

八 HDFS HA 高可用

8.1 HA 概述

- 1) 所谓 HA (high available), 即高可用 (7*24 小时不中断服务)。
- 2) 实现高可用最关键的策略是消除单点故障。HA 严格来说应该分成各个组件的 HA 机制: HDFS 的 HA 和 YARN 的 HA。

- 3) Hadoop2.0 之前, 在 HDFS 集群中 NameNode 存在单点故障 (SPOF)。

- 4) NameNode 主要在以下两个方面影响 HDFS 集群

NameNode 机器发生意外, 如宕机, 集群将无法使用, 直到管理员重启

NameNode 机器需要升级, 包括软件、硬件升级, 此时集群也将无法使用

HDFS HA 功能通过配置 Active/Standby 两个 nameNodes 实现在集群中对 NameNode 的热备来解决上述问题。如果出现故障, 如机器崩溃或机器需要升级维护, 这时可通过此种方式将 NameNode 很快的切换到另外一台机器。

8.2 HDFS-HA 工作机制

- 1) 通过双 namenode 消除单点故障

8.2.1 HDFS-HA 工作要点

- 1) 元数据管理方式需要改变:

内存中各自保存一份元数据;

Edits 日志只有 Active 状态的 namenode 节点可以做写操作;

两个 namenode 都可以读取 edits;

共享的 edits 放在一个共享存储中管理 (qjournal 和 NFS 两个主流实现);

- 2) 需要一个状态管理功能模块

实现了一个 zkfailover, 常驻在每一个 namenode 所在的节点, 每一个 zkfailover 负责监控自己所在 namenode 节点, 利用 zk 进行状态标识, 当需要进行状态切换时, 由 zkfailover 来负责切换, 切换时需要防止 brain split 现象的发生。

- 3) 必须保证两个 NameNode 之间能够 ssh 无密码登录。

- 4) 隔离 (Fence), 即同一时刻仅仅有一个 NameNode 对外提供服务

8.2.2 HDFS-HA 自动故障转移工作机制

前面学习了使用命令 `hdfs haadmin -failover` 手动进行故障转移，在该模式下，即使现役 NameNode 已经失效，系统也不会自动从现役 NameNode 转移到待机 NameNode，下面学习如何配置部署 HA 自动进行故障转移。自动故障转移为 HDFS 部署增加了两个新组件：ZooKeeper 和 ZKFailoverController (ZKFC) 进程。ZooKeeper 是维护少量协调数据，通知客户端这些数据的改变和监视客户端故障的高可用服务。HA 的自动故障转移依赖于 ZooKeeper 的以下功能：

1) 故障检测：集群中的每个 NameNode 在 ZooKeeper 中维护了一个持久会话，如果机器崩溃，ZooKeeper 中的会话将终止，ZooKeeper 通知另一个 NameNode 需要触发故障转移。

