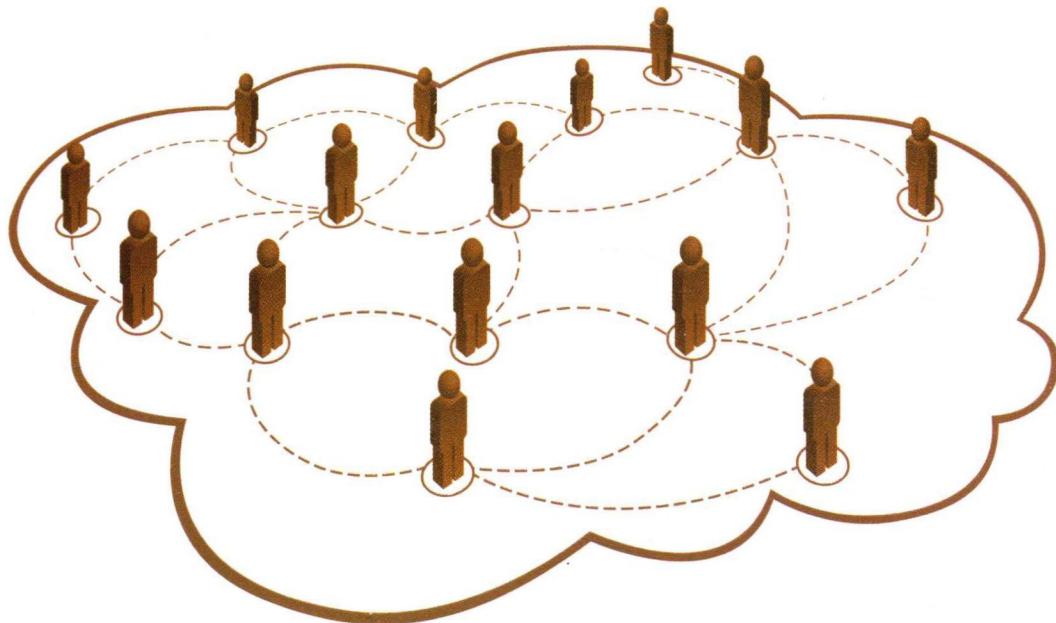


TURING

图灵程序
设计丛书

图解OpenFlow

[日] 晃通 宫永直树 岩田淳 著
李战军 薛文玲 译



189张图表轻松入门

从工作原理到应用实例
一本书掌握OpenFlow协议



中国工信出版集团



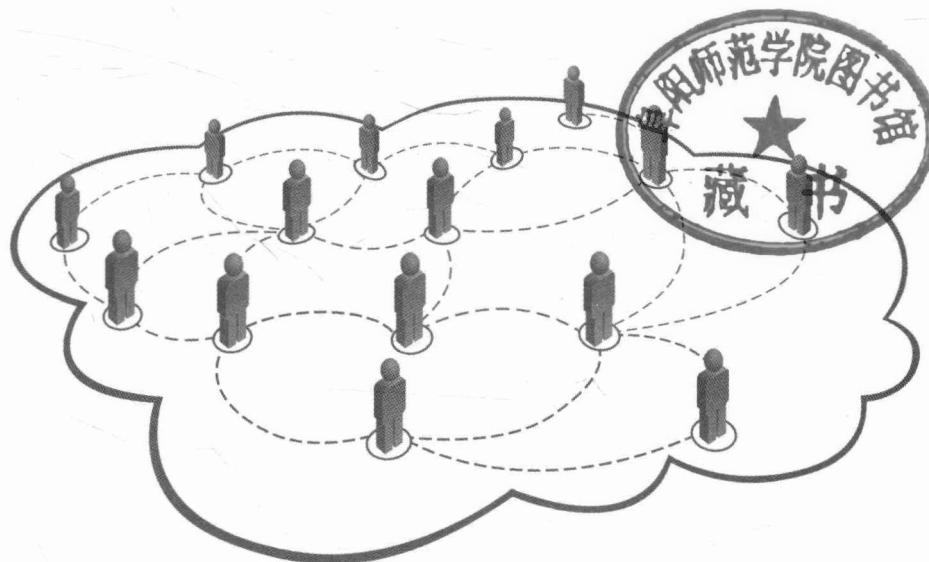
人民邮电出版社
POSTS & TELECOM PRESS

TURING

图灵程序
设计丛书

图解 OpenFlow

[日] 晃通 宫永直树 岩田淳 著
李战军 薛文玲 译



人民邮电出版社
北京

图书在版编目(CIP)数据

图解 OpenFlow / (日) 晃通, (日) 宫永直树, (日) 岩田淳著; 李战军, 薛文玲译. -- 北京: 人民邮电出版社, 2016.1

(图灵程序设计丛书)

ISBN 978-7-115-41125-9

I . ①图… II . ①晃… ②宫… ③岩… ④李… ⑤薛… III. ①计算机网络—图解 IV. ①TP393-64

中国版本图书馆CIP数据核字(2015)第275480号

Original Japanese edition

Mastering TCP/IP OpenFlow-Hen By Akimichi, Naoki Miyanaga and Atsushi Iwata

Copyright © 2013 by Akimichi, Naoki Miyanaga and Atsushi Iwata Published by Ohmsha, Ltd.

This Simplified Chinese Language edition published by Post & Telecom Press Copyright © 2016

All rights reserved.

本书中文简体字版由 Ohmsha, Ltd. 授权人民邮电出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。

内 容 提 要

本书从OpenFlow的概要入手，以OpenFlow 1.0规范为核心，详细讲解了OpenFlow的基本机制和结构。为了加深理解，还介绍了实际的设置示例，包括OpenFlow中的LLDP和网络虚拟化等。另外，还讲解了从OpenFlow 1.0到OpenFlow 1.3.2之间版本的主要变化，以及OpenFlow的注意事项和未来的发展动向。本书适合所有网络开发和管理人员阅读。

◆ 著 [日] 晃通 宫永直树 岩田淳

译 李战军 薛文玲

责任编辑 乐 馨

执行编辑 高宇涵

