

# 基于 Rocks cluster 的集群高性能计算系统的搭建

杨寅

(上海工程技术大学, 上海 201600)

**摘要:**高性能计算系统在国民经济的许多领域有非常重要的应用。全世界高性能计算机500TOP甚至是一个国家综合国力的体现。基于Linux的Rocks cluster系统提供了一种可行的高性能集群计算机的搭建方案。介绍了Rocks cluster的特点,提供了Rocks cluster服务器的详细部署方法,并给出了集群管理和监控的方案。旨在提高集群系统的易用率,为高性能计算系统集群应用系统提供了一个易部署、高效且稳定的构建方法。

**关键词:**高性能计算系统;集群;HPC;Linux;Rocks Cluster

**中图分类号:** TP315 **文献标识码:** A

**文章编号:** 1009-3044(2022)32-0068-03

**DOI:** 10.14004/j.cnki.ckt.2022.1985

开放科学(资源服务)标识码(OSID):



## 1 概述

### 1.1 高性能计算机集群

超级计算机于20世纪60年代问世,几十年后,矢量计算机成为20世纪90年代的主导设计。拥有数万个现成处理器的大规模并行超级计算机现在已经成为常态。中国的超级计算发展也十分迅速,截至目前已部署天河、神威·太湖之光等高性能计算机。

高性能计算集群(High Performance Computing Cluster),通常是指使用多个处理器或以几台计算机组成的计算系统,用于各种领域的大型计算任务,包括量子力学、气候研究、分子建模和密码分析。可以由多台超级计算机组成,也可以通过由网络连接的PC机组成的集群。由于在某些廉价而通用的计算平台上运行并行计算集群可以提供极佳的性能价格比,所以近年来这种解决方案越来越受到用户的青睐<sup>[1]</sup>。

### 1.2 高性能计算系统

现如今,大多数现代超级计算机使用Linux内核作为操作系统,但每台超算的设计人员都会对他们使用的系统做了特定更改,并且没有行业标准,主要原因是不同计算机硬件架构的差异。例如IBM Blue Gene超级计算机在计算节点上使用CNK操作系统、天河超级计算机使用Slurm,以及OpenHPC。然而。对于非专业人员来说,集群的建立,安装、管理和监控一直是一致和困扰的问题。

Rocks集群是一个基于SMP的分布式架构开发的Linux系统。作为业内一流的免费集群操作系统,Rocks集群具有低成本、易扩展、结构精简、多软件可继承等特点,向广泛的科学应用提供集群计算能力,

提高研究者们的工作效率。

本文选用Rocks6.2版本以最便捷的方式搭建高性能计算集群平台,低成本建立计算机集群,为普通用户大规模计算任务提供了快捷有效的解决方案。

## 2 Rocks 集群结构

Rocks Cluster在软件结构上由FRONTEND(前端节点)和NODE(计算节点)组成。一般来说,前端节点可通过局域网交换机连接到多个计算节点,前端节点可配置一个网络适配器连接到公共网络(外网),通过前端节点上的外网端口,计算节点可以同时实现对公共网络(外网)的访问。Rocks Cluster的具体结构如图1所示:

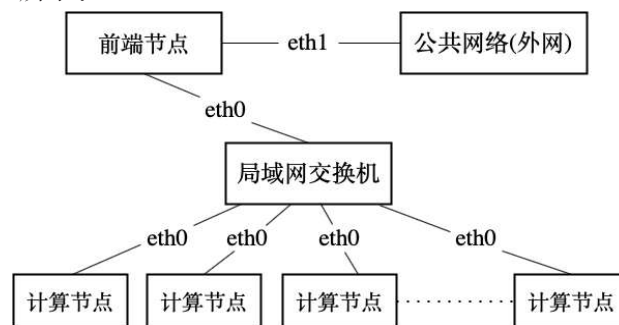


图1 Rocks Cluster 的具体结构

另外,Rocks系统中的模块化Roll机制也是其关键组成部分之一,它可以为特定的专业领域定制Roll包。在使用GEON Roll构建集群前端节点时,所有计算节点都将安装Geo专用软件。通过这种机制,我们可以轻松地将某些专业软件安装到Rocks集群系统中。

Rocks是一份完整的光盘机群解决方案,它面向

收稿日期:2022-03-06

基金项目:上海市大学生创新创业计划(S202110856037)

作者简介:杨寅(2000—),男,贵州人,主要研究方向为大数据处理。

x86 及 IA64 的 Red Hat Linux COTS 机群。组建一套 Rocks cluster 并不需要任何机群方面的经历,实际上,机群架构师将能找到一条灵活的并且标题化的方式来重新设计整个软件栈,而这对大多数用户而言则适当地隐藏了<sup>[2]</sup>。尽管 Rocks 包括了在任何机群软件结构中都应指望的工具(PBS、Maui、GM 支持、Ganglia 等),它的安装简易性则是独一无二的<sup>[3]</sup>。

