

自主实现 SDN虚拟网络与企业私有云

YY游戏云平台组 著

这项全新的远大战略旨在把强大得超乎想象的计算能力分布到众人手中



自主实现 SDN虚拟网络与企业私有云

YY游戏云平台组 著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

本书深入浅出地介绍了 YY 游戏云平台团队在云计算领域的心路历程和实践经验,不仅总结了多年来在 OpenStack 驱动的 Cloud 1.0 上开发和使用的经验教训,而且花大幅笔墨深入讲解了自主研发的私有云平台 Cloud 2.0 的设计和实现。

本书内容全面而详尽,依次讲解了 Cloud 2.0 的选型思路、基于 VXLAN 技术的 VPC 网络架构设计和实现、业务层架构设计和实现、基于 libvirt 的虚拟计算实践、基于 Ceph 的虚拟存储实践、云数据源等产品的架构选型及实现、容量管理的详细思路,以及云平台方方面面的测试,是云计算领域实践类书籍中不可多得的一本好书。

本书干货众多,不仅适合初入云计算领域的读者阅读,更适合开发人员参考学习和实践。本书各章均可独立成册,读者可以根据自己的需要来阅读。

未经许可,不得以任何方式复制或抄袭本书之部分或全部内容。
版权所有,侵权必究。

图书在版编目(CIP)数据

自主实现 SDN 虚拟网络与企业私有云/YY 游戏云平台组著. —北京:电子工业出版社,2017.4
ISBN 978-7-121-31015-7

I. ①自… II. ①Y… III. ①计算机网络—网络结构 IV. ①TP393.02

中国版本图书馆 CIP 数据核字(2017)第 043502 号

策划编辑:张春雨

责任编辑:葛娜

印刷:北京天宇星印刷厂

装订:北京天宇星印刷厂

出版发行:电子工业出版社

北京市海淀区万寿路 173 信箱

邮编:100036

开本:787×980 1/16 印张:19

字数:411 千字

版次:2017 年 4 月第 1 版

印次:2017 年 4 月第 1 次印刷

定 价:69.00 元

凡所购买电子工业出版社图书有缺损问题,请向购买书店调换。若书店售缺,请与本社发行部联系,联系及邮购电话:(010) 88254888, 88258888。

质量投诉请发邮件至 zltz@phei.com.cn, 盗版侵权举报请发邮件至 dbqq@phei.com.cn。

本书咨询联系方式:010-51260888-819 faq@phei.com.cn。

推荐序一

在群星璀璨的信息技术领域，云计算无疑是最耀眼的一颗星星。从 2006 年春天亚马逊上线第一个云计算产品 Amazon S3 以来，云计算市场和技术都取得飞速发展。2015 年全球云计算市场规模达到 522 亿美元，预计到 2020 年将达到 1435 亿美元。

在全球主流技术公司纷纷拥抱云计算之时，YY 游戏早在 2013 年就搭建了自己的私有云平台，利用云计算驱动游戏业务。一些月流水过亿的知名游戏，都曾经运行在我们的私有云平台上。

YY 游戏云平台团队在过去几年对云计算的研发、运维中，积累了非常多的经验，并在此基础上，自主实现了一套私有云平台，以更好地满足公司业务日新月异发展所带来的不断挑战。今后，云平台不止服务游戏，还将服务虎牙直播，利用云平台的分布式计算能力，带给用户更好的体验和感受。

欢聚时代是一个技术驱动的公司，将近 60% 员工为研发人员。不管是游戏平台，还是虎牙直播，我们在技术研发上的投入都领先于同行业。希望 YY 游戏云平台继续加强发展，为云计算自主创新领域树立一个标杆。

本书是团队的智慧结晶，用户在阅读本书后，将对云计算本质有一个全新的理解。

董荣杰

欢聚时代执行副总裁

推荐序二

随着近年来云计算应用技术的快速发展和逐渐成熟，企业如何上云、如何用云是摆在 CIO 面前的重要课题。各种解决方案有其特定应用场景，企业要综合考虑业务情况、资金投入、技术能力等多种因素进行选择——是削足适履，对自身 IT 架构进行调整去适应相应平台？还是量体裁衣，根据自身业务定制自有云平台？无论采用何种方案，都需要从业务角度出发，对所涉及的相关技术有足够深刻的理解。

