



OpenStack System Architecture and Practice

OpenStack系统架构 设计实战

..... 陆平 赵培 左奇 等编著

详细介绍OpenStack技术架构和核心模块
深入解析OpenStack各模块的设计思想、实现方案和部署方案
介绍OpenStack在大数据服务、数据库服务等PaaS领域的实现方案



机械工业出版社
China Machine Press



OpenStack System Architecture and Practice

OpenStack系统架构 设计实战

..... 陆平 赵培 左奇 等编著

详细介绍OpenStack技术架构和核心模块

深入解析OpenStack各模块的设计思想、实现方案和部署方案

介绍OpenStack在大数据服务、数据库服务等PaaS领域的实现方案



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

OpenStack 系统架构设计实战 / 陆平等编著. —北京: 机械工业出版社, 2016.7
(云计算与虚拟化技术丛书)

ISBN 978-7-111-54333-6

I. O… II. 陆… III. 计算机网络 IV. TP393

中国版本图书馆 CIP 数据核字 (2016) 第 167999 号

OpenStack 系统架构设计实战

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 王颖 张梦玲

责任校对: 殷虹

印刷: 三河市宏图印务有限公司

版次: 2016 年 8 月第 1 版第 1 次印刷

开本: 186mm × 240mm 1/16

印张: 18

书号: ISBN 978-7-111-54333-6

定价: 69.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88379426 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzit@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

本书编委会

主 编：陆 平

副主编：赵 培 左 奇

编 委：董振江 邓芳伟 张 晗

杨 勇 彭 涛 王 蔚

推 荐 序 *Foreword*

2009年第一届中国云计算大会的盛况仿佛还在眼前，转眼间7年快过去了。在这7年间，云计算在中国从萌芽到发展，如今云计算的浪潮正在影响着数据中心、应用系统的建设，甚至无时无刻影响着人们的生活。随着云计算技术的成熟，运营商、互联网公司、政府企业都纷纷在自身的IT建设中使用了云计算。

云计算的来龙去脉是什么？为什么工业界需要云计算？其背后的技术背景、相关公司、非营利化开源组织、商业利益集团在云计算方面的策略是什么？随着云计算技术的成熟，企业如何部署自己的云计算，选用什么样的云计算、云平台来搭建IT系统？另外，云计算领域出现了许多分支——公有云、私有云、混合云等，其各自的布局和未来发展如何？

这些都是相关领域的国家技术发展政策研究人员、企业CIO/CTO、高级研发人员、高校研究人员必须了解和能够回答的问题。本书和《云计算基础架构及关键应用》^①，深入浅出地解释了上述问题，是难得的好书。

云计算相关的技术书籍已经有一些，且各有亮点，如虚拟化、OpenStack、KVM等方面，都有大量的参考书籍。这两本书的特点是从云计算技术及应用全貌进行完整介绍，兼顾了系统与细节：包括云计算的虚拟资源层、IaaS云管理层、PaaS等各个平台服务层次；详细介绍了KVM、Xen、Docker、OpenStack、Cloud Foundry、Ceph、SDN等云计算的关键技术；介绍了计算、存储、网络虚拟化的技术发展和应用；介绍了NFV、公有云、私有云、混合云的架构、部署和应用场景。这两本书的作者是长期从事云计算的一线研发专家，他们从云计算的关键技术着手，同时站在云计算提供者、使用者、IT建设

^① 《云计算基础架构及关键应用》已由机械工业出版社出版，书号是9787111531760。——编辑注

决策者的多方角度来考虑云计算的应用场景和技术，非常难能可贵。

本人近年来一直从事分布式存储编码与系统的研发，部分理论成果——基于 G2 域的二进制纠删码 (Binary Reed-Solomen, BRS) 成功融入中兴通讯公司的大数据存储系统中。这一合作的过程使得我对这两本书的作者的专业水平有了更深入的了解，因此强烈推荐计算机、网络、系统和相关专业的研发人员阅读此书。我相信，通过对这两本书的阅读，大家将会进一步加强对云计算的全面认识，综合理解 SDN、NFV、云存储、云计算的部署和运维，全面掌握云计算的整体技术知识。