3 安装 Rocks 集群

在硬件方面,主节点为一台 IBM System x3650 M4 服务器,两颗 CPU:2.6GHz Intel Xeon E5-2670,64GB 内存,2.4TB 硬盘。

系统软件方面,在 Rocks cluster 的官方网站 <http://www.rocksclusters.org> 下载 rocks 6.2 ISO 安装包,安装包内提供的工具详情如下表:

表 1 Rocks ISO 安装包工具

卷名	描述	功能
Area51	系统安全相关的服务和应用程序	分析集群上文件和内核的完整性
Base	基础卷	Rocks 基本安装包
SGE	分布式资源管理软件	保证集群内的资源得到有效利用
ZFS-linux	ZFS 设备驱动程序	用于构建可靠的存储系统
Fingerprint	应用程序依赖性软件	分析、管理与跟踪应用程序依赖关系
Ganglia	集群监控软件	监控系统性能
HPC	性能峰值测试	并行运行的已配置软件工具的应用
HTCondor	高吞吐量计算	计算系统吞吐量
Kernel	核心卷	Rocks 启动引导包
KVM	系统虚拟化模块	让进程有专属的内核和用户模式
OS	CentOS 安装卷	安装 Linux 系统
Perl	实际抽取与回报语言	从 CPAN 模块创建基于 Rocks 的 RPM
Java、Python	编程语言	编写计算机程序

以上准备工作完成后,安装内容如下:

3.1 RAID 磁盘阵列设置

Redundant Array of Independence Disks(RAID)磁盘阵列是一种数据存储虚拟化技术,它将多个物理磁盘驱动器组件按不同的方式组合起来形成一个或多个逻辑单元,以实现数据冗余、提升性能。下面是常见 RAID 方案的对比:

表 2 RAID 方案对比

RAID 等级	RAID0	RAID1	RAID3	RAID5	RAID6	RAID10
别名	条带	镜像	专用奇偶校验条带	分布奇偶校验条带	双重奇偶校验条带	镜像加条带
容错性	无	有	有	有	有	有
冗余类型	无	有	有	有	有	有
热备份选择	无	有	有	有	有	有
读性能	高	低	高	高	高	高
随机写性能	高	低	低	一般	低	一般
连续写性能	高	低	低	低	低	一般
需要磁盘数	n≥2	2n(n≥1)	n≥3	n≥3	n≥4	2n≥n≥4
可用容量	100%	50%	(n-1)/n	(n-1)/n	(n-2)/n	50%

磁盘阵列卡进行 RAID1+0 模式后,磁盘容量减半获得容错能力<sup>[4]</sup>。综合可靠性与性能,选用 RAID 10 作为系统磁盘阵列。设置步骤如下:

(1) 服务器启动后,等待进入阵列卡提示出现,按“Ctrl + H”进入 RAID 卡选项,点击“start”进入配置界面。

(2) 在配置界面中,进入“Configuration Wizard”选项,选择“new configuration”创建新配置,确认后系统将清除磁盘内所有数据。

(3) 在下一个界面中选择“manual configuration”进行手动分配磁盘空间,根据 RAID 10 的配置方式,将四个硬盘为一组添加到“Drives groups”,按此方式添加两组。

(4) 将制作好的两组硬盘放入“Span”,在“virtual disk definition”配置界面选择 RAID Level 为 RAID 10,其他参数保持默认。更新数据“Update Size”后开始磁盘阵列创建。

(5) 保存设置并返回主配置界面,可以看到阵列设置为 RAID 10,VD 状态为最佳,两组共 8 块硬盘状态在线。

3.2 Rocks 主节点安装

(1) 插入 Kernel/Boot Roll 光盘,从光盘启动服务器,在安装选项界面输入 Build。

(2) 选择安装的必要功能,包括 area51、kernel、base、OS、SGE、HPC、Ganglia、Java、Python 等。

(3) 输入集群信息,集群名称保持默认,输入接口 eth0 与 eth1 的外网地址、私网 IP、外网网关、外网 DNS 等。

(4) 配置系统 root 的密码,设置时区为 Asia/Shanghai,硬盘自动分区。

(5) 系统主节点开始自动安装,安装完成后会自动重启,至此主节点安装完成。

### 3.3 Rocks 计算节点安装

Rocks 最大的优点在于把你所需要的软件在系统安装时都一次安装好,不需要安装者更多地配置<sup>[5]</sup>。计算节点可以从主节点获取所需文件,自行安装。