责任印制 杨林杰

◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号

邮编 100164 电子邮件 315@ptpress.com.cn

网址 <http://www.ptpress.com.cn>

北京天宇星印刷厂印刷

◆ 开本: 787×1092 1/16

印张: 14.5

字数: 276千字 2016年1月第1版

印数: 1~4 000册 2016年1月北京第1次印刷

著作权合同登记号 图字: 01-2013-8604号

定价: 49.00元

读者服务热线: (010)51095186转600 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京崇工商广字第0021号

前言

OpenFlow 是一种可实现多种功能的技术，甚至可以认为该技术能在网络上实现任何功能。当然，OpenFlow 有时也受到规范或设备的限制，但只要稍下点功夫仍可以实现相当多的功能。但是，反过来也可以说，如果技术人员不进行非常细致的设置，则不能实现任何功能。OpenFlow 本身只不过是能够实现细致设置的“框架”。事实上，未进行任何设置的 OpenFlow 网络甚至连自学习桥接器（L2 交换机）这样普通的网络功能都无法实现。

您是否听过“OpenFlow 是可编程的”？OpenFlow 网络中的各种设置，并不是通过转发数据包的 OpenFlow 交换机来完成的，而是通过统一控制多台 OpenFlow 交换机的 OpenFlow 控制器来完成。以往的网络设备都是各自通过自律分散协作来构建网络的。与此相对，在 OpenFlow 中，为了能让 OpenFlow 控制器对网络内各设备的详细动作等进行综合管理，而对其进行各种设置，这种行为就好像对网络整体进行“编程”一样。而这也正是人们认为 OpenFlow “可编程”的原因。有人说，如果将 OpenFlow 比作一门编程语言，那么它就是汇编语言。但笔者认为，从某种意义上讲，这种说法并不正确。汇编语言对硬件的依赖程度很高，而 OpenFlow 并不是单一供应商建立的独有规范，它能进行各种设备之间的互通。从这种意义上而言，笔者更喜欢“OpenFlow 是无标准库的 C 语言”这一比喻。C 语言本身虽然是非常基础的编程语言，但如果完全不使用标准库及 POSIX 等标准 API，创建任何内容都很困难。未经任何设置的 OpenFlow 因存在标准而能够互通，也正反映了这一点。笔者推测，将来在实际环境中使用 OpenFlow 时，很多用户并不是直接使用 OpenFlow，而是像使用 C 语言的标准库那样，会在整合的、完善的环境中使用 OpenFlow。