北京大学信息工程学院
北京大学大数据技术研究院存储编码及系统实验室
深圳市云计算重点实验室

前 言 *Preface*

早在 20 世纪 90 年代，云计算就已作为一种全新的技术模型被提出，但直到 2007 年，才因 Google、亚马逊等云计算先驱将其付诸于商业实践并获得丰厚利润，从而得到业界的广泛重视。与互联网、物联网等技术一样，云计算是电子信息技术和信息社会的需求发展到一定阶段的必然产物。从 2007 年至今，云计算已经成为人们进行信息交互与存储的重要方式，云管理平台也成为大数据处理和深度挖掘的主要平台。

高盛研究公司在 2015 年的一份报告中指出，花在云计算基础建设以及云管理平台上的费用在 2013 ~ 2018 年的年均增长率为 30%，而整个 IT 行业的预计增长率仅为 5%。面对这个蓬勃发展的市场，许多咨询公司和研究机构都对云计算市场有着不同的预测，但是他们都一致认为，在全球范围内，云计算的发展正在加速。在巨大需求的刺激下，云计算核心得到快速发展，商业云计算与开源云计算技术在竞争中共同推进，而云计算与行业结合，也形成了形态各异、特色鲜明的电子政务云、教育云、医疗云、金融云、环保云、旅游云等云计算服务，云计算大数据的发展空间则更加广阔。

中兴通讯公司在云计算方面有多年的技术积累和应用实践。本书结合云计算最新技术趋势和中兴通讯公司的长期实践，对云计算技术提出系统性的阐述，对云计算实践提供了思路和建议。本书首先从云计算的需求和现状出发，分析目前云计算面临的问题，针对这些问题分析了 IaaS 云管理平台、IaaS 云平台部署，并对 PaaS（平台即服务）等概念进行了充分的探索和讨论。

本书结构

本书由 12 章组成。

第 1 章对各主流云管理平台进行介绍及对比，并对 OpenStack 平台进行了重点介绍。

第2章系统地介绍了 Nova 的各个子模块，以及 Nova 的基本运行原理。第3章重点介绍了 OpenStack 项目的存储管理 (Cinder) 模块，描述了 Cinder 的架构、API、主要功能和工作流程。第4章从网络虚拟化主要面临的问题出发，讨论了 Neutron 架构及其具体功能。第5章从 Ceilometer 的体系架构谈起，首先简单介绍了 Ceilometer 的起源和几个重要概念，之后介绍 Ceilometer 的架构及关键组件，使读者对 Ceilometer 有整体的了解，然后详细剖析 Ceilometer 的数据采集机制，包括计量数据采集、计量数据转换和发布、计量数据存储，并介绍 Ceilometer 的二次开发。第6章重点讲述 OpenStack 中编排子系统 (Heat) 的相关概念、架构及其实现，并分别介绍 Heat 模板、Heat 资源类型以及 Heat 引擎，结合典型的场景，对基于 Heat 的业务弹性伸缩流程进行深入的分析与阐述。第7章从介绍裸机管理的 PXE、IPMI 通用技术开始，对 Ironic 的架构、基本运行原理和流程，以及 Ironic 的完全安装、简化安装等进行介绍。第8章系统地介绍 OpenStack 的消息总线及其基本运行原理，让读者对 OpenStack 的消息队列协议以及常用的消息队列方案有一个全面了解。第9章通过对 Sahara 使用模式、架构的介绍，及其与 Amazon、VMware 解决方案的对比，让读者对大数据即服务的概念建立清晰的认识。第10章对 Trove 总体架构、主要功能和 API 接口、安装和配置，以及二次开发进行全面的介绍。第11章通过对 Keystone 的介绍，阐明 OpenStack 作为云管理平台，如何应对云计算带来的包括虚拟化安全、数据安全、身份和访问管理安全等新的安全挑战。第12章分别从使用场景、逻辑架构等视角对当前 OpenStack 社区中比较活跃的孵化项目，例如，消息队列服务 (Zaqar)、共享文件系统服务 (Manila)、DNS 管理服务 (Designate)、密钥管理服务 (Barbican)、容器管理服务 (Magnum) 进行介绍。

本书适合高校计算机相关专业学生、云计算研究人员、云计算开发者和工程技术人员阅读参考。由于作者水平所限，书中难免存在一些谬误和不足之处，敬请读者批评指正。本书在写作过程中得到了很多领导和同事的大力支持，在此一并表示谢意。

