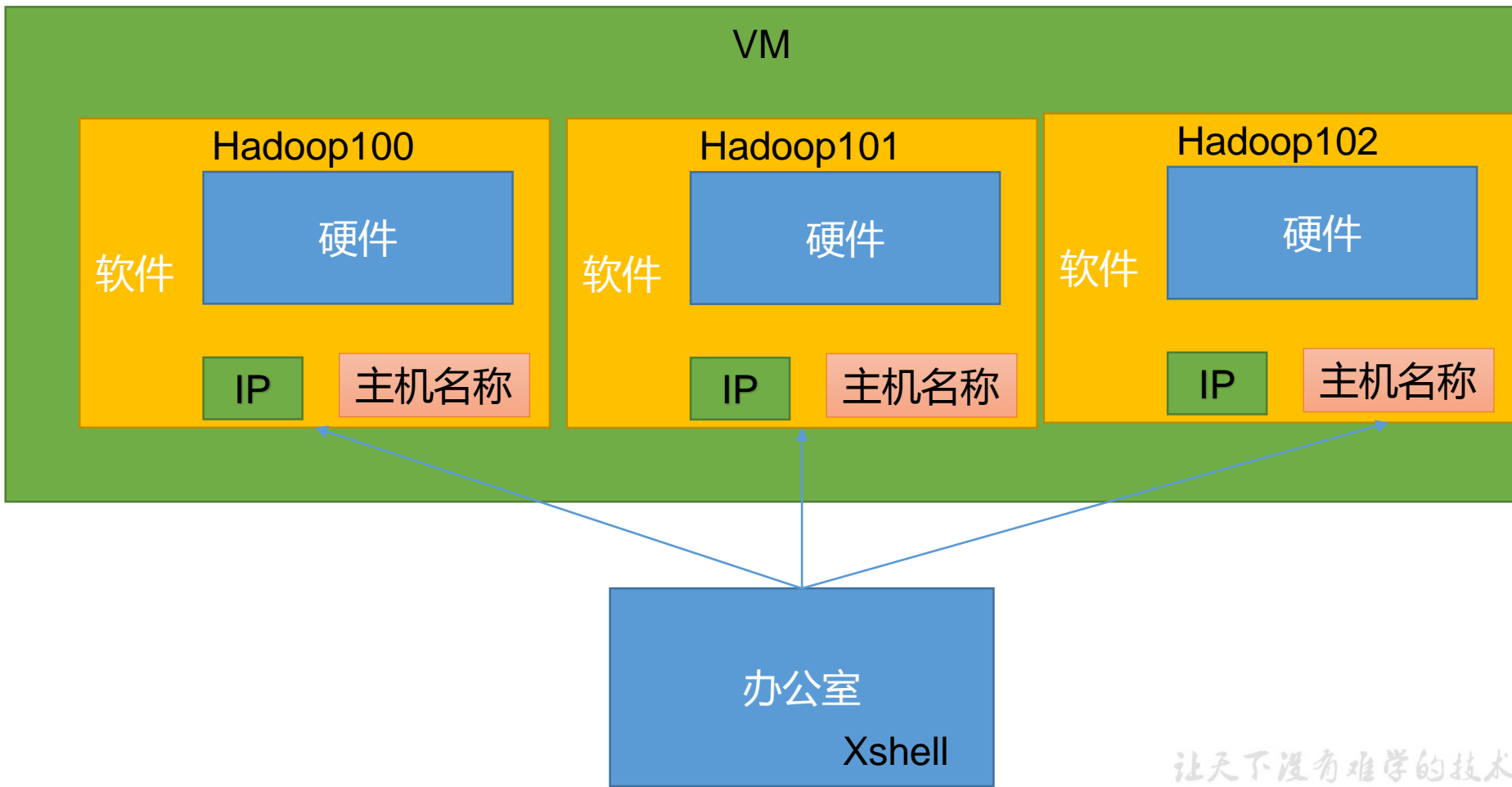
192.168.10.105 hadoop100

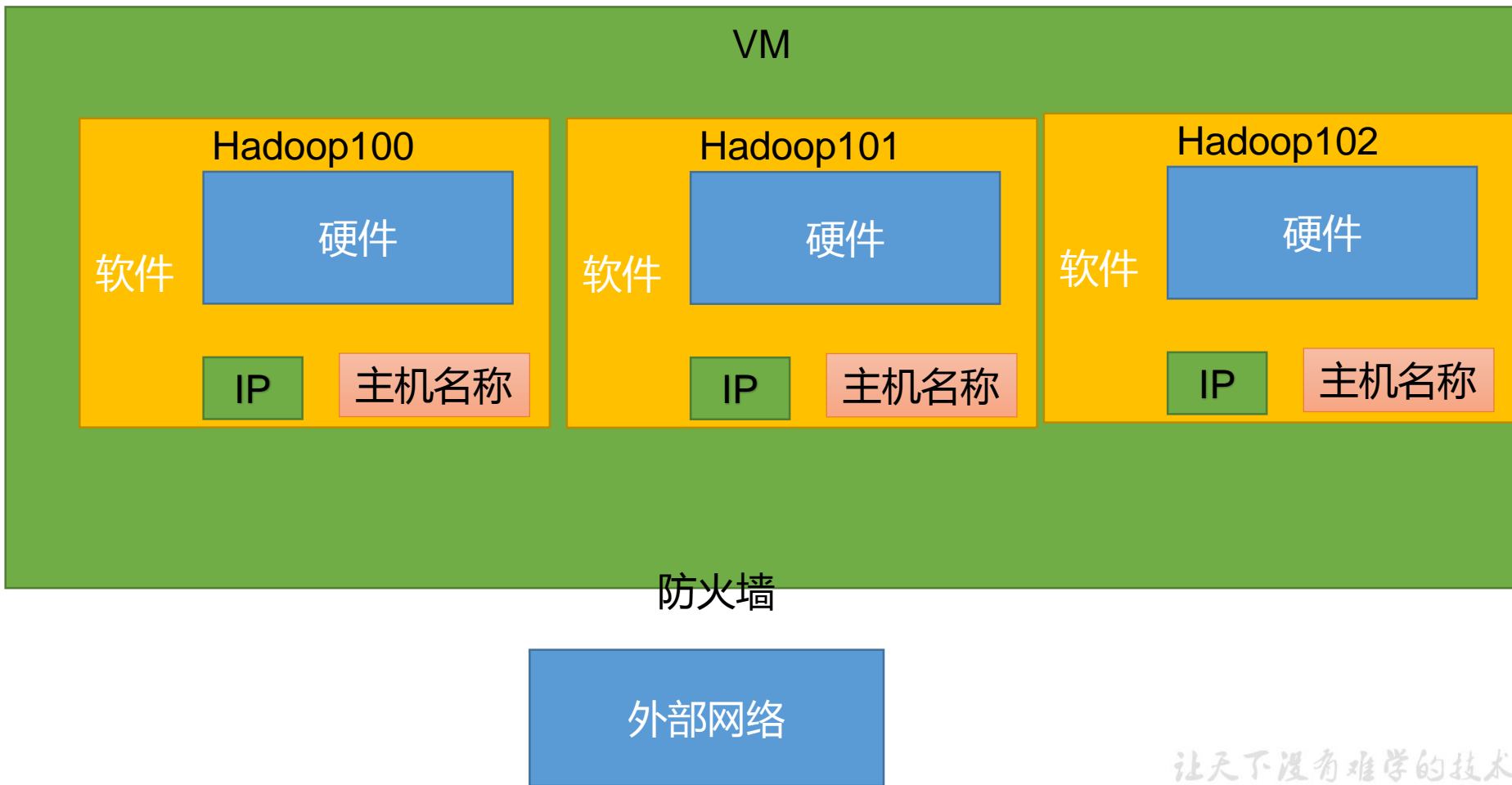
hadoop100

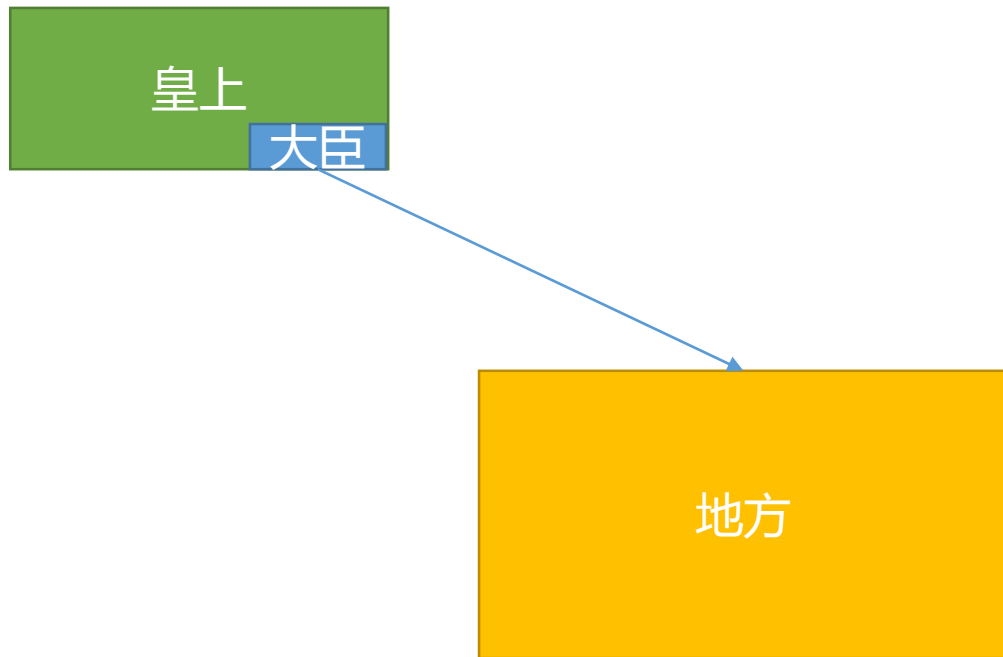
hadoop100

hadoop100

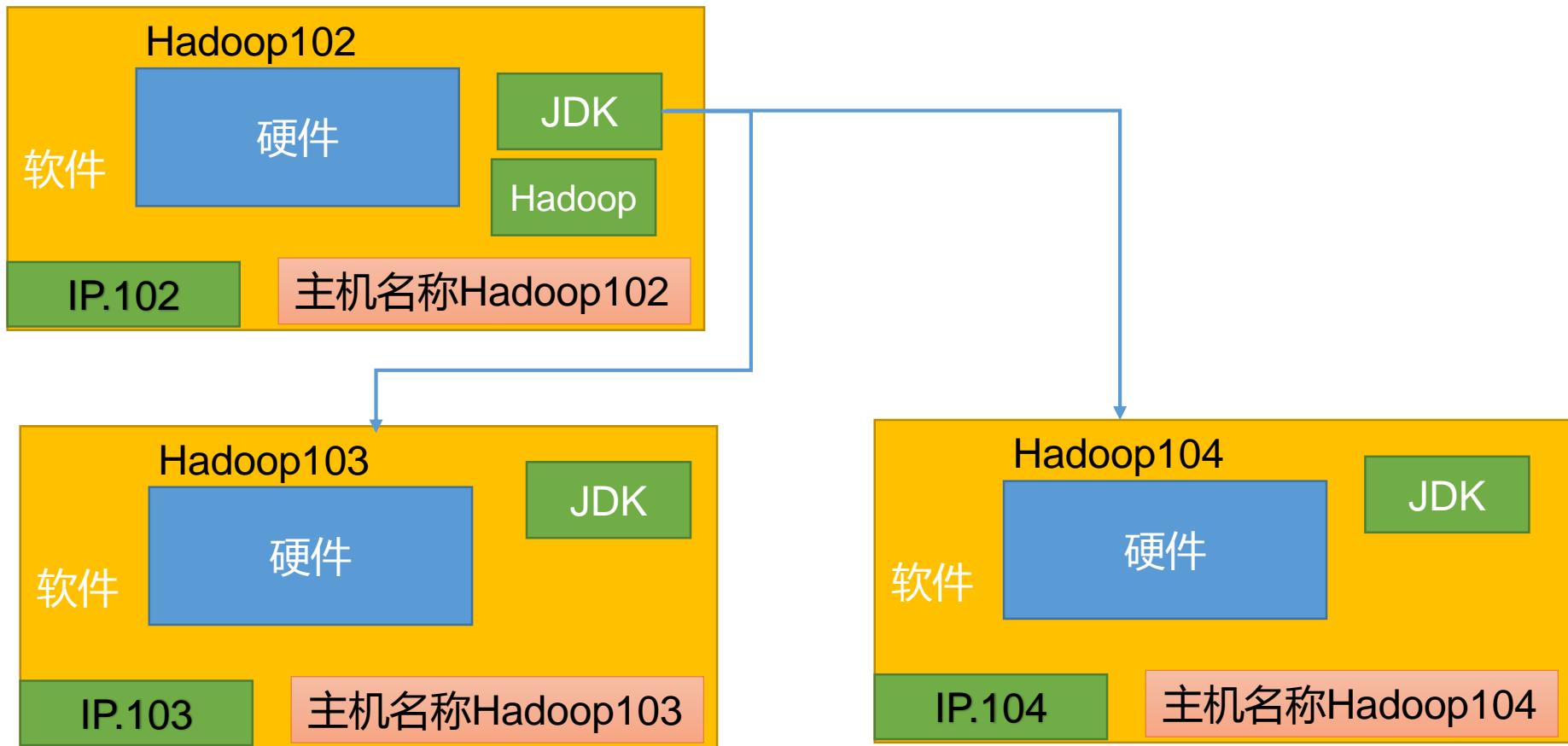
hadoop100













本地
hadoop100