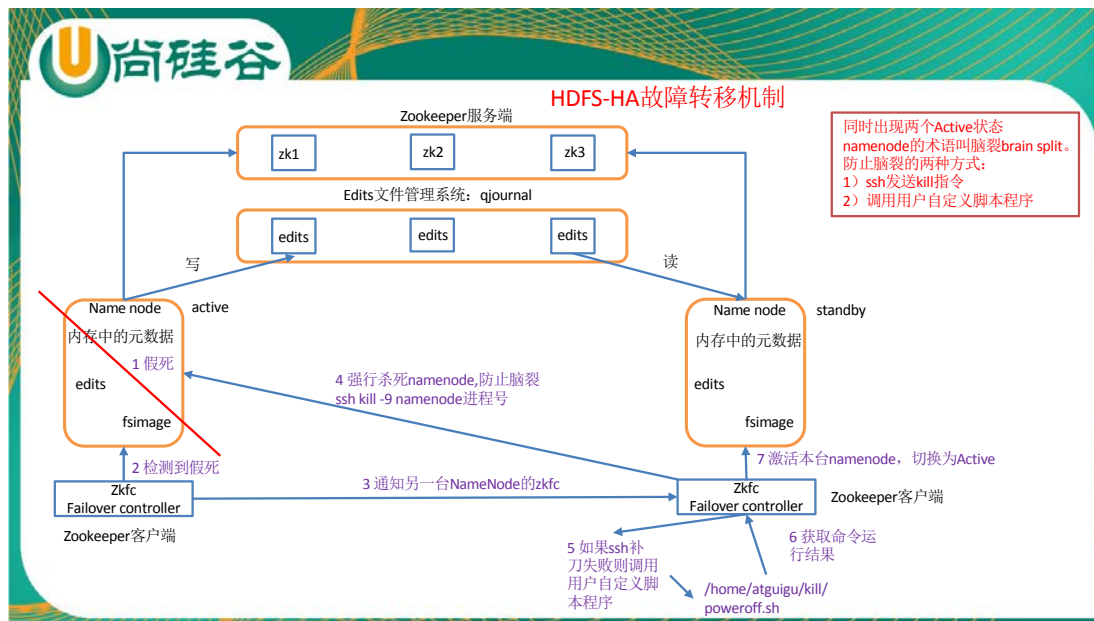
2) 现役 NameNode 选择：ZooKeeper 提供了一个简单的机制用于唯一的选择一个节点为 active 状态。如果目前现役 NameNode 崩溃，另一个节点可能从 ZooKeeper 获得特殊的排外锁以表明它应该成为现役 NameNode。

ZKFC 是自动故障转移中的另一个新组件，是 ZooKeeper 的客户端，也监视和管理 NameNode 的状态。每个运行 NameNode 的主机也运行了一个 ZKFC 进程，ZKFC 负责：

1) 健康监测：ZKFC 使用一个健康检查命令定期地 ping 与之在相同主机的 NameNode，只要该 NameNode 及时地回复健康状态，ZKFC 认为该节点是健康的。如果该节点崩溃，冻结或进入不健康状态，健康监测器标识该节点为非健康的。

2) ZooKeeper 会话管理：当本地 NameNode 是健康的，ZKFC 保持一个在 ZooKeeper 中打开的会话。如果本地 NameNode 处于 active 状态，ZKFC 也保持一个特殊的 znode 锁，该锁使用了 ZooKeeper 对短暂节点的支持，如果会话终止，锁节点将自动删除。

3) 基于 ZooKeeper 的选择：如果本地 NameNode 是健康的，且 ZKFC 发现没有其它的节点当前持有 znode 锁，它将为自己获取该锁。如果成功，则它已经赢得了选择，并负责运行故障转移进程以使它的本地 NameNode 为 active。故障转移进程与前面描述的手动故障转移相似，首先如果必要保护之前的现役 NameNode，然后本地 NameNode 转换为 active 状态。



8.3 HDFS-HA 集群配置

8.3.1 环境准备

- 1) 修改 IP
- 2) 修改主机名及主机名和 IP 地址的映射
- 3) 关闭防火墙
- 4) ssh 免密登录
- 5) 安装 JDK，配置环境变量等

8.3.2 规划集群

hadoop102	hadoop103	hadoop104
NameNode	NameNode	
JournalNode	JournalNode	JournalNode
DataNode	DataNode	DataNode
ZK	ZK	ZK
	ResourceManager	
NodeManager	NodeManager	NodeManager

8.3.3 配置 Zookeeper 集群

0) 集群规划

在 hadoop102、hadoop103 和 hadoop104 三个节点上部署 Zookeeper。

1) 解压安装

(1) 解压 zookeeper 安装包到/opt/module/目录下

```
[atguigu@hadoop102 software]$ tar -zxvf zookeeper-3.4.10.tar.gz -C /opt/module/
```