作者

2016年5月

目 录 Contents

推荐序	
前 言	
第 1 章 云管理平台概述	1
1.1 主流云管理平台对比	2
1.2 OpenStack 简介	6
1.2.1 OpenStack 设计原理和体系 结构	7
1.2.2 OpenStack 社区和项目开发 流程	11
1.2.3 OpenStack 应用现状与发展 趋势	12
1.3 OpenStack 入门体验	15
1.3.1 初探 OpenStack	15
1.3.2 创建 OpenStack 虚拟机实例	17
1.3.3 创建虚拟机的流程概述	19
第 2 章 计算管理 (Nova)	23
2.1 概述	23
2.2 逻辑架构	24
2.3 物理架构	24
2.4 对主流 Hypervisor 的支持架构	26
2.5 与 VMware 的对接	27
2.6 支持的 Hypervisor	28
2.7 Nova 关键组件	28
2.7.1 API 服务 (nova-api)	28
2.7.2 消息队列 (AMQP)	29
2.7.3 nova-compute	32
2.7.4 nova-cell	35
2.7.5 nova-conductor	36
2.7.6 nova-scheduler	37
2.7.7 nova-volume	39
2.7.8 nova-network	40
2.8 nova-objectstore	61
2.9 OpenStack 使用流程	62
2.9.1 初始化数据库与 IP 池	62
2.9.2 创建用户与项目	62
2.9.3 使用 euca2tools 工具	63
2.9.4 创建镜像	65
2.9.5 创建虚拟机	68
2.10 K 版本新特性	70
2.11 小结	73

第 3 章 存储管理 (Cinder)	75	4.1.8 Overlay 网络.....	96
3.1 Cinder 的架构.....	75	4.1.9 Network NameSpace	97
3.2 Cinder API	76	4.1.10 NAT 地址转换.....	97
3.3 cinder-scheduler.....	77	4.2 Neutron 的由来	98
3.4 cinder-volume.....	79	4.2.1 nova-network 的问题	99
3.5 cinder-backup.....	80	4.2.2 Neutron 项目要解决的问题	100
3.6 Cinder 对存储设备及 Ceph 的 支持.....	81	4.3 Neutron 的架构	100
3.7 Nova 与 Cinder 的交互流程分析.....	81	4.3.1 Neutron API	102
3.8 Cinder 功能及典型工作流程	84	4.3.2 Neutron 插件及代理介绍	103
3.8.1 cinder-api 服务启动流程	84	4.3.3 ML2	104
3.8.2 cinder-scheduler 服务启动 流程	85	4.3.4 Neutron 核心数据模型	105
3.8.3 cinder-volume 服务启动流程	86	4.3.5 Neutron 消息交互	106
3.8.4 cinderclient 部分创建流程.....	86	4.3.6 租户网络与提供商网络	106
3.9 Glance.....	88	4.3.7 OpenStack 网络部署架构.....	107
3.10 K 版本的存储管理新功能.....	89	4.3.8 业务处理流程简述	108
3.10.1 Glance 新功能	89	4.4 K 版本新功能.....	112
3.10.2 Cinder 新功能.....	90	4.5 小结.....	113
3.11 小结.....	90	第 5 章 计量与监控 (Ceilometer)	114
第 4 章 网络管理模块 (Neutron)	92	5.1 Ceilometer 的体系架构	114
4.1 网络基本概念.....	93	5.2 Ceilometer 计量数据采集机制	116
4.1.1 L2 与 L3.....	93	5.2.1 概述	116
4.1.2 交换机与路由器	93	5.2.2 计量数据采集	117
4.1.3 防火墙	94	5.2.3 计量数据转换和发布.....	122
4.1.4 负载均衡	94	5.2.4 计量数据存储	123
4.1.5 DHCP 服务.....	94	5.3 Ceilometer 告警.....	124
4.1.6 子网和 ARP.....	94	5.4 Ceilometer API 服务器.....	126
4.1.7 VLAN	95	5.5 Ceilometer 的二次开发	127
		5.5.1 Notification Listener 插件 开发	128