数据存储在linux本地
测试偶尔用一下

伪分布
hadoop101

数据存储在HDFS
公司中比较差钱 2m 16g 1台 =》 2m -1T

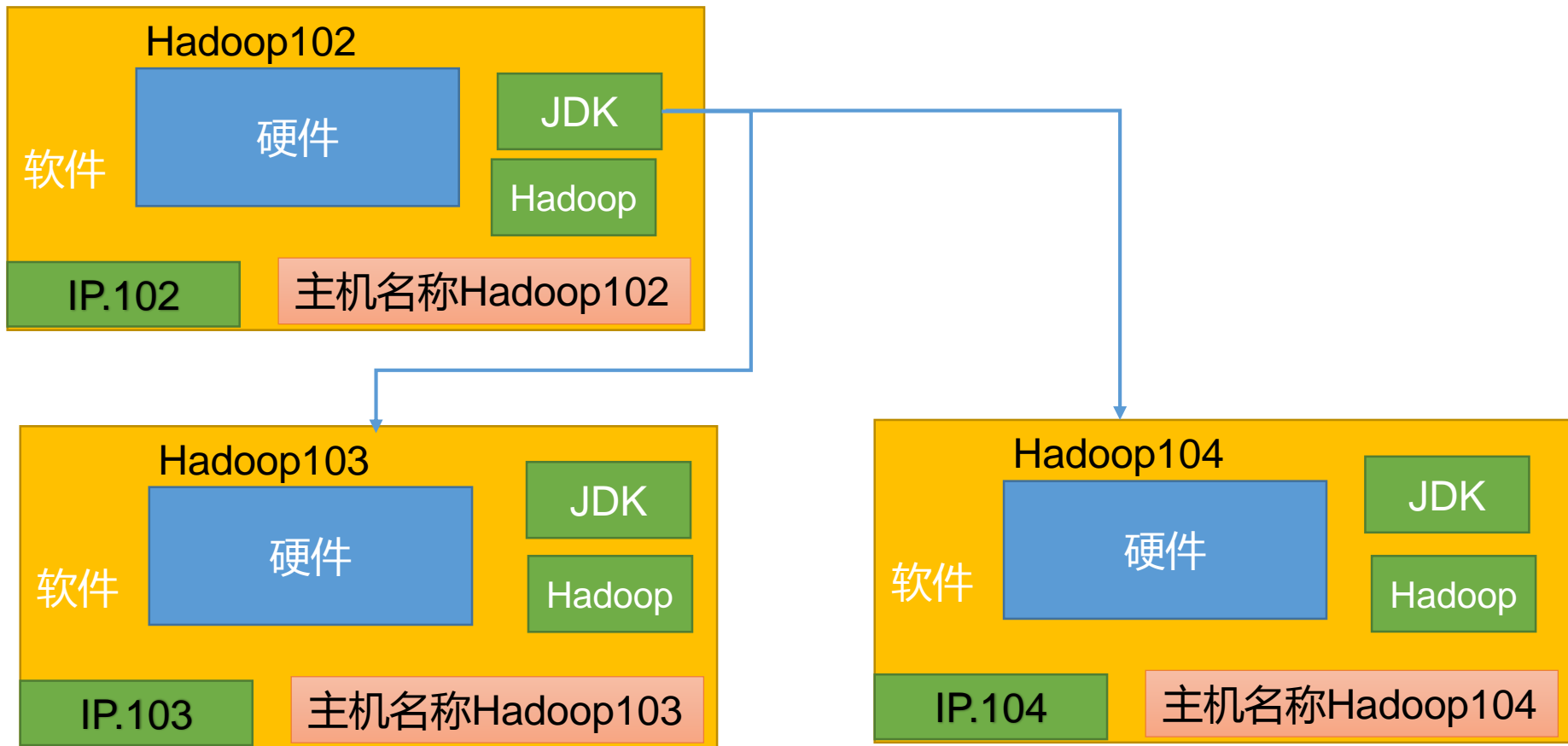
hadoop102

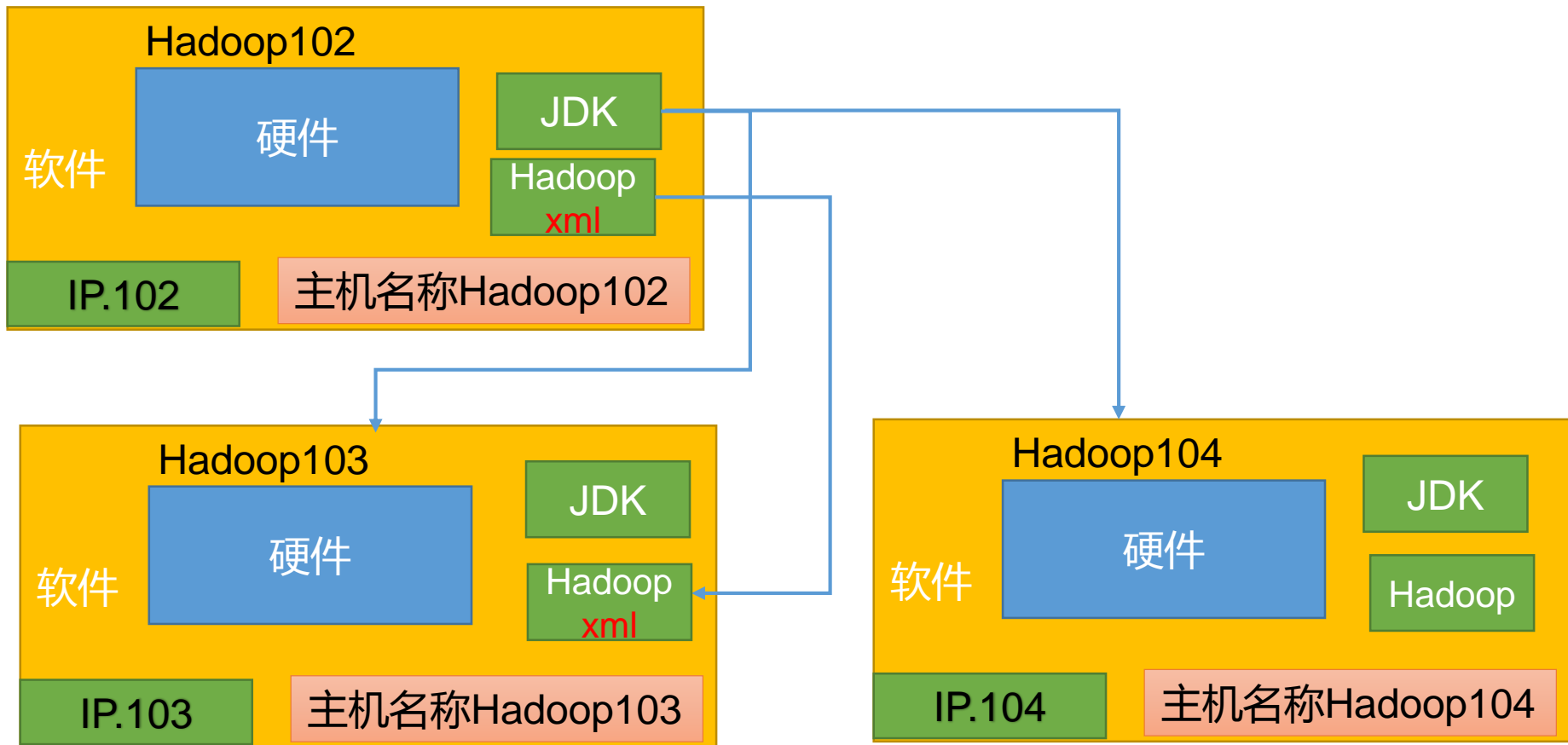
hadoop103

hadoop104

完全分布式 数据存储在HDFS/多台服务器工作

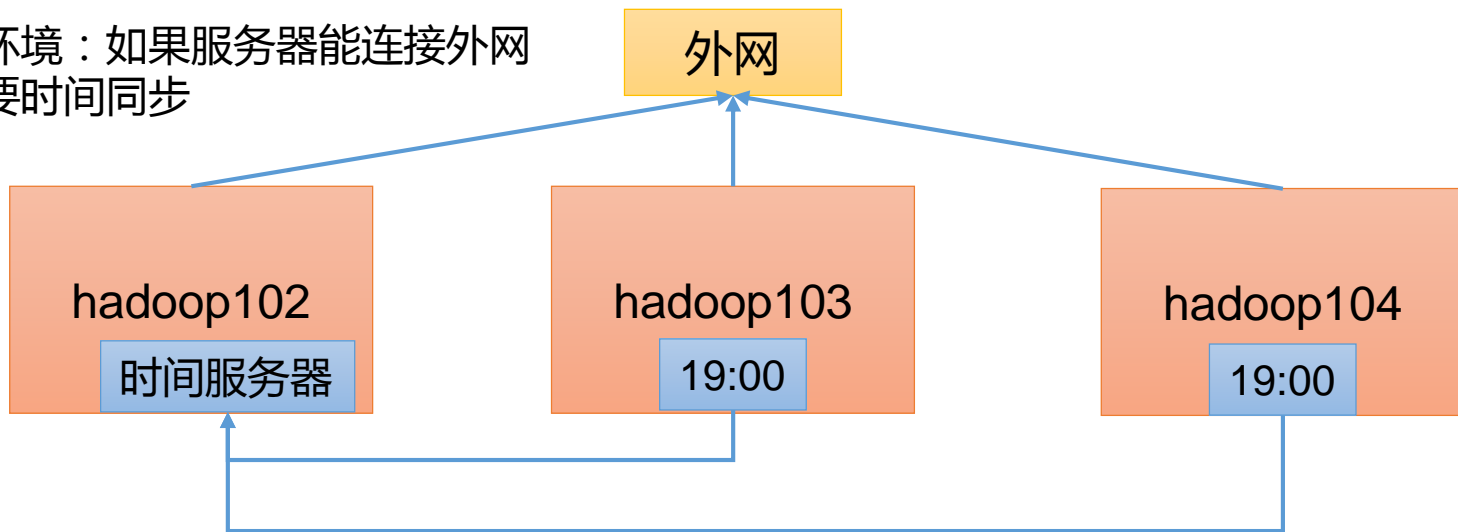
企业里面大量使用



