

# Hadoop 管理员的十个最佳实践



## 前言

接触 Hadoop 有两年的时间了，期间遇到很多的问题，既有经典的 NameNode 和 JobTracker 内存溢出故障，也有 HDFS 存储小文件问题，既有任务调度问题，也有 MapReduce 性能问题。遇到的这些问题有些是 Hadoop 自身的缺陷（短板），有些则是使用的不当。

在解决问题的过程中，有时需要翻源码，有时会向同事、网友请教，遇到复杂问题则会通过 mail list 向全球各地 Hadoop 使用者，包括 Hadoop Committer（Hadoop 开发者）求助。在获得很多人帮助后，自己将遇到问题和心得整理成文，希望本文可以对那些焦头烂额的 Hadoop 新手们有所帮助，少走笔者的弯路。

PS. 本文基于 Cloudera CDH 3u4（同 Apache Hadoop 1.0）编写。相关推荐配置为官方推荐值或者笔者经验数值，它不是绝对的，可能会因为不同的应用场景和硬件环境有所出入。

相关厂商内容



## 1. 选择 Cloudera CDH 部署你的 Cluster

### 动机

大多数管理员都是从 Apache Hadoop 开始学习。笔者最开始也使用 Apache 版本 Hadoop 进行开发和部署工作，但接触到 Cloudera CDH 后，我发现它可以使管理员的工作更简单，不仅可以获得最新的特性和 Bug 修复，有时也会带来令人惊喜的性能改善。

CDH 为什么更好？笔者罗列了以下几点：

1. CDH 基于稳定版 Apache Hadoop，并应用了最新 Bug 修复或者 Feature 的 Patch。Cloudera 常年坚持季度发行 Update 版本，年度发行 Release 版本，更新速度比 Apache 官方快，而且在实际使用过程中 CDH 表现无比稳定，并没有引入新的问题。
2. Cloudera 官方网站上安装、升级文档详细，省去 Google 时间。
3. CDH 支持 Yum/Apt 包，Tar 包，RPM 包，Cloudera Manager 四种方式安装，总有一款适合您。官方网站推荐 Yum/Apt 方式安装，笔者体会其好处如下：
  1. 联网安装、升级，非常方便。当然你也可以下载 rpm 包到本地，使用 Local Yum 方式安装。
  2. 自动下载依赖软件包，比如要安装 Hive，则会级联下载、安装 Hadoop。
  3. Hadoop 生态系统包自动匹配，不需要你寻找与当前 Hadoop 匹配的 Hbase，Flume，Hive 等软件，Yum/Apt 会根据当前安装 Hadoop 版本自动寻找匹配版本的软件包，并保证兼容性。
  4. 自动创建相关目录并软链到合适的地方（如 conf 和 logs 等目录）；自动创建 hdfs, mapred 用户，hdfs 用户是 HDFS 的最高权限用户，mapred 用户则负责 mapreduce 执行过程中相关目录的权限。

推荐指数：★★★★

推荐理由：获取最新特性和最新 Bug 修复；安装维护方便，节省运维时间。

## 2. Hadoop 集群配置与管理

安装和维护 Hadoop 集群涉及大量的管理工作，包括软件安装，设备管理（crontab、iptables 等）、配置分发等。

对于小型集群软件分发和节点管理可以使用 PDSH 这款软件，它可以通过免密钥的 SSH 将文件分发到目标服务器，以及为一组目标设备发送命令并获得反馈。如果是大型集群或者硬件配置差别很大的集群，推荐使用 puppet 这样的工具帮助你维护配置文件，或者通过 Cloudera Manager 以 GUI 的方式的管理集群（注意：Cloudera Manager 不是开源软件，免费版最多支持 50 个节点）。

推荐指数：★★★★

推荐理由：提高运维效率

## 3. 开启 SecondaryNameNode

SecondaryNameNode（下称 SNN）的主要功能是工作是帮助 NameNode（下称 NN）合并编辑日志，然后将合并后的镜像文件 copy 回 NN，以减少 NN 重启时合并编辑日志所需的时间。SNN 不是 NN 的热备，但是通过以下步骤可以实现将 SNN 切

换为 NN 的目的。首先，SNN 节点上导入从 NN Copy 过来的镜像文件，然后修改 SNN 机器名和 IP 与 NN 一致，最后重启集群。

特别注意的是 SNN 的内存配置要与 NN 一致，因为合并编辑日志的工作需要将 metadata 加载到内存完成。另外，不仅仅是 SNN，任何保存 NN 镜像的节点都可以通过上面步骤变为 NN，只是 SNN 更适合罢了。

推荐指数：★★★★

推荐理由：减少 NN 重启导致集群服务中断时间；NN 节点故障后，SNN 充当 NN 角色

## 4. 使用 Ganglia 和 Nagios 监控你的集群

当运行一个大型 mapreduce 作业时，我们通常非常关心该作业对 TaskTracker（下称 TT）CPU、内存、磁盘，以及整个网络的带宽情况，这时候就需要 Ganglia 这个工具为我们生成相关图表来诊断、分析问题。

Ganglia 可以监控集群状态，但当你的服务器 down 机或者某个 TT 挂掉，它却无法通知到你，这时我们可以使用 Nagios 这款告警软件，它可以配置邮件告警和短息告警。通过编写 plugins，可以实现自己的监控功能。我们的集群目前做了如下监控：