5.5.2	Pollster 插件开发	130	6.4.7	模板依赖	156
5.5.3	Discovery 插件开发	131	6.5	Heat 资源类型	156
5.5.4	Compute Agent Inspector 插件 开发	132	6.5.1	资源类型的使用	156
5.5.5	Publisher 插件开发	133	6.5.2	资源类型的实现	157
5.6	OpenStack 组件计量	134	6.6	Heat 引擎	158
5.6.1	Nova 计量	134	6.7	典型场景分析	160
5.6.2	Glance 计量	138	6.8	K 版本新特性	163
5.6.3	Cinder 计量	138	6.9	与 AWS CloudFormation 的对比	164
5.6.4	Swift 计量	139	6.10	小结	165
5.6.5	Neutron 计量	139	第 7 章 裸机管理 (Ironic)		167
5.6.6	Keystone 计量	140	7.1	裸机管理通用技术	167
5.6.7	Heat 计量	141	7.2	Ironic 介绍	169
5.6.8	Ironic 计量	141	7.3	Ironic 架构	169
5.6.9	Ceph 计量	141	7.4	基本运行原理和流程	171
5.7	K 版本新功能	142	7.5	Ironic 安装	173
5.8	Ceilometer 对接外部系统	143	7.5.1	完全安装	173
5.9	OpenStack 监控	144	7.5.2	简化安装	190
5.10	小结	148	7.6	K 版本新功能	192
第 6 章 编排 (Heat)		150	7.7	小结	194
6.1	Heat 概述	150	第 8 章 消息总线		195
6.2	Heat 架构	150	8.1	概述	195
6.3	Heat API	151	8.2	AMQP 消息队列协议	196
6.4	Heat 模板	152	8.3	OpenStack 支持的消息总线类型	198
6.4.1	模板结构	152	8.4	小结	200
6.4.2	输入参数	153	第 9 章 OpenStack 大数据 服务 (Sahara)		201
6.4.3	资源	153	9.1	Sahara 概述	201
6.4.4	资源依赖	154			
6.4.5	输出参数	154			
6.4.6	模板执行	155			

9.1.1 Sahara 的定位	201	11.2.2 启动	235
9.1.2 Sahara 的发展历程	202	11.2.3 用户认证和令牌获取	236
9.1.3 Sahara 的主要特点	203	11.2.4 签名证书生成	238
9.2 Sahara 的使用模式	205	11.2.5 多级 Keystone 架构	240
9.3 Sahara 的架构	207	11.2.6 Keystone 与现有用户安全 认证系统的对接	241
9.3.1 Sahara 外部架构	207	11.3 K 版本新特性	242
9.3.2 Sahara 内部架构	208	11.4 基于可信计算的云安全体系	242
9.4 Sahara 与 EMR、Serengeti 的 对比	210	11.4.1 可信计算平台	242
9.4.1 Sahara 与 Amazon EMR 的 对比	210	11.4.2 OpenStack 中的可信计算池	244
9.4.2 Sahara 与 VMware Serengeti 的 对比	211	11.5 小结	246
9.5 K 版本新特性	213	第 12 章 OpenStack 孵化项目简介	248
9.6 小结	214	12.1 消息队列服务 (Zaqar)	249
第 10 章 OpenStack 数据库服务 (Trove)	215	12.1.1 概述	249
10.1 Trove 概述	215	12.1.2 使用场景	251
10.2 Trove 总体构架	216	12.1.3 逻辑架构	251
10.3 Trove 主要功能和 API 接口	218	12.1.4 本节小结	252
10.4 Trove 的安装和配置	221	12.2 共享文件系统服务 (Manila)	253
10.5 Trove 创建实例过程	225	12.2.1 概述	253
10.6 Trove 二次开发	226	12.2.2 使用场景	254
10.7 小结	229	12.2.3 逻辑架构	255
第 11 章 OpenStack 安全方案	230	12.2.4 本节小结	259
11.1 OpenStack 安全概述	230	12.3 DNS 管理服务 (Designate)	259
11.2 Keystone	231	12.3.1 概述	259
11.2.1 Keystone 介绍	232	12.3.2 使用场景	260
		12.3.3 逻辑架构	260
		12.3.4 本节小结	262
		12.4 密钥管理服务 (Barbican)	263
		12.4.1 概述	263

12.4.2 使用场景	263	12.5.2 使用场景	268
12.4.3 逻辑架构	264	12.5.3 逻辑架构	268
12.4.4 本节小结	266	12.5.4 本节小结	271
12.5 容器管理服务 (Magnum)	266		
12.5.1 概述	266	参考文献	273

云管理平台概述

云服务的核心在于服务所运行的技术平台，云平台在计算、存储和网络等方面为云服务提供支撑，为用户提供所需要的 IT 资源。云管理平台允许开发者或将写好的程序放在“云”里运行，或使用“云”提供的服务，或两者皆有。至于这种平台的名称，现在可以听到不止一种，比如按需平台（on-Demand Platform）、平台即服务（Platform as a Service, PaaS）等。但无论它叫什么，这种新的支持应用的方式无疑有着巨大的潜力。云平台是云计算的重要组成部分，如图 1-1 所示。云平台将虚拟化的计算资源、存储资源、网络资源统一管理，并面向用户提供服务，形成了云服务。可以说没有云平台，就没有云计算和云服务。



图 1-1 云平台基础架构

1.1 主流云管理平台对比

目前，业界中有 4 种有影响力的主流开源软件平台，分别是 OpenStack、CloudStack、Eucalyptus、OpenNebula。同时 VMware 作为云的商业软件提供商，也有很大的影响力。本章节对这几个云平台做简单的对比。

如表 1-1 所示，对 4 种开源云平台从背景、架构、商业模式等方面进行了全面的对比，希望读者有更全面的认识。

表 1-1 4 种开源云平台对比

	OpenStack	CloudStack	Eucalyptus	OpenNebula
项目背景	为 Rackspace 与 NASA 共同发起的开源项目，DELL、Citrix、思科和国内众多厂家也对其做出了重要的贡献	源于 2008 年成立的 VM Ops 公司，2012 年 4 月加入 Apache 基金会，此前采用 GPLv3 授权协议	为加利福尼亚大学圣芭芭拉学院研究项目，2009 年成立公司，实现商业化运营，仍对开源项目进行维护和开发	为由欧洲研究学会发起的虚拟基础设施的计划，2008 年发布首个开放源代码版本，2010 年起大力推进开源社区的建设
架构概述	以 Python 语言编写，包含 Nova、Neutron、Cinder 等多个主要模块，提供虚拟计算、网络、存储资源管理，提供类似 Amazon 的云 IaaS 服务	使用 Java 语言编写，采用框架+插件的系统架构，通过不同插件来提供对不同虚拟化技术的支持。CloudStack 提供 3 种管理云资源的途径：Web 界面、命令行和全功能 RESTful API	开发语言为 Java、C/C++，包括云控制器 (CLC)、Walrus、集群控制器 (CC)、存储控制器 (SC) 和节点控制器 (NC) 5 个主要组件，是 Amazon EC2 的一个开源实现	采用驱动层、核心层、工具层三层架构，驱动层负责虚拟机的创建、启动和关闭，监控虚拟机运行状态；核心层负责对虚拟机、存储、网络等进行管理；工具层通过命令行、浏览器和 API 方式提供用户交互接口
授权协议	Apache2.0 授权协议	GPLv3 授权协议，2012 年宣布加入 Apache 基金会，使用 Apache2.0 授权协议	社区版采用 GPLv3 授权协议，企业版使用自定义的商业授权协议	Apache2.0 授权协议
虚拟化支持	VMware\Xen\KVM\Power-VM\Hyper-V	VMware\Xen\KVM\Oracle-VM\Hyper-V	VMware\Xen\KVM	VMware\Xen\KVM
商业模式	免费使用	社区版免费使用；企业版提供增强功能和技术支持	社区版免费使用；企业版按处理器核数收费	社区版免费使用；企业版将社区版打包，提供补丁等程序的访问权限，以订阅的方式提供服务
总结	拥有超高的社区开发人气和庞大的生态系统，OpenStack 已经发布到 Kilo 版本（简称为 K 版本）。企业可很容易地将数据和应用迁移到公共云中。主流的 Linux 操作系统都支持 OpenStack。OpenStack 在可扩展性上有优势	在进入 Apache 阵营之前，CloudStack 在商业领域进行了长期的积累，帮助了近百个大规模生产平台，并实现了数十亿美元的运营收入；也提供友好的用户界面和丰富的功能，用户体验好，安装简单	从大学发源，有浓厚的研究风格，全面兼容 Amazon API，已经拥有虚拟化环境的用户能够使用它增强自己的虚拟化环境	项目启动早，一直处于稳步发展状态。社区规模较小，主要参与者为支持和参与该项目的企业人员，以及少量用户，用户能够获取的技术支持和交流空间有限

VMware 与 OpenStack 相比，两者在设计原则、商业模式等方面都有所不同，导致其在架构、功能、实施和维护性方面有一定的差异。VMware 软件套件是以虚拟化技术为核心，自底向上的架构，下端边界为虚拟机管理器。像 VMware 的 vSphere 和 vCloud director 产品都依赖于免费的 ESX(i) 虚拟机管理器，ESX(i) 虚拟机管理器能提供非常优秀的部署架构。VMware 的软件套件也是经过全面测试的，并且都有单一部署框架。总的来说，VMware 的产品由于其架构的健壮性，很多高规格用户在多数据中心规模的环境中都会使用。而 OpenStack 作为一个开源系统，没有任何一家单独的公司控制 OpenStack 的发展路线。OpenStack 本身是年轻的，但是却具有巨大的市场动力，与此同时，很多大公司都在支持 OpenStack 的发展。有了如此多公司的资源投入，OpenStack 的发展将是多元化的。

从具体功能来看，VMware 的核心功能是 VMware vMotion。它是 vSphere DRS、DPM 和主机维护三大功能的集合，同时 VMware 具有 FT 高容错、跨数据中心的容灾迁移等特色功能。OpenStack 也支持虚拟机的动态迁移，KVM 动态迁移允许一个虚拟机由一个虚拟机管理器迁移到另一个，说详细一点，你可以基于共享存储来来回回地将一台虚拟机在 AMD 架构主机与 Intel 架构主机上进行迁移。OpenStack 目前并不支持 FT 高容错等特色功能，但 OpenStack 的优势在于开放的架构以及对广大的 IT 设备厂家硬件的支持，各个厂家可以基于 OpenStack 的架构，开发出很多特色功能。FT 功能的实用性存在问题，也并不能保证备机状态的完整性。因此从应用角度我们可以看到，在功能的支持和功能细节方面，OpenStack 相比 VMware 还是有差距的。但是 OpenStack 还是有优势的，因为与 VMware 昂贵的价格相比，OpenStack 免费、开放的优势明显。对于 VMware 高投入带来的功能，OpenStack 可以免费提供给客户大部分。从 VMware 在功能方面的领先优势可以看出，VMware 还在继续研发除了 vMotion、高可用、容错以外其他的新功能，以保护它们的虚拟；OpenStack 一方面跟随 VMware 的脚步，另一方面也投入精力在支持更多的硬件厂商解决方案上。

关于 OpenStack、OpenNebula、Eucalyptus、CloudStack 社区的活跃度，这里借用蒋清野先生提供的一组对比数据，以对这些社区的活跃度进行分析和比较，如图 1-2 和图 1-3 所示。

图 1-2 和图 1-3 分别是上述 4 种开源项目的相关社区每个月所产生的讨论主题数量和帖子数量。可以看出：

1) 从 2012 年开始，与 OpenStack 和 CloudStack 相关的主题讨论数量在同一水平上，与 Eucalyptus 和 OpenNebula 相关的主题讨论数量在同一水平上。

2) 从 2012 年开始，与 OpenStack 和 CloudStack 相关的主题讨论数量远大于与 Eucalyptus 和 OpenNebula 相关的主题讨论数量。