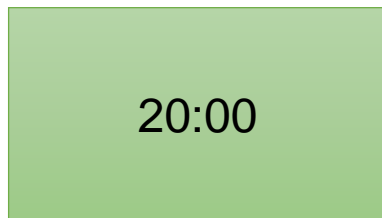




生产环境：如果服务器能连接外网
不需要时间同步

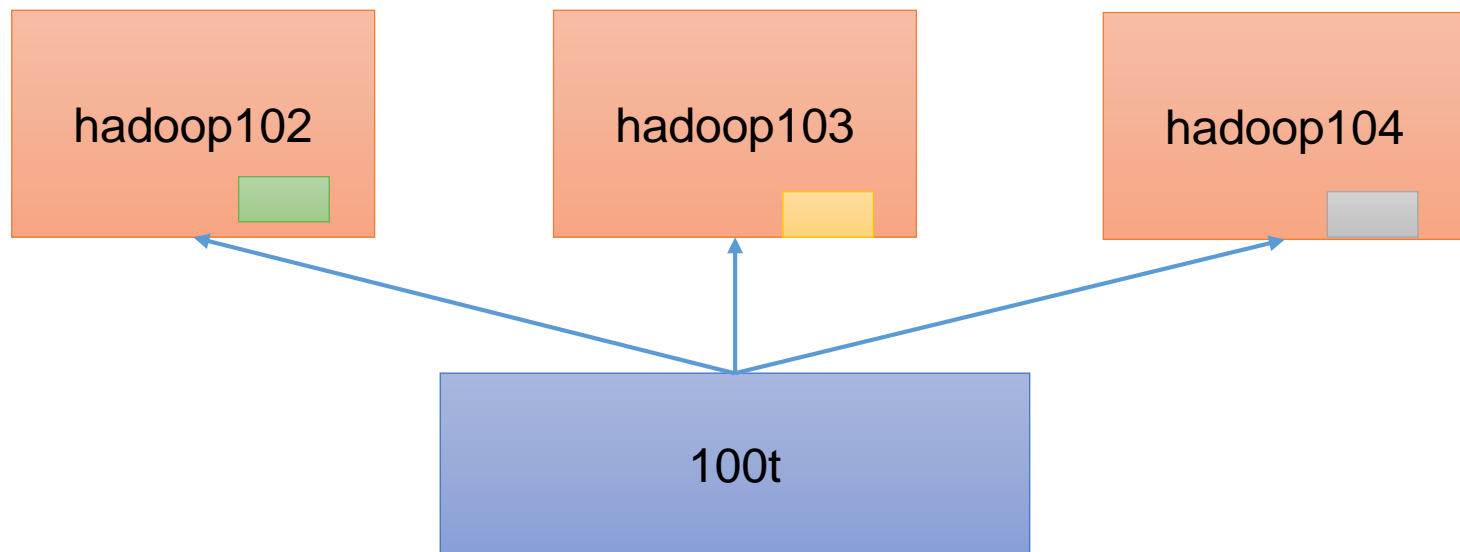


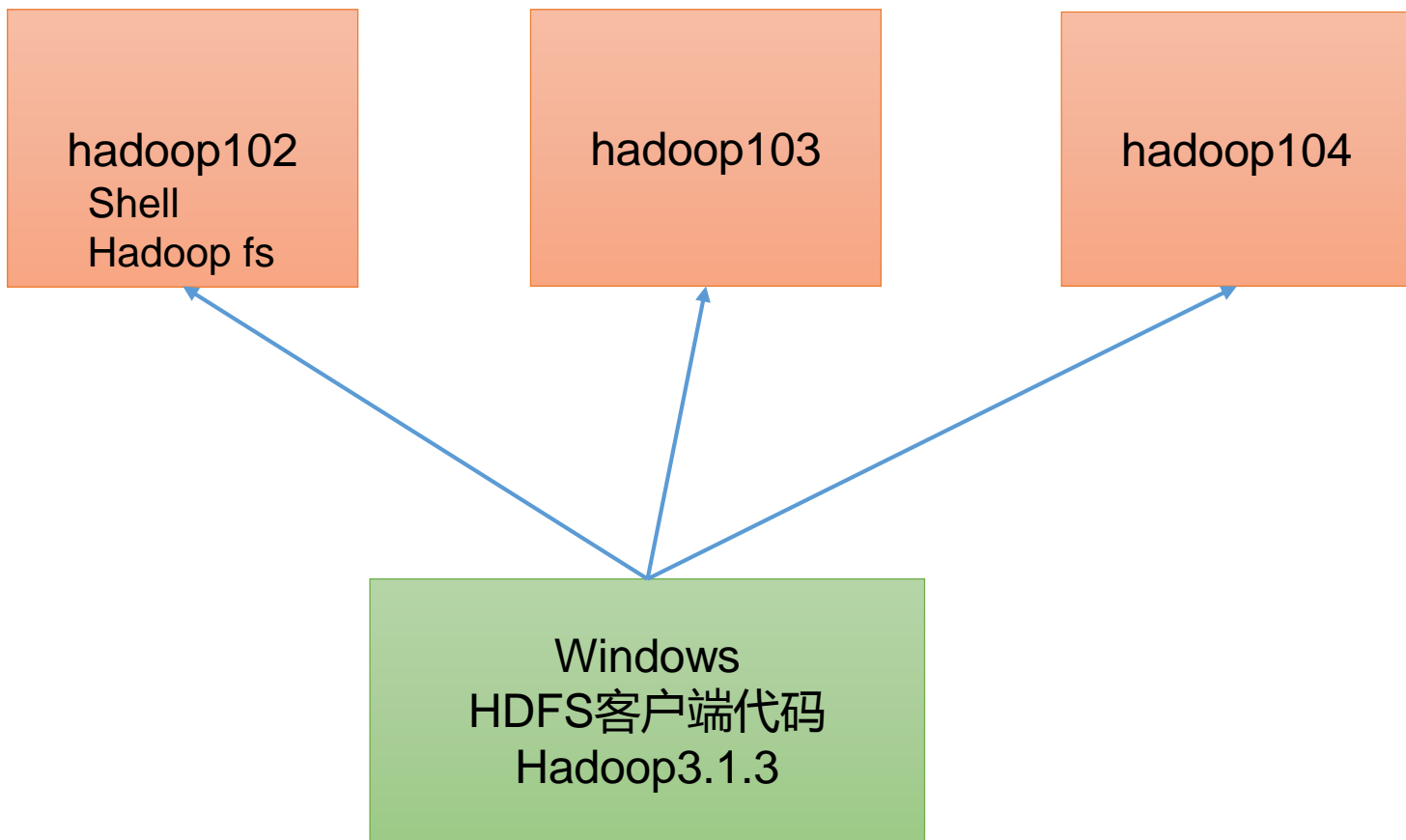
生产环境：如果服务器
能连接不了外网
需要时间同步





HDFS







hadoop102

NameNode

hadoop103

hadoop104

内存

内存：

好处：计算快
坏处：可靠性差