云平台所涉及的技术和运营复杂度超过了以往大多数 IT 系统，不少自建私有云或行业云的企业大都经历了如下两个阶段。

一是战略规划和选型。

从战略层面看，企业选择自建私有云或构建更大的行业云，往往都是从自身业务发展和战略角度出发的，尤其是业务对 IT 有重度依赖的企业。一方面，这些企业对产品技术的要求和复杂度超过了市场上公有云的能力，不希望在基础设施架构方面受制于人；另一方面，IT 基础架构本身就是核心竞争力，如果在此方面能有所突破，就可构筑本行业的技术壁垒。

从产品技术层面看，云平台技术选型中的计算、网络、存储、安全都是要重点考虑的核心因素。以 OpenStack 为主流的开放平台极大地加速了技术选型的过程，并带来了丰富的应用、社区资源以及技术人才。云平台所涉及的技术体系广泛，一般企业很难有完整的驾驭能力，比较实际的做法是采用以开源平台为基础、自有研发为核心、专业厂商为关键技术补充的模式。

二是上线和迭代。

云平台上线之后，企业根据自身业务特点，一般会将部分业务迁移至云平台进行试运行，从技术、运营方面不断磨合，为后续更大的业务上线做积累储备。但业务上线之后，一般都会

面临诸多不适应或者难题，如开源软件的本身质量问题、功能缺陷、网络性能瓶颈、磁盘 IO 性能瓶颈、运维复杂度高、安全问题等，这些都对提供高品质的业务基础团队构成了挑战。

云计算大量使用了虚拟化、分布式等新技术，在新技术和传统 IT 产品融合的过程中会遇到诸多问题。比如在大规模云平台中，当防火墙很难满足性能要求时如何扩展；如何利用现有 SDN 交换机的性能优势，为租户网络提供高性能的东西向网关；在虚拟网络中，云网络虚拟化技术使得传统的监控产品如 NPM 很难放置探针点，如何进行快速运维诊断分析，这些在技术方案中都需要考虑，并且必须有相应的解决方案。

云计算所采用的技术对支撑的软硬件系统既有要求降低的部分，如服务器、存储等，采用大量性价比更高的通用服务器，更大发挥软件的价值；也有要求增加的部分，如 SDN 交换机、NFV 功能、数据分析产品，即要考虑传统 IT 产品如何能以服务方式交付到云环境中。在对厂商降低依赖的同时要求平台运维者能够深刻理解平台运行原理，在硬件故障、软件故障、外部因素（如网络攻击）、内部因素（如人为误操作）等任何情况下都能做到自主可控，才可能有 SLA 的保证。

在和 YY 云平台团队多次深入交流中，感受最深刻的是他们对系统架构和技术方向的深刻把握，这一方面来自于运营 RiseCloud（升龙云）云平台多年的技术积累，另一方面来自于重负载的 YY 游戏业务需要持续催生、迭代出最适合的技术方案。RiseCloud 是一个游戏行业云，从技术角度看，游戏业务对平台的要求均衡且很高。RiseCloud 能够承载大量的重载游戏，从性能、服务品质和产品功能成熟度看，该平台已经完全是一个成熟的商业化平台。

YY 云平台团队对技术的开放态度很令人钦佩，一个规模不大的技术团队，能够自研 RiseCloud 并将其运营起来，是因为他们能与业界、社区广泛合作，吸纳诸多领先技术和理念。同时他们能够把积累多年的经验写成案例共享出来，也非常难能可贵。本书涉及了 RiseCloud 完整的系统设计理念和核心技术实现方案，相信是对开发、运营平台多年来的经验、教训的总结。尤其是在 SDN 的相关实践方面，书中的案例非常典型，可以用来解决 OpenStack 网络方面的诸多问题。这些对于关注大规模云计算平台体系架构、技术选型、功能实现、平台运维的读者来说，将获益良多。

张天鹏

云杉网络 CTO

推荐序三

收到《自主实现 SDN 虚拟网络和企业私有云》书稿后，只是目录，就吸引我迫不及待地一口气通读了一遍。第一感觉是国内搞云计算的朋友有福了，你们以后必定要经常翻阅此书，因为这样做可以节省大量的摸索时间。

