

英特尔® Apache Hadoop* 软件发行版高可用性 操作手册

版本 2.3

2013 年 3 月



免责声明和法律信息

本文件中包含关于英特尔产品的信息。本文件不构成对任何知识产权的授权，包括明示的、暗示的，也无无论是基于禁止反言的原则或其他。除英特尔产品销售的条款和条件规定的责任外，英特尔不承担任何其他责任。英特尔在此作出免责声明：本文件不构成英特尔关于其产品的使用和 / 或销售的任何明示或暗示的保证，包括不就其产品的 (i) 对某一特定用途的适用性、(ii) 适销性以及 (iii) 对任何专利、版权或其他知识产权的侵害的承担任何责任或作出任何担保

除非经过英特尔的书面同意认可，英特尔的产品无意被设计用于或被用于以下应用：即在这样的应用中可因英特尔产品的故障而导致人身伤亡。

英特尔有权随时更改产品的规格和描述而无需发出通知。设计者不应信赖任何英特尔产品所不具有的特性，设计者亦不应信赖任何标有“保留权利”或“未定义”说明或特性描述。对此，英特尔保留将来对其进行定义的权利，同时，英特尔不应为其日后更改该等说明或特性描述而产生的冲突和不相容承担任何责任。此处提供的信息可随时改变而无需通知。请勿根据本文件提供的信息完成一项产品设计。

本文件所描述的产品可能包含使其与宣称的规格不符的设计缺陷或失误。这些缺陷或失误已收录于勘误表中，可索取获得。

在发出订单之前，请联系当地的英特尔营业部或分销商以获取最新的产品规格。

索取本文件中或英特尔的其他材料中提的、包含订单号的文件的复印件，可拨打 1-800-548-4725，或登陆 <http://www.intel.com/design/literature.htm>。

英特尔处理器标号不是性能的指标。处理器标号仅用于区分同属一个系列的处理器的特性，而不能够用于区分不同系列的处理器。**详情请登陆：**
http://www.intel.com/products/processor_number

Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance. Buyers should consult other sources of information to evaluate the performance of systems or components they are considering purchasing. For more information on performance tests and on the performance of Intel products, visit [Intel Performance Benchmark Limitations](#).

结果基于模拟测算得出，仅作参考之用。结果通过系统模拟器或模型测算得出。任何系统硬件、软件的设计或配置的不同均可能影响实际性能。

Intel, 英特尔® Apache Hadoop® 软件发行版, Intel® Distribution, Intel® Manager for Apache Hadoop® software, Intel® Manager 是英特尔在美国和 / 或其他国家的商标。

* 其他的名称和品牌可能是其他所有者的资产。

英特尔公司 2013 年版权所有。所有权保留。



文档修订记录

日期	修订	描述
2013 年 3 月	001	英特尔® Apache Hadoop* 软件发行版 v2.3 第一版文档
2013 年 3 月	002	英特尔® Apache Hadoop* 软件发行版 v2.3 文档更新



目录

1.0	简介	2
1.1	支持的高可用性种类	2
1.2	高可用性是如何工作的	2
1.3	支持的操作系统	3
2.0	设置高可用性的要求和推荐	4
2.1	总体要求	4
2.2	推荐配置	5
2.3	限制	5
3.0	使用 Intel® Manager 来配置高可用性	6
4.0	高可用性维护	9
4.1	监控高可用性	9
4.2	更改高可用性配置	9
4.3	基本 DRBD 操作	10
4.3.1	资源概述	10
4.3.2	检查 DRBD 状态	10
4.3.3	常用命令	12
4.4	Pacemaker 的基本操作	12
4.4.1	资源状态	12
4.4.2	操作资源	13
4.4.3	操作集群节点	13
4.4.4	维护模式	13
4.4.5	手工更改 CRM 配置	13
4.5	处理硬盘故障	13
4.6	节点故障后的 NameNode 恢复	14
4.7	脑裂恢复	14
4.8	运行脚本排除故障	15
5.0	卸载高可用性	16



1.0 简介

Apache Hadoop* 集群中通常有多个用户长时间地运行多个作业。这些作业生成的数据分析具有商业上的重要性，可帮助公司节省大笔开支或产生收入。因此，集群的高可用性至关重要，几分钟、几小时或几天的宕机可能花费大量的金钱。