磁盘：

好处：可靠性高
坏处：计算慢

fsImage 存储数据

Edits 追加

内存 + 磁盘 =》效率低

fsImage 存储数据（如果是随机读写效率

$a = 10 \quad a + 10 \Rightarrow a = 20$ ）

Edits 追加 =>

$a + 10$

$a - 30$

$a * 20$



静态
数据



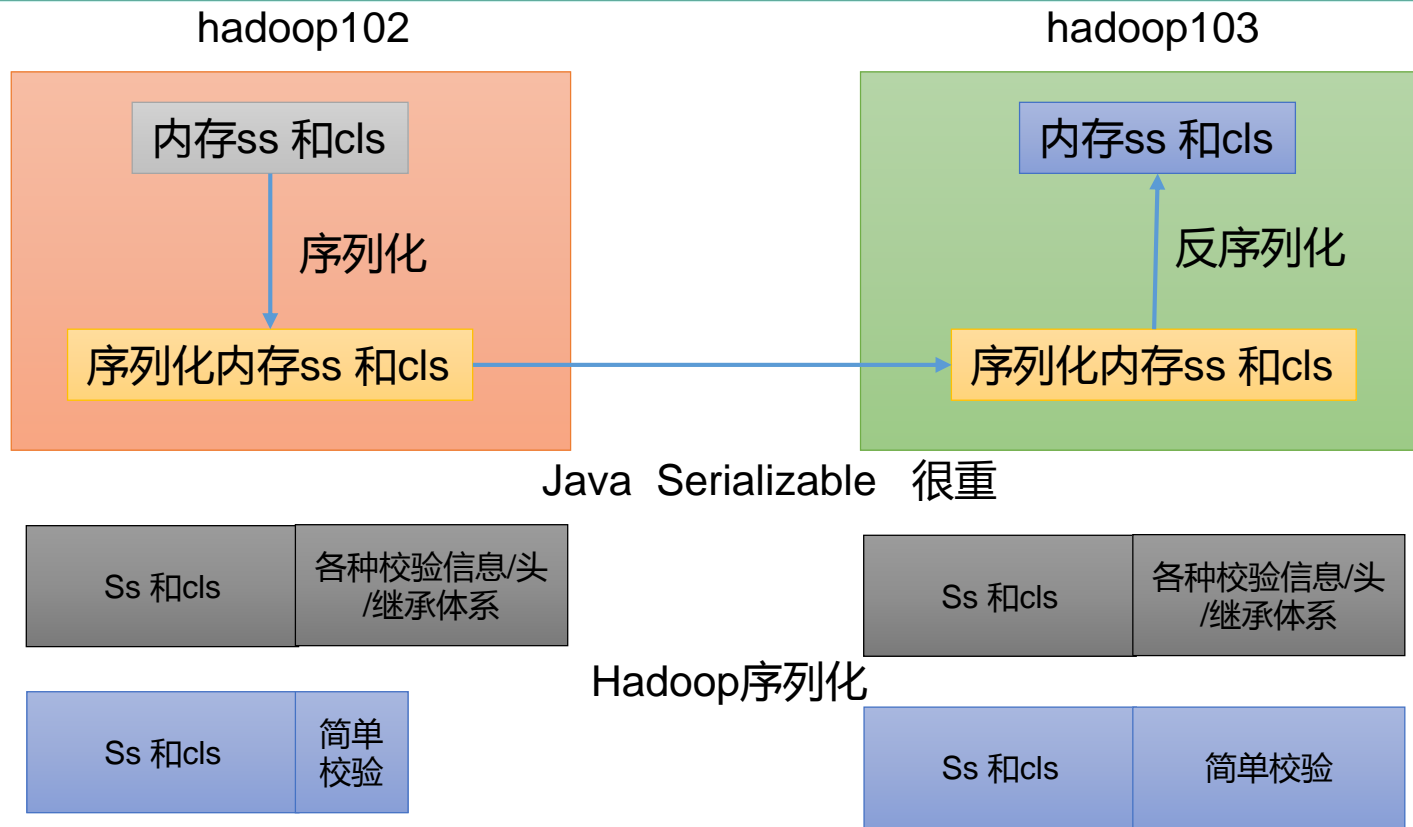
MapReduce: 自己处理业务相关代码+ 自身的默认代码

优点:

- 1、易于编程。 用户只关心，业务逻辑。 实现框架的接口。
- 2、良好扩展性: 可以动态增加服务器，解决计算资源不够问题
- 3、高容错性。任何一台机器挂掉，可以将任务转移到其他节点。
- 4、适合海量数据计算 (TB/PB) 几千台服务器共同计算。

缺点:

- 1、不擅长实时计算。 Mysql
- 2、不擅长流式计算。 Sparkstreaming flink
- 3、不擅长DAG有向无环图计算。 spark



紧凑：存储空间少
快速：传输速度快
互操作性：



序列化

sumFlow

downFlow

upFlow



反序列化

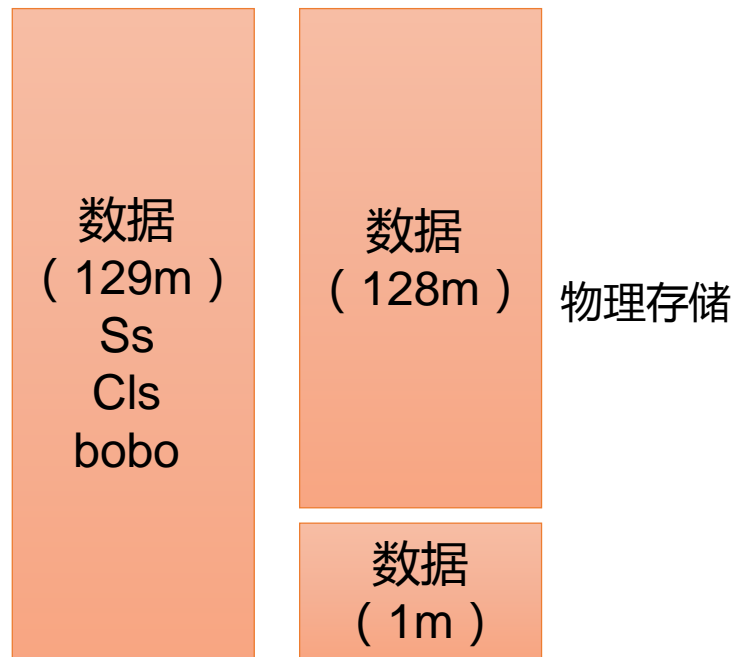
upFlow

downFlow

sumFlow

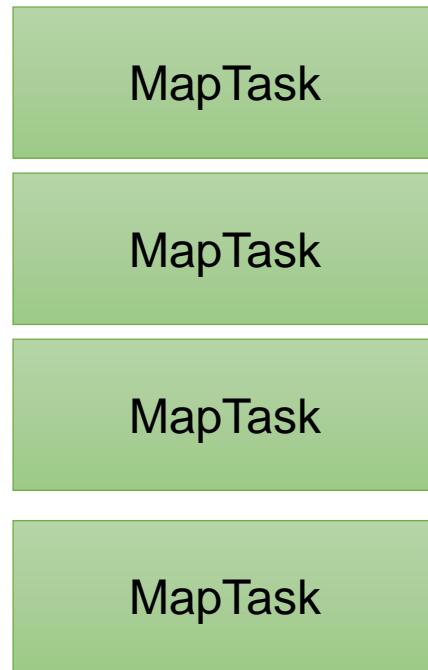
Map

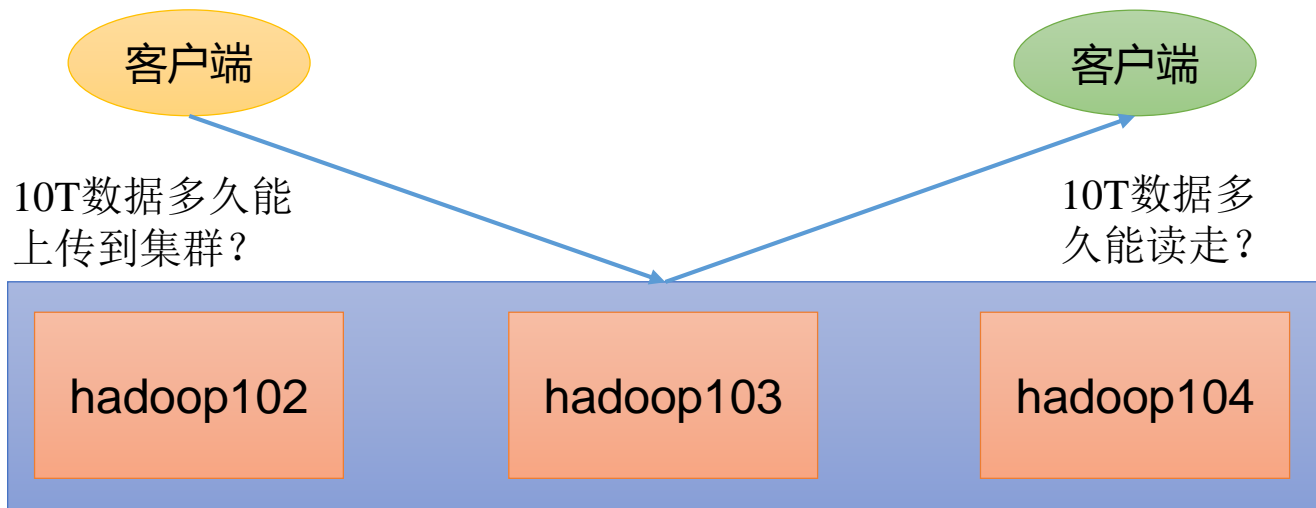
<key,val, key, val>

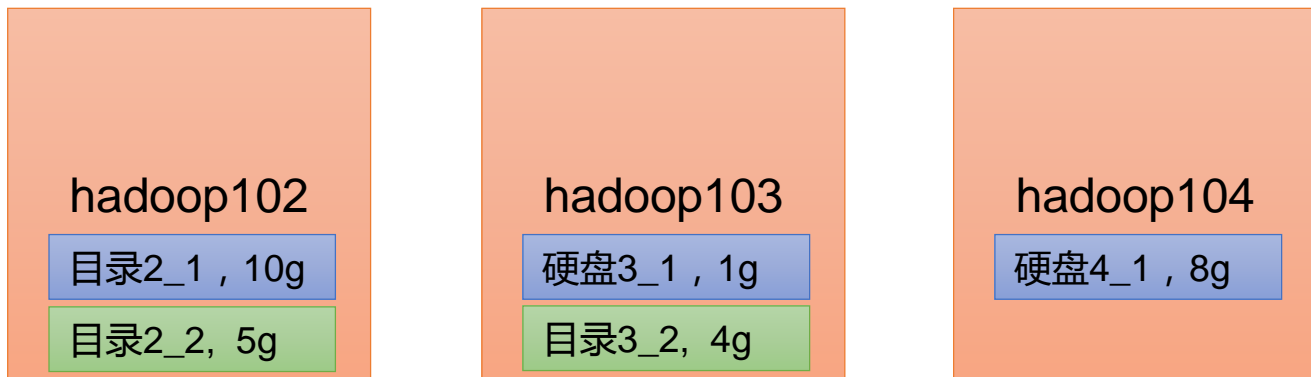


0-128索引 是一片
128-129索引 是一片
逻辑存储

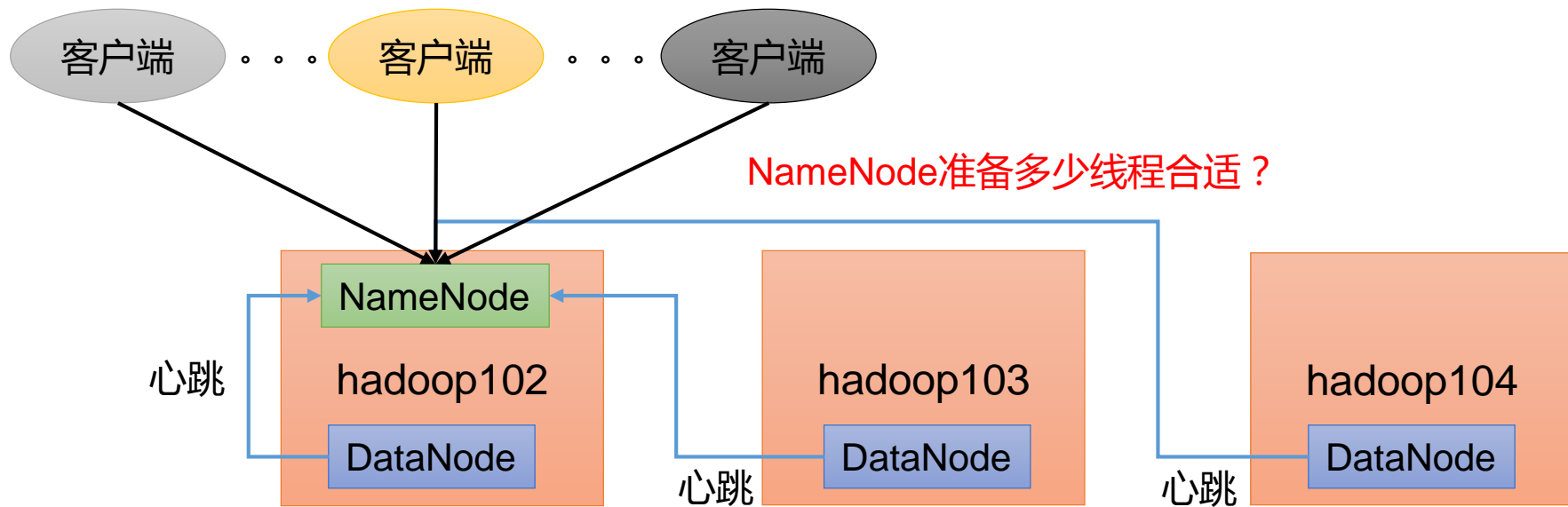
MapTask个数，决定了我的并行度







每个目录存储的数据不一样





hadoop102

硬盘2_1 , 10g

硬盘2_2, 0g

hadoop103

硬盘3_1 , 1g

hadoop104

硬盘4_1 , 8g



hadoop102

硬盘2_1, 10g

hadoop103

硬盘3_1, 1g

hadoop104

硬盘4_1, 8g



hadoop102

mrAPPM
aster

hadoop103

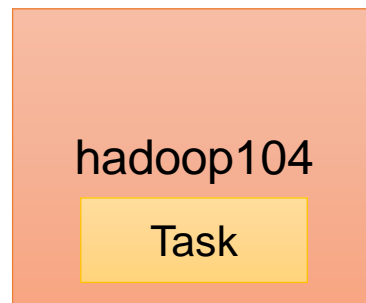
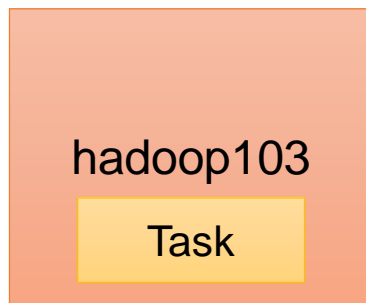
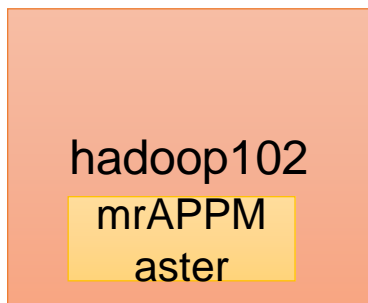
Task

hadoop104

Task

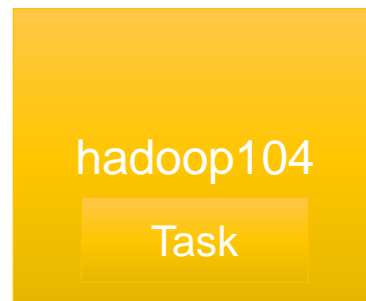
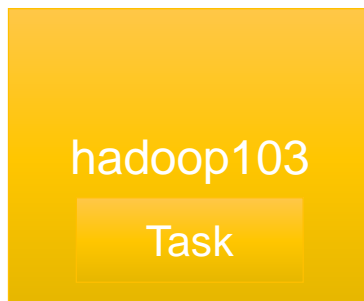
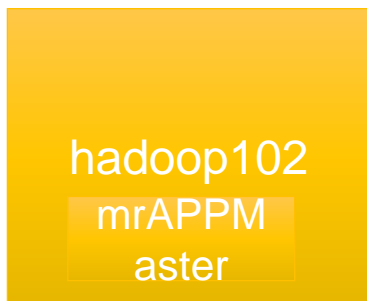