■ 本书所讲解的 OpenFlow 版本

本书以 OpenFlow 1.0 为主，介绍了 OpenFlow 1.3.2 以下的版本。以 OpenFlow 1.0 为主的原因是，OpenFlow 1.0 清晰地反映了 OpenFlow 的初期设计思想。因此要想理解 OpenFlow，从 OpenFlow 1.0 入手会比较容易明白。

在笔者开始写作时，OpenFlow 交换机的硬件在实现时大多支持 OpenFlow 1.0，基本上不支持 OpenFlow 1.1 以上的版本。但是，随着本书的写作接近尾声，基于 OpenFlow 1.3 的产品也已上市，今后支持 OpenFlow 1.3 以上版本的硬件将越来越多。因此，介绍时以哪个版本为主，这个问题一直困扰着笔者。

考虑再三，笔者还是决定以 OpenFlow 1.0 为主来说明 OpenFlow，另外也介绍了本书写作阶段的最新版本 OpenFlow 1.3.2 以下版本的发展过程。虽然也考虑过从一开始就直接介绍 OpenFlow 1.3，但如此一来，就很难理解有些功能为什么要添加了。实际上，在 OpenFlow 1.0 之前，“无需设计新的硬件，只对现有硬件更新其软件”的思想就已经确立了，不过为了实现 OpenFlow 的旧版本中不能实现的功能，OpenFlow 1.3 仍旧进行了扩展。

另外，本书还根据需要对协议相应的变更进行了说明，同时用三章的篇幅来分别解说 OpenFlow 1.1、OpenFlow 1.2 和 OpenFlow 1.3。

■ 本书的结构

本书力求避免介绍特定的实现及产品，只是单纯地介绍作为协议的 OpenFlow 本身。另外，还会尽可能地介绍用户在使用协议的过程中容易陷入的误区。

第 1 章介绍了 OpenFlow 及 OpenFlow 相关的情况。第 2 章介绍了 OpenFlow 1.0 的机制及格式。

为加深大家对 OpenFlow 的理解，第 3 章至第 6 章介绍了 OpenFlow 的设置示例。

第 3 章介绍了如何在 OpenFlow 中应用能够实现网络拓扑检测功能的 LLDP；第 4 章介绍了使用 OpenFlow 实现各种 L2 网络功能的方法；第 5 章介绍了基于 OpenFlow 的网络虚拟化，以及基于 VMware 等的服务器虚拟化和 OpenFlow 之间的关系；第 6 章则通过其他各种用例对 OpenFlow 进行了介绍。通过这些设置示例，应该可以将 OpenFlow 的实际情况介绍给大家。

此后的章节介绍了 OpenFlow 1.0 以上的版本及今后的发展动向。

第 7 章、第 8 章、第 9 章分别介绍了 OpenFlow 1.1、1.2、1.3 的变化。如果将 OpenFlow 1.3 及其之前版本的所有变化全部汇总到一处，可能会因不同点太多而导致混乱。所以，笔者将各个版本分别单列一章，以便于大家理解各版本中的变更内容及各版本升级的设计思想。

第 10 章介绍了实际应用 OpenFlow 时的注意事项。第 11 章介绍了为实现标准化时的重要因素——互通性而采取的措施，还介绍了 OpenFlow 规范中没有直接包含的相关讨论等。