系统管理员面对的问题主要是 JobTracker 和 Primary NameNode 的单点故障。如果其中一个服务失败，则在问题解决前集群功能将不可用。而且，这些故障可能需要花费大量的时间和人力去解决，这将导致长时间宕机，这对公司业务尤其是关键业务来说是不可接受的。

要解决这些问题，英特尔® Apache Hadoop* 软件发行版支持 JobTracker 和 Primary NameNode daemon 的高可用性 (High Availability) 功能。如果其中一个 daemon 发生故障，则备用 daemon 将立即托管，集群可继续处理数据，与此同时，系统管理员可进行故障排除解决问题。

本文档解释如何为英特尔® 发行版设置高可用性。

1.1 支持的高可用性种类

高可用性功能支持 JobTracker 和 Primary NameNode 的 active-standby 配置。这表示每个 Primary NameNode 在另一个节点上有一个完全冗余的对象，它只有当关联的 Primary NameNode 发生故障时会在线。Primary NameNode 的冗余对象被称为 *Standby NameNode*。JobTracker 的冗余对象被称为 *Backup JobTracker*。Primary NameNode 的服务和 Standby 节点上的服务完全一致。因此，如果 Primary NameNode 发生故障，则 standby 节点上接手托管的服务被命名为 *hadoop-namenode*。如果 JobTracker 发生故障，则 Backup JobTracker 节点上接手托管的服务被命名为 *hadoop-jobtracker*。

高可用性配置包括以下四类节点：Standby NameNode、Backup JobTracker、Primary NameNode 和 JobTracker 节点。

如果 Primary NameNode 服务不可用，则高可用性功能将自动探测故障，并切换到 *Standby NameNode*。如果 JobTracker 服务不可用，则高可用性功能将自动探测故障，并切换到 *Backup JobTracker*。如果 Primary NameNode 和 JobTracker 同时发生故障，则高可用性功能将切换到相应的 Standby 进程。

发生故障时，通常会产生短时间的服务中断，无需人工干预会自动恢复。此外，如果在运行作业时 JobTracker 发生故障，期间所有作业将丢失。

1.2 高可用性是如何工作的

高可用性功能包含以下服务：

- Distributed Replicated Block Device (DRBD) — 是一个用软件实现的、无共享的、服务器之间镜像块设备内容的存储复制解决方案。也就是说，这是位于不同机器的二个磁盘分区的基于 TCP 的磁盘冗余阵列 (RAID)。
- Pacemaker 是一个集群资源管理 (CRM) 的框架，它能自动启动、停止、监控和迁移资源。
- Corosync 是 Pacemaker 能够使用的集群通讯层。

在为高可用性配置 DRBD 时，存在主要和次要设备。



主要设备存在于 Primary NameNode。这一设备是存储 NameNode 的 fs_image 的逻辑磁盘分区。次要设备存在于 Standby NameNode。这一设备是一个大小和主要设备一致的逻辑磁盘分区。当数据块放入主要设备时，它们会被自动复制到次要设备上。数据复制是指次要设备仅包含主要设备的数据，但主要设备不从次要设备上获得数据的一种方法。

Pacemaker 负责探测 Active 节点的故障，并启动切换到合适的 Standby 节点。服务也管理和监控以下发生故障时切换所需资源：

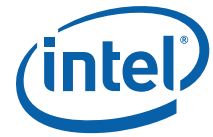
- 在 Active 和 Standby 节点之间浮动的虚拟 IP (VIP) 地址。虚拟 IP 用于客户机（比如 DataNodes 和 TaskTracker）和 active 节点的通讯。
- Primary NameNode、Standby NameNode、JobTracker 和 Backup JobTracker。
- DRBD

Pacemaker 本身是一个集群，它包含了 Apache Hadoop* 集群中的多个节点。Pacemaker 服务根据 quorum 做出判断。也就是说，当 Pacemaker 集群中的多数节点同意故障切换已发生、而且哪个 Standby 节点应成为 Active 节点后，则 Pacemaker 判断故障已被探测到，且切换已启动。