拍摄快照
Ss 18岁



岁月杀猪刀

Ss 50岁





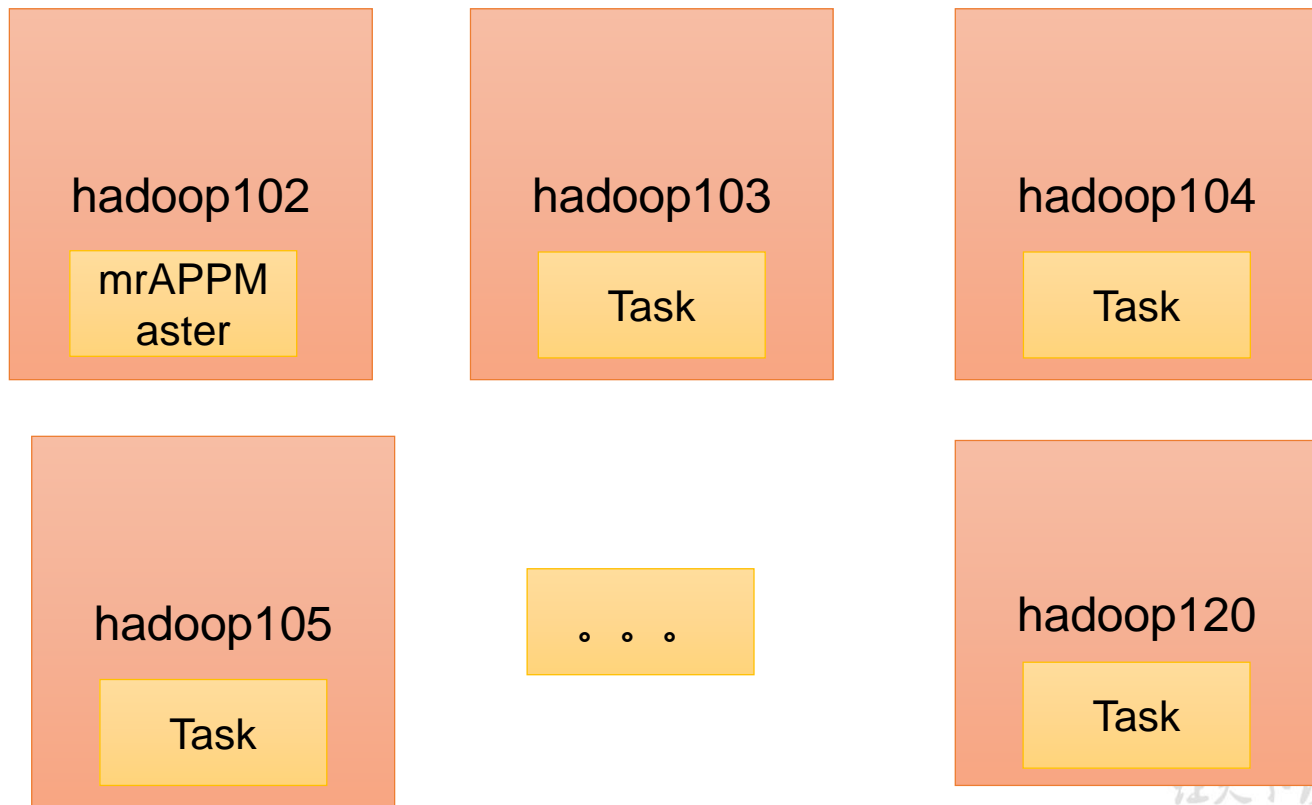
多个队列

11.11 6.18

10个队列

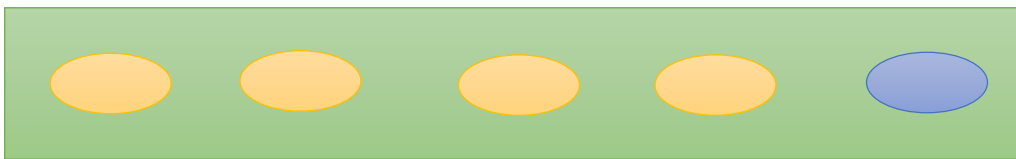
1 / 2 3

降级





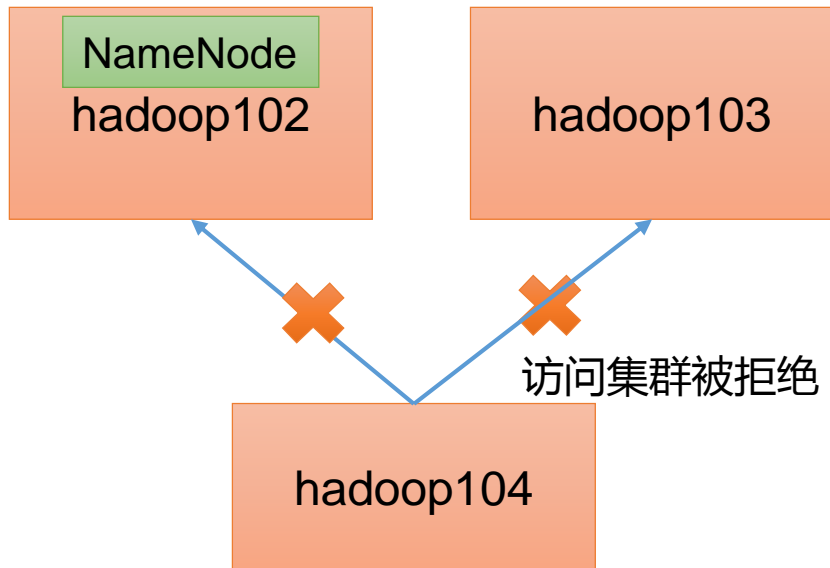
有5个任务



任务1

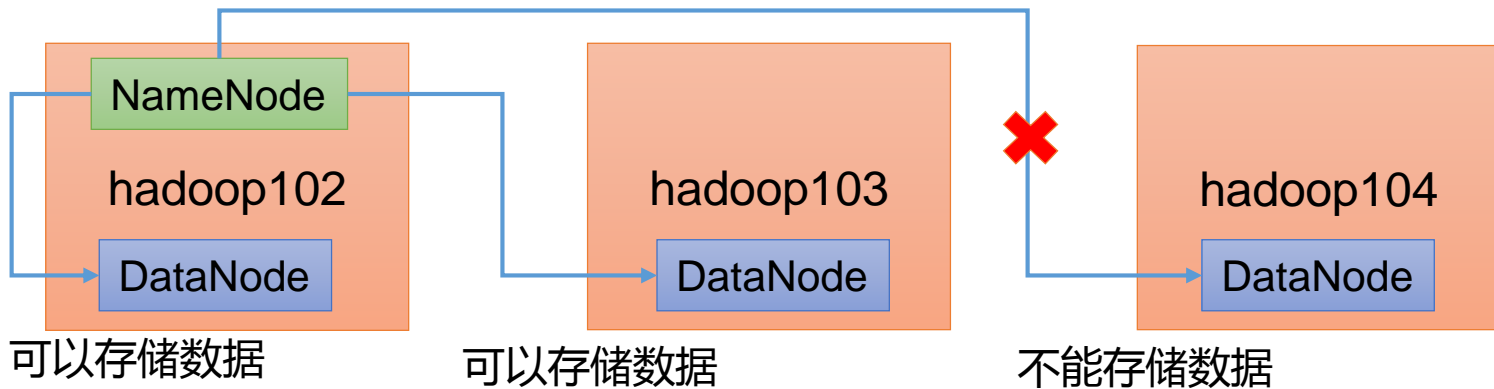


白名单：hadoop102、hadoop103



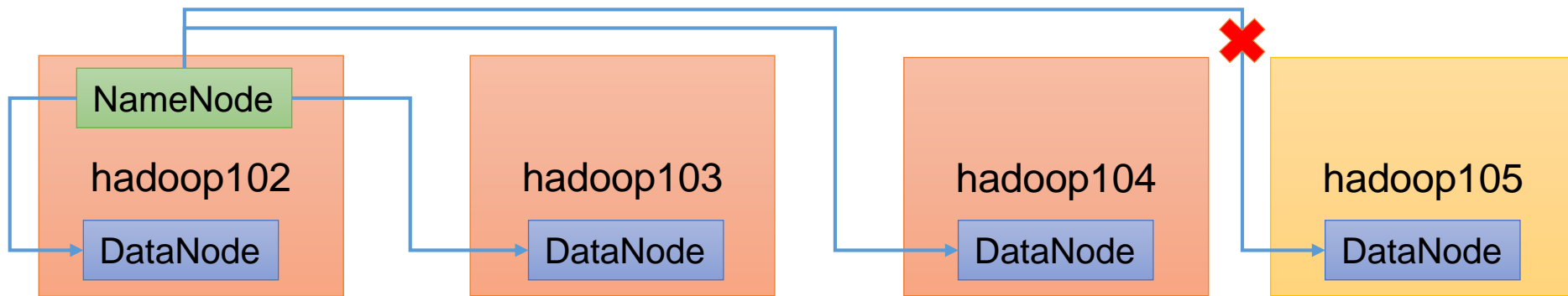


白名单：hadoop102、hadoop103