(2) 在/opt/module/zookeeper-3.4.10/这个目录下创建 zkData

```
mkdir -p zkData
```

(3) 重命名/opt/module/zookeeper-3.4.10/conf 这个目录下的 zoo_sample.cfg 为 zoo.cfg

```
mv zoo_sample.cfg zoo.cfg
```

2) 配置 zoo.cfg 文件

(1) 具体配置

```
dataDir=/opt/module/zookeeper-3.4.10/zkData
```

增加如下配置

```
#####cluster#####
```

```
server.2=hadoop102:2888:3888
```

```
server.3=hadoop103:2888:3888
```

```
server.4=hadoop104:2888:3888
```

(2) 配置参数解读

Server.A=B:C:D。

A 是一个数字，表示这个是第几号服务器；

B 是这个服务器的 ip 地址；

C 是这个服务器与集群中的 Leader 服务器交换信息的端口；

D 是万一集群中的 Leader 服务器挂了，需要一个端口来重新进行选举，选出一个新的 Leader，而这个端口就是用来执行选举时服务器相互通信的端口。

集群模式下配置一个文件 myid，这个文件在 dataDir 目录下，这个文件里面有一个数据就是 A 的值，Zookeeper 启动时读取此文件，拿到里面的数据与 zoo.cfg 里面的配置信息比较从而判断到底是哪个 server。

3) 集群操作

(1) 在/opt/module/zookeeper-3.4.10/zkData 目录下创建一个 myid 的文件

```
touch myid
```

添加 myid 文件，注意一定要在 linux 里面创建，在 notepad++ 里面很可能乱码

(2) 编辑 myid 文件

```
vi myid
```

在文件中添加与 server 对应的编号：如 2

(3) 拷贝配置好的 zookeeper 到其他机器上

```
scp -r zookeeper-3.4.10/ root@hadoop103.atguigu.com:/opt/app/
```

```
scp -r zookeeper-3.4.10/ root@hadoop104.atguigu.com:/opt/app/
```

并分别修改 myid 文件中内容为 3、4

(4) 分别启动 zookeeper

```
[root@hadoop102 zookeeper-3.4.10]# bin/zkServer.sh start
```

```
[root@hadoop103 zookeeper-3.4.10]# bin/zkServer.sh start
```

```
[root@hadoop104 zookeeper-3.4.10]# bin/zkServer.sh start
```

(5) 查看状态

```
[root@hadoop102 zookeeper-3.4.10]# bin/zkServer.sh status
```

JMX enabled by default

Using config: /opt/module/zookeeper-3.4.10/bin/../conf/zoo.cfg

Mode: **follower**

```
[root@hadoop103 zookeeper-3.4.10]# bin/zkServer.sh status
```

JMX enabled by default

Using config: /opt/module/zookeeper-3.4.10/bin/../conf/zoo.cfg

Mode: **leader**

```
[root@hadoop104 zookeeper-3.4.5]# bin/zkServer.sh status
```

JMX enabled by default

Using config: /opt/module/zookeeper-3.4.10/bin/../conf/zoo.cfg

Mode: **follower**

8.3.4 配置 HDFS-HA 集群

1) 官方地址: <http://hadoop.apache.org/>

2) 在 opt 目录下创建一个 ha 文件夹

```
mkdir ha
```

3) 将/opt/app/下的 hadoop-2.7.2 拷贝到/opt/ha 目录下

```
cp -r hadoop-2.7.2/ /opt/ha/
```

4) 配置 hadoop-env.sh

```
export JAVA_HOME=/opt/module/jdk1.8.0_144
```

5) 配置 core-site.xml

```
<configuration>
  <!-- 把两个 NameNode 的地址组装成一个集群 mycluster -->
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://mycluster</value>
  </property>

  <!-- 指定 hadoop 运行时产生文件的存储目录 -->
  <property>
    <name>hadoop.tmp.dir</name>
    <value>/opt/ha/hadoop-2.7.2/data/tmp</value>
  </property>
</configuration>
```