如果 Primary NameNode 发生故障，则 Pacemaker 会将之切换到 Standby NameNode。因为 Standby NameNode 可访问 DRBD 次要设备，且这一设备包含和 Primary NameNode 的 fs_image 完全相同的复制品，存储在 HDFS 的数据不会丢失，客户机仍可通过 Standby NameNode 访问这些数据。

1.3 支持的操作系统

Red Hat Enterprise Linux 6.3 和 CentOS 6.3 操作系统支持高可用性。



2.0 设置高可用性的要求和推荐

2.1 总体要求

在配置英特尔® Apache Hadoop* 软件发行版高可用性前，你需要了解或操作以下事项：

- 你需要参照英特尔® Apache Hadoop* 软件发行版操作管理手册创建和配置一个不带高可用性的集群。
- 集群中的每个节点必须符合英特尔® Apache Hadoop* 软件发行版操作管理手册中所述的要求：
 - 硬件和软件要求
 - 内存要求
 - 磁盘分区要求
 - 集群和网络拓扑要求
 - 端口要求
- 确定要成为 Standby NameNode 和 Backup JobTracker 的节点。
- 高可用性设置后，HDFS 必须被重新格式化。也就是说，在设置高可用性之前任何存在于 HDFS 的数据都将在设置高可用性后丢失。
- 确定当故障发生进行切换时，从 Primary NameNode 浮动到 Standby NameNode 的虚拟 IP 和主机名。虚拟 IP 必须是静态的，网络内的其他节点不能使用这一主机名。
- 确定当故障发生进行切换时，从 JobTracker 浮动到 Backup JobTracker 的虚拟 IP 和主机名。虚拟 IP 必须是静态的，网络内的其他节点不能使用这一主机名。
- 高可用性配置中的所有节点（Primary NameNode、Standby NameNode、JobTracker 和 Backup JobTracker）必须被分配静态 IP 地址。
- 存储 fs_image 的 Primary NameNode 的逻辑磁盘分区的大小必须等同于系统的内存数量。比如，如果 Primary NameNode 内存为 64GB，则该逻辑磁盘分区必须为 64GB。
- 存储 fs_image 的 Primary NameNode 的逻辑磁盘分区必须是 DRBD 主要设备。除了 fs_image，其他数据不能存储在逻辑磁盘分区。
- Standby NameNode 必须有一个大小和 Primary NameNode 上的 DRBD 主要设备相同的逻辑磁盘分区。逻辑磁盘分区必须是 DRBD 次要设备。
- DRBD 复制通过 TCP 端口 7799 来使用 7788，每个资源的侦听都分开使用端口。DRBD 对每个配置的资源使用 TCP 连接。要让 DRBD 正常工作，Standby NameNode 和 Primary NameNode 上的防火墙不能阻止和这些端口的 TCP 通讯。
- 除了 DRBD，没有其他服务或进程可通过 TCP 端口 7799 来使用 7788。
- 在配置高可用性之前，Pacemaker、Corosync 和 DRBD 必须没有被安装在集群中的任何节点上。安装包必须只从英特尔® 发行版提供的 repository 进行安装。
- Standby NameNode 和 Primary NameNode 必须有以下内核版本 2.6.32-279.el6.x86_64。
- 高可用性功能仅在内部部署的（on-premise）集群中工作。这一功能不被云支持，比如 Amazon EC2。



2.2 推荐配置

- 当通过交换机进行 DRBD 复制时，推荐以 active-backup 模式使用冗余组件和绑定的驱动。
- 不能通过路由器进行 DRBD 复制。
- 推荐通过网卡进行 DRBD 复制，可专门用于 DRBD 网络通讯。连接必须是直接的、背对背的千兆以太网通讯。
- 高可用性集群通讯可以是单播或多播。高可用性组件仅在多播模式下测试过，强烈推荐你在集群通讯中使用单播模式。
- 为避免某个机架成为单点故障，高可用性配置中的任一节点（Primary NameNode、Standby NameNode、JobTracker 和 Backup JobTracker）不能和高可用性配置中的其他节点位于同一机架上。