黑名单：hadoop105





hadoop102

硬盘2_1, 10g

hadoop103

硬盘3_1, 1g

hadoop104

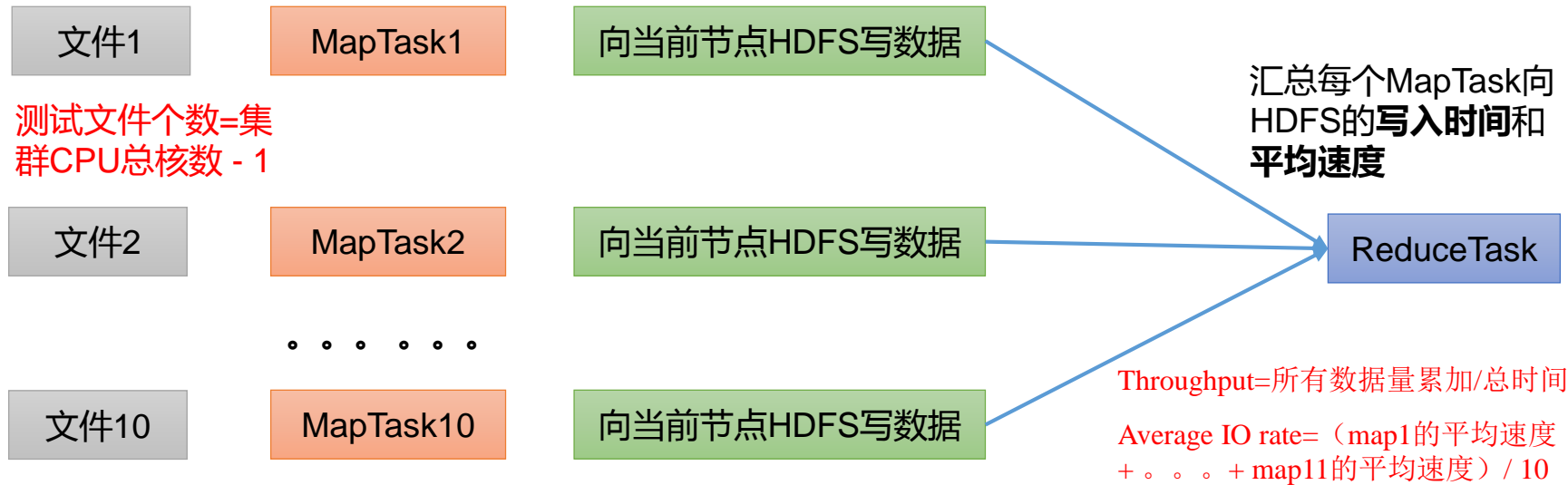
硬盘4_1, 8g

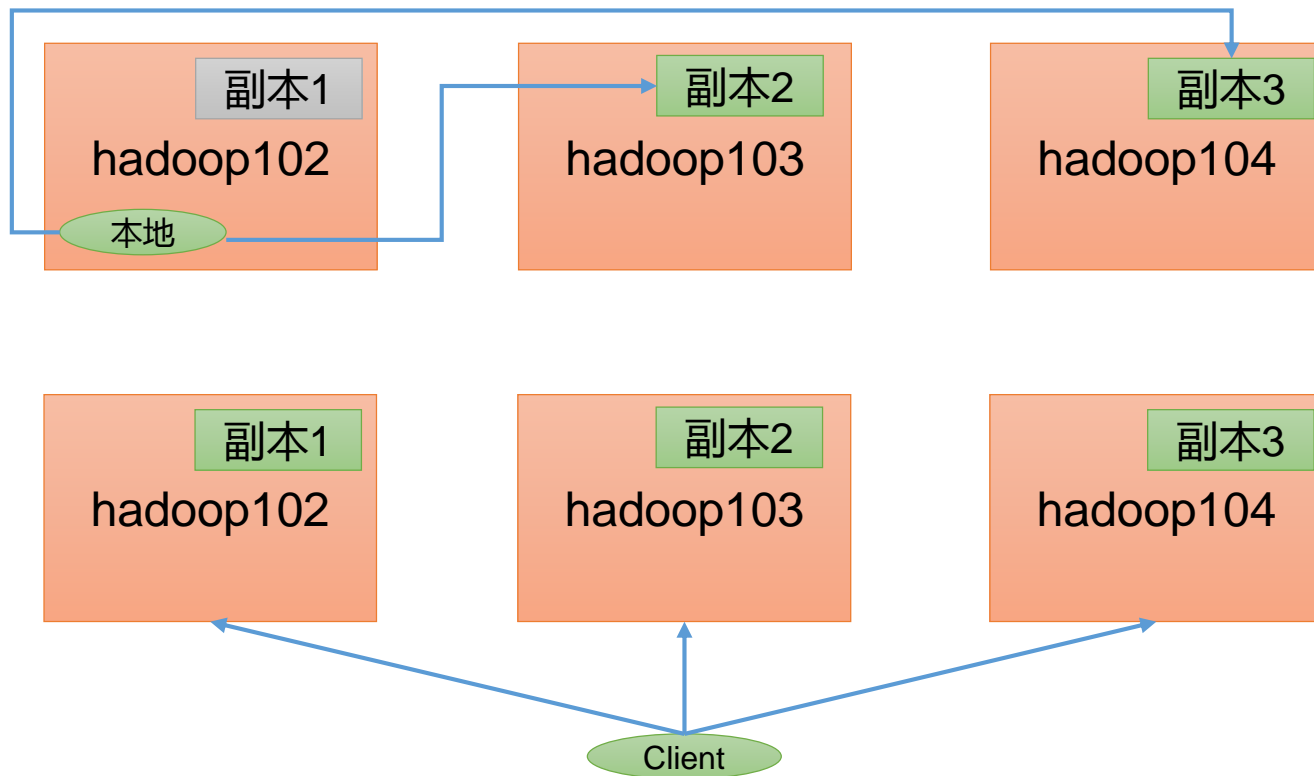
hadoop105

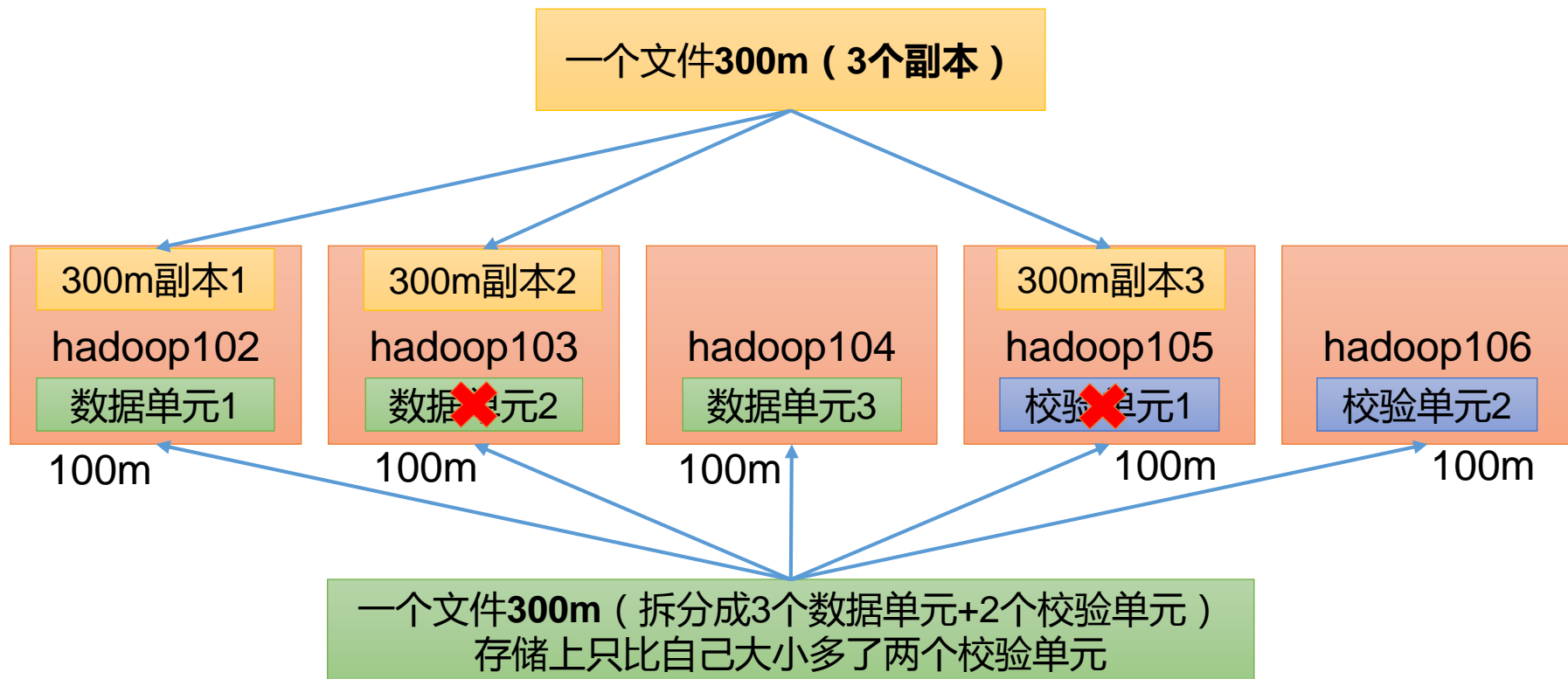
硬盘5_1, 0g



记录每个Map的**写时间**和**平均速度**



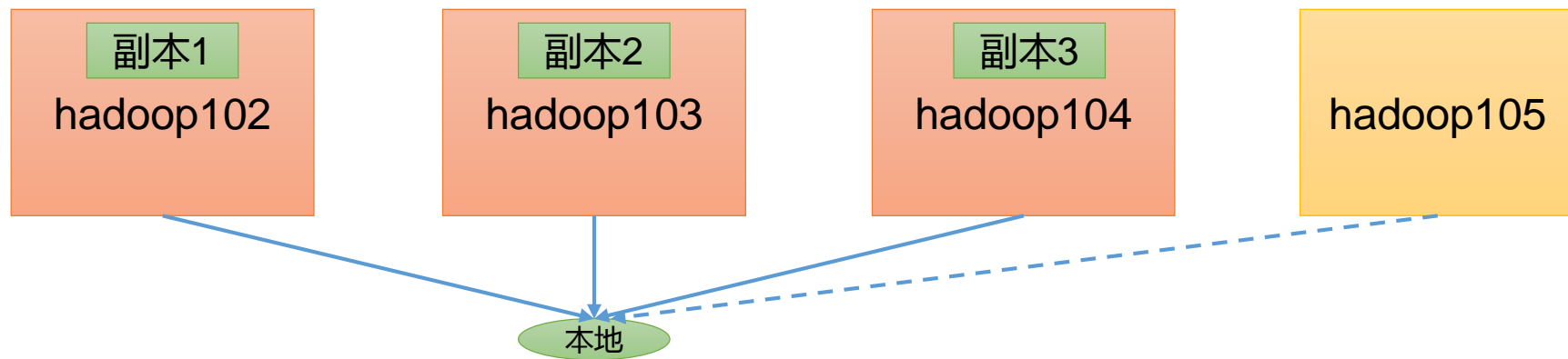






纠删码策略是给具体一个路径设置。