6) 配置 hdfs-site.xml

```
<configuration>
  <!-- 完全分布式集群名称 -->
  <property>
    <name>dfs.nameservices</name>
    <value>mycluster</value>
  </property>

  <!-- 集群中 NameNode 节点都有哪些 -->
  <property>
    <name>dfs.ha.namenodes.mycluster</name>
    <value>nn1,nn2</value>
  </property>

  <!-- nn1 的 RPC 通信地址 -->
  <property>
    <name>dfs.namenode.rpc-address.mycluster.nn1</name>
    <value>hadoop102:9000</value>
  </property>

  <!-- nn2 的 RPC 通信地址 -->
  <property>
    <name>dfs.namenode.rpc-address.mycluster.nn2</name>
    <value>hadoop103:9000</value>
  </property>
</configuration>
```

```
</property>

<!-- nn1 的 http 通信地址 -->
<property>
  <name>dfs.namenode.http-address.mycluster.nn1</name>
  <value>hadoop102:50070</value>
</property>

<!-- nn2 的 http 通信地址 -->
<property>
  <name>dfs.namenode.http-address.mycluster.nn2</name>
  <value>hadoop103:50070</value>
</property>

<!-- 指定 NameNode 元数据在 JournalNode 上的存放位置 -->
<property>
  <name>dfs.namenode.shared.edits.dir</name>
  <value>qjournal://hadoop102:8485;hadoop103:8485;hadoop104:8485/mycluster</value>
</property>

<!-- 配置隔离机制，即同一时刻只能有一台服务器对外响应 -->
<property>
  <name>dfs.ha.fencing.methods</name>
  <value>sshfence</value>
</property>

<!-- 使用隔离机制时需要 ssh 无密钥登录-->
<property>
  <name>dfs.ha.fencing.ssh.private-key-files</name>
  <value>/home/atguigu/.ssh/id_rsa</value>
</property>

<!-- 声明 journalnode 服务器存储目录-->
<property>
  <name>dfs.journalnode.edits.dir</name>
  <value>/opt/ha/hadoop-2.7.2/data/jn</value>
</property>

<!-- 关闭权限检查-->
<property>
  <name>dfs.permissions.enable</name>
  <value>false</value>
</property>
```

```
<!-- 访问代理类: client, mycluster, active 配置失败自动切换实现方式-->
<property>
  <name>dfs.client.failover.proxy.provider.mycluster</name>
  <value>org.apache.hadoop.hdfs.server.namenode.ha.ConfiguredFailoverProxyProvider</value>
</property>
</configuration>
```

7) 拷贝配置好的 hadoop 环境到其他节点

8.3.5 启动 HDFS-HA 集群

1) 在各个 JournalNode 节点上, 输入以下命令启动 journalnode 服务:

```
sbin/hadoop-daemon.sh start journalnode
```

2) 在[nn1]上, 对其进行格式化, 并启动:

```
bin/hdfs namenode -format
```

```
sbin/hadoop-daemon.sh start namenode
```

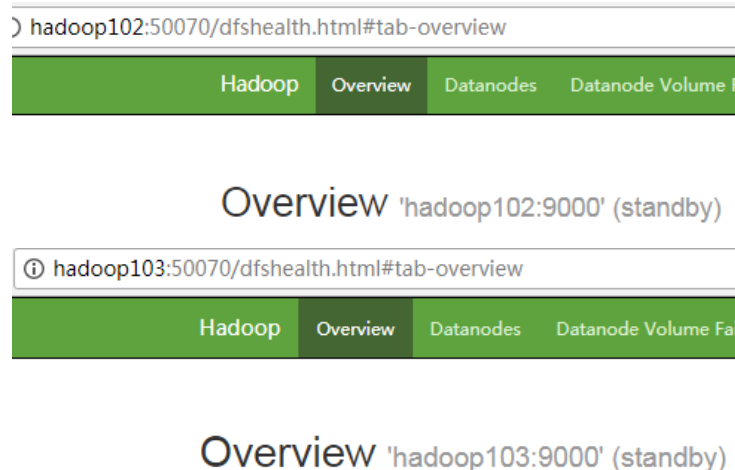
3) 在[nn2]上, 同步 nn1 的元数据信息:

```
bin/hdfs namenode -bootstrapStandby
```