2.3 限制

如果高可用性已配置，则集群受到以下限制：

- 不支持 Kerberos 验证。
- 不支持基于角色的访问控制和访问控制列表。



3.0 使用 Intel® Manager 来配置高可用性

要配置高可用性，在 Intel® Manager for Apache Hadoop* software 中执行以下步骤：

1. 以管理员权限的用户登录 Intel® Manager。
2. 如果有任何服务还在运行中，你必须先停止它们。
3. 点击**配置向导**链接。
4. 在第 1 步页面，勾选**高可用性**选项。

第1步

配置新的集群

集群名称：

选择集群中将会使用的组件，包括HDFS，MapReduce，HBase，Hive，Sqoop，Pig，Flume 和 Oozie；另外高可用性组件将会使用两台主备机器来保证集群的高可用性。

集群组件：

☒ HDFS：HDFS是一个分布式的文件系统。

☒ MapReduce：MapReduce是一种用于分布式系统的并行计算框架。

☒ ZooKeeper：ZooKeeper是一个针对大型分布式系统的可靠协调系统。

☒ HBase：HBases是基于HDFS的分布式的，可伸缩的，版本化的数据库系统。

☒ Hive：Hive是基于Hadoop的数据仓库工具。

☒ Sqoop：Sqoop是用于结构化数据存储和Hadoop之间的数据传输的工具。

☒ Pig：Pig是一个基于Hadoop的大规模数据分析平台。

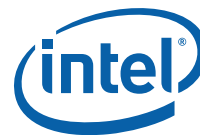
☒ Flume：Flume是一个分布式的、可靠的、和高可用的海量日志聚合的系统。

☒ Oozie：Oozie是一个用于管理和调度Hadoop任务的工作流引擎。

☒ 高可用性：集群中将会有一台备份机器来保证高可用性。

下一步

取消



5. 在配置向导中，在每个页面中点击**下一步**，直至页面显示**配置集群高可用性节点使用的通信模式**。

第5步

配置集群高可用性节点使用的通信模式

通信模式：

组播地址：

6. 在**通信模式**下拉菜单，选择 **Corosync 使用组播模式与其他节点通信** 选项。
7. 点击**下一步**。

第6步

HDFS高可用性配置

集群使用主备两台机器来保证集群的高可用性，当主机宕机时，备机会切换成主机接替原先主机的工作，从而保证集群正常运行不受影响。您需要指定主备两台机器所使用的虚拟主机名和IP地址；另外还需要指定这两台机器留作高可用性数据备份用的DRBD分区。

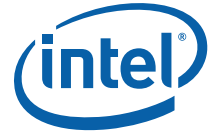
虚拟主机名： Namenode的主机名，它和被主备两台机器共用的一个虚拟主机名。

虚拟IP地址： 虚拟主机名对应的IP地址。确保虚拟IP地址没有被局域网内其他机器使用。

Primary NameNode： Primary NameNode DRBD分区：

Standby NameNode： Standby NameNode DRBD分区：

8. 在 HDFS 高可用性配置页面，执行以下步骤。
 - a. 在**虚拟主机名**栏内，输入当故障发生机器切换时从 active 节点浮动到 standby 节点的主机名。这是和 Primary NameNode 通讯的客户机（比如 DataNode）所使用的主机名。
 - b. 在**虚拟 IP 地址**栏内，输入当故障发生机器切换时从 active 节点浮动到 standby 节点的 IP 地址。这是和 Primary NameNode 通讯的客户机（比如 DataNode）所使用的 IP 地址。
 - c. 在 **Primary NameNode** 下拉菜单中，选择将会成为 Primary NameNode 节点的主机名。
 - d. 在 **Primary NameNode DRBD Partition** 下拉菜单中，选择会成为 DRBD 主要设备的机器上的设备，以及 fs_image 存储的位置。分区大小必须和机器的内存数量相同。
- 注释：** DRBD 分区下拉菜单仅显示设备的逻辑分区。如果设备没有被逻辑分区，则设备不在下拉菜单中显示。
- e. 在 **Standby NameNode** 下拉菜单中，选择将会成为 Standby NameNode 节点的主机名。



- f. 在 **Standby NameNode DRBD 分区** 下拉菜单中，选择会成为 DRBD 次要设备的机器上的设备，以及从 Primary NameNode 复制 fs_image 后存储的位置。分区大小必须和 Primary NameNode 的内存数量相同。
9. 在要为 Primary Namenode 指定的虚拟主机名上执行 nslookup 或 dig 命令。DNS 必须返回一个有效的条目，返回的 IP 地址必须和第 6 步页面中指定的匹配，而且网络中其它任何机器都不能和该主机名绑定。
10. 在 HDFS 高可用性配置页面，点击**下一步**。

第7步

MapReduce高可用性配置

集群使用主备两台机器来保证MapReduce的高可用性，当主Jobtracker宕机时，备机会切换到主机接替原先主机的工作，从而保证集群正常运行不受影响。您需要指定主备两台Jobtracker所使用的虚拟主机名和IP地址。

虚拟主机名： Jobtracker的主机名，它和被主备两台机器共用的一个虚拟主机名。

虚拟IP地址： 虚拟主机名对应的IP地址。**确保虚拟IP地址没有被局域网内其他机器使用。**

Jobtracker：

Backup Jobtracker：

11. 在 MapReduce 高可用性配置页面，执行以下步骤。
 - a. 在**虚拟主机名**栏内，输入当故障发生机器切换时从 active 节点浮动到 standby 节点的主机名。这是和 JobTracker 通讯的客户机（比如 TaskTracker）所使用的主机名。
 - b. 在**虚拟 IP 地址**栏内，输入当故障发生机器切换时从 active 节点浮动到 standby 节点的 IP 地址。这是和 JobTracker 通讯的客户机（比如 TaskTracker）所使用的 IP 地址。
 - c. 在 **JobTracker** 下拉菜单中，选择将会成为 JobTracker 的主机名。
 - d. 在 **Backup JobTracker** 下拉菜单中，选择将会成为 Backup JobTracker 的主机名。
12. 在要为 JobTracker 指定的虚拟主机名上执行 nslookup 或 dig 命令。DNS 必须返回一个有效的条目，返回的 IP 地址必须和第 7 步页面中指定的匹配，而且网络中其它任何机器都不能和该主机名绑定。
13. 点击每个页面的**下一步**直至配置向导完成。
14. 在集群配置菜单，点击**集群节点**选项。
15. 在集群节点子页面，点击**配置所有节点**链接。
16. 在配置所有节点过程完成后，点击**格式化集群**链接。
17. 在警告对话框，点击**确定**。

警告： 在 HDFS 重新格式化的过程中，每个 DRBD 主要设备上的块被复制到 DRBD 次要设备上。因此，取决于 DRBD 主要设备的容量大小，重新格式化可能需要半小时或更长时间。

4.0 高可用性维护

4.1 监控高可用性

你可以在**集群概况 > 高可用概述 > HA 概况**中查看 Intel® Manager for Apache Hadoop* software 高可用性的状态。



The screenshot shows the Intel Hadoop Management Platform interface. The main content area is titled "HA 概况" (HA Overview). It includes a sidebar on the left with navigation options like "集群概况" (Cluster Overview), "HDFS概述" (HDFS Overview), "MapReduce概述" (MapReduce Overview), "ZooKeeper概述" (ZooKeeper Overview), and "HBase概述" (HBase Overview). The main content area is divided into three sections:

- HA 概况** (HA Overview): Shows Current DC (ip-10-0-20-17.ec2.internal), Partition (WITH Quorum), HA Nodes (3(2Online, 1Standby)), and HA Resources (7(6运行, 1停止)).
- HA 节点状况** (HA Node Status): A table showing the status of HA nodes.
- HA 详细状态信息** (HA Detailed Status Information): A list of nodes and their current status, with links to view the current status.

服务器	状态	服务器	状态
ip-10-0-20-17.ec2.internal	Standby	ip-10-0-20-19.ec2.internal	Online
ip-10-0-20-18.ec2.internal	Online		

你可查看以下信息：

- HA 概况：显示当前 DC，分区，高可用性节点和高可用性资源
- HA 节点状况：显示高可用性节点主机或 IP 地址及其状态
- HA 详细状态信息：点击**查看当前状态**以查看高可用性节点状态。

4.2 更改高可用性配置

配置高可用性后，你可在 Intel® Manager for Apache Hadoop* software 的**集群配置 > 高可用性 > 简要配置**中重新配置。

你可以在配置页面进行以下操作：

- 储存：保存高可用性的新设置
- 重置：将新的设置重置为你在配置向导中设置的原设置。
- 展开所有：展开所有设置。
- 收起所有：按类别收起设置。



4.3 基本 DRBD 操作

4.3.1 资源概述

在 DRBD 中，资源是一个概括性的术语，特指某复制的数据集的所有方面。这些包括资源名称、数量和 DRBD 设备。

DRBD 资源配置文件保存在 `/etc/drbd.d` 目录下，目录名称和资源名称保持一致。比如，你可以输入以下命令查看定义为 `r0` 资源的配置内容：

```
vi /etc/drbd.d/r0.res
```

以下为内容示例：

```
resource r0
{
    device /dev/drbd0;

    on xtt-portal
    {
        disk          /dev/sdb1;
        address        10.239.47.81:7789;
        meta-disk       internal;
    }

    on xmlqa-clv9
    {
        disk          /dev/sdb1;
        address        10.239.47.35:7789;
        meta-disk       internal;
    }
}
```

这一示例对 DRBD 做出如下配置：

- 高可用性集群由二个节点组成，`xtt-portal` 和 `xmlqa-clv9`。
- 我们有一个名称为 `r0` 的资源，它使用 `/dev/sdb1` 作为底层存储，并已配置内部数据元。
- 资源使用 TCP 端口 7789 进行网络连接，并分别和 IP 地址 10.239.47.81 及 10.239.47.35 绑定。

4.3.2 检查 DRBD 状态

要通过 `/proc` 伪文件系统获取已存在的 DRBD 资源信息，输入以下命令：

```
cat /proc/drbd
```

本章节中的表 1、2、3、4 可帮助你理解 DRBD 状态的含义。

？ 1. DRBD 信息的重要缩写

缩写	含义	命令	状态
cs	连接状态 (Connection State)	drbdadm cstate <resource>	参见表 2
ro	角色 (Roles)	drbdadm role <resource>	本地资源角色 / 远程资源角色 参见表 3
ds	磁盘状态 (Disk State)	drbdadm dstate <resource>	本地资源状态 / 远程资源状态 参见表 4
p	复制协议 (Replication Protocol)		A、B 或 C

表 2. 主要连接状态

状态	含义
StandAlone	无可用的网络配置。资源尚未连接，或者因为管理原因连接中断，或者因为验证失败或脑裂连接未成功。
WFConnection	节点在等待对等节点在网络上可见。
WFReportParams	TCP 连接已建立，这一节点在等待对等节点发出的第一个网络数据包。
Connected	DRBD 连接已建立，数据镜像状态为活动。这是正常的工作状态。

表 3. 角色和含义

角色	含义
Primary	资源目前处于 Primary 角色，可进行读写。除非你启用了双主（dual-primary）模式，否则二个节点中只有一个可被赋予该角色。
Secondary	资源目前处于 secondary 角色。可以正常接收从对等节点发出的更新（除非处于 disconnected 模式），但不能被读写。其中一个或二个节点都可成为该角色。
Unknown	资源当前角色为 unknown。本地资源角色则永不会有这一状态。仅对等节点的资源角色会显示这一状态，而且仅在 disconnected 模式时显示。

表 4. 磁盘状态和含义

磁盘状态	含义
Diskless	没有本地块设备被分配到 DRBD 驱动。这可能表示资源从未被附加给它的备用设备，也就是说，它可能通过 drbdadm 命令被人工分离，或因为底层 I/O 错误被分离。
Inconsistent	数据不一致。这一状态在创建新的资源时同时在二个节点上立即产生（在初始 full sync 前）。而且，这一状态可在同步过程中的其中一个节点（同步对象）上找到。
Outdated	资源数据一致，但已过时。

表 4. 磁盘状态和含义

DUnknown	如果没有可用的网络连接，对等磁盘将处于这一状态。
Consistent	没有网络连接时数据是一致的。当连接建立后，数据状态将变为 UpToDate 或 Outdated。
UpToDate	数据状态为一致并已更新。这是正常的状态。

输入以下命令可得到更简洁的概述：

```
drbd-overview
```

4.3.3 常用命令

表 5 列出了常用的 DRBD 命令及其功？

表 5. 常用命令

命令	功能
drbdadm up <resource>	在主机上启用 DRBD 资源
drbdadm down <resource>	禁用 DRBD 资源
drbdadm primary <resource>	将 DRBD 资源转换为 Primary 模式
drbdadm secondary <resource>	将 DRBD 资源转换为 Secondary 模式
drbdadm attach <resource>	将资源附加到备用设备
drbdadm detach <resource>	将 DRBD 从备用设备分离
drbdadm connect <resource>	启用 DRBD 资源的网络连接
drbdadm disconnect <resource>	禁用 DRBD 资源的网络连接
drbdadm pause-sync <resource>	中断一个正在运行的再同步进程
drbdadm resume-sync <resource>	恢复再同步进程

4.4 Pacemaker 的基本操作

本章节介绍一些在高可用性维护中可能经常会用到的 Pacemaker 基本操作。你也可输入以下命令获得更多帮助：

```
crm help <command>
```

4.4.1 资源状态

在 Pacemaker 中，资源包括：fs_hadoop、ms_drbd_hadoop (drbd_hadoop)、NameNode、ip_hadoop、hive、mysqld、hive_metastore 和 jobtracker。

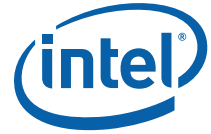
如果你仅需要知道所有资源的当前状态概况，输入以下命令：

```
crm resource show
```

要知道更详细的集群当前状态概况，包括节点、失败计数和非活动状态资源的信息，输入以下命令：

```
crm_mon -l -rf
```

如果你想要一个永久视图而不是一次性视图，去掉命令中的 -l，这样当有最新事件时，你的视图将被更新。



4.4.2 操作资源

要启动一个已停止的资源，输入以下命令：

```
crm resource start <resource>
```

要停止一个资源，输入以下命令：

```
crm resource stop <resource>
```

4.4.3 操作集群节点

如果你想要让集群中某个节点成为 *Standby* 模式，输入以下命令：

```
crm node standby <node>
```

要让集群中某个节点成为 *Online* 模式，以允许 Pacemaker 在节点上启动资源，输入以下命令：

```
crm node online <node>
```

如果你在想要的节点上执行命令，去掉命令中的 *<node>*。

4.4.4 维护模式

如果你要对集群进行维护，并且希望 Pacemaker 暂时不干预集群的工作，输入以下命令启用维护模式：

```
crm configure property maintenance-mode="true"
```

要重新启用 Pacemaker 的集群管理，输入以下命令关闭维护模式：

```
crm configure property maintenance-mode="false"
```

4.4.5 手工更改 CRM 配置

要编辑当前集群配置，在配置过程中输入以下命令启动 CRM shell：

```
crm configure
edit
```

4.5 处理硬盘故障

下述步骤适用于你直接在物理硬盘上运行 DRBD 的情形。

1. 使用以下命令查看你配置的 I/O 错误处理策略：

```
vi /etc/drbd.d/global_defaults.conf
```

策略将在以下 *<strategy>* 部分显示：

```
disk {
    on-io-error <strategy>;
    ...
}
```

如果策略是 *detach*，你可以跳过此步骤。如果策略是 *pass_on*，你需要手工将 DRBD 从你的硬盘分离。

```
drbdadm detach <resource>
```

2. 替换某个有故障的磁盘的过程包括创建一个新的数据元并重新附加资源。如果你使用内部数据元，输入以下命令：



```
drbdadm create-md <resource>
```

```
drbdadm attach <resource>
```

3. 如果你使用外部数据元，DRBD 将不能独立判断硬盘是否已交换，因此，你需要输入以下命令执行更多步骤：

```
drbdadm inAvalidate <resource>
```

4.6 节点故障后的 NameNode 恢复

故障排除后，active 状态的 NameNode 可能已更改。你可以手工恢复主节点，执行以下步骤：

1. 在要成为新的主节点的节点上执行以下命令：

```
service drbd start
```

```
service corosync start
```

```
service pacemaker start
```

```
crm node online
```

2. 找到之前的主节点，执行以下命令：

```
crm node standby
```

3. 要查看 DRBD 资源：

```
crm resource status ms_drbd_hadoop
```

4. 要查看主节点的 DRBD 分区是否已成功挂载在 `/hadoop/drbd` 目录：

```
crm resource status fs_hadoop
```

5. 要查看虚拟 IP 地址是否已绑定到主节点：

```
crm resource status ip_hadoop
```

6. 最后，在当前的 secondary NameNode 上执行以下命令：

```
crm node online
```

4.7 脑裂恢复

脑裂是由于集群节点之间的网络连接暂时失败造成的，原因可能是集群管理软件的干预，或人工操作失误，双方节点在连接失败后都转换成 primary 角色。DRBD 在连接重新可用时会探测到脑裂和对等节点。

交换初始 DRBD 握手协议

执行以下命令：

```
vi /var/log/messages
```

以下消息将显示：

```
Split-Brain detected, dropping connection!
```

这表示存在脑裂故障。Primary NameNode 和 Standby NameNode 上的数据不再同步。你可在 Standby NameNode 上输入以下命令，以检查 DRBD 服务状态：

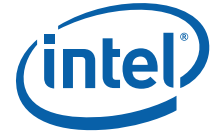
```
service drbd status
```

在此情形下，你必须完成以下步骤，以删除不一致的数据：

1. 删除不一致的数据，在 Standby NameNode 上输入以下命令：

```
drbdadm disconnect r0
```

```
drbdadm secondary r0
```



```
drbdadm -- --discard-my-data connect r0
```

2. 将 Primary NameNode 和资源连接，在 Primary NameNode 上输入以下命令：

```
drbdadm connect r0
```

3. 在 Standby NameNode 上启动 DRBD 服务：

```
service drbd start
```

4. 确认 Standby NameNode 上的 DRBD 服务状态为启动：

```
service drbd status
```

4.8 运行脚本排除故障

英特尔® Apache Hadoop* 软件发行版提供用于高可用性故障排除的一些脚本文件。你可以在以下目录中找到这些脚本 `/usr/lib/intelcloudui/tools/`。

在 CentOS 操作系统上：

- `./ha/CentOS`
- `./ha/CentOS/operate_single_resource.sh`
- `./ha/CentOS/operate_single_node.sh`
- `./ha/CentOS/clean_cluster_settings.sh`
- `./ha/CentOS/clean_single_settings.sh`

用于 DRBD 故障排除：

- `./drbd/recover-drbd.sh`

`operate_single_resource.sh` 和 `operate_single_node.sh` 脚本可在单点资源或节点上恢复高可用性设置，而卸载脚本无此功能。输入以下命令以运行脚本：

```
sh <script file name>
```

`clean_cluster_settings.sh` 和 `clean_single_settings.sh` 脚本可清除某个集群或单个节点的高可用性设置，而卸载脚本无此功能。输入以下命令运行脚本：

```
sh <script file name>
```

你可使用 `recover-drbd.sh` 脚本手工恢复 DRBD。如果 Primary NameNode 发生故障进行机器切换时，DRBD 分区可能会损坏。你可以使用脚本来恢复 DRBD。英特尔® Apache Hadoop* 软件发行版需要自动探测这一情形何时发生并自动运行脚本以清理损坏的节点。你必须首先确定一个处于良好状态的 DRBD 分区。输入以下命令运行脚本：

```
sh recover-drbd.sh
```



5.0 卸载高可用性

要卸载高可用性组件，在 Intel[®] Manager 中执行以下步骤。

1. 确认集群中所有的服务都已停止。如果没有，现在停止这些服务。
2. 点击**配置向导**按钮开始集群配置向导。
3. 在**第 1 步**页面，不勾选**高可用性**选项。
4. 根据配置向导完成各步骤。
5. 执行部署服务属性的步骤。
6. 在集群中所有节点上执行以下命令：
 - `rm -rf /etc/corosync/*`
 - `rm -rf /var/lib/heartbeat/crm/*`
 - `rm -rf /hadoop/ hadoop_image_local/*`