1. NameNode、JobTracker 内存
2. DataNode 和 TaskTracker 运行状态
3. NFS 服务状态
4. 磁盘使用情况
5. 服务器负载状态

推荐指数：★★★★

推荐理由：Ganglia 可以帮你记录集群状态，方便诊断问题；Nagios 可以再遇到问题时第一时间通知你。

## 5. 设置好内存至关重要

Hadoop 集群安装完毕后，第一件事就是修改 bin/hadoop-env.sh 文件设置内存。主流节点内存配置为 32GB，典型场景内存设置如下

NN: 15-25 GB

JT: 2-4GB

DN: 1-4 GB

TT: 1-2 GB, Child VM 1-2 GB

集群的使用场景不同相关设置也有不同，如果集群有大量小文件，则要求 NN 内存至少要 20GB，DN 内存至少 2GB。

推荐指数：★★★★★

推荐理由：几个组件中 NN 对内存最为敏感，它有单点问题，直接影响到集群的可用性；JT 同样是单点，如果 JT 内存溢出则所有 MapReduce Job 都无法正常执行。

## 6. 管理员玩转 MapReduce

Hadoop 原生 MapReduce 需要 Java 语言编写，但是不会 Java 也没问题，通过 Hadoop streaming 框架管理员可以使用 Python, Shell, Perl 等语言进行 MapReduce 开发，但更简单的办法是安装和使用 Hive 或者 Pig。

推荐指数：★★★

推荐理由：减少运维时间，快速响应各种 ad-hot 需求和故障诊断。

## 7. NameNode HA

前面已经说过，NN 是整个集群可能出现的单点故障。

Hadoop 通过在 `hdfs.site.xml` 文件的 `dfs.name.dir` 属性指定保持的 metadata 路径，如果希望保持到多个路径，可以使用逗号分割配置多个路径。

```
<property>
  <name>dfs.name.dir</name>
  <value>/data/cache1/dfs/nn,/data/cache2/dfs/nn</value>
</property>
```

Hadoop 官方推荐配置为 metadata 配置多个 path，其中包含一个 NFS 的路径。但根据笔者一次集群严重故障经验，即使这样，还是导致了所有镜像文件损坏，包括 SNN 上的镜像文件，所以定期备份一个可用的副本还是很有必要的。

推荐指数：★★★★★

推荐理由：Cloudera3uX 和 Apache1.0 的 NN 单点问题是大家最头痛问题之一，多些准备，少许痛苦。

## 8. 使用 firewall 阻止坏人进入

Hadoop 的安全控制非常简单，只包含简单的权限，即只根据客户端用户名，决定使用权限。它的设计原则是：“避免好人做错事，但不阻止坏人做坏事”。

如果你知道某台 NN 的 IP 和端口，则可以很轻松获取 HDFS 目录结构，并通过修改本机机器用户名伪装成 HDFS 文件所属 owner，对该文件进行删除操作。

通过配置 kerberos，可以实现身份验证。但很多管理员使用更简单有效的办法——通过防火墙对访问 IP 进行控制。

推荐指数：★★★★★

推荐理由：安全无小事，防范于未然。

## 9. 开启垃圾箱(trash)功能

### 动机

我曾经犯下一个错误，在我加班非常累，大脑稍有混乱的时候，不小心删除执行了一个命令“`hadoop fs -rmr /xxx/xxx`”，没有删除提示，几 TB 的数据，一下子就没有了。简直让我崩溃，后悔莫及。这时你多希望有个时间机器可以让 HDFS 恢复到删除前的状态。

trash 功能就是这个时间机器，它默认是关闭的，开启后，被你删除的数据将会 mv 到操作用户目录的“.Trash”文件夹，可以配置超过多长时间，系统自动删除过期数据。这样一来，当操作失误的时候，可以把数据 mv 回来。开启垃圾箱步骤如下：

vi core-site.xml ，添加下面配置，value 单位为分钟。

```
<property>
  <name>fs.trash.interval</name>
  <value>1440</value>
</property>
```

笔者在 CDH3u4 下不用重启 Namenode 就可以生效。开启垃圾箱后，如果希望文件直接被删除，可以在使用删除命令时添加“-skipTrash”参数，如下：

```
hadoop fs -rm -skipTrash /xxxx
```

推荐指数：★★★★★

推荐理由：想要时间机器吗？

## 10. 去社区寻找帮助

Hadoop 是一个非常优秀的开源项目，但它仍存有很多尚未解决的问题，诸如，NN, JT 单点问题，JT 挂死问题，Block 在小文件下汇报效率低下等问题。此时可以通过如下渠道找到可以帮助你的人，笔者几次集群严重故障都是通过 Cloudera 公司的 google user group 直接获得几位 committer 的帮助。通常前一天提问，第二天就会有反馈。下面是两个能够帮助的你的社区，当然你也可以帮助其他人：

Apache hadoop 的 mail list：

[http://hadoop.apache.org/mailling\\_lists.html](http://hadoop.apache.org/mailling_lists.html)

Cloudera CDH google group：

<https://groups.google.com/a/cloudera.org/forum/#!forum/cdh-user>

推荐指数：★★★★★

推荐理由：没有人比软件作者更熟悉 Hadoop 本身，去社区求助，帮你解决很多自己无法跨越的问题。

## Cloudera 简介：

公司是一家 Hadoop 软件服务公司，提供免费软件 CDH 和 Cloudera Manager Free Edition，同时提供 Hadoop 相关资讯、培训、技术支持等服务。Hadoop 创始人 Dong Cutting 在该公司任架构师，同时该公司拥有多名 Apache Committer。