默认只开启对RS-6-3-1024k策略的支持，如要使用别的策略需要提前启用。





SS

正在看的、经常看的、不经常看的、永久保存的视频，怎么存？

hadoop102

硬盘5_1，内存镜像

hadoop103

硬盘2_1，固态

hadoop104

硬盘3_1，机械

hadoop105

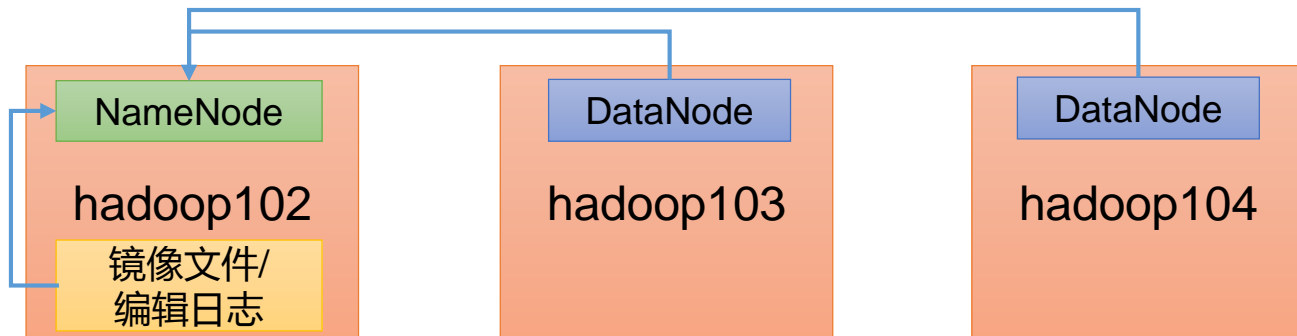
硬盘4_1，破旧



1) 安全模式：文件系统只接受读数据请求，而不接受删除、修改等变更请求

2) 进入安全模式场景

- NameNode在加载镜像文件和编辑日志期间处于安全模式；
- NameNode再接收DataNode注册时，处于安全模式





100个1k文件块和100个128m的文件块，占用NN内存大小一样