4) 启动[nn2]:

```
sbin/hadoop-daemon.sh start namenode
```

5) 查看 web 页面显示



6) 在[nn1]上, 启动所有 datanode

```
sbin/hadoop-daemons.sh start datanode
```

7) 将[nn1]切换为 Active

```
bin/hdfs haadmin -transitionToActive nn1
```


8) 查看是否 Active

```
bin/hdfs haadmin -getServiceState nn1
```

8.3.6 配置 HDFS-HA 自动故障转移

1) 具体配置

(1) 在 hdfs-site.xml 中增加

```
<property>
  <name>dfs.ha.automatic-failover.enabled</name>
  <value>true</value>
</property>
```

(2) 在 core-site.xml 文件中增加

```
<property>
  <name>ha.zookeeper.quorum</name>
  <value>hadoop102:2181,hadoop103:2181,hadoop104:2181</value>
</property>
```

2) 启动

(1) 关闭所有 HDFS 服务:

```
sbin/stop-dfs.sh
```

(2) 启动 Zookeeper 集群:

```
bin/zkServer.sh start
```

(3) 初始化 HA 在 Zookeeper 中状态:

```
bin/hdfs zkfc -formatZK
```

(4) 启动 HDFS 服务:

```
sbin/start-dfs.sh
```

(5) 在各个 NameNode 节点上启动 DFSZK Failover Controller, 先在哪台机器启动, 哪个机器的 NameNode 就是 Active NameNode

```
sbin/hadoop-daemon.sh start zkfc
```

3) 验证

(1) 将 Active NameNode 进程 kill

```
kill -9 namenode 的进程 id
```

(2) 将 Active NameNode 机器断开网络

```
service network stop
```

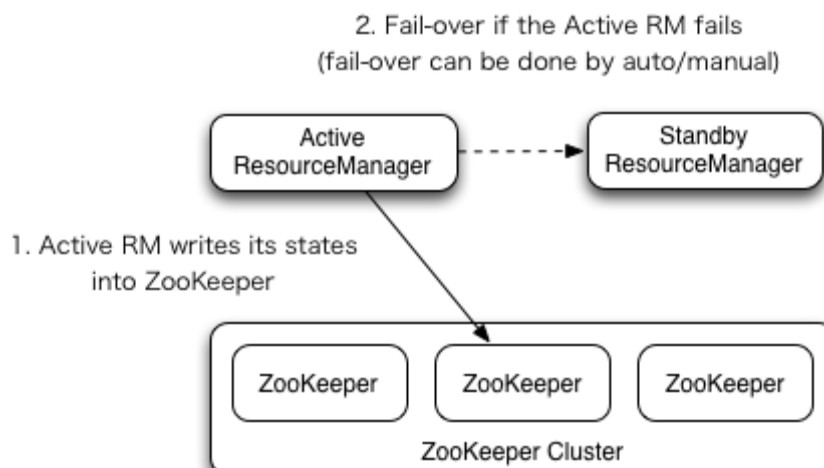
8.4 YARN-HA 配置

8.4.1 YARN-HA 工作机制

1) 官方文档:

<http://hadoop.apache.org/docs/r2.7.2/hadoop-yarn/hadoop-yarn-site/ResourceManagerHA.html>

2) YARN-HA 工作机制



8.4.2 配置 YARN-HA 集群

0) 环境准备

- (1) 修改 IP
- (2) 修改主机名及主机名和 IP 地址的映射
- (3) 关闭防火墙
- (4) ssh 免密登录
- (5) 安装 JDK，配置环境变量等
- (6) 配置 Zookeeper 集群