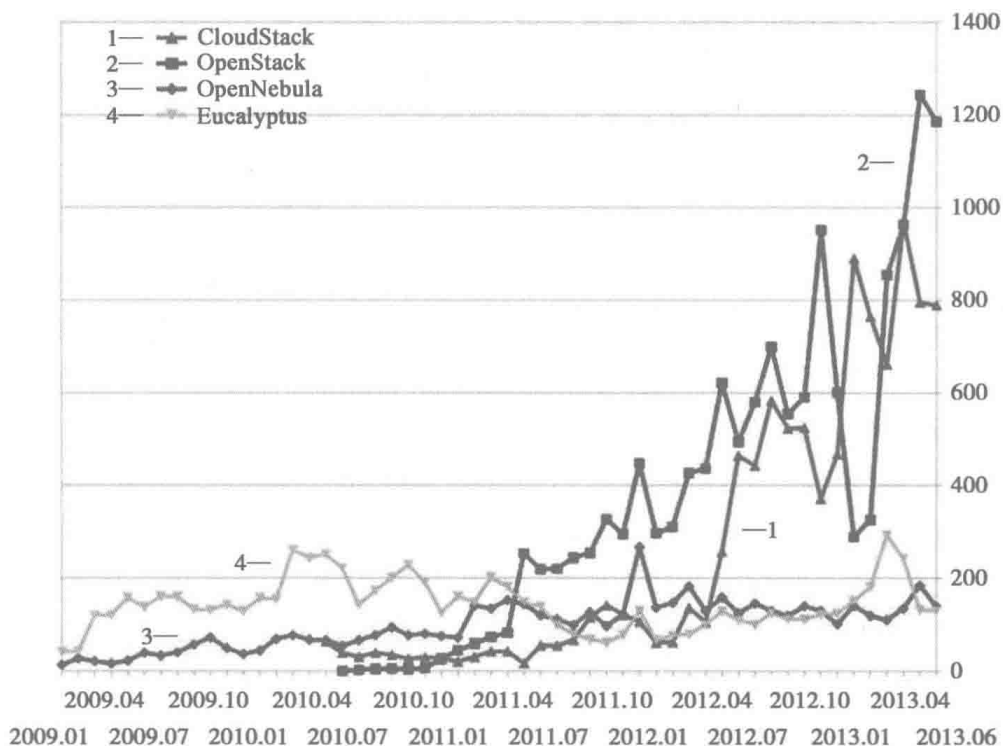


图 1-2 4 种开源云平台论坛主题数量对比

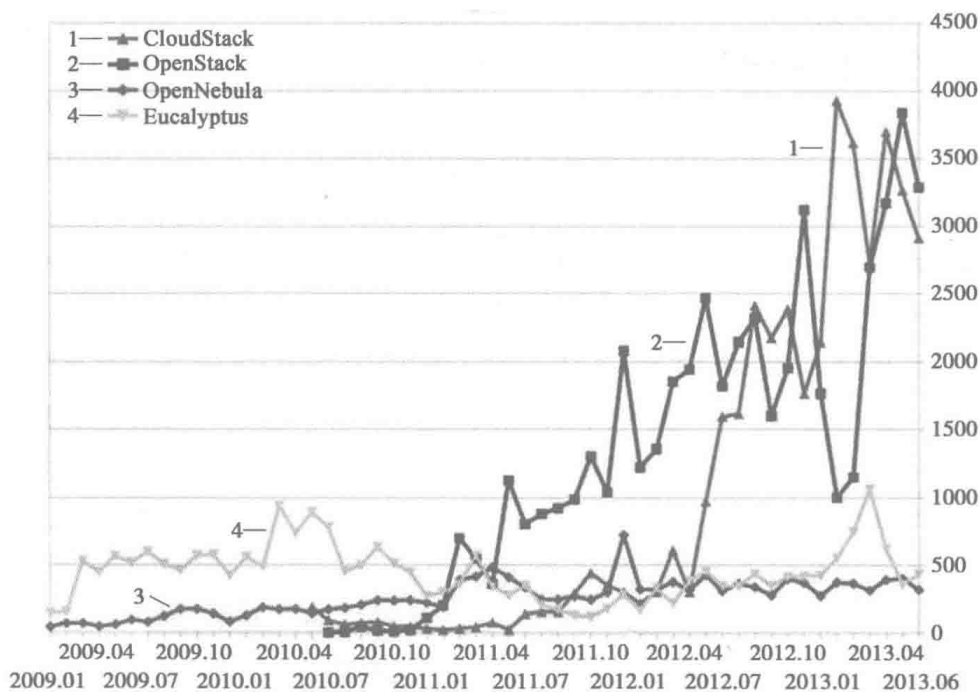


图 1-3 4 种开源云平台论坛帖子数量对比