

虚拟化技术 原理与实现

The Principles & Practice of Virtualization Technology

广小明 胡杰 陈龙 郭京 等编著



· 十二五国家重点图书出版规划项目 ·
云计算实践指南丛书

转型时代丛书
中国电信北京研究院 专家奉献

The Principles & Practice of Virtualization Technology

虚拟化技术 原理与实现

广小明 胡杰 陈龙 郭京 侯光华 司伟 顾茜 编著

电子工业出版社
Publishing House of Electronics Industry
北京·BEIJING

内 容 简 介

本书对云计算中关键技术之一的虚拟化技术进行了深入的分析,从 x86 计算机体系结构以及操作系统的工作原理出发,介绍了虚拟化技术原理以及业界主流虚拟化软件产品,并以 Xen、KVM 开源软件为例分析了虚拟化软件的架构及其实现方法,最后对虚拟化软件管理接口的工作原理以及实现方法进行了全面的梳理。

本书注重技术理论与应用实践的紧密结合,可供从事云计算技术研究开发、设备制造、咨询设计、工程建设、运营维护与管理的技术人员和管理人员阅读,也可供高等院校通信工程专业、计算机专业师生参考,还可作为 IT 培训机构的培训参考书。

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

虚拟化技术原理与实现 / 广小明等编著. —北京: 电子工业出版社, 2012.10
(转型时代丛书)

ISBN 978-7-121-18528-1

I. ①虚… II. ①广… III. ①虚拟处理机 IV. TP338

中国版本图书馆 CIP 数据核字 (2012) 第 222700 号

策划编辑: 刘 皎

责任编辑: 李利健

印 刷: 三河市双峰印刷装订有限公司

装 订: 三河市双峰印刷装订有限公司

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: 720×1000 1/16

印张: 19.25 字数: 330 千字

印 次: 2012 年 10 月第 1 次印刷

印 数: 4000 册 定价: 59.00 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。

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。

服务热线: (010) 88258888。

转型时代丛书

指导委员会

主任委员：吴基传

副主任委员：杨杰

委员：陈俊亮 李未 韦乐平

邬贺铨 张继平（按拼音顺序排序）

编委会

主任：李志刚

副主任：侯春雨 赵慧玲

委员：毕奇 朱健 野永东 谢朝阳

陈自清 杨峰义 王晓平 张成良



转型时代丛书

总序

“变化，无论是突如其来的，还是循序渐进的，有时都会淘汰你认为理所当然的一切。”


——《转型》

21世纪以来，信息化更加快速而深刻地改变着这个世界，大到全球经济社会的发展格局，小到每个人的日常工作生活。许多国家把数字化、信息化、智能化作为国家战略的关键主题，把信息基础设施建设作为后金融危机时代振兴经济的重要手段。同样，我国“十二五”规划也把全面提高信息化水平，特别是加快建设下一代国家信息基础设施、推动信息化和工业化深度融合、推进经济社会各领域信息化作为重要工作列入其中。

信息通信产业中新技术、新业务的不断快速发展不仅催生着新的经济增长点，造就了谷歌、Twitter、腾讯等一个又一个明星企业，引领整个行业及社会经济的发展方向，更重要的是它对人们生产、生活产生了深刻而久远的影响。我们的生产资料不仅仅是机器，还有计算机、手机、互联网；我们通过点击“百度”打开未知世界，通过“淘宝”购买商品，利用手机登录“Facebook”去了解彼此、评论时政，所有这一切都表明信息通信产业正在更广，更深地影响着我们每一个人，互联网/移动互联网已成为像水、电一样的生产、生活“必需品”。

环顾全球，整个信息通信产业正在朝着宽带化、移动化、智能化发展，特别是3G的普及和LTE的逐步成熟使得移动互联网一跃成为整个行业中最前沿、最具革命性的领域。智能管道、物联网、下一代互联网、云计算等一个个新的理念、新的信息服务模式正在席卷全球成为新热点。而这一切变化，都将对从事信息服务的企业（包括电信运营商）带来前所未有的机遇和挑战。适者生存法则同样适用于多变的企业生态系统。无论是百年老店，还是创业新秀，只有顺应信息化时代发展潮流，重新审视并及时调整企业的商业模式，抓住信息化带来的重大机遇，才能在变化中顺势前进。

鉴于此，这套“转型时代丛书”既有对智能宽带网络、移动互联网、云计算等新技术、新网络的研究和实践总结，也有对商业模式、营销变革等现代管理中关键问题的长期探索。相信此系列书籍能帮助您了解趋势，廓清谜团，抓住机会，与信息化时代共同成长。



2012年7月

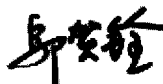
序

在科学技术发展到 21 世纪的今天，科技领域伴随着人类社会生活多样化正面临着新一轮的变革。经历了 20 世纪计算机和互联网两轮发展浪潮之后，互联网和移动互联网已经真正融入到社会生活的方方面面，随之而来的传统计算和服务模式也发展到一个转折点，我们不断面临着来自资源的无限需求和有限供给、总体社会资源高效利用和粗放管理等诸多挑战。云计算作为一种新型资源分配和使用模式，正是顺应了这种社会发展的要求，它为“移动为端，云为中心”的移动互联网业务体系打下了基础，起到了信息化社会发展的中坚作用。

在云计算的 SaaS、PaaS 和 IaaS 三个层面中，IaaS 承担着底层基础作用，而虚拟化技术又是 IaaS 技术的核心，它提供资源的多种颗粒度分配、动态可扩展和回收等手段。随着云计算在电信网的应用不断深入，作为全球最大的宽带运营商的中国电信组建了云计算研究中心，研究开发虚拟化技术及其在电信网的应用，重点针对开源的虚拟机软件 KVM 和开源的虚拟机监视器软件 Xen 的技术和配套的管理技术进行攻关，他们将研究的体会总结在本书中，与从事这方面技术研究或希望了解其应用的技术人员分享，共同推动云计算技术在开发网络业务方面发挥更大的作用。虚拟化技术及其在云计算中的应用还处于发展阶段，将它们应用到电信网在全球也处于起步阶段。

虚拟化技术在电信网领域的应用带来很多机遇的同时也面临很大的挑战，希望通过本书的抛砖引玉，能引起相关专业人员更多的关注，也希望更多的有志之士投入这一领域，在实践中发展和创新。

中国工程院院士



前 言

随着云计算热潮的兴起，构成其中关键技术之一的虚拟化技术再次成为业内外关注的焦点。但是与之前有所不同，虚拟化技术不再仅限于计算机从业者范围，而是走向更为广阔的融合宽带网络在内的广义的资源基础设施层面。更为重要的是，受虚拟化概念的影响，正在形成一种新的资源体系，动态组合调整所需的计算、存储和网络资源以适应最终应用的需求。这一变化直接带动创新的服务模式、快速部署和资源弹性扩展等一系列优势。

中国电信作为基础运营商，一直密切关注云计算技术和业务的发展，成立了专业化的研究中心对云计算相关核心技术进行研究，并将研究与中国电信自身业务服务的优势相结合。中国电信拥有世界上最大的宽带网络和国内最丰富的数据中心资源，如何将传统的物理资产转化成逻辑可管理的虚拟资源以提高资源的利用率，降低运营和维护成本，满足国家绿色环保的要求是研究重点之一。

当前基于 x86 的虚拟化技术已经成熟，但是其核心技术一直掌握在少数虚拟化软件厂家手中，阻碍了虚拟化技术的规模应用。开源技术 Xen 和 KVM 的出现以及该类技术在亚马逊成功的商用，给虚拟化技术的应用注入了强心剂，越来越多的企业和个人参与到开源技术的研究和开发过程中。随着开源软件版本的不断升级，其与管理平台的接口也日渐成熟和标准化，一些知名的平台开发厂商也开始支持这些开源技术，为开源技术的商用奠定了良好的基础。

本书作者均来自中国电信云计算研究中心，他们以饱满的工作热情参与到虚拟化开源技术的研发过程中，秉承开源的精神。希望自己的研究成果不再局限于中国电信研究中心内部，而且像云计算资源一样借助图书这个介质为更多的研究人员和开发者服务，让更多的人加入到开源技术的研究和开发过程中。

本书的主要内容包括 4 个部分。

第一部分为虚拟化技术原理篇。介绍虚拟化技术，包括它的发展历程、基本原理、技术架构和一些主流的虚拟化产品等。

第二部分为 Xen 虚拟化技术篇。介绍 Xen 软件模块结构、工作原理及流程，并借助源码分析 Xen 最核心的软件模块 Hypervisor 的工作原理、流程及实现。

第三部分为 KVM 虚拟化技术篇。从 KVM 的体系结构出发，分别从日常使用、KVM 内核代码和 qemu-kvm 用户态代码三个角度进行分析，情景化地剖析了 KVM 虚拟化技术的实现方案，包括重要的数据结构、处理流程等。

第四部分为虚拟化软件开放接口篇。从典型 Xen 管理接口入手，通过源码分析介绍了 Xen 管理接口的工作原理与实现方式，并推而广之。考虑到未来需要对 Xen、KVM 等虚拟化技术构建的异构资源池进行统一管理，需要提供统一标准化管理接口，本书最后对 libvirt 虚拟化控制中间件进行了介绍。

本书涉及的知识面较广，由于笔者的知识水平和认知的局限，书中难免有纰漏之处，恳请各位专家和读者不吝赐教。

作 者

目 录

第一篇 云计算与虚拟化技术

第 1 章 虚拟化技术基本原理	2
1.1 云计算与虚拟化技术	3
1.2 x86 和非 x86 体系结构基础	4
1.2.1 x86 的发展历程	4
1.2.2 x86-64	6
1.2.3 x86 内存架构	7
1.2.4 x86-64 的基本模式	23
1.2.5 x86-64 的寄存器组	25
1.2.6 中断与异常	26
1.2.7 I/O 架构	31
1.2.8 DMA	32
1.2.9 时钟	33
1.3 操作系统与虚拟化	34
1.3.1 操作系统	35
1.3.2 进程	35
1.3.3 系统虚拟化	38
1.3.4 系统虚拟化的发展历程	39
1.3.5 可虚拟化条件	41
1.3.6 虚拟化的原理与分类	43
1.4 VMM 技术架构分类	45
1.4.1 Hypervisor 模型	46
1.4.2 宿主 (Hosted) 模型	47
1.4.3 混合模型	48
1.5 本章小结	49

第 2 章 虚拟化实现技术架构	50
2.1 处理器虚拟化实现技术.....	52
2.1.1 Intel VT.....	53
2.1.2 AMD SVM.....	55
2.1.3 vCPU.....	55
2.2 中断虚拟化实现技术.....	56
2.3 内存虚拟化实现技术.....	58
2.3.1 影子页表.....	60
2.3.2 Intel EPT.....	65
2.3.3 AMD NPT.....	67
2.4 I/O 设备虚拟化实现技术.....	68
2.4.1 Intel VT-d.....	69
2.4.2 DMA 重映射.....	70
2.4.3 I/O 页表.....	73
2.4.4 AMD IOMMU.....	74
2.5 网络虚拟化技术.....	76
2.6 时间虚拟化技术.....	79
2.6.1 操作系统和客户机的时间概念.....	79
2.6.2 客户机时间概念的实现.....	82
2.7 主流虚拟化产品及其特点.....	84
2.7.1 Xen.....	84
2.7.2 VMware.....	86
2.7.3 Hyper-V.....	87
2.7.4 KVM.....	88
2.8 本章小结.....	90

第二篇 Xen 虚拟化技术

第 3 章 Xen 软件系统原理	92
3.1 Xen 软件模块结构.....	93
3.1.1 Xen Hypervisor.....	93
3.1.2 特权虚拟域 0 (Dom0).....	94
3.1.3 独立设备驱动域 (IDD).....	95
3.1.4 非特权虚拟域 U(DomU).....	96
3.1.5 硬件虚拟域 (HVM).....	96
3.2 Xen 系统启动工作原理及流程.....	96
3.2.1 系统引导过程.....	97
3.2.2 Hypervisor 启动与初始化过程.....	98

3.2.3	Dom0 启动过程.....	99
3.2.4	DomU 的启动.....	99
3.3	Xen CPU 虚拟化工作原理.....	100
3.3.1	x86 体系虚拟化存在的问题.....	100
3.3.2	CPU 虚拟化——半虚拟化（又称为泛虚拟化）.....	102
3.3.3	CPU 虚拟化技术——硬件虚拟化技术支持的全虚拟化.....	103
3.4	Xen 内存虚拟化工作原理.....	105
3.4.1	内存虚拟化——直接模式.....	106
3.4.2	内存虚拟化——影子模式.....	107
3.5	I/O 虚拟化工作原理.....	108
3.5.1	半虚拟化 I/O.....	108
3.5.2	全虚拟化 I/O.....	109
3.6	Xen 虚拟机（DomU）生命周期管理.....	110
3.7	本章小结.....	112

第 4 章 Xen Hypervisor 技术实现..... 113

4.1	Xen Hypervisor 关键技术概述.....	114
4.2	Hypercall.....	114
4.2.1	Hypercall 的实现机制.....	115
4.2.2	自定义 Hypercall 的方法.....	118
4.2.3	应用程序使用 Hypercall 的方法.....	120
4.3	事件通道.....	121
4.3.1	事件通道的初始化.....	121
4.3.2	事件通道的绑定.....	122
4.3.3	发送事件通知.....	136
4.3.4	事件通知的处理.....	138
4.4	数据共享.....	142
4.4.1	授权表（Grant table）.....	142
4.4.2	XenStore 和 XenBus.....	146
4.4.3	分离设备驱动.....	149
4.5	本章小结.....	154

第三篇 KVM 虚拟化技术

第 5 章 qemu-kvm 虚拟化解决方案..... 156

5.1	概述.....	157
-----	---------	-----

5.2	内核模块组成概述	158
5.2.1	KVM 的内核模块结构	158
5.2.2	Linux 内核源码中的 KVM	160
5.3	KVM 所提供的 API	162
5.3.1	KVM API 纵览	162
5.3.2	system ioctls 调用	163
5.3.3	vm ioctl 系统调用	164
5.3.4	vcpu ioctl 系统调用	165
5.4	KVM 内核模块重要的数据结构	168
5.4.1	KVM 结构体	168
5.4.2	kvm_vcpu 结构体	169
5.4.3	kvm_x86_ops 结构体	169
5.4.4	KVM API 中重要的结构体	171
5.5	KVM 内核模块重要流程的分析	173
5.5.1	初始化流程	173
5.5.2	虚拟机的创建	175
5.5.3	vCPU 的创建	177
5.5.4	vCPU 的运行	180
5.6	qemu-kvm 软件架构分析	184
5.6.1	QEMU 的三种运行模式	184
5.6.2	libvirt 和 virt-manager	185
5.6.3	KVM 的调试接口	186
5.7	本章小结	187
第 6 章 qemu-kvm 原理与分析		188
6.1	QEMU 软件架构	189
6.1.1	qemu-kvm 的配置与编译	189
6.1.2	qemu-kvm 的架构与配置	190
6.2	QEMU 组件	190
6.2.1	模块模型	190
6.2.2	libkvm 模块	193
6.2.3	virtio 组件	196
6.3	基于 KVM 的 QEMU PC Emulator	199
6.3.1	KVM 中的 Machine 模块	199
6.3.2	基于 KVM 加速支持的 CPU 虚拟化模块	207
6.3.3	虚拟机的内存管理	216
6.3.4	I/O 管理	223

6.4 本章小结..... 225

第四篇 虚拟化软件开放接口

第 7 章 Xen API 接口技术及实现..... 228

- 7.1 Xen Management API 接口概述..... 229
- 7.2 XML-RPC 工作原理..... 230
 - 7.2.1 XML-RPC 概述..... 231
 - 7.2.2 XML-RPC 请求..... 232
 - 7.2.3 XML-RPC 响应..... 234
- 7.3 Xen Managemnet API 的实现..... 236
 - 7.3.1 C 语言和 Python 语言的扩展与嵌入..... 236
 - 7.3.2 Xen Management API 类的定义..... 237
 - 7.3.3 Xen Management API 处理流程分析..... 238
- 7.4 本章小结..... 242

第 8 章 libvirt 虚拟化控制中间件..... 243

- 8.1 libvirt 概述..... 244
 - 8.1.1 libvirt 简介及使用样例..... 244
 - 8.1.2 基于 libvirt 所开发的开源应用..... 245
 - 8.1.3 安装与配置..... 245
- 8.2 libvirt 架构与开发..... 247
 - 8.2.1 libvirt 架构说明..... 247
 - 8.2.2 libvirt API 控制接口..... 250
 - 8.2.3 libvirt 的主机域管理..... 254
 - 8.2.4 libvirt 的网络架构..... 254
 - 8.2.5 libvirt 的存储管理..... 256
- 8.3 基于 libvirt 的 XML 配置解析..... 256
 - 8.3.1 XML 配置格式简析..... 256
 - 8.3.2 针对 Xen 的 libvirt 配置详解..... 264
 - 8.3.3 针对 KVM/QEMU 的 libvirt 配置详解..... 271
- 8.4 本章小结..... 282

参考文献..... 283

图目录

图 1-1 x86 的发展历程 4

图 1-2 线性地址空间构造 9

图 1-3 分段机制流程分布 11

图 1-4 段选择符的结构 12

图 1-5 段寄存器的构造 13

图 1-6 段描述符的结构 14

图 1-7 通过段选择符索引段描述符表 16

图 1-8 分页机制流程分布 17

图 1-9 未启用 PAE 的 4KB 页——二级页表 18

图 1-10 启用 PAE 的 4KB 页——三级页表 20

图 1-11 四级页表结构 23

图 1-12 APIC 系统架构 28

图 1-13 中断门的格式 30

图 1-14 陷阱门的格式 30

图 1-15 DMA 传输示意图 33

图 1-16 系统虚拟化结构 38

图 1-17 系统虚拟化的发展历程 39

图 1-18 虚拟环境的组成 41

图 1-19 Hypervisor 模型的 VMM 46

图 1-20 宿主模型的 VMM 47

图 1-21 混合模型的 VMM 48

图 2-1 Intel VT 的组成 51

图 2-2 VT-x 的基本思想 54

图 2-3 物理平台的中断架构 57

图 2-4 虚拟机的中断架构 57

图 2-5 内存虚拟化示意图 59

图 2-6 影子页表的作用 60

图 2-7 客户机操作系统页表与影子页表 61

图 2-8 EPT 原理图	66
图 2-9 传统分页技术下的地址转换	67
图 2-10 嵌套分页技术下的地址转换	68
图 2-11 使用 VT-d 后访问内存架构	69
图 2-12 BDF 结构	70
图 2-13 根条目的结构	71
图 2-14 上下文条目的结构	72
图 2-15 根条目表和上下文条目表构成的两级结构	73
图 2-16 DMA 重映射的 4KB 页面地址转换过程	74
图 2-17 IOMMU 技术示意图	75
图 2-18 虚拟化网卡的基本原理	77
图 2-19 SR-IOV 原理图	78
图 2-20 操作系统的时钟概念	80
图 2-21 客户机时间概念 1	
图 2-22 客户机时间概念 2 (客户机时间与实际时间统一)	81
图 2-23 客户机被调度出去情况下时间概念的实现	83
图 2-24 客户机被调度出去情况下中断注入的微观示意图	84
图 2-25 Xen 架构图	85
图 2-26 微软 Hyper-V 架构图	88
图 2-27 KVM 架构图	89
图 3-1 Xen 软件体系结构图	93
图 3-2 Xm、Xend、Xen Hypervisor 调用关系图	94
图 3-3 Xen 系统启动流程图	97
图 3-4 x86 体系、RISC 体系指令示意图	102
图 3-5 Xen 系统特权级分布图 (IA32 保护模式)	103
图 3-6 基于 VT-x 技术的虚拟机生命周期	104
图 3-7 直接模式下客户机操作系统对页表项的读/写操作	106
图 3-8 虚拟机页表与影子页表的对应关系	108
图 3-9 虚拟 I/O 驱动分离模型	109
图 3-10 Xen 网络 I/O 全虚拟化模型	110

图 3-11 Xen 虚拟机状态转移图	111
图 4-1 Hypercall 的作用	114
图 4-2 分离设备驱动处理 I/O 请求流程图	150
图 5-1 KVM 和 Xen 虚拟化方案比较	157
图 5-2 qemu-kvm 虚拟机的运行时进程信息	158
图 5-3 KVM 的用户空间访问接口	162
图 5-4 KVM 模块初始化阶段	173
图 5-5 KVM 的初始化流程	174
图 5-6 QEMU 的三种模块架构	184
图 5-7 基于 libvirt 和 virt-manager 的虚拟机管理工具	186
图 5-8 KVM 调试信息图	186
图 6-1 qemu-kvm 的默认参数	189
图 6-2 QEMU 系统中模块的相互关系	192
图 6-3 设备层半虚拟化与全虚拟化架构图	196
图 6-4 virtio 驱动程序架构	197
图 6-5 KVM 的 CPU 执行架构	207
图 6-6 KVM 的内存映射原理	216
图 6-7 Guest OS 内存访问流程示意图	217
图 6-8 EPT 内存管理模式	217
图 6-9 qemu-kvm 的 I/O 处理流程	224
图 7-1 Xen Management API 分层架构图	229
图 7-2 XML-RPC 远程调用图	231
图 7-3 Xen Management API 函数调用流程图	241
图 8-1 没有使用 libvirt 的虚拟机管理方式	247
图 8-2 使用 libvirt 的虚拟机管理方式	248
图 8-3 使用 libvirt 的远程虚拟机管理方式	249