1) 规划集群

hadoop102	hadoop103	hadoop104
NameNode	NameNode	
JournalNode	JournalNode	JournalNode
DataNode	DataNode	DataNode
ZK	ZK	ZK
ResourceManager	ResourceManager	

NodeManager

NodeManager

NodeManager

2) 具体配置

(1) yarn-site.xml

```
<configuration>

  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>

  <!-- 启用 resourcemanager ha-->
  <property>
    <name>yarn.resourcemanager.ha.enabled</name>
    <value>true</value>
  </property>

  <!-- 声明两台 resourcemanager 的地址-->
  <property>
    <name>yarn.resourcemanager.cluster-id</name>
    <value>cluster-yarn1</value>
  </property>

  <property>
    <name>yarn.resourcemanager.ha.rm-ids</name>
    <value>rm1,rm2</value>
  </property>

  <property>
    <name>yarn.resourcemanager.hostname.rm1</name>
    <value>hadoop102</value>
  </property>

  <property>
    <name>yarn.resourcemanager.hostname.rm2</name>
    <value>hadoop103</value>
  </property>

  <!-- 指定 zookeeper 集群的地址-->
  <property>
    <name>yarn.resourcemanager.zk-address</name>
    <value>hadoop102:2181,hadoop103:2181,hadoop104:2181</value>
  </property>


```

```
<!--启用自动恢复-->
<property>
  <name>yarn.resourcemanager.recovery.enabled</name>
  <value>true</value>
</property>

<!--指定 resourcemanager 的状态信息存储在 zookeeper 集群-->
<property>
  <name>yarn.resourcemanager.store.class</name>
  <value>org.apache.hadoop.yarn.server.resourcemanager.recovery.ZKRMStateStore<
    /value>
</property>

</configuration>
```

(2) 同步更新其他节点的配置信息

3) 启动 hdfs

(1) 在各个 JournalNode 节点上, 输入以下命令启动 journalnode 服务:

```
sbin/hadoop-daemon.sh start journalnode
```

(2) 在[nn1]上, 对其进行格式化, 并启动:

```
bin/hdfs namenode -format
```

```
sbin/hadoop-daemon.sh start namenode
```

(3) 在[nn2]上, 同步 nn1 的元数据信息:

```
bin/hdfs namenode -bootstrapStandby
```

(4) 启动[nn2]:

```
sbin/hadoop-daemon.sh start namenode
```

(5) 启动所有 datanode

```
sbin/hadoop-daemons.sh start datanode
```

(6) 将[nn1]切换为 Active

```
bin/hdfs haadmin -transitionToActive nn1
```

4) 启动 yarn

(1) 在 hadoop102 中执行:

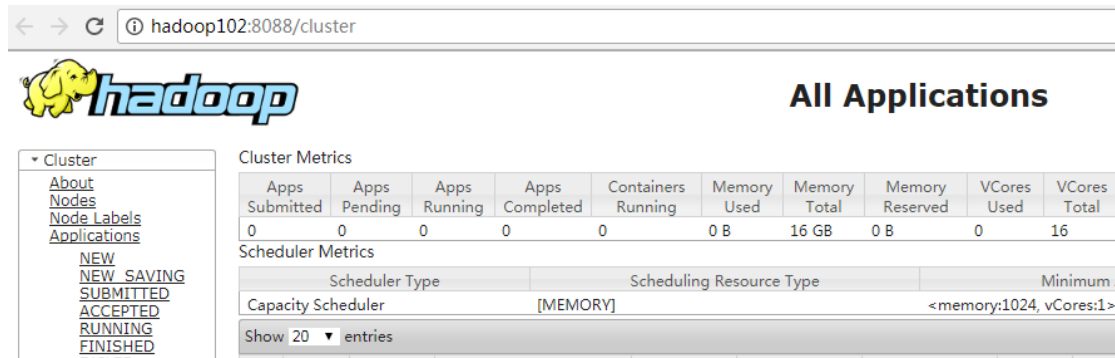
```
sbin/start-yarn.sh
```

(2) 在 hadoop103 中执行:

```
sbin/yarn-daemon.sh start resourcemanager
```

(3) 查看服务状态

```
bin/yarn rmadmin -getServiceState rml
```



The screenshot shows the Hadoop web interface at `hadoop102:8088/cluster`. It displays the Hadoop logo and the title "All Applications". On the left, there is a sidebar with navigation links: Cluster, About, Nodes, Node Labels, Applications, NEW, NEW SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, and FAILED. The main content area shows "Cluster Metrics" and "Scheduler Metrics".

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total
0	0	0	0	0	0 B	16 GB	0 B	0	16

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum
Capacity Scheduler	[MEMORY]	<memory:1024, vCores:1>

Below the scheduler metrics, there is a "Show 20 entries" button.

8.5 HDFS Federation 架构设计

1) NameNode 架构的局限性

(1) Namespace（命名空间）的限制

由于 NameNode 在内存中存储所有的元数据（metadata），因此单个 namenode 所能存储的对象（文件+块）数目受到 namenode 所在 JVM 的 heap size 的限制。50G 的 heap 能够存储 20 亿（200million）个对象，这 20 亿个对象支持 4000 个 datanode，12PB 的存储（假设文件平均大小为 40MB）。随着数据的飞速增长，存储的需求也随之增长。单个 datanode 从 4T 增长到 36T，集群的尺寸增长到 8000 个 datanode。存储的需求从 12PB 增长到大于 100PB。

(2) 隔离问题

由于 HDFS 仅有一个 namenode，无法隔离各个程序，因此 HDFS 上的一个实验程序就很有可能影响整个 HDFS 上运行的程序。

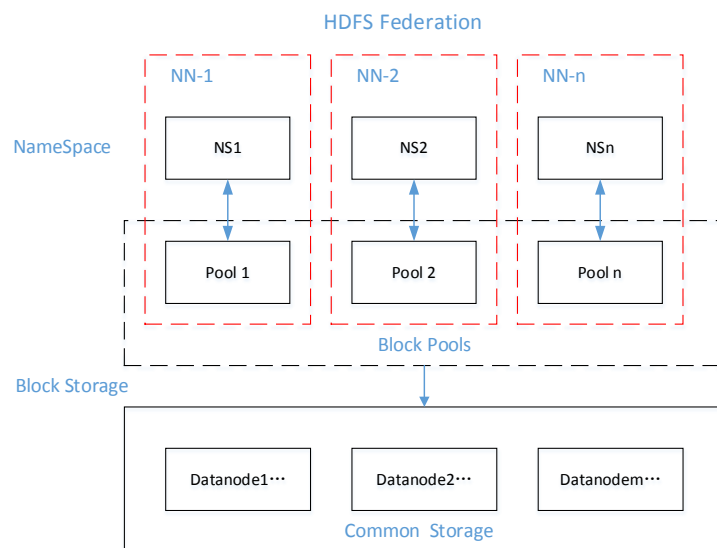
(3) 性能的瓶颈

由于是单个 namenode 的 HDFS 架构，因此整个 HDFS 文件系统的吞吐量受限于单个 namenode 的吞吐量。

2) HDFS Federation 架构设计

能不能有多个 NameNode

NameNode	NameNode	NameNode
元数据	元数据	元数据
Log	machine	电商数据/话单数据



3) HDFS Federation 应用思考

不同应用可以使用不同 NameNode 进行数据管理

图片业务、爬虫业务、日志审计业务

Hadoop 生态系统中，不同的框架使用不同的 namenode 进行管理 namespace。（隔离性）