■ 致谢

对于本书的内容，大久保修一先生（SAKURA 互联网株式会社）、菊池之裕先生（博科通讯系统株式会社）、安田丰先生（京都产业大学）、@tsubo 等提供了很多宝贵意见，在此深表感谢。

日本电气株式会社的芦原浩司先生、高飞先生、铃木一哉先生、须尧一志先生、高宫安仁先生及 TechVan 株式会社的山梨晓先生回答了笔者各种各样的问题，并对本书的内容提供了宝贵的意见，在此深表感谢。

如果没有大家的协助，不可能完成本书的写作。大家从专业角度给出了各种建议及感想，谢谢大家！

如果对 OpenFlow 感兴趣的读者能从本书中获取自己所需的知识，笔者将深感荣幸。

2013 年 7 月

笔者

目录

第1章 OpenFlow 概要 1

1.1	OpenFlow 的发展历程	2
1.1.1	OpenFlow 的历史	2
1.1.2	ONF	2
1.2	有效运用现有硬件，实现高效设计	4
1.3	所谓 OpenFlow，具体是指什么	5
1.3.1	OpenFlow 的网络构成示例	5
1.3.2	控制面和数据面的分离	5
1.3.3	控制面的构建方法	6
1.3.4	数据面的构建方法	7
1.3.5	OpenFlow 控制器和 OpenFlow 通道	9
1.3.6	流表	10
1.4	控制器和交换机的基本动作	11
1.4.1	OpenFlow 交换机的初始动作	11
1.4.2	Proactive 模式设置	11
1.4.3	Reactive 模式设置	11
1.5	OpenFlow 的动作示例	14
1.5.1	动作示例 1	14
1.5.2	动作示例 2	15
1.5.3	动作示例 3	15
1.5.4	其他动作示例	16
1.6	不断变化的 OpenFlow 规范	17
1.6.1	OpenFlow 1.0 版本以后快速变化的规范	17
1.6.2	实现 OpenFlow 1.3 标准化后	17
1.6.3	本书中对 OpenFlow 1.1 以上版本的说明	18

第 2 章**OpenFlow 1.0 的机制**19

2.1	OpenFlow 1.0 中的流表和流表项	20
2.2	流表项	22
2.2.1	头字段	22
2.2.2	计数器	25
2.2.3	行动	26
2.3	行动	27
2.3.1	Forward 行动	27
2.3.2	Drop 行动	28
2.3.3	Enqueue 行动 (可选)	29
2.3.4	Modify-Field 行动 (可选)	29
2.4	控制器和交换机之间的消息	32
2.4.1	安全通道的建立	32
2.4.2	消息格式	32
2.4.3	安全通道的建立和初始设置	33
2.4.4	握手	35
2.4.5	Flow-Mod 消息	39
2.4.6	Packet-In 消息	45
2.4.7	Packet-Out 消息	46
2.4.8	Port-Status 消息	47
2.4.9	Flow-Removed 消息	48
2.4.10	Error 消息	49
2.4.11	Barrier 消息	50
2.4.12	Echo 消息	51
2.5	OpenFlow Switch Errata Version 1.0.1	53
2.6	新规范变化很大，需要注意	54

第 3 章**LLDP 和 OpenFlow**55

3.1	LLDP 和 OpenFlow	56
3.2	LLDP 的机制	58
3.2.1	在 LLDP 中使用的目标以太网地址	58
3.2.2	LLDP 中使用的 3 种组播以太网地址	59
3.2.3	LLDPDU 格式	60

3.3	OpenFlow 中有效使用 LLDP 的方法示例	62
3.3.1	事先准备.....	62
3.3.2	通过 Packet-Out 消息发送 LLDP 帧.....	63
3.3.3	通过 Packet-In 消息将 LLDP 帧发送至 OpenFlow 控制器.....	64
3.4	基于迪杰斯特拉算法的路径计算	66

