



英特尔® Hadoop 发行版
版本 2.2
新手指南

Contents

1. 简介.....	1
1.1 文档目的	1
1.2 产品简介	1
1.3 集群结构	2
1.4 构造集群的主要步骤	3
2. 系统要求.....	4
2.1 硬件要求	4
2.2 软件要求	5
2.3 网络要求	5
2.3.1 链路聚合	5
3. 规划 Hadoop 集群.....	7
4. 对集群中所有节点进行操作系统安装.....	8
4.1 安装操作系统	8
4.2 磁盘分区	8
5. 在管理节点安装英特尔® Hadoop 发行版	10
6. 集群配置.....	12
6.1 登录并接受用户许可协议	12
6.2 集群安装配置向导	12
6.3 输入许可证	26
6.4 配置节点	27
6.5 启动集群	29
6.6 手动配置部分组件	33



1. 简介

1.1 文档目的

本文档用于指导英特尔® Hadoop 发行版初级用户安装、部署、验证和开始使用英特尔® Hadoop 发行版。

1.2 产品简介

英特尔® Hadoop 发行版供给客户稳定可靠易用的 HDFS、MapReduce、HBase 和 Hive 等框架的商业版本。具体套件及其特点如表 1.1 所示。

商业套件	功能和特点
分布式协同工作系统 (Zookeeper)	<ul style="list-style-type: none">● 高效的选举算法● 确保分布式系统一致性● 保证集群数据及配置同步● 实现统一命名服务
分布式文件系统 (HDFS)	<ul style="list-style-type: none">● 可自我修复的高带宽集群文件存储系统● 高可扩展性，无需停机无缝动态扩容● 高容错性，数据自动复制和校验● 改进的可靠性和扩展性
分布式数据库 (HBase)	<ul style="list-style-type: none">● 分布式、面向列、多维度的数据库系统。● 数据自动切分和分布存储● 高可扩展性，无宕机线性扩容● 高性能并发读写
分布式计算框架 (Map/Reduce)	<ul style="list-style-type: none">● 高度并行和可扩展的分布式批处理计算框架● 高容错能力，支持任务自动迁移和重试● 公平调度算法，支持任务抢占，兼顾长短任务● 调度任务到最近的数据节点，有效降低网络带宽● 灵活的资源分配和调度，达到资源利用最大化
分布式数据仓库 (Hive)	<ul style="list-style-type: none">● 高性能分布式海量数据仓库● 强大的查询与分析功能● 类 SQL 查询语言
	<ul style="list-style-type: none">● 分布并行运行分析任务

1. 简介

分布式数据分析平台 (Pig)	<ul style="list-style-type: none"> 把类 SQL 数据分析请求转换为经优化的 Map/Reduce 运算
机器学习类库 (Mahout)	<ul style="list-style-type: none"> 可扩充能力的机器学习类库 实现了一些可扩展的机器学习领域经典算法 有效地使用 Map/Reduce 实现高性能计算
分布式海量日志处理系统 (Flume)	<ul style="list-style-type: none"> 高效采集、聚合、迁移海量日志数据 高可靠性及可管理性 可扩展性，三层架构水平扩展
结构化数据源转移系统 (Sqoop)	<ul style="list-style-type: none"> 分布并行导入关系型数据库中数据

表 1.1 英特尔® Hadoop 发行版提供商业套件的功能和特点

1.3 集群结构

集群由管理节点、Hadoop 集群以及客户端组成。图 1.1 描述了集群的结构。

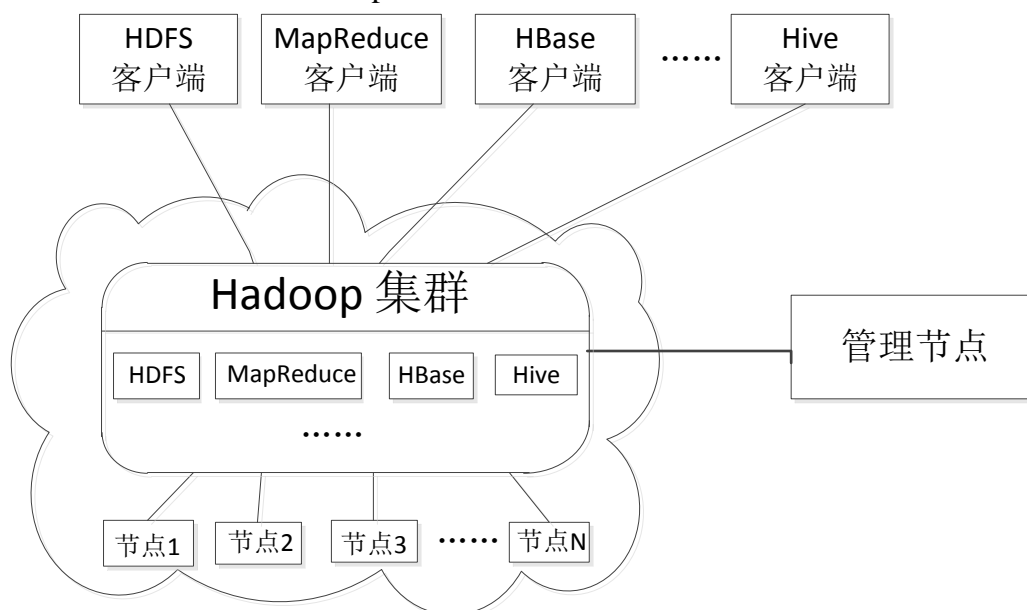


图 1.1 集群结构

客户端包括 HDFS 客户端，MapReduce 客户端，HBase 客户端，Hive 客户端等。这些客户端能被运行在一个或多个服务器上。

Hadoop 集群拥有的组件包括 Zookeeper, HDFS, MapReduce, HBase, Hive 等。这些组件需要一个以上的服务器以实现其功能。而每一个服务器可以运行一个或一个以上组件服务。

管理节点能够监视和管理 Hadoop 集群和客户端。

1.4 构造集群的主要步骤

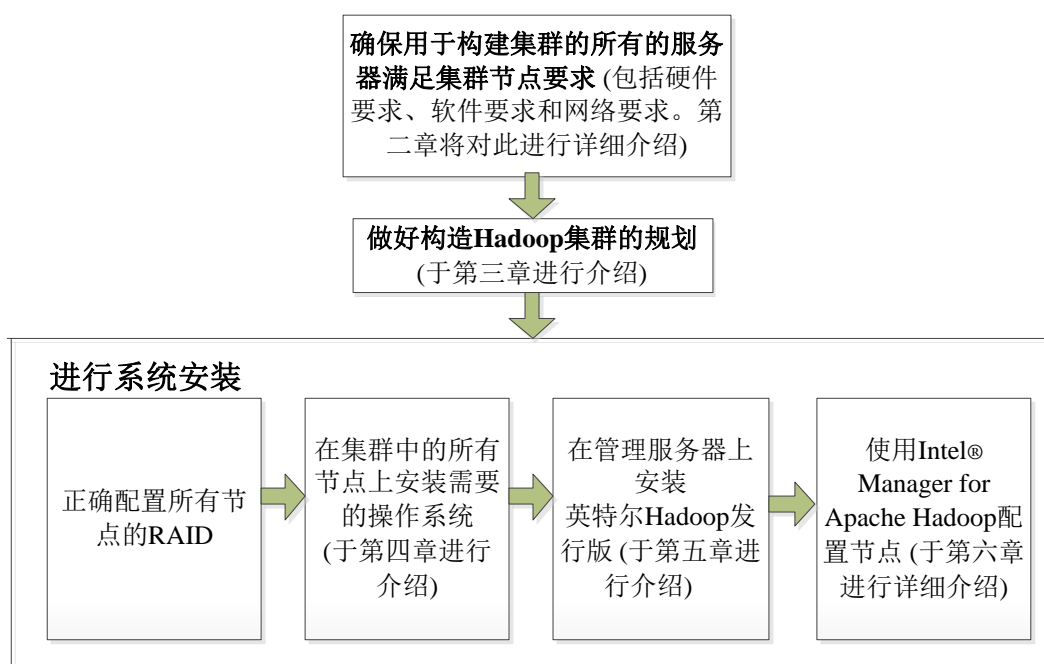
创建一个集群前，首先必须保证将用于构造集群的服务器满足一些要求。这些要求包括硬件要求、软件要求以及网络要求。第二章将对这些要求进行详细介绍。

满足了这些要求后，就要做好相关的构造 Hadoop 集群的规划。第三章介绍了在系统安装前必须做好的规划。

规划 Hadoop 集群后，就可以开始进行系统安装了。系统安装包括以下步骤：

1. 正确配置所有节点的 RAID。
2. 在集群中所有的节点（包括管理节点及 Hadoop 集群中所有的节点）上安装所需要的操作系统。这将在第四章中进行介绍。
3. 在管理节点上安装英特尔® Hadoop 发行版。第五章将对此进行详细介绍。
4. 使用管理节点上的 Web 用户界面——Intel® Manager for Apache Hadoop——来把所有的节点加到集群中、对这些节点部署相关软件和配置，并运行服务。第六章将给出相关的详细步骤。

图 1.2 描述了构造一个集群的主要步骤。



2. 系统要求

本章介绍了集群中不同节点必须满足的硬件要求、软件要求和网络要求。

2.1 硬件要求

服务器运行英特尔® Hadoop 发行版至少需要英特尔® 至强处理器，推荐使用双路 4 核英特尔® 处理器。

服务器运行英特尔® Hadoop 发行版的最低内存要求为 16GB 内存。在此基础上，不同服务器角色和服务类型有着各自的内存要求，如表 2.1 所示。推荐内存配置是针对服务器上运行的服务种类，对表 2.1 中的相关内存要求进行叠加。

服务器角色及服务类型	内存要求
管理节点	8GB
Hadoop 集群:	
• MapReduce Job Tracker	2GB
• MapReduce Task Tracker	2GB
• MapReduce Slots on Task Tracker	512MB * slot 数量
• HDFS NameNode	16GB
• HDFS Secondary NameNode	16GB
• HDFS DataNode	2GB
• ZooKeeper	4GB
• HBase Master Server	2GB
• HBase Region Server	16GB
• Hive Server	2GB
客户端	8GB

表 2.1 基于服务器角色及服务类型的运行英特尔® Hadoop 发行版内存要求

用户可以通过简单地叠加相应服务需要的内存要求来计算推荐的内存要求。比如一个服务器计划运行如下服务：HDFS DataNode, MapReduce TaskTracker 和 HBase Region Server。同时计划的 slot 数量（包括 map slots 和 reduce slots）为 16。这样，对于这个服务器的推荐内存为：2GB + 2GB + 512MB*16 + 16GB = 28GB。

对于所有集群中的服务器（除了主命名节点和从命名节点），推荐在物理硬盘中不要使用 RAID。但在 RAID 无法被移除的情况下，每一个物理硬盘可以被



设为一个单独的 RAID 0。

例子: IBM 服务器 LSI MegaRaid 的 RAID 设置

首先, 每个物理硬盘应该被指定为一个 Drive Group
然后, 每个 Drive Group 应该被指定为一个 Span
最后, 每个 Span 应该被制定为一个基于 RAID 0 的一个 Virtual Drive

对于主命名节点和从命名节点, 推荐在命名节点数据目录所在的分区使用 RAID 1 或 RAID 5。

2.2 软件要求

本节介绍安装英特尔® Hadoop 发行版所需要的软件环境。
支持的操作系统包括:

1. Red Hat Enterprise Linux 5.7
2. CentOS 5.7
3. Red Hat Enterprise Linux 6.1, 6.2, 6.3
4. Oracle Enterprise Linux 6.1, 6.2, 6.3
5. CentOS 6.1, 6.2, 6.3
6. SUSE* Linux Enterprise Server 11 sp1

在把一台服务器加入集群前, 必须先在该服务器上安装相应的操作系统, 并且须要确保 opnssh-server 在运行。如果 opnssh-server 没有运行, 请在 Hadoop 集群中的所有节点中安装 openssh-server 包被确保该服务运行。

2.3 网络要求

安装英特尔® Hadoop 发行版对网络的最低要求为千兆以太网。有关网络方面的其他要求见下。

2.3.1 链路聚合

当一台机器上有多个网络适配器时, 用户可以在安装英特尔® Hadoop 发行版之前对其进行链路聚合配置以提高网络带宽。



3. 规划 Hadoop 集群

在安装英特尔® Hadoop 发行版前，必须制定一些相关的规划。具体如下：

1. 决定将使用哪些 Hadoop 组件。这些组件包括 Zookeeper, HDFS, MapReduce, HBase, Hive, HA 等。
2. 决定集群大小（服务器数量）和服务器硬件配置。
3. 决定集群物理布局。确定将使用多少机架和每个机架上有多少机器。
4. 决定集群网络。包括：
 - ① 决定网络带宽和交换机背板带宽。决定交换机型号。
 - ② 决定如何连接到交换机。必须知道需要用到哪些以太网端口和是否需要绑定。
 - ③ 确定每台机器的 IP 地址和主机名。决定如何分配 IP（使用 DHCP 或静态分配）。决定如何解析主机名（使用 DNS 或/etc/hosts）。如果使用/etc/hosts，管理节点将负责更新集群中每台机器的/etc/hosts。
 - ④ 决定如何进行时间同步。管理节点将负责所有服务器上的时间的同步，但您需要决定是否使用外部的 NTP 服务。如果不使用外部 NTP 服务，虽然集群中所有服务器的时间是相同的，但这个时间有可能不是标准时间，这有可能导致当 Hadoop 集群与外部连接时错误的产生。
5. 决定哪台机器（当需要 HA 时为机器对）是 NameNode 和 JobTracker（当需要 MapReduce 时）。对于较大的集群，还需要两台机器来分别充当 NameNode 和 JobTracker。此时，如需 HA 的话，就需要 4 台机器。
6. 决定哪台机器为管理节点。
7. 大致决定由哪些机器用于构建 Hadoop 集群、哪些机器用于构建客户端。



4. 对集群中所有节点进行操作系统安装

在安装英特尔® Hadoop 发行版之前，集群中的所有节点必须满足 2.2 节所列举的所有要求。

4.1 安装操作系统

使用者可以使用两种方式来安装集群中的服务器的操作系统，单独安装方式和 PXE 安装方式。

单独安装方式是使用 Red Hat Enterprise Linux for Servers、CentOS 64、SUSE Linux Enterprise Server 或 Oracle Enterprise Linux 的安装光盘在每台服务器上独立安装操作系统。

PXE 安装方式是使用基于英特尔® 服务器解决方案的 Preboot Execution Environment 技术。集群中的服务器通过网络从 PXE 服务器上自动下载操作系统镜像进行安装。本文档暂不对 PXE 安装方式详细步骤进行介绍。

4.2 磁盘分区

在硬盘分区时需要遵守以下几点原则：

1. 至少要分出 boot, swap 和加载于 “/” 的系统分区。
2. 包含有操作系统的 root 目录的根分区需要 100GB 以上空间。推荐此分区使用 ext4 文件系统。
3. 推荐把每个物理磁盘挂载为在 /mnt/disknn (nn 为 1 至 2 位的数字) 上不同的挂载点。DataNode 上每个这样的目录会被管理节点自动配置为 HDFS 的数据目录。建议使用 ext4 文件系统。
4. HDFS DataNode 的数据目录不能放在 root 分区，以避免空间不足和 IO 竞争。同时也建议不要将他们放在 root 分区所在的系统盘以避免 IO 竞争。但是当磁盘空间不足时，可以在系统盘的剩余空间中创建一个新的分区，通过#2 所描述的方式挂载来作为数据目录。

例子：有一台机器有两个 500GB 大小的磁盘，请如下将其进行分区：

驱动	大小	挂载点	类型
/dev/sda1	100GB	/	ext4
/dev/sda2	400GB	/mnt/disk1	ext4
/dev/sdb1	500GB	/mnt/disk2	ext4



4. 对集群中所有节点进行操作系统安装

5. 如果需要高可用性，则需要对这些用于高可用性的机器配置备用分区，这个分区必须略大于系统内存，并且将该分区设置成**非自动挂载(删除/etc/fstab 中对应的条目)**。同时，两台机器上的备用分区大小必须一致。如果您需要配置高可用性，建议在操作系统安装完成之后再划分该备份分区，防止该分区被设置成自动挂载。

例子：有两台用于高可用性配置的机器各拥有两个 500GB 大小的硬盘，当前的内存为 48GB，但它有可能变为 64GB。所以他们的备用分区应该至少为 64GB。请如下对其进行分区：

驱动	大小	挂载点	类型
/dev/sda1	100GB	/	ext4
/dev/sda2	336GB	/mnt/disk1	ext4
/dev/sda3	64GB		
/dev/sdb1	500GB	/mnt/disk2	ext4

其中，/dev/sda1 被作为系统分区，/dev/sda3 被作为备用分区，备用分区不需要挂载，将通过 DRBD 服务来挂载。

5. 在管理节点安装英特尔® Hadoop 发行版

英特尔® Hadoop 发行版需要通过安装光盘在管理节点上进行安装。具体步骤如下：

1. 插入光盘，在终端窗口输入以下命令进入安装目录并开始安装。

```
./install
```

2. 检查时间、日期和时区。如果这些信息是正确的，选择 yes 继续。

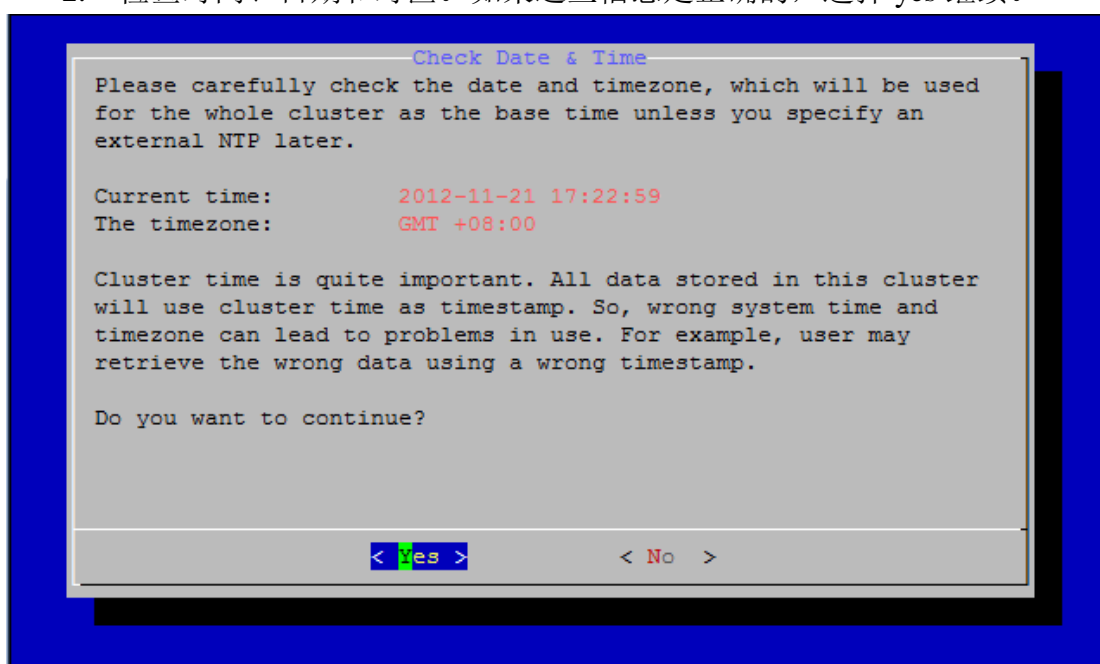


图 5.1 检查时间、日期和时区

3. 检查主机名与全称域名是否相同，如果正确，选择 Yes 继续，否则请根据提示进行修改。如果无法解析出本机 hostname，请修改/etc/hosts 根据 IP 添加本机 hostname 条目。

```
127.0.0.1    localhost.localdomain localhost
::1         localhost localhost.localdomain localhost6 localhost6.localdomain6
192.168.1.71 intelidh-01
192.168.1.73 intelidh-03
192.168.1.74 intelidh-04
192.168.1.75 intelidh-05
192.168.1.76 intelidh-06
192.168.1.77 intelidh-07
```

图 5.2 根据本机 IP 修改 hostname

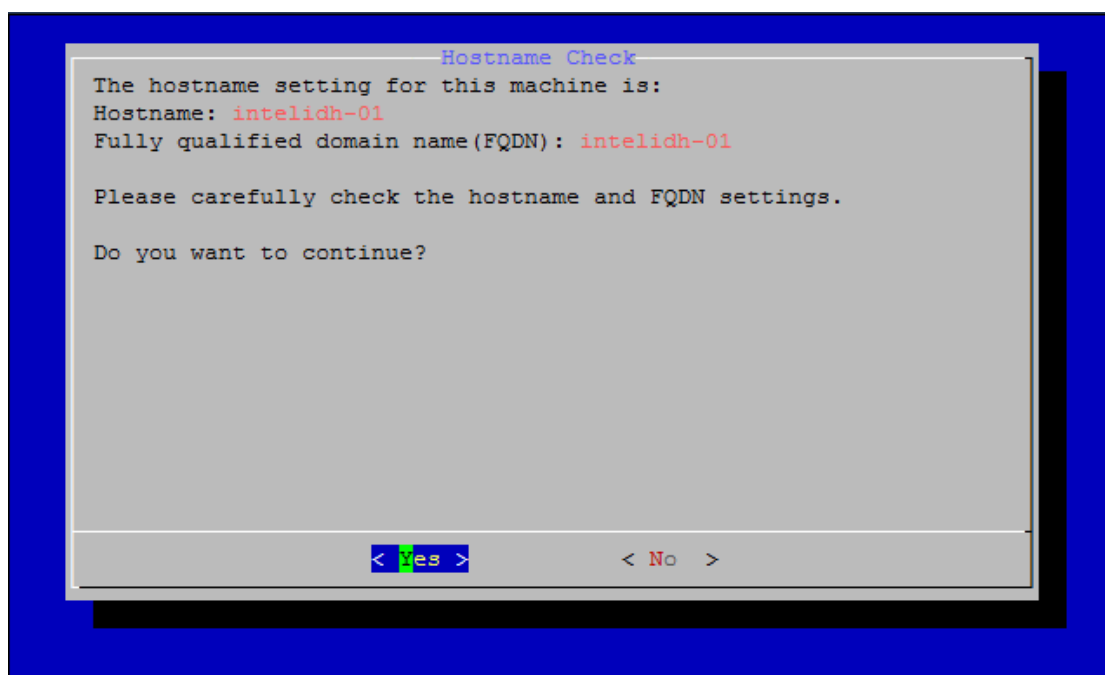


图 5.3 检查主机名

4. 开始部署 Intel® Manager for Apache Hadoop, 完成后选择 Continue 继续。

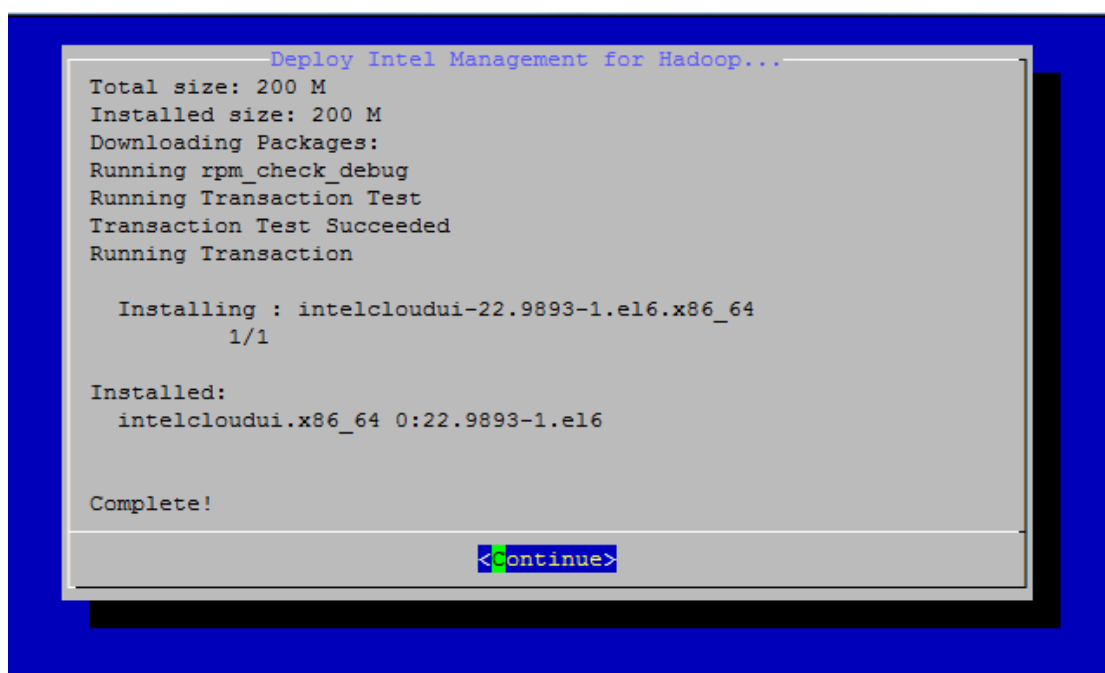


图 5.4 部署 Intel® Manager for Apache Hadoop

5. 如果电脑只有一个网络接口, 您将会跳过这一步。如果电脑有多个网络接口, 请选择 Intel® Manager for Apache Hadoop 将使用的网络接口, All 表示绑定所有接口。用户 (而不是集群中的服务器) 通过这个网络接口访问 Intel® Manager for Apache Hadoop。一般不需要考虑这个区别。但在内外网分离, 并且



5. 对管理节点进行英特尔® Hadoop 发行版安装

管理节点正好在内外网连接的网关上时要注意这个区别。

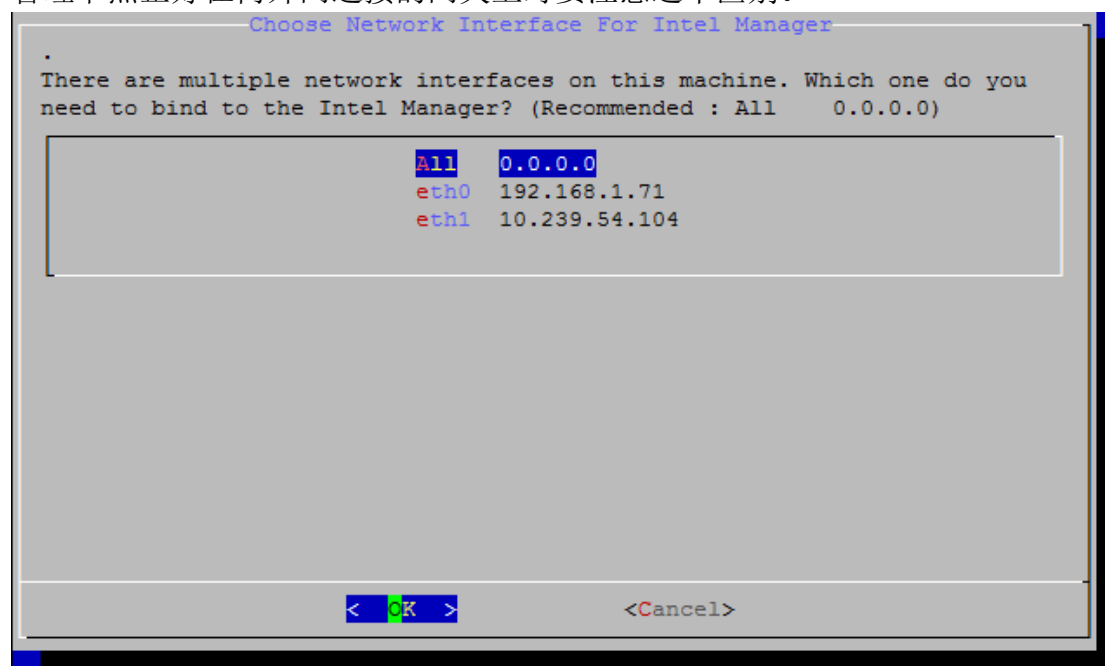


图 5.5 选择 Intel® Manager for Apache Hadoop 使用的网络接口

6. 确认绑定的网络接口

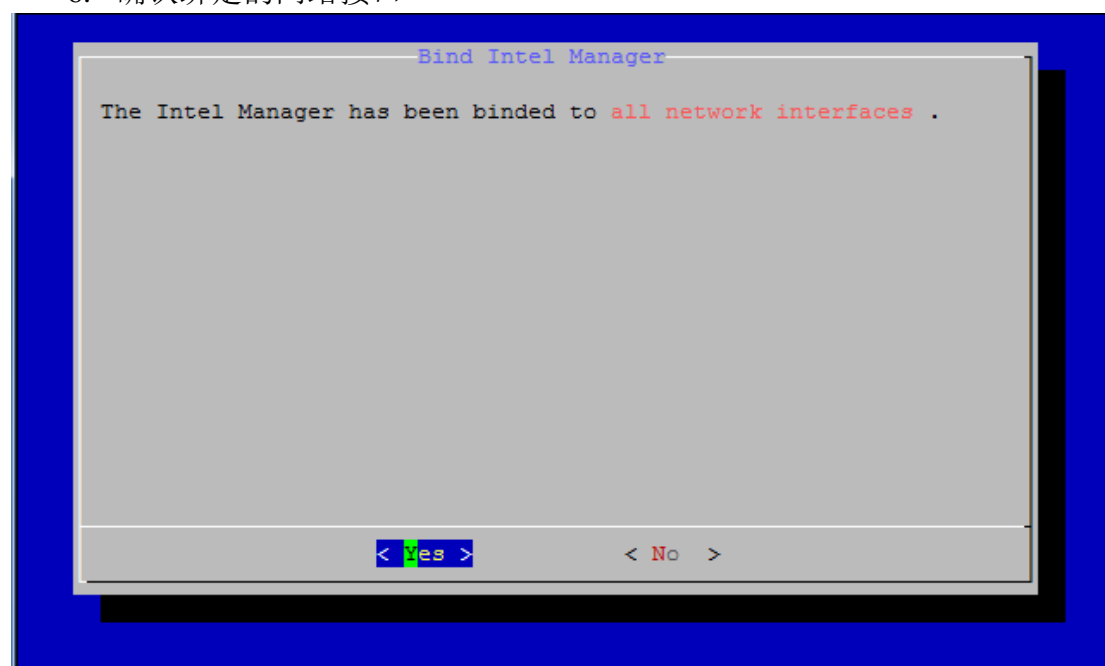


图 5.6 确认 Intel® Manager for Apache Hadoop 使用的网络接口

7. 安装英特尔® Hadoop 发行版和进行集群管理需要配置 Linux OS Package 资源库，选择是否清理节点上已存在的条目。

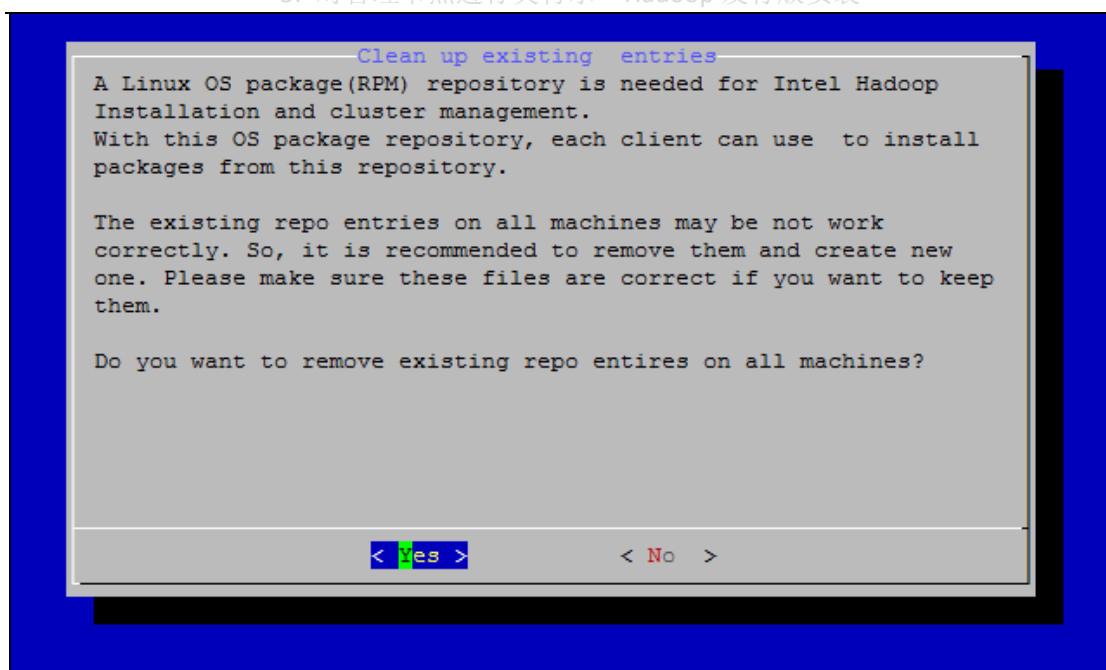


图 5.7 确认清理 Linux OS Package 资源库的现有条目

8. 安装英特尔® Hadoop 发行版和进行集群管理需要一个 Linux 系统的软件包资源库。您可以选择在本地主机上创建一个资源库，或使用一个存在的资源库。如果您选择在本地主机上创建，您则需要一个 Linux 系统的安装 DVD 光盘或者 ISO 文件来创建资源库。如果您选择使用一个存在的资源库，须要输入已存在 Linux 系统资源库的 URL。

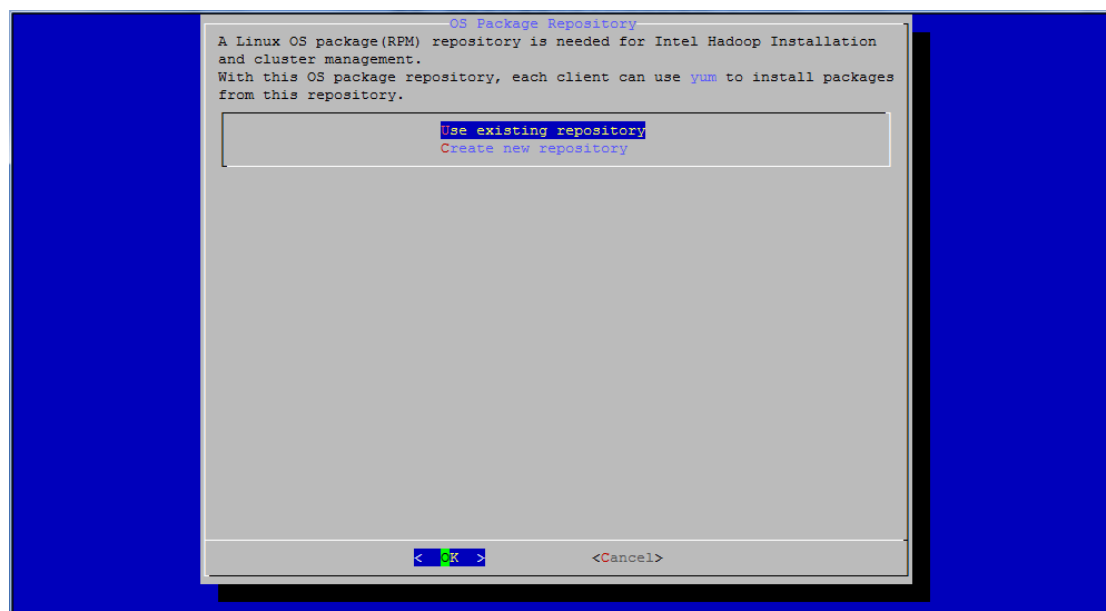


图 5.8 使用 Linux 系统安装 DVD 创建 yum 目录

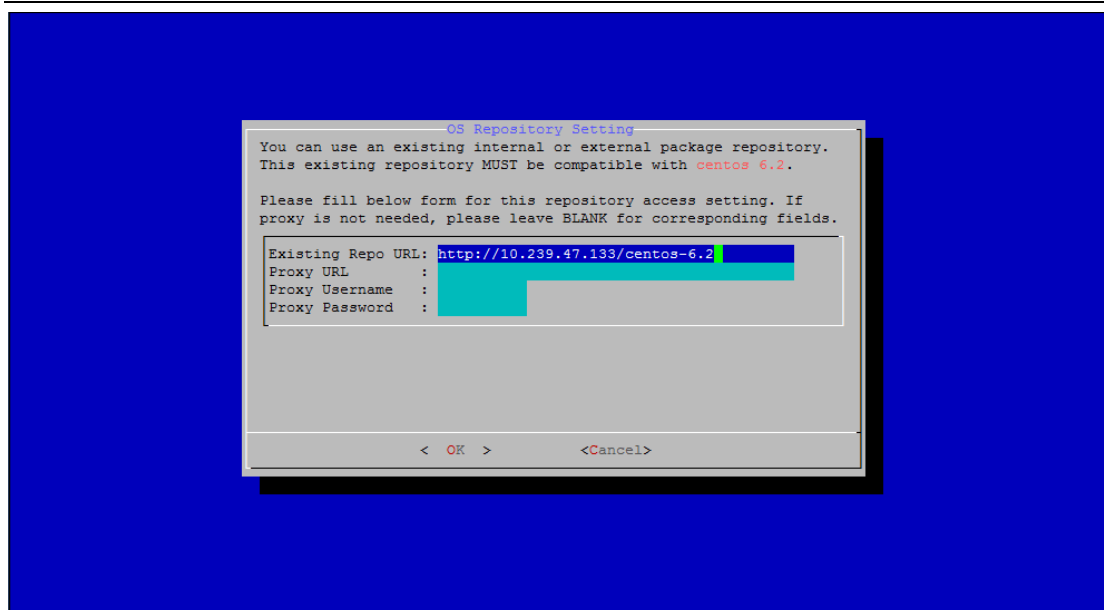


图 5.9 使用一个存在的目录

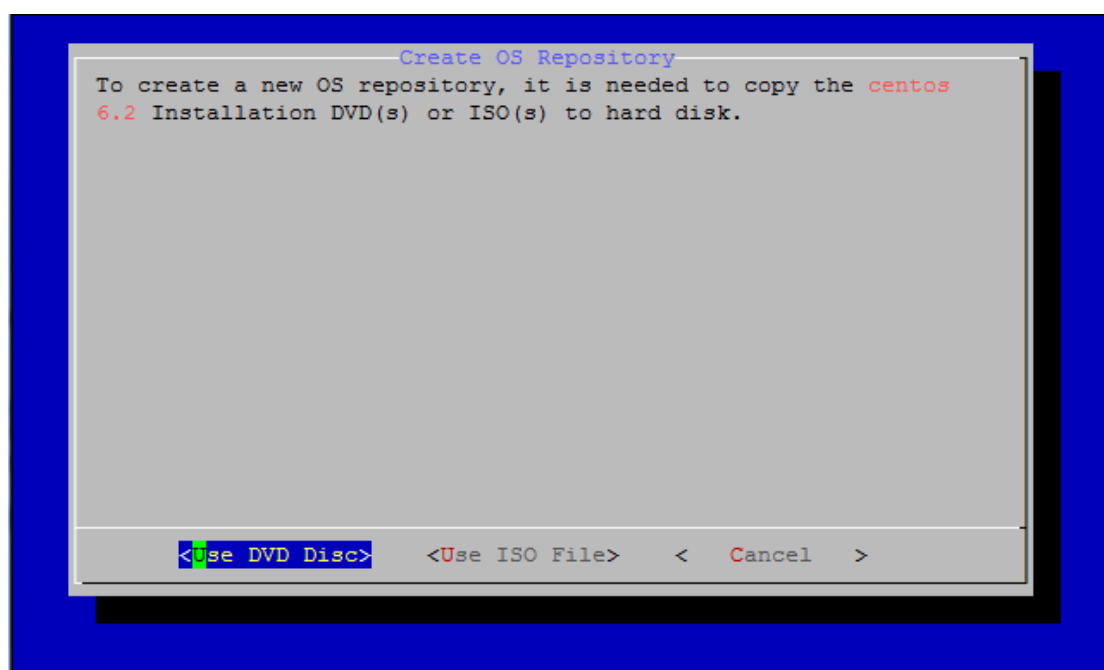


图 5.10 创建一个新的资源库目录

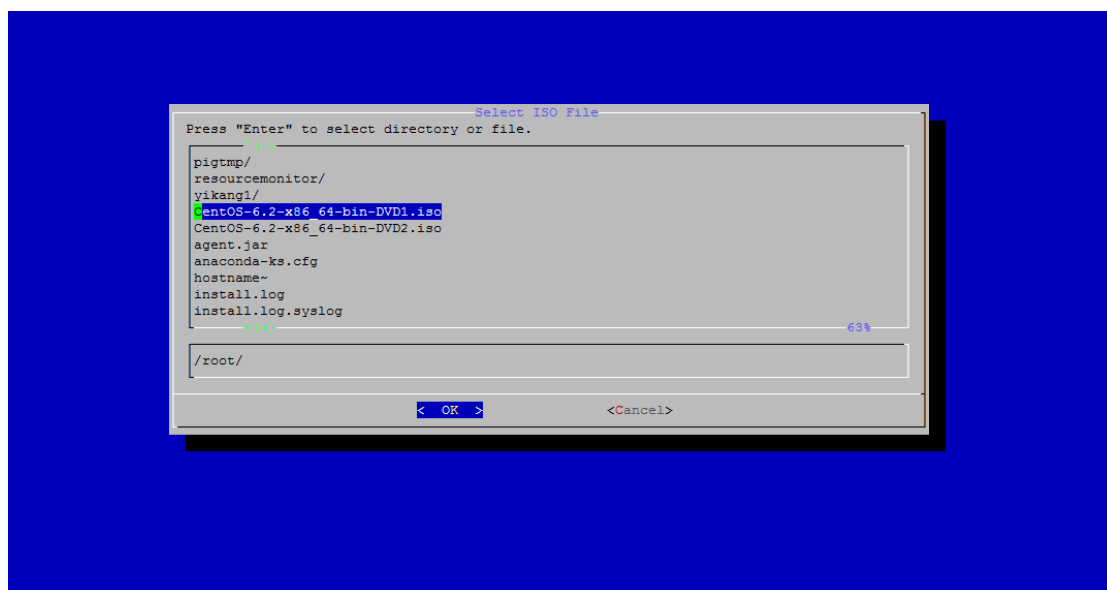


图 5.11 指定资源库 ISO 文件

如需要从多个 ISO 中复制文件，请选择 Yes 并指定新的 ISO 文件。

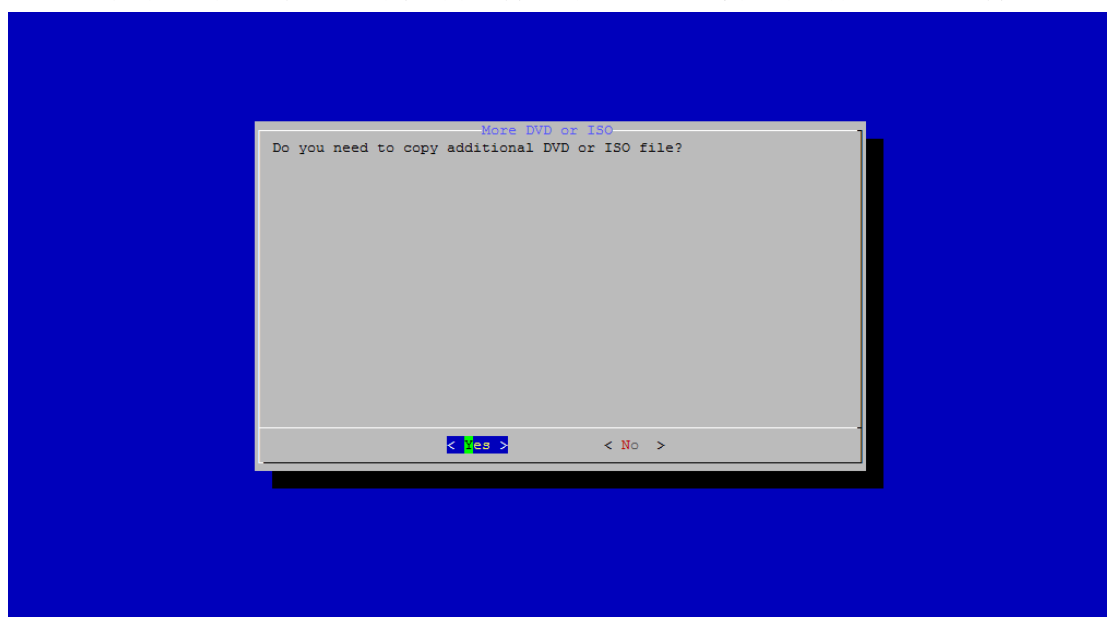


图 5.12 指定多个 ISO 文件

9. 安装英特尔® Hadoop 发行版和进行集群管理需要在本机上配置一个资源库地址以加速集群中其他节点的安装。如果有多个网络接口，您须要指定一个网络接口来进行。必须确保集群中的所有服务器都可以连接到这个网络接口。

5. 对管理节点进行英特尔® Hadoop 发行版安装

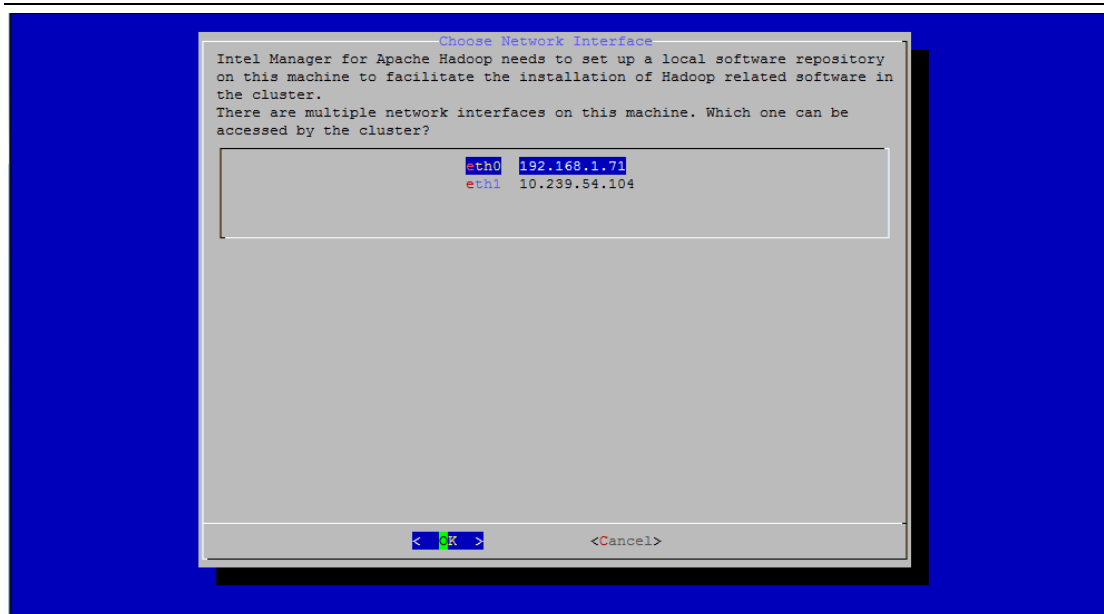


图 5.13 指定本地软件资源库使用的网络接口

10. 安装程序会自动安装软件包，Intel® Manager for Apache Hadoop 并配置 Puppet 服务器，点击 Continue 继续。

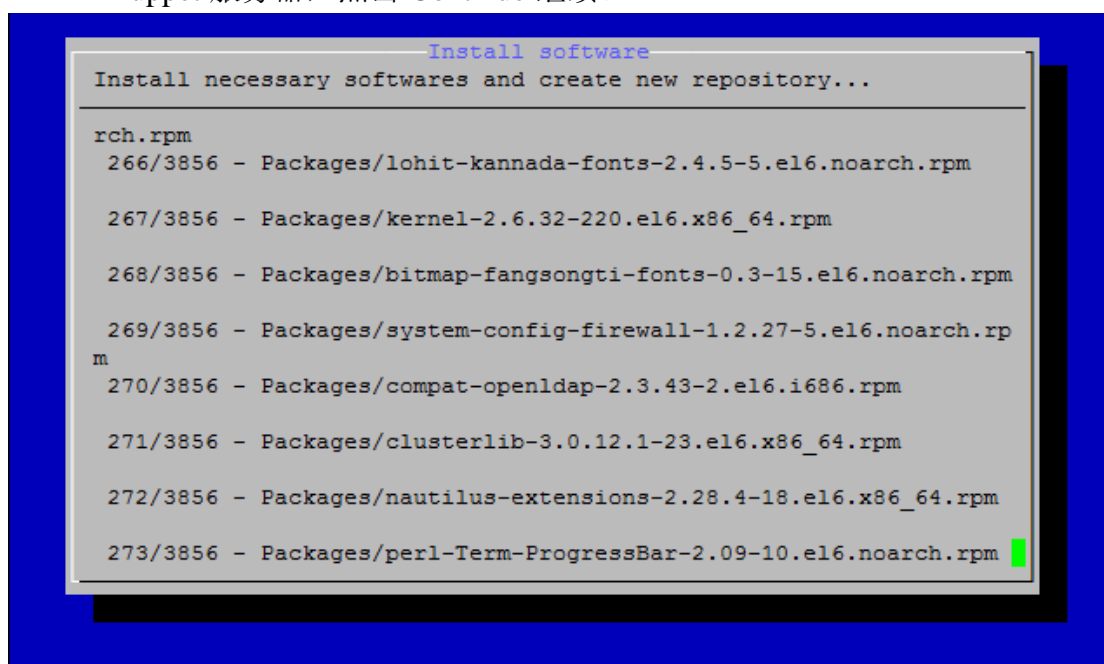


图 5.14 安装相关软件包

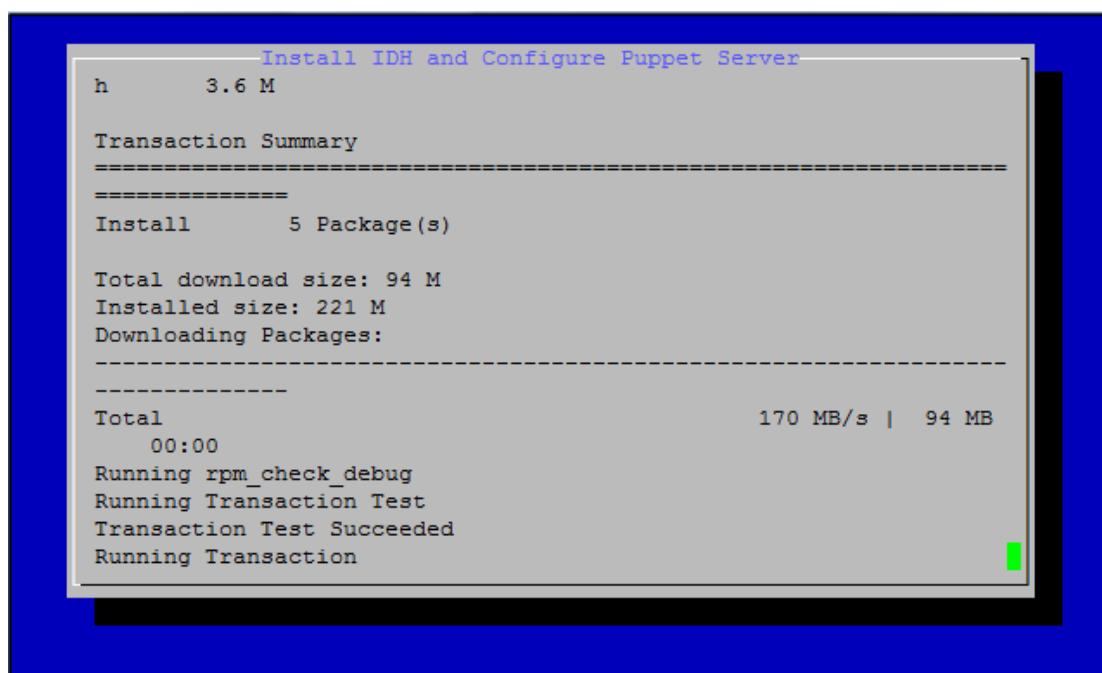


图 5.15 安装 IDH 并配置 Puppet 服务器

11. 当 FINISH 出现在终端窗口时, 英特尔® Hadoop 发行版已经被成功安装。

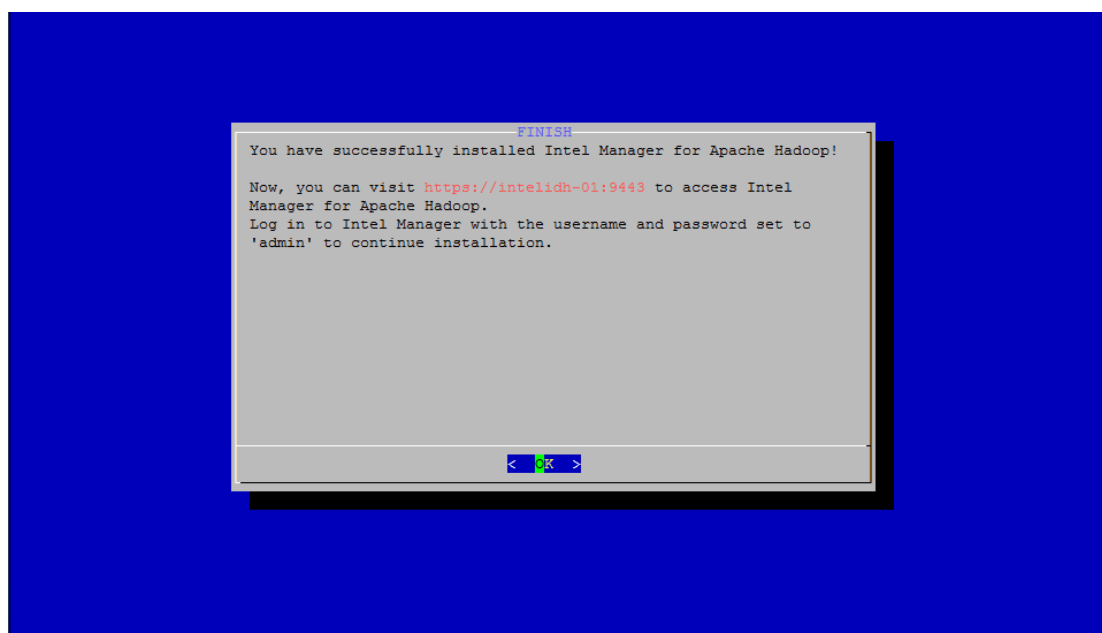


图 5.16 成功安装英特尔® Hadoop 发行版

6. 集群配置

6.1 登录并接受用户许可协议

安装完英特尔® Hadoop 发行版后，管理员可以通过管理界面——Intel® Manager for Apache Hadoop 来完成配置。建议使用 Chrome 或 Firefox 5 以上版本，并使用窗口大小为 1024 * 768 以上的浏览器来开启管理界面。

通过浏览器访问英特尔® Hadoop 发行版管理界面地址 <https://管理节点地址:9443>，输入用户名和密码登录后，英特尔® 软件最终用户许可协议将会出现。请接受许可协议。



图 6.1 英特尔® 软件最终用户许可协议界面

6.2 集群安装配置向导

接受英特尔® 软件最终用户许可协议后，将出现欢迎使用英特尔® Hadoop 管理中心的提示窗口。选择“是的，我想立刻使用配置向导帮助我配置”，如下图所示，并点击“确认”按钮。

第一步，在配置新的集群界面中输入集群名称和集群中所要安装的组件。“高可用性”仅限试用版与商业版用户使用。如果您须要配置集群的高可用性，请参阅《英特尔® Hadoop 发行版 2.2 HA 操作手册》。点击“下一步”继续。

第1步

配置新的集群

集群名称：

Cluster

选择集群中将会使用的组件，包括HDFS，MapReduce，HBase，Hive，Sqoop，Pig 和Flume；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：HBase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。

☐ 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

图 6.2 配置新集群

第二步，首先选择网络环境。然后您可以添加或删除节点。如果配置集群中的节点可以通过用主机名互相访问，用户必须配置有效的 DNS 服务器或已在集群中的所有机器上配置好了/etc/hosts 文件。否则请选择不能通过主机名访问，Intel® Manager 会相应为您配置/etc/hosts 文件。

注：用户必须指定合法的 **hostname**（如以字母开头），保证 **hostname** 解析出的 **IP** 地址不能为 127.0.0.1。

13

第2步

指定集群节点以及网络环境

网络环境：集群中的节点能通过主机名互相访问(通过配置好的DNS服务器或/etc/hosts文件)

节点名称	节点IP	状态
inteliidh-01	192.168.1.71	已连通

添加节点
删除节点

上一步
下一步
取消

图 6.3 指定集群节点以及网络环境

如果您点击的“添加节点”的按钮，添加机器的窗口就会出现。有两种方法可以添加节点——单个增加或批量增加。如果选择了批量增加，如图所示，在起始 IP 地址和结束 IP 地址栏目里输入需要搜索的 IP 地址段，并输入 root 密码。在界面中点击开始查找按钮，进入检测机器界面。

添加机器

添加机器。可以添加单台机器，也可以添加指定IP地址范围内的机器，需要提供root用户的凭证。

添加方式：批量添加

起始IP地址：192.168.1.71

结束IP地址：192.168.1.77

Root用户密码：●●●●●●

开始查找
取消

图 6.4 添加机器

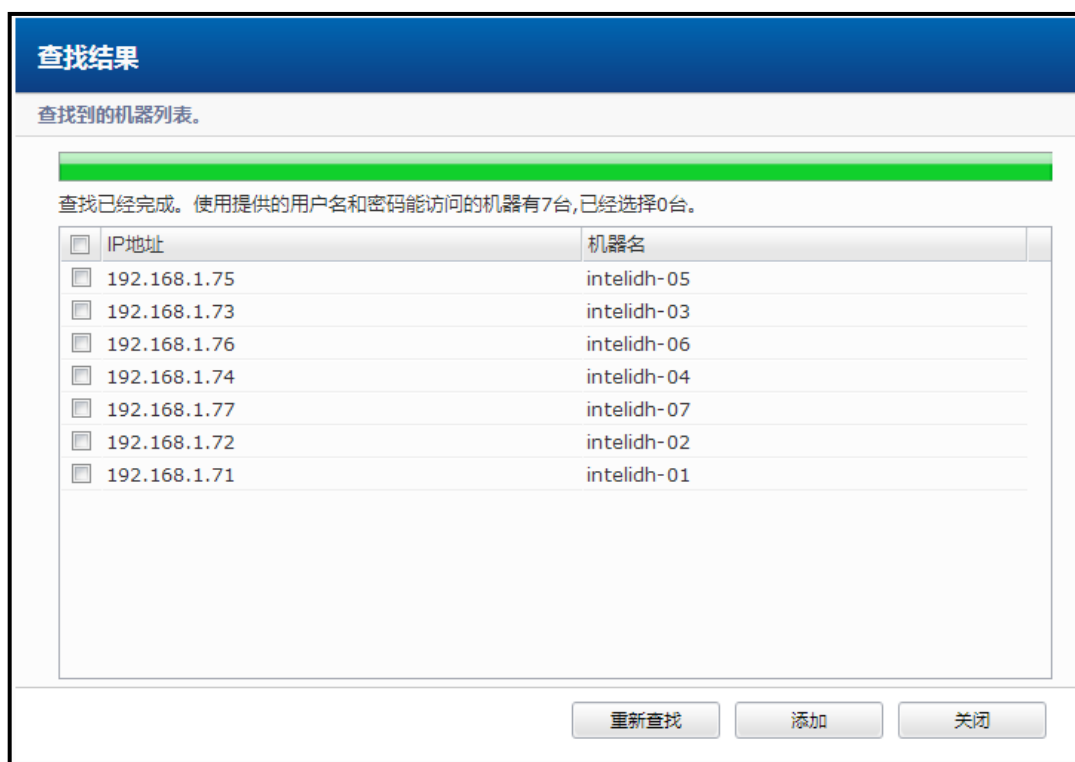


图 6.5 查找结果

当查找结束后，请勾选需要加入集群的机器，并点击“添加”。这时会出现一个如下图所示的窗口让您确认您的选择。点击“确认”继续。



图 6.6 确认对所选择的节点安装软件

如果有些节点的状态为“未连通”，请稍微等待一会。节点将自动被连接，其状态将由“未连通”变为“正在连接”，最后变为“已连通”，如下图所示。点击“完成”继续。

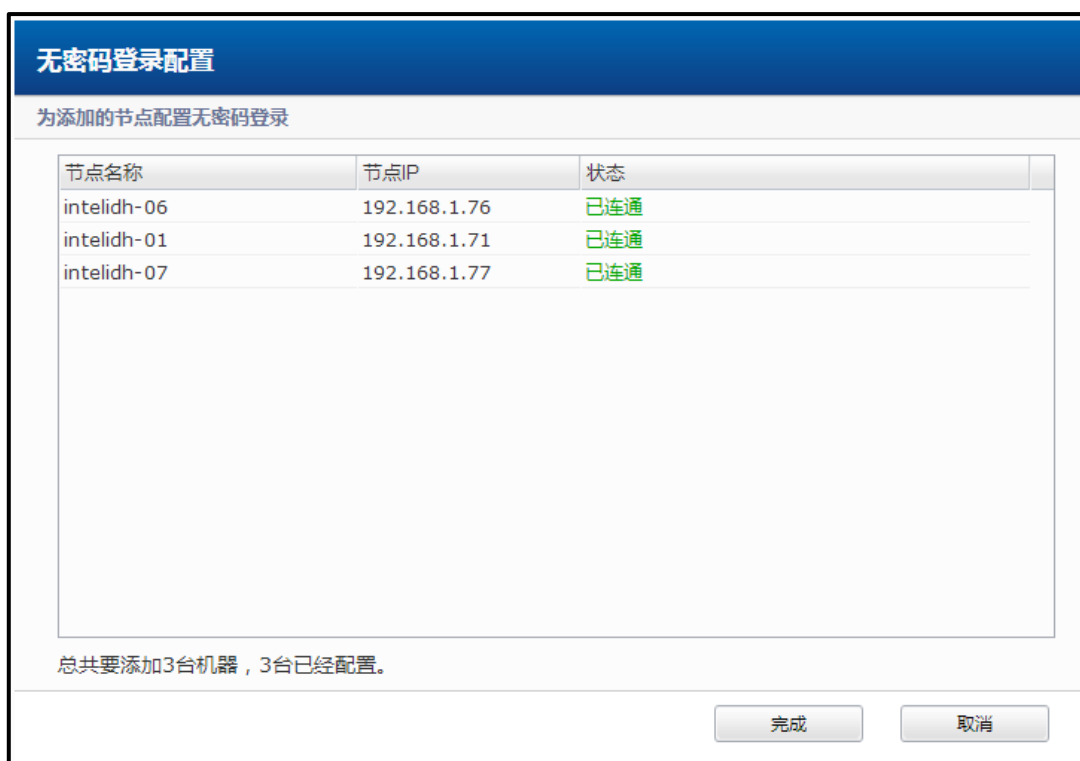


图 6.7 “已连通”的状态

所选节点将被添加到集群中，点击“下一步”继续。

第2步

指定集群节点以及网络环境

网络环境：

集群中的节点能通过主机名互相访问(通过配置好的DNS服务器或/etc/hosts文件)

节点名称	节点IP	状态
intelidh-01	192.168.1.71	已连通
intelidh-06	192.168.1.76	已连通
intelidh-07	192.168.1.77	已连通

添加节点

删除节点

上一步

下一步

取消

图 6.8 节点成功添加

第三步，配置集群的机柜并点击“下一步”继续。

第3步

设置集群的机柜

选择机柜： /Default

未分配节点

No items to show.

>

>>

<

<<

已有节点

intelidh-06

intelidh-07

intelidh-01

添加机柜

编辑机柜

删除机柜

上一步

下一步

取消

图 6.9 设置集群的机柜

第四步，选择安全策略并点击“下一步”继续。

第4步

配置集群节点认证协议

设置集群中节点使用的安全策略。

安全策略： 集群使用简单安全策略

上一步

下一步

取消

图 6.10 选择安全策略

第五步，如果有节点的状态为“未安装”，如下图所示，点击“安装未成功安装的节点”并点击“确认”按钮。

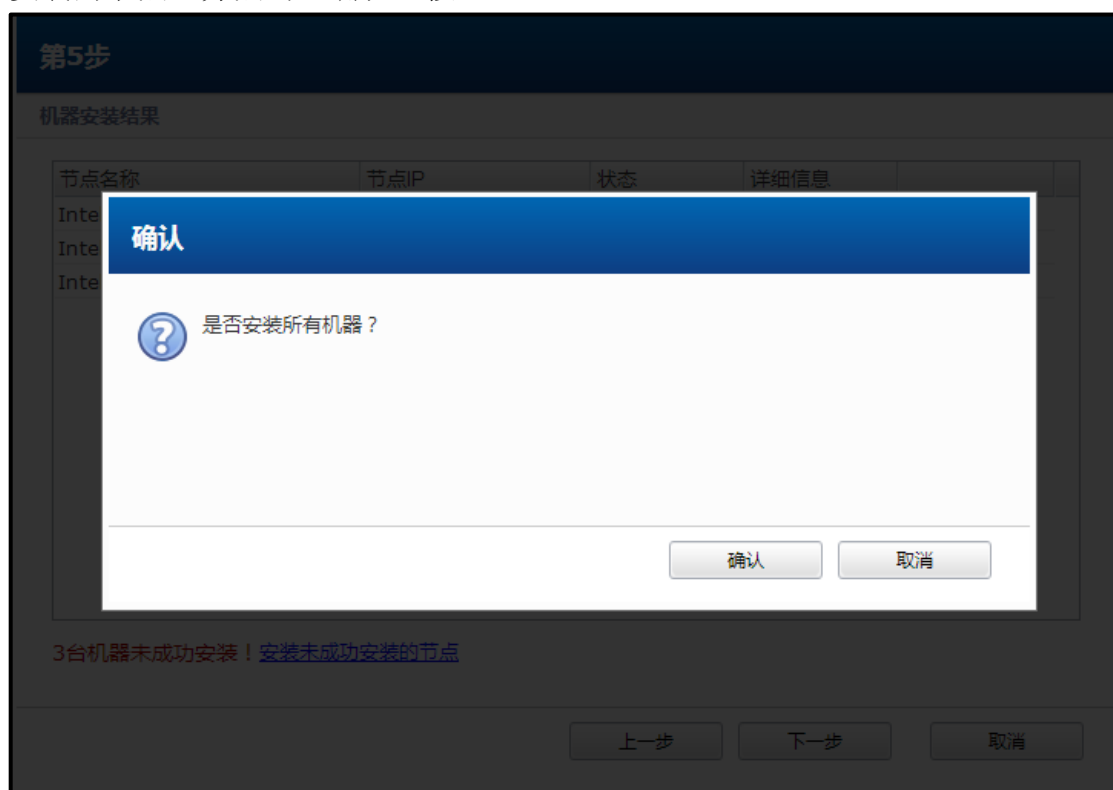


图 6.11 “未安装”的状态

由于 NTP 服务器所需同步周期一般为 5 分钟，请您耐心等待直到软件安装全部完成，安装完成后请点击“确认”并将弹出的命令行窗口关闭。

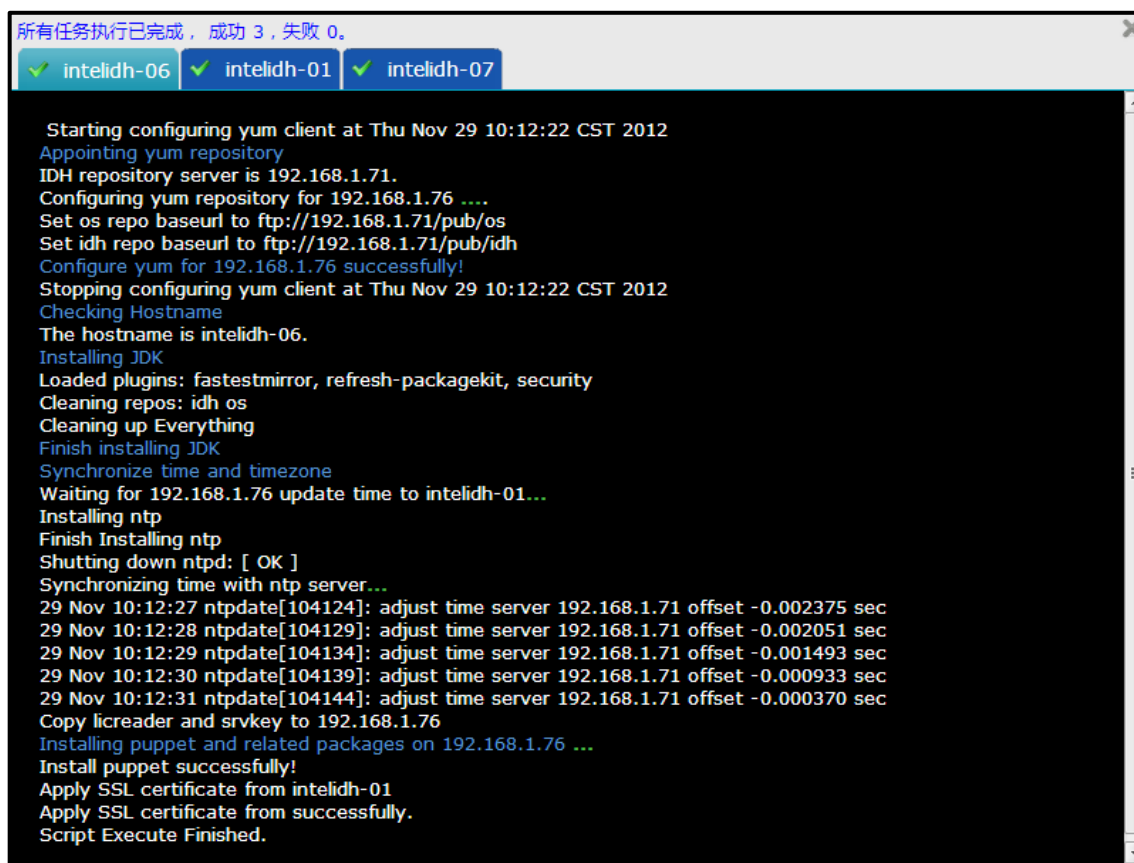


图 6.12 对所选择的节点安装软件

如果所有节点的状态都为“成功”，点击“下一步”继续。只有安装成功的机器菜会被成功添加到集群中。

第5步

机器安装结果

节点名称	节点IP	状态	详细信息
intelidh-01	192.168.1.71	成功	详细信息
intelidh-06	192.168.1.76	成功	详细信息
intelidh-07	192.168.1.77	成功	详细信息

3台机器安装成功。

上一步

下一步

取消

图 6.13 “成功” 状态

第六步，开始进行集群拓扑配置。进入 HDFS 组件控制节点的配置界面。分别一台服务器选择作为主命名节点，从命名节点。其中必须选择服务器作为主命名节点。而从命名节点是可选的。选择结束后，点击“下一步”继续。

第6步

HDFS组件控制节点的配置

Primary NameNode :

intelidh-01

(*)必填，集群中必须包含一个Primary NameNode。

Secondary NameNode :

选填，Secondary Namenode可以备份Primary NameNode的元数据。

上一步

下一步

取消

图 6.14 配置 HDFS 组件控制节点

第七步，进入 MapReduce 组件控制节点的配置界面。这里必须选择一台服务器作为 MapReduce 的任务分配器，而在上一步配置 HDFS 组件控制节点中，您选择的主命名节点将被默认作为任务分配器。选择结束后点击“下一步”继续。

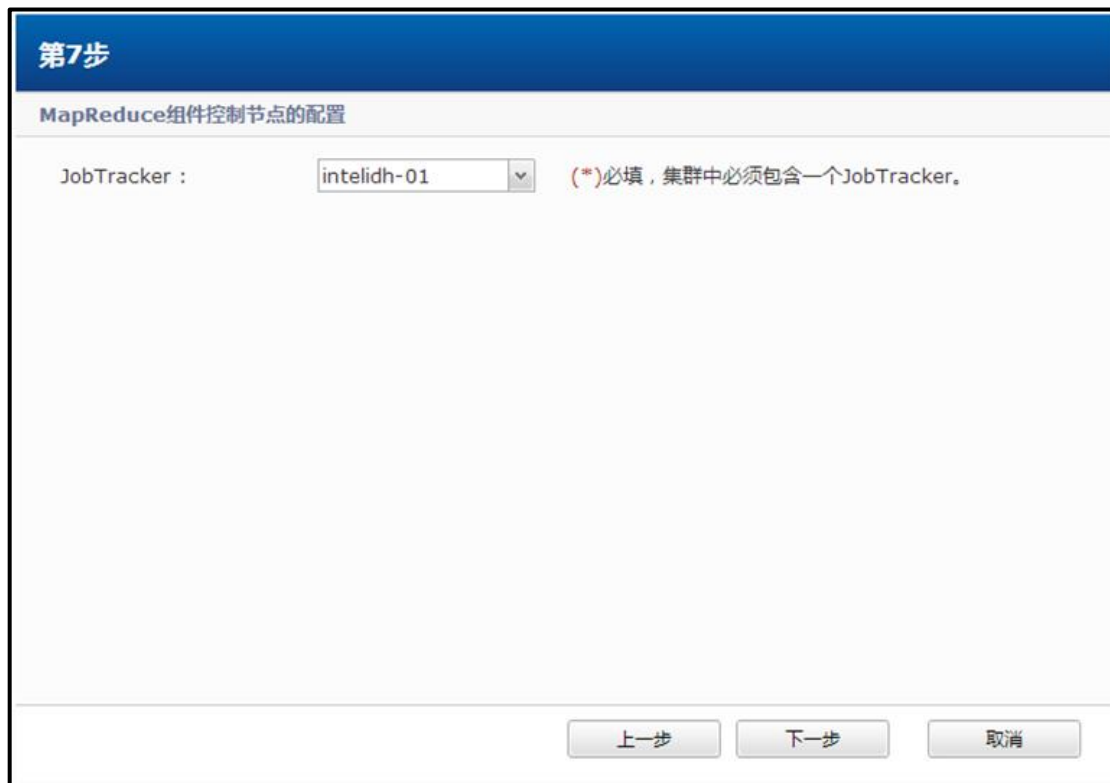


图 6.15 MapReduce 组件控制节点配置界面

第八步，进入 Zookeeper 组件控制节点配置界面。这里可以选择 ZooKeeper 节点，建议使用奇数并且数量至少为 3。

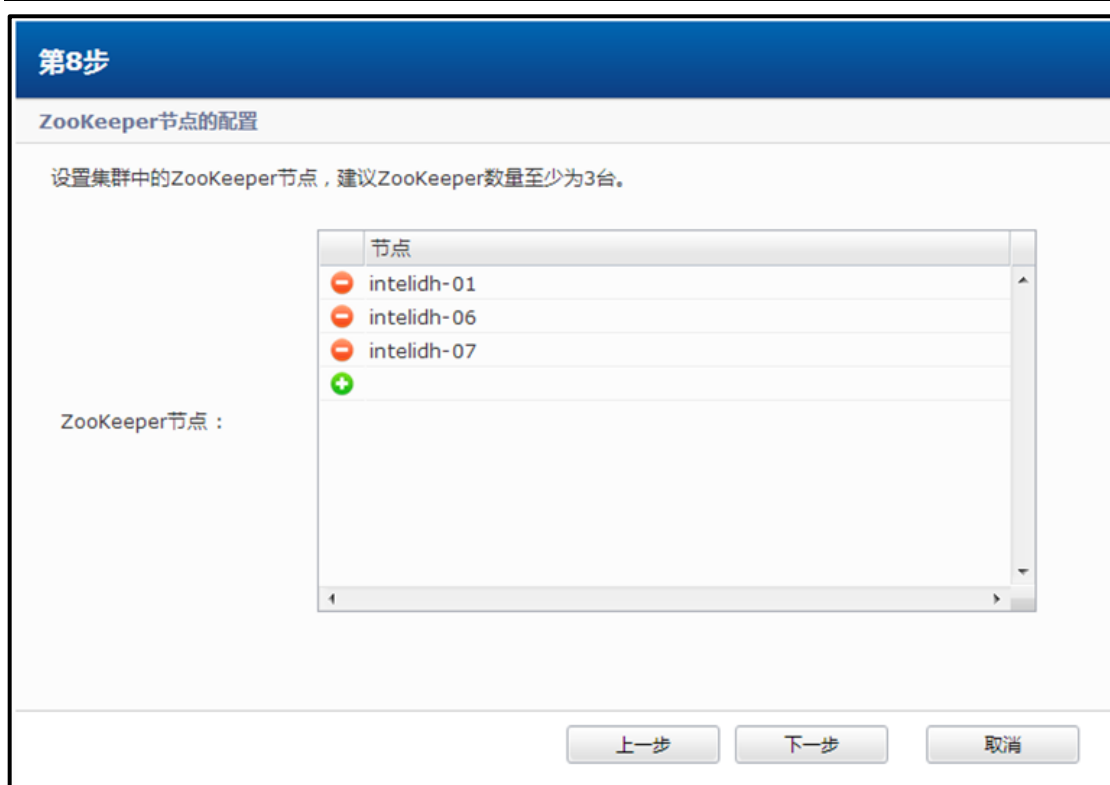


图 6.16 HBase 组件 Zookeeper 节点配置界面

第九步,选择结束后点击“下一步”继续,选择 HMaster 节点,默认与 Zookeeper 节点一致。

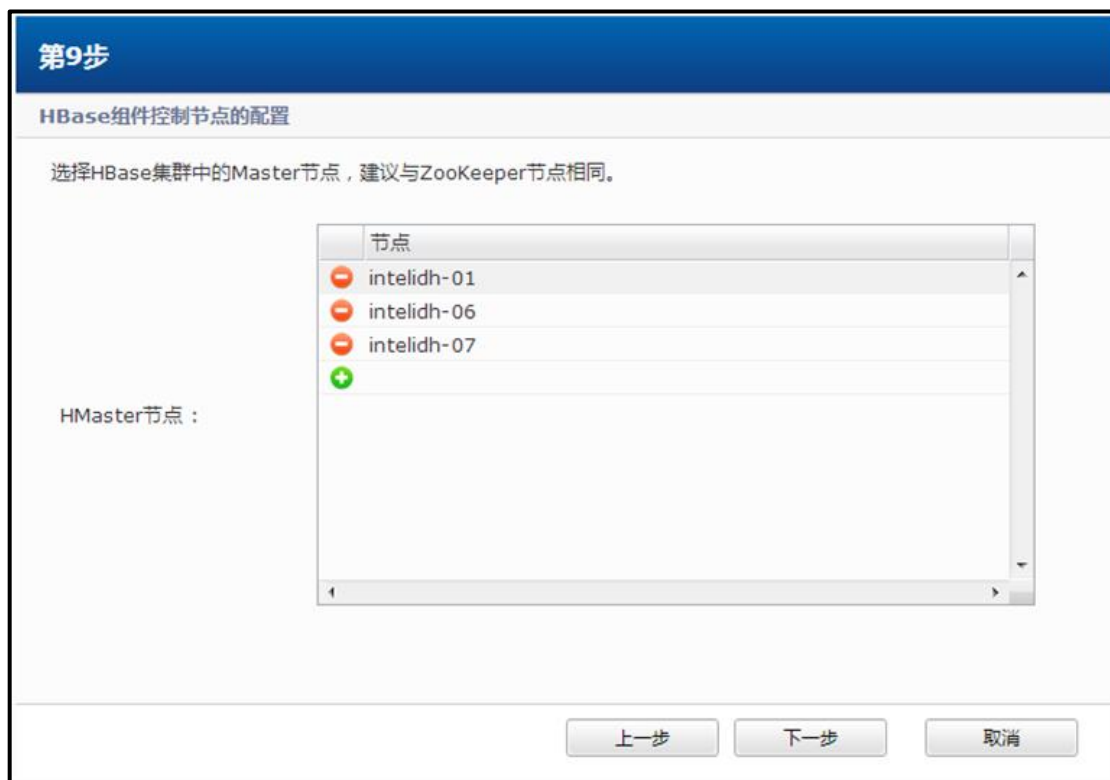


图 6.17 HMaster 节点配置界面节点配置界面

第十步，进入 Hive 组件控制节点配置界面。这里可以选择 Hive 服务所安装的服务器。选择结束后点击“下一步”继续。

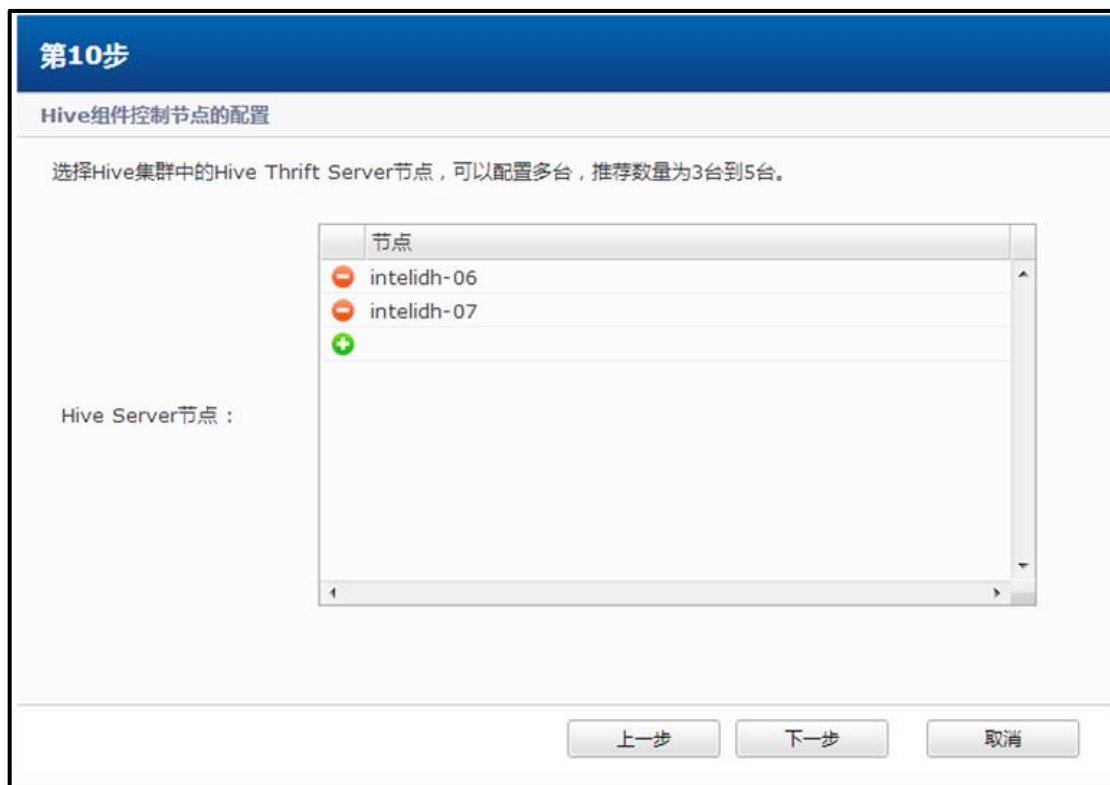


图 6.18 Hive 组件控制节点配置界面

第十一步，确认集群拓扑主要节点配置完成，点击“完成”关闭向导。

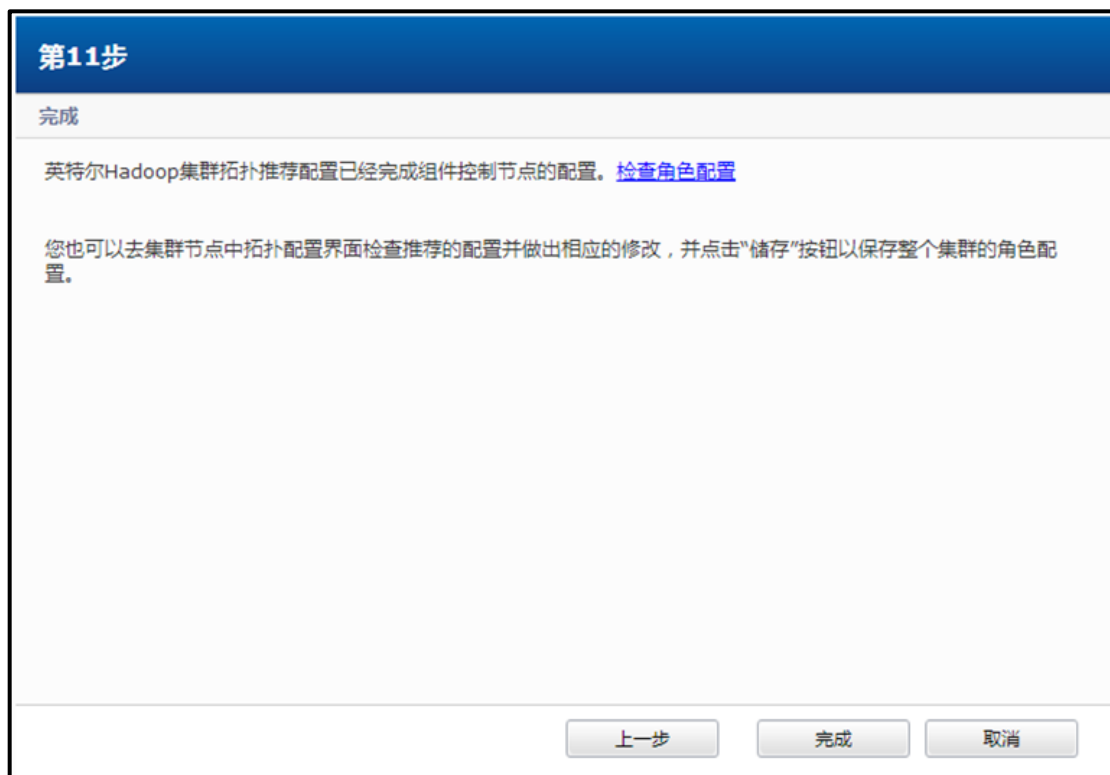


图 6.19 完成集群拓扑配置

您可以检测角色配置确认之前的配置情况。

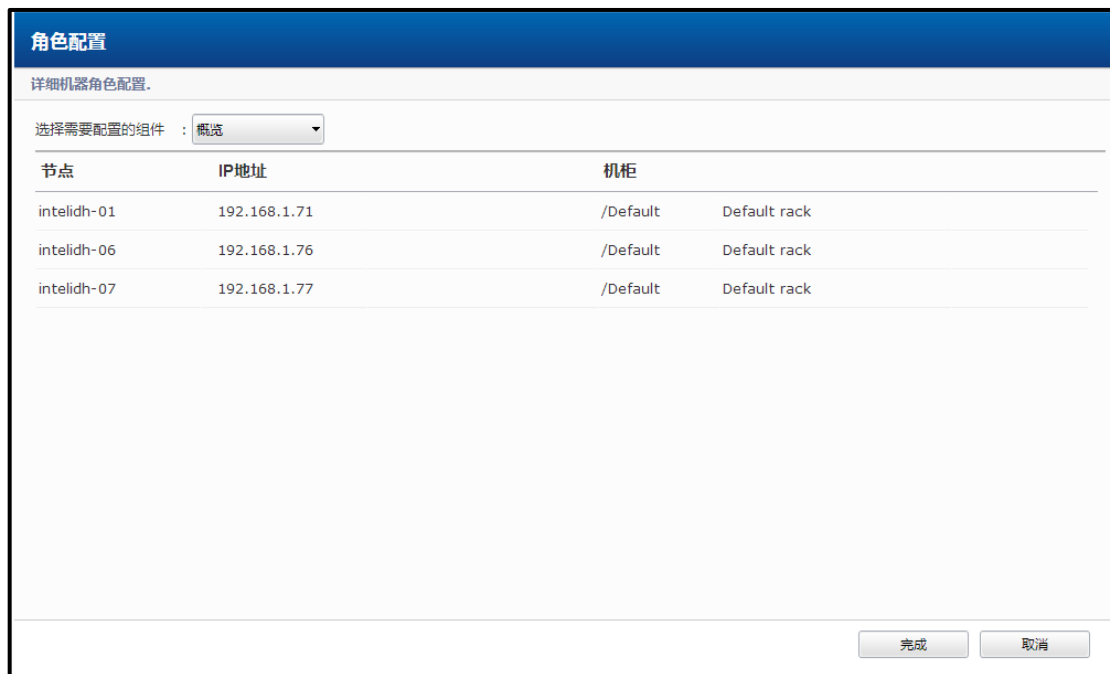


图 6.20 检查集群角色配置

点击确认来储存来保存机器角色配置，如果是第一次安装，请点击确认进行“格式化集群并进行配置”，否则在节点配置中点击“配置所有节点”选项。

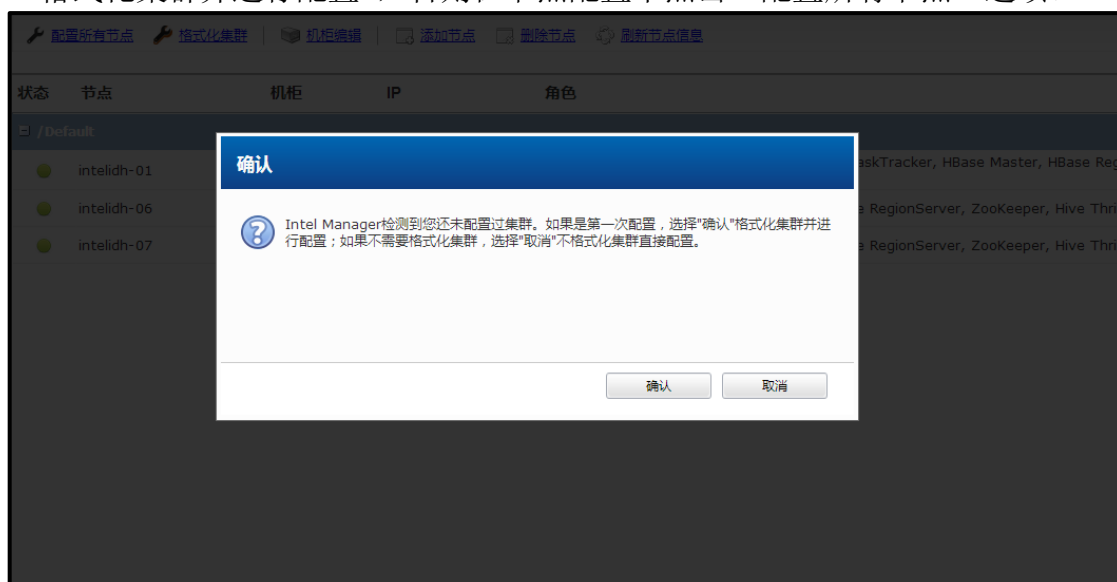


图 6.21 确认格式化集群配置

至此，集群安装配置已全部完成。



图 6.22 完成集群安装配置

6.3 输入许可证

英特尔® Hadoop 发行版需要输入许可证才可进行配置并启动运行完整功能。免费版用户默认拥有 *免费版的许可证*，可以直接进入 6.5 章节。如果您拥有 *试用版或商业版的许可证*，通过点击管理主界面左侧菜单栏中的许可证管理菜单进入许可证管理界面。

点击下图所示位置的“刷新许可证信息”，在许可证界面中会列出所有具有许可证验证功能的服务器列表。



图 6.23 许可证管理界面

双击列表中的任何服务器，进入许可证输入主界面，如下图所示。将许可证文件中的内容输入许可证文件输入框，然后点击上传按钮。

许可证信息

查看许可证信息与状态，可以上传新的许可证。

许可证信息

服务器ID : figYvpfMtbvhusxg8LvtTbK/pqa/
许可证类型 :
包含组件 :

升级许可证

许可证文件 :

上传

取消

图 6.24 上传许可证文件

然后点击图许可证管理主界面左上角的部署许可证按钮，将许可证部署到机器中的相应服务器中。

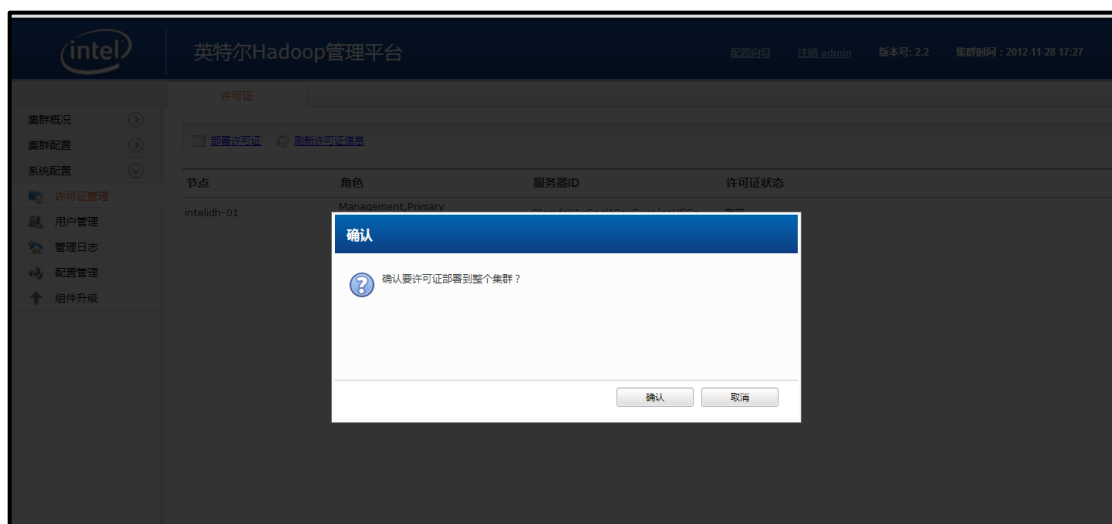


图 6.25 许可证部署

6.4 配置节点

在输入许可证后，请在管理主界面左侧菜单标签集群配置中选择集群节点，并点击“配置所有节点”，系统会同步在后台实施集群中各个服务器的安装配置

工作。

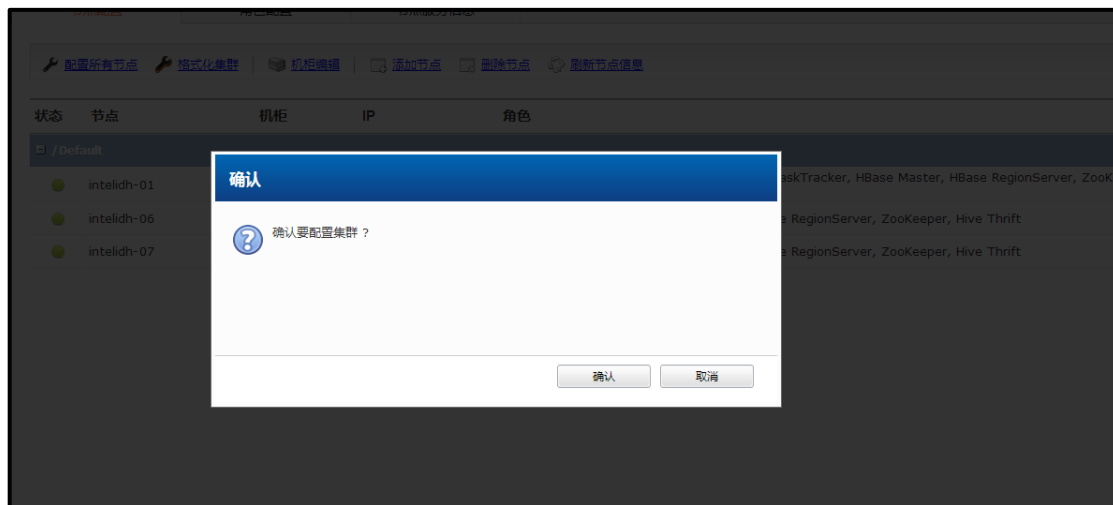


图 6.26 配置所有节点

如下图所示，在启动过程界面的标题栏中可以查看启动服务器的总个数和完成个数。

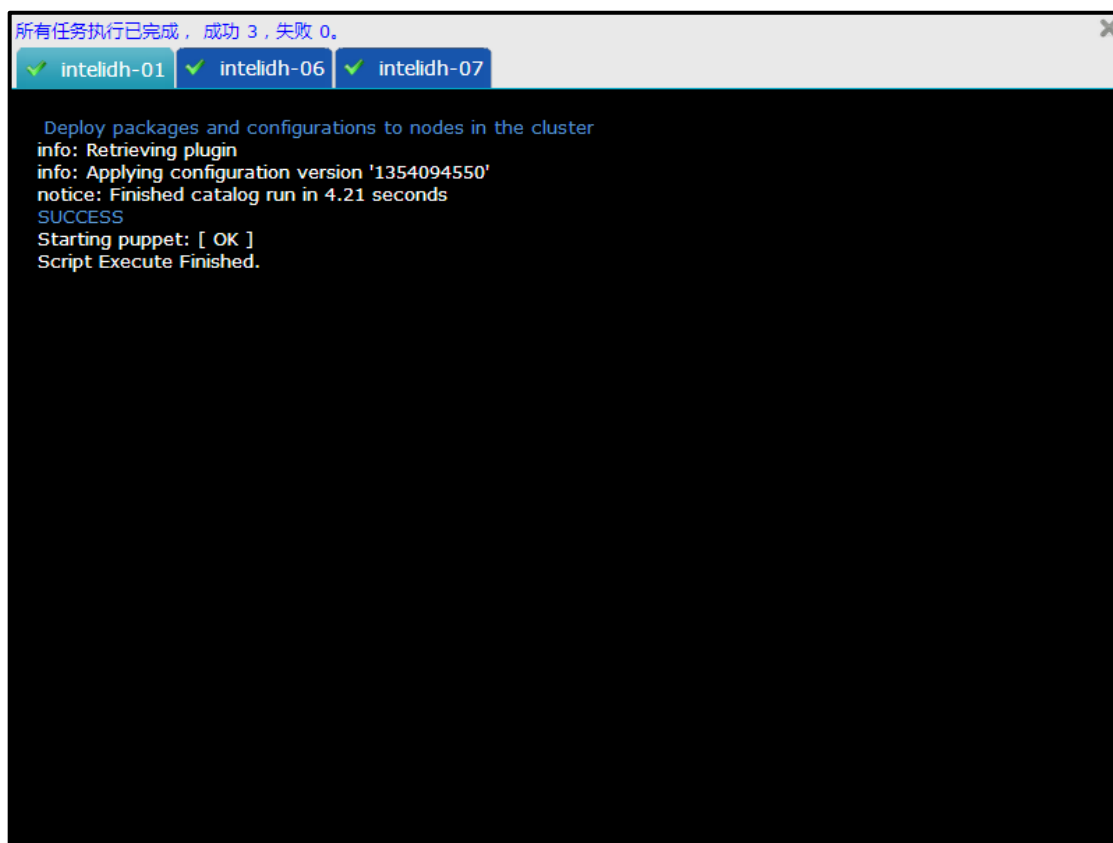


图 6.27 配置节点的运行界面

当所有服务器都启动完成后，点击启动过程界面右上角的关闭按钮关闭窗口。

6.5 启动集群

在运行中的组件分别为 Zookeeper, HDFS, MapReduce, HBase 和 Hive。除 HDFS 外，所有组件都只有两种状态“运行中”和“未运行”，在界面的最右侧有操作列表，可以通过点击按钮来对集群单一组件进行“启动/停止”操作。



图 6.28 控制面板

单一组件的启动必须满足如下要求：

- Zookeeper 启动之前不需要依赖另外组件；
- HDFS 启动之前不需要依赖另外组件；
- MapReduce 启动之前，需要确保 HDFS 处于运行状态下；
- HBase 启动之前，需要确保 Zookeeper, HDFS 处于运行状态下；
- Hive 启动之前，需要确保 HDFS, MapReduce 以及 HBase 处于运行状态下。

所以如需启动集群，需要严格按照启动顺序：Zookeeper, HDFS, MapReduce, HBase, Hive。点击控制面板界面上相应控件旁的启动按钮，自上而下顺序启动每个服务。系统会显示每个服务启动的进度，如下图所示。

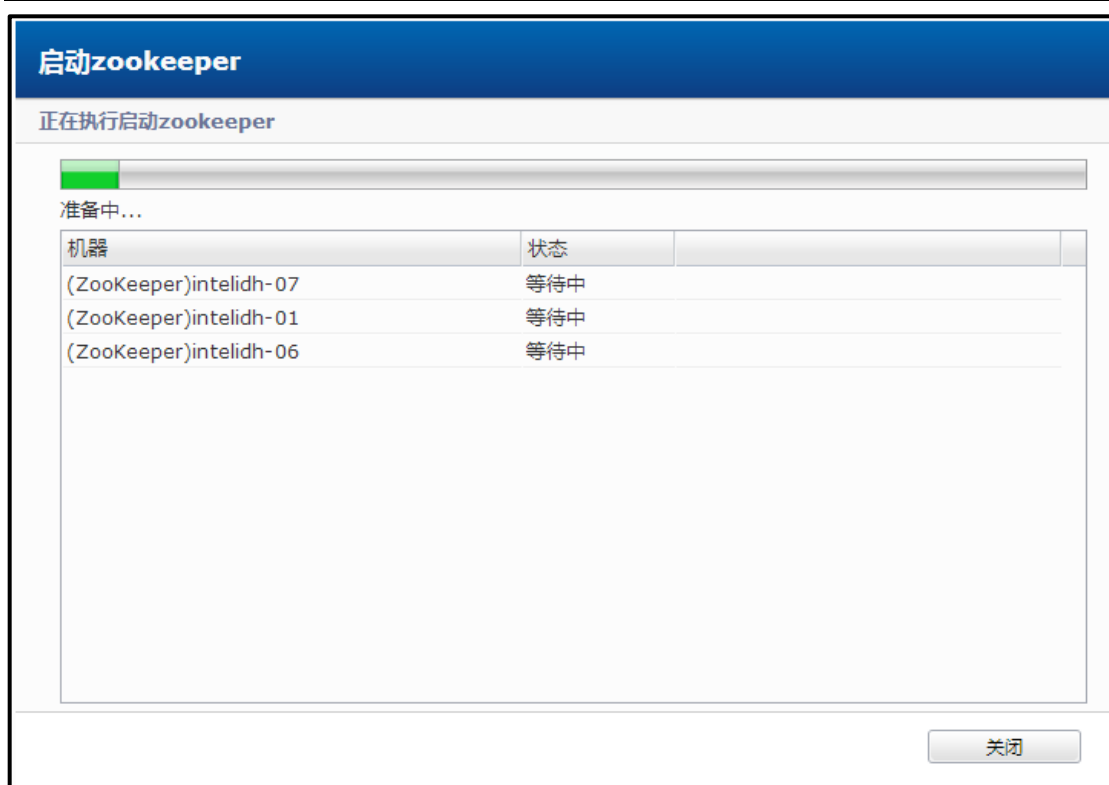


图 6.29 启动 Zookeeper 的进度

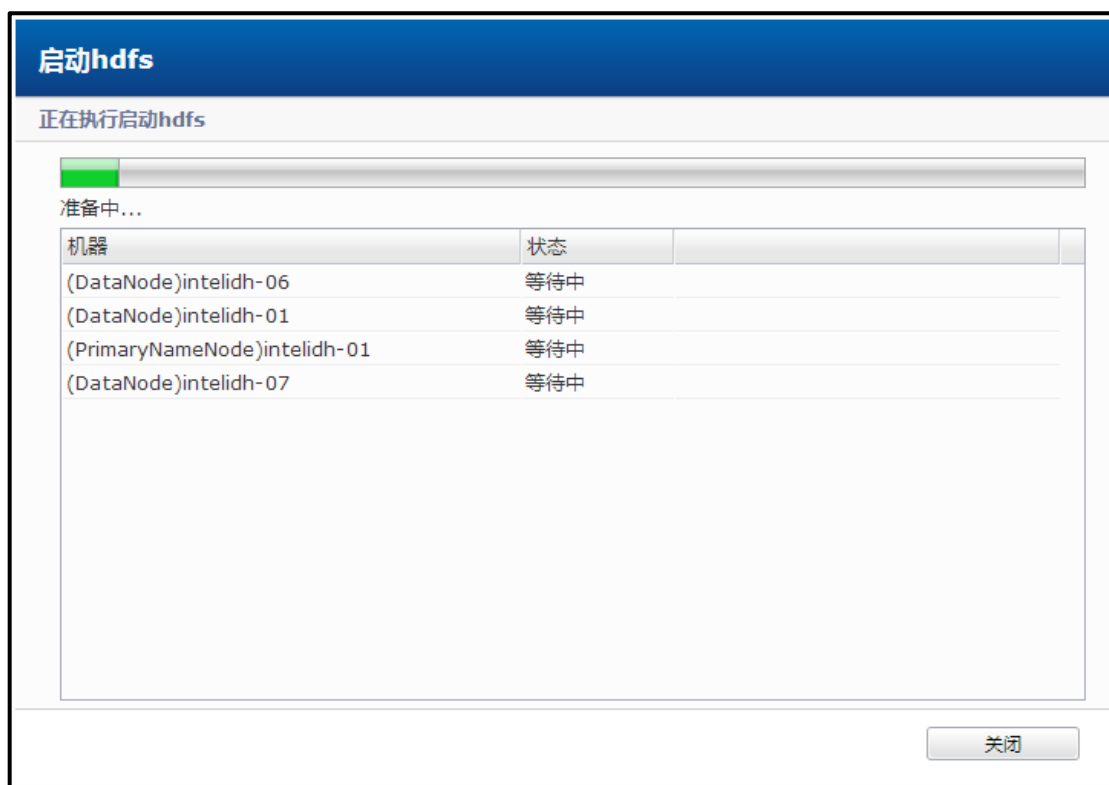


图 6.30 启动 HDFS 的进度

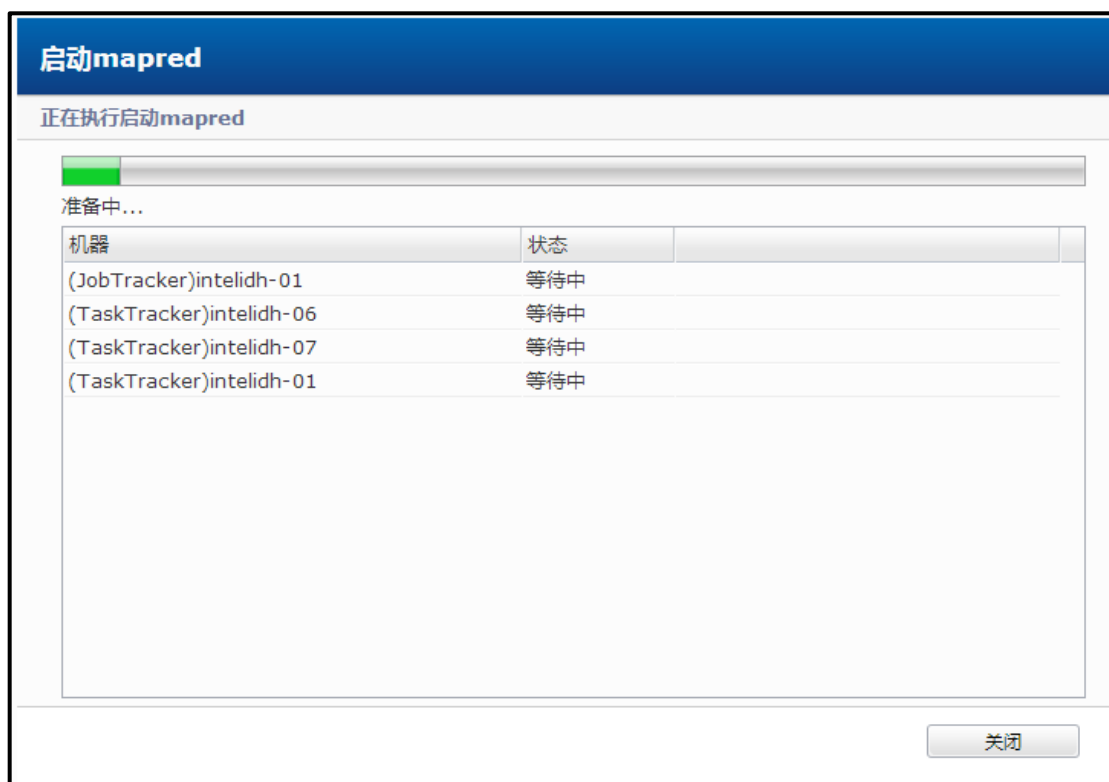


图 6.31 启动 MapReduce 的进度

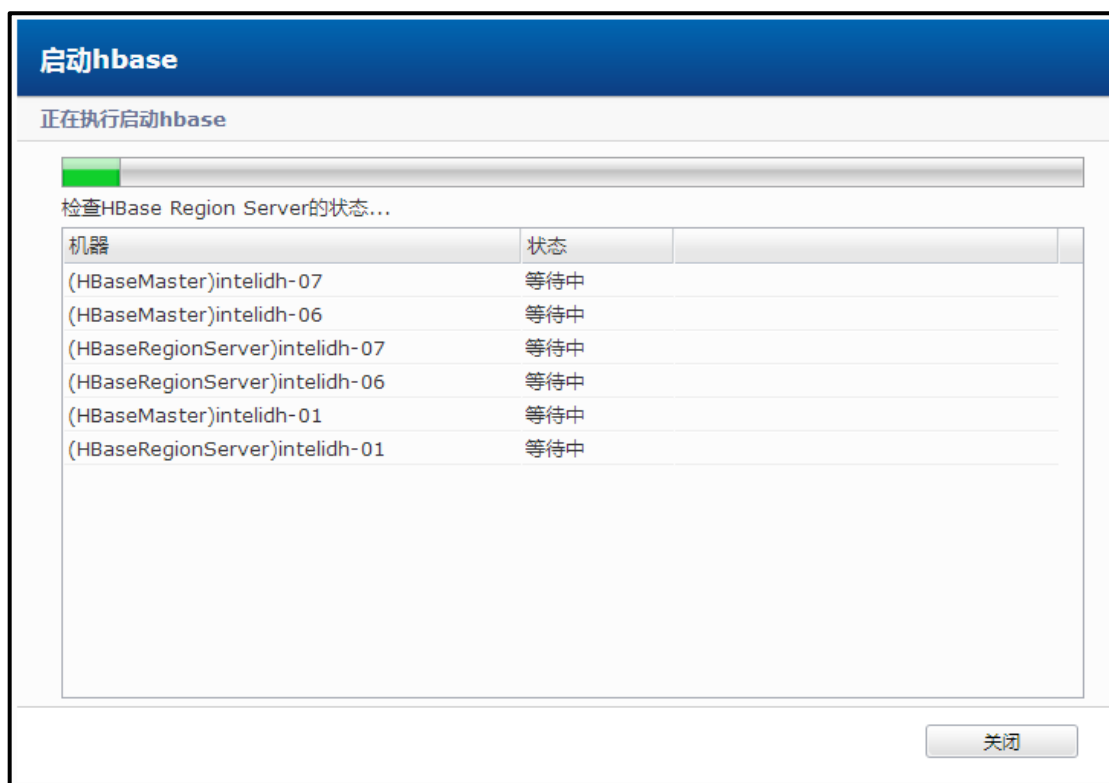


图 6.32 启动 HBase 的进度

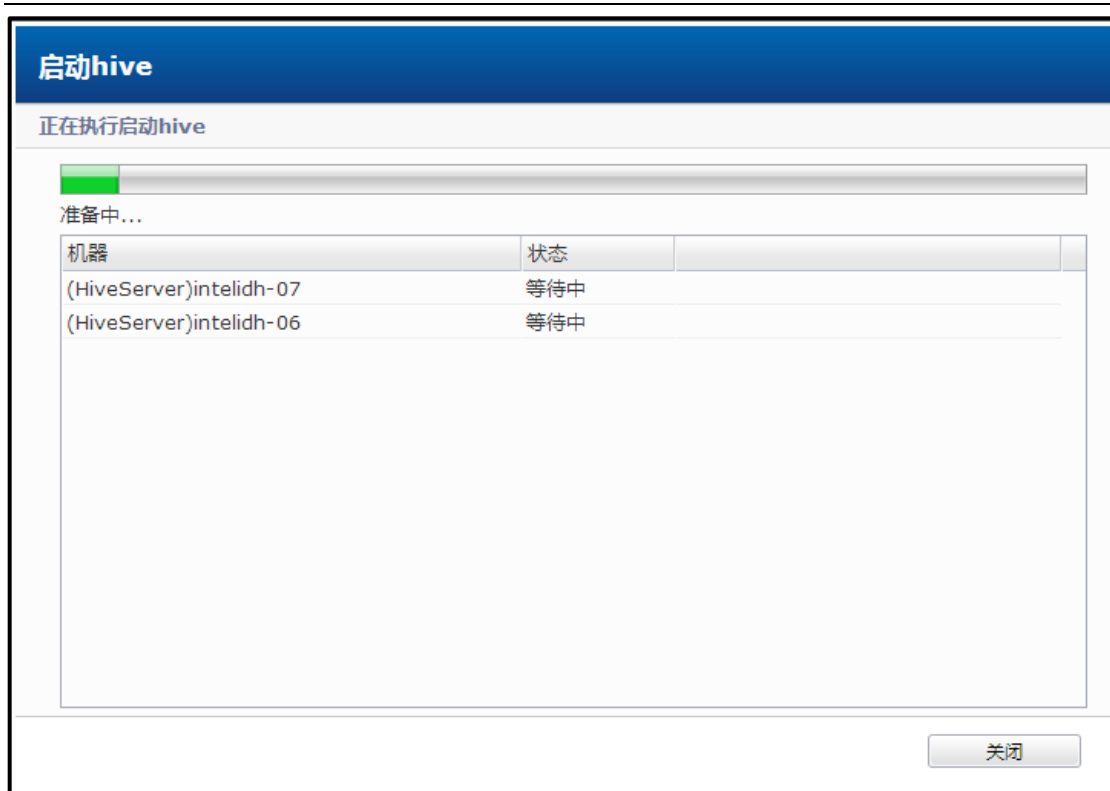


图 6.33 启动 Hive 的进度

启动完成后，系统会显示所有服务已经在运行中，证明系统安装成功。如下图所示。

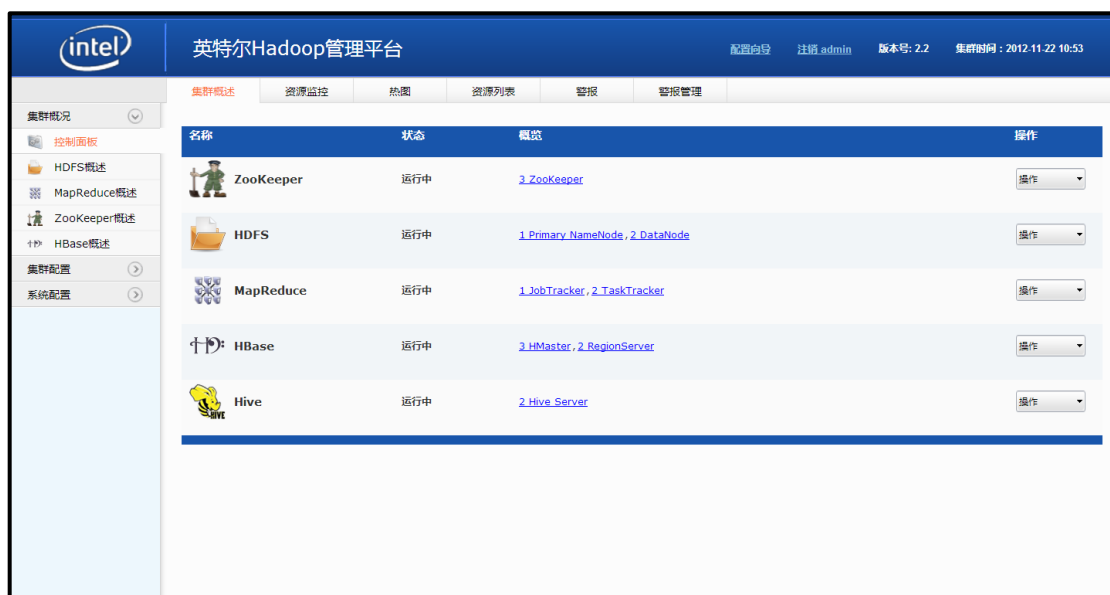


图 6.34 所有服务已经在运行

6.6 手动配置部分组件

在配置完所有节点后，如果在配置时选择了 Sqoop, Pig 和 Flume，组件将被自动安装。但您必须手动对其进行配置。