(1) 以 root 用户登录主管理节点,打开终端,执行 `# insert-ethers` 命令,选择计算节点属性 Compute。

(2) 通过交换机连接主节点与计算节点,将子节点服务器启动, BIOS 设置引导顺序为: CD、PXE(网络引导)、Hard Disk。主节点识别成功后,子节点将自行安装。安装完成后,子节点自动重新启动并加入计算机群。

(3) 更改计算节点的根密码以提高 Rocks 系统的安全性。执行命令如下: `# rocks set host sec_attr compute attr=root_pw`,系统确认用户自定义密码后即设置完成。

(4) 至此 Rocks 集群搭建完毕,可通过命令 `# rocks list host` 查看集群主节点及计算节点状态。

### 3.4 Rocks 系统网络配置

节点设置完成后,此时系统仍然无法连接互联网,不能及时进行远程调试和提交运行任务。经调试,系统还需要修改以下参数:

(1) 使用命令 `# ifconfig` 查看网络配置,再通过 `vim` 修改 `ONBOOT=yes`,将 Rocks 系统的双网卡 `eth0` 和 `eth1` 网络连接设置为开机自动启用。

(2) 修改完成后使用命令 `# service network restart` 重启网络。

(3) 至此,网络配置完成。可以使用 `# ping IP` 查看网络状态。

## 4 Rocks 集群管理和监控

### 4.1 集群用户管理

Rocks 的用户管理方式是使用 411 安全信息服务来代替传统 NIS(Network Information Service)服务。这个系统设计的目的是提高拓展性、安全性与意外处理能力,并且 411 服务与 NIS 服务的使用和配置方式与 NIS 服务几乎没有区别<sup>[6]</sup>。

### 4.2 集群任务管理

Rocks 安装包内包括 Sun Grid Engine(SGE)分布式资源管理软件,保证集群内的资源得到有效利用并且不需要 NFS 进行操作。使用它可以极大提升用户提交计算任务的效率,并更加合理分配和提高整个集群的计算资源<sup>[7]</sup>。

### 4.3 集群监控系统

Ganglia 是一个开源的分布式监控系统,由美国加利福尼亚大学发起,用于集群和网格这样的高性能计算系统。Ganglia 偏向于操作系统低层一点的监控,主要是收集集群中的各个机器的 CPU 使用情况、内存使用情况、磁盘 IO、网络 IO、磁盘容量等<sup>[8]</sup>。它基于针对集群整体的层次结构设计,利用诸如用于数据表示的 XML、用于紧凑、可移植的数据传输的 XDR 以及用于数据存储和可视化的 RRDtool 工具。Ganglia 使用精心设计的数据结构和算法来实现高效的节点调度与高并发性。

Rocks 系统安装包中已包含 ganglia 卷,我们仅需要使用命令 `# rocks enable roll ganglia` 激活 roll 中的 ganglia 安装程序,并执行语句 `# cd /export/rocks/install` 切换到对应安装包目录进行安装。

安装完成后,可以在浏览器中进入 `http://localhost/ganglia` 生成系统实时状态显示,向用户提供直观的 Web 界面。

## 5 总结

本文基于 Rocks cluster 系统搭建一组高性能计算集群。介绍了当今高性能计算机集群与系统的现状,针对 Rocks 集群结构展开说明,详细描述基于 Rocks 构建高性能计算集群系统的过程,并且讨论服务器 RAID 磁盘阵列的设置过程,最后对 Rocks 计算节点进行配置与部署,为用户提供一个快速易用的集群搭建方案。

### 参考文献:

- [1] 张澜宇,邹溪. PC 实验室的高性能运算平台的实施[J]. 实验室科学,2016,19(5):53-56,62.
- [2] Rocks Cluster Distribution: Users Guider[M]. UC Regents, 2011.
- [3] 林先津. 桌面虚拟化技术在分布式设备管理中的研究与应用[J]. 实验技术与管理,2013,30(4):68-70.
- [4] 赖策,祝元仲. 虚拟机基于不同存储池模式下的磁盘性能测试分析[J]. 轻工科技,2020,36(7):90-91.
- [5] 张洋,陈文波,李廉. 基于 Rocks 的高性能集群平台搭建与应用[J]. 计算机工程与科学,2008,30(11):137-139,143.
- [6] 徐明,罗明宇,蔡文举,等. Rocks 集群应用系统发布技术研究[J]. 自动化技术与应用,2015,34(3):38-41.
- [7] 侯建军. 数字电子技术基础[M]. 北京:高等教育出版社,2003.
- [8] 狄晓娇. 企业级 Hadoop 平台实现的相关技术[J]. 中国新通信,2016,18(4):89-90.

【通联编辑:王力】