YY 的私有云规模非常大，游戏业务对云的考验也非常大，在计算能力、磁盘 IO 能力、网络质量方面都有非常高的要求。尤其是端游，比手游要求更高，压力更大，所以一般在公有云上都很少跑端游。但是 YY 的私有云不光跑了大量的页游、手游，还跑了一定数量的端游，证明 YY 的私有云切切实实经受住了实践的检验。

《自主实现 SDN 虚拟网络和企业私有云》分享了 YY 私有云建设的思路、方法及大量实践。更难得的是，本书不是一个人而是一个团队的经验，含金量非常高，章章都是干货，尤其是 YY 在私有云网络建设方面独特的解决方案和实践。

肖力
云技术社区创始人

推荐序四

YY 是国内最早做 PC 端直播的，而且在很长一段时间内很火爆。本书写得很好，从最初的调研云平台，到使用过程中的踩坑总结都不乏亮点。YY 能够不随波逐流，摒弃目前流行的 OpenStack，还是很有勇气的。还是那句话：只有适合自己的才是最好的。Ceph 是目前主流的开源分布式存储软件。看到越来越多的云服务商和企业用户开始用 Ceph，作为国内较早布道、推广 Ceph 的人来说很是高兴和激动。总体来说，YY 云平台沉着稳重又不缺乏创新和突破，而出自平台建设者之手的本书我觉得也是国内 IT 企业建设私有云最值得借鉴的一本。

耿航

Ceph 中国社区联合创始人 & XSKY 市场技术专家

推荐序五

在 SDN 盛行的当下，专注于 SDN 与私有云实践类的书籍甚少。

《自主实现 SDN 虚拟网络和企业私有云》一书能及时填补空白，让广大从业者犹如久干逢雨露。

书中介绍了 YY 私有云的搭建思路、技能选型以及技术细节，为在 SDN 和私有云前沿探索的战士们指路、点灯，可谓功德无量。

胡凯

bilibili 运维负责人

前言

本书由 YY 游戏云平台团队编写。YY 游戏在平台技术上保持开放的心态，勇于采用新兴的技术来促进业务发展、降低成本。从 2013 年起，YY 大部分游戏都运行在私有云平台上。这个私有云平台也经历了两代的发展，从第一代基于 OpenStack 的私有云，到第二代完全自主实现的 SDN 和私有云，历经磨难，实现了质的改进。

在 OpenStack 如日中天时期，我们也选择它来部署第一代私有云。但是很快发现 OpenStack 在实际运行中存在诸多问题，这些问题表明 OpenStack 在没有足够的开发、运维力量的情况下，不适合大规模的生产应用。当然，正是在使用 OpenStack 的过程中，我们学习了它的一些设计思想，例如虚拟计算的调度、虚拟网络的实现等，这些经验为我们研发第二代云平台做了有效铺垫。

从 2015 年 8 月起至 2016 年 6 月止，历经将近一年的时间，我们自主研发和上线了第二代私有云平台。第二代私有云平台完全抛弃了 OpenStack，它的体系架构、业务组件、功能模块全部由我们自己实现，其中涉及对虚拟计算、虚拟存储、虚拟网络的深入分析和调研。尤其是虚拟网络，我们跟华为、华三、云杉等公司进行了大量交流，最终确定了适合自己的 SDN 产品和方案。

在第一代云平台向第二代云平台转型过程中，以及在第二代云平台的研发、运维中，团队进行了大量实践，并积累了许多宝贵经验。团队成员深入理解云计算本质，取其精华为己所用，构建了一个稳固的私有云平台，并成功服务于页游、端游、手游的运营。

今天，我们把建设私有云平台的实战经验以写书的方式分享出来。阅读完本书后，你会发现私有云并没有那么复杂，也并非高不可攀。只要搞懂了它的体系架构，以及每个子系统（计算、存储、网络）的实现方式，再结合企业自身的业务特点，自主构建一个私有云平台不是很难的事。

本书对于对云计算感兴趣、有计划建设私有云的用户，格外具有参考、借鉴意义。

本书内容由 YY 游戏云平台组的开发、运维、QA 工程师共同完成，包括：

- 风河，编写了第 1、2 章背景介绍部分；
- 张春，编写了第 3 章虚拟网络部分；
- 何招武、张兴平，编写了第 4 章云平台业务部分；
- 朱辉，编写了第 5 章虚拟计算部分；
- 戚昱，编写了第 6 章虚拟存储部分；
- 张博，编写了第 7 章云数据库部分；
- 刘亚丹，编写了第 8 章云平台容量管理部分；
- 黄书明、廖志委、任方超，编写了第 9 章测试与安全部分。

上述同学在各自领域都是资深工程师，感谢他们在工作之外的辛勤付出，才有了本书。同时感谢 QA 组的马飞同学在本书编写过程中的项目管理工作。也感谢电子工业出版社的编辑们对本书的排版、校对工作。

如果读者对书中内容有疑问或建议，请发邮件反馈给我们：g-yygame-booking@yy.com。

轻松注册成为博文视点社区用户 (www.broadview.com.cn)，您即可享受以下服务：

- **提交勘误：**您对书中内容的修改意见可在【提交勘误】处提交，若被采纳，将获赠博文视点社区积分（在您购买电子书时，积分可用来抵扣相应金额）。
- **与作者交流：**在页面下方【读者评论】处留下您的疑问或观点，与作者和其他读者一同学习交流。

页面入口：<http://www.broadview.com.cn/31015>

二维码：



目录

第 1 章 绪论	1
1.1 云计算发展趋势	1
1.2 YY 游戏使用云平台的经验	3
1.3 云计算随想	5
第 2 章 选型思路	8
2.1 为什么放弃 OpenStack	8
2.2 Cloud 2.0 研发思路	9
2.3 发展规划：基于云的 VDC 实现	14
第 3 章 Cloud 2.0 虚拟网络实现	17
3.1 Cloud 1.0 的实践经验	17
3.1.1 Neutron 与企业私有云	18
3.1.2 问题与不足	19
3.1.3 拥抱 SDN	22
3.2 虚拟网络架构	26
3.2.1 Overlay 网络模型	26
3.2.2 虚拟网络架构	33
3.2.3 网络设备技术要点	37
3.3 虚拟网络实现	41
3.3.1 Underlay 网络	41
3.3.2 Overlay 网络——VXLAN VPC	44

3.3.3	SDN 的核心——控制器实现.....	49
3.3.4	服务如臂使指——北向接口 API.....	54
3.3.5	网络触手可及——南向控制协议.....	57
3.3.6	SDN 与 NFV.....	61
3.4	虚拟网络业务.....	62
3.4.1	Underlay 网络配置流程.....	62
3.4.2	云主机创建流程.....	65
3.4.3	云主机迁移流程.....	66
第 4 章	云平台业务.....	67
4.1	业务组件.....	67
4.2	业务架构.....	69
4.3	安全子系统.....	70
4.3.1	用户安全.....	70
4.3.2	组件安全.....	72
4.3.3	技术实现.....	75
4.4	调度子系统.....	90
4.4.1	流程引擎设计.....	90
4.4.2	统一实现审计.....	94
4.4.3	乐观锁加记录锁的并发控制.....	94
4.4.4	异步线程池管理.....	96
4.4.5	基于 Redis 实现的分布式锁.....	97
4.4.6	基于 Hibernate 与 Spring JDBC 的灵活持久层.....	100
4.5	云控制台.....	101
4.5.1	为什么选择 AngularJS.....	102
4.5.2	开发心得总结.....	105
第 5 章	虚拟计算.....	116
5.1	虚拟化概述.....	116
5.2	KVM/QEMU/libvirt 浅析.....	117
5.2.1	KVM 简介.....	117
5.2.2	KVM 与 QEMU.....	118
5.2.3	libvirt 介绍.....	119

5.3	KVM 虚拟化环境安装.....	120
5.3.1	APT 源安装.....	120
5.3.2	源码编译安装.....	120
5.4	使用 qemu-img 管理虚拟机磁盘镜像.....	121
5.4.1	qemu-img 基本命令.....	122
5.4.2	在宿主机上如何挂载镜像文件.....	124
5.5	使用 libvirt 管理 KVM 虚拟机.....	124
5.5.1	libvirt Java API 的使用.....	124
5.5.2	虚拟机 XML 配置文件详解.....	133
5.5.3	virsh 常用命令.....	141
5.6	实战系列.....	142
5.6.1	使用 Cloudinit 实现虚拟机启动初始化.....	142
5.6.2	在线更改虚拟机内存大小.....	146
5.6.3	热添加虚拟机 CPU.....	147
5.6.4	如何限制虚拟机磁盘 I/O.....	148
5.6.5	在虚拟机内部如何正确自动挂载磁盘.....	148
5.6.6	虚拟机如何使用 Ceph 块设备.....	149
5.6.7	libvirt hook 机制.....	151
5.6.8	如何支持使用 virsh 控制台登录虚拟机.....	152
5.6.9	虚拟机如何通过 OpenvSwitch 接入网络.....	153
5.6.10	宿主机如何通过 qemu-guest-agent 与虚拟机通信.....	154
5.6.11	虚拟机的迁移.....	156
第 6 章	虚拟存储.....	161
6.1	概念和术语.....	162
6.2	硬件配置.....	162
6.2.1	Ceph 网络配置.....	163
6.2.2	服务器配置.....	165
6.3	软件配置.....	166
6.4	部署.....	168
6.4.1	设置副本分布到不同的机架上.....	168
6.4.2	创建 Pool.....	172
6.4.3	测试 PG 副本的机架分布性.....	173

6.5	监控	173
6.5.1	监控层次	173
6.5.2	与 Zabbix 监控系统集成	174
6.5.3	告警条件	174
6.5.4	监控面板	174
6.6	性能测试和调优	176
6.6.1	块设备性能测试	176
6.6.2	调优	177
6.7	维护操作	183
6.7.1	使用 systemctl 管理 Ceph 进程	183
6.7.2	OSD 机器重启	184
6.7.3	扩容	184
6.7.4	升级 Ceph 软件版本	184
6.8	故障定位和处理	185
6.8.1	查看集群状态	185
6.8.2	日志	186
6.8.3	MON	187
6.8.4	OSD	188
6.8.5	PG	192
6.8.6	实际运维中的问题	196
第 7 章	云数据库	199
7.1	云数据库服务功能介绍	199
7.2	1.0 版本三层架构	199
7.2.1	Manager 开放 API 的认证方式	201
7.2.2	Manager 节点的调度策略	202
7.2.3	Manager 节点的配额管理	203
7.2.4	Backend 节点的资源	203
7.3	云数据库 2.0 版本架构演化及改进	214
7.3.1	云数据库 1.0 版本的问题与不足	215
7.3.2	云数据库 2.0 版本的改进	215
7.3.3	CloudMySQL 2.0 Agent 的设计模型	217
7.3.4	CloudMySQL 2.0 KVM 配置	221

7.4	小结	225
第 8 章	云平台容量管理	226
8.1	容量管理概述	226
8.2	容量管理真实案例	227
8.3	容量管理特点和成熟度	229
8.3.1	ITIL 容量管理特点	229
8.3.2	云平台容量管理特点	229
8.3.3	云平台容量管理成熟度模型	230
8.4	容量管理组件	231
8.5	容量管理总体模型	231
8.5.1	容量与其他系统的关系	233
8.5.2	网络模型	234
8.5.3	计算模型	236
8.5.4	存储模型	237
8.5.5	云 MySQL 模型	239
8.6	容量预测	241
8.7	容量管理可视化	242
8.8	容量管理之资源采购	243
8.9	小结	243
第 9 章	云平台测试	244
9.1	云平台测试策略	244
9.1.1	云平台项目测试特性	244
9.1.2	测试方法与策略	245
9.1.3	测试环境	248
9.2	云计算测试	249
9.2.1	功能测试	249
9.2.2	自动化测试	252
9.3	云数据库测试	257
9.3.1	云 MySQL 测试	257
9.3.2	云 Redis 测试	260
9.4	云网络测试	262

9.4.1	虚拟网络测试.....	262
9.4.2	云平台业务的网络功能测试.....	273
9.4.3	迁移后的网络测试.....	277
9.5	云平台安全测试.....	280
9.5.1	API 安全测试.....	280
9.5.2	云 Redis 安全测试.....	283
附录 A	缩略词表.....	286