第 4 章 通过实现 L2 交换机的功能来学习 OpenFlow 69

4.1	通过具体网络设备的实现理解 OpenFlow	70
4.2	中继器 HUB	71
4.2.1	该示例中的网络构成.....	71
4.2.2	通过 Proactive 模式设置实现.....	71
4.2.3	将所有数据包 Packet-In 至 OpenFlow 控制器的方法.....	73
4.3	自学习桥接器	77
4.3.1	该示例中的网络构成.....	77
4.3.2	使用 OpenFlow 1.0 挑战自学习桥接器.....	78
4.3.3	监控 ARP 并创建流表项.....	83
4.3.4	如果将 PCA 和 PC B 对调，结果会怎样	89
4.3.5	通过发送源和目标以太网地址的配对进行管理的方法	90
4.3.6	在 OpenFlow 1.1 以上版本中实现自学习桥接器的方法	90
4.4	Tagged VLAN	92
4.4.1	该示例中的网络构成.....	92
4.4.2	实现 Tagged VLAN 的设置 (OpenFlow 交换机 1).....	93
4.4.3	实现 Tagged VLAN 的设置 (OpenFlow 交换机 2).....	96
4.4.4	该示例中的注意事项	98

第 5 章 OpenFlow 与虚拟化 99

5.1	服务器虚拟化和网络虚拟化	100
5.1.1	服务器虚拟化	100
5.1.2	动态迁移	101
5.1.3	多租户	102
5.1.4	网络虚拟化	103
5.2	基于 OpenFlow 的网络虚拟化的实现方法示例	105
5.2.1	使用 VLAN ID 的方法	105

5.2.2 使用物理端口 / 逻辑端口的方法	105
5.2.3 OpenFlow 控制器的实现要点	106
5.2.4 其他方法	106
5.3 FlowVisor	107
5.3.1 FlowVisor 概要	107
5.3.2 FlowVisor 和 OpenFlow 控制器之间的 OpenFlow 通道	107
5.3.3 FlowVisor 的串联	108
5.3.4 FlowVisor 和虚拟网络	109
5.3.5 设置 FlowVisor 时的注意事项	110

第 6 章

通过用例考察 OpenFlow

111

6.1 使用以太网地址的用户管理	112
6.2 ECMP	114
6.2.1 该示例中的网络构成	114
6.2.2 通过发送源地址区分时	115
6.2.3 通过 TCP 端口号区分时	116
6.2.4 轮询方式	116
6.3 简易负载均衡	117
6.4 选择性端口映射	118
6.4.1 单纯的端口映射	118
6.4.2 仅映射特定的 TCP 端口	119
6.4.3 OpenFlow 1.1 的“组”和映射	120
6.4.4 从多个 OpenFlow 交换机持续进行选择性映射并转发至监控设备	120
6.5 重定向至安全产品	121
6.6 与虚拟路由近似的动作 (多层交换机)	122
6.6.1 该示例中的网络构成	122
6.6.2 同一子网内的数据包转发处理	123
6.6.3 经过路由器的数据包转发处理	123
6.6.4 作为虚拟路由器响应 ARP 请求	124
6.6.5 虚拟路由器使用 ARP 解决以太网地址	125
6.6.6 TTL 的处理	126

第7章

OpenFlow 1.1

127

7.1	OpenFlow 1.1 中的变更要点	128
7.2	匹配字段的变更	129
7.3	多流表规范的变更 (流水线处理)	130
7.3.1	流水线处理	130
7.3.2	元数据	134
7.3.3	OpenFlow 1.1 中的自学习桥接器的实现手法	135
7.4	指令	136
7.4.1	何谓指令	136
7.4.2	行动、行动集、行动列表、指令的区别	136
7.4.3	对行动的变更	137
7.5	组	139
7.5.1	组表	139
7.5.2	组表项	139
7.5.3	组类型	140
7.5.4	组的组	144
7.6	虚拟端口的扩展	145
7.7	TTL 字段操作	146
7.7.1	copy TTL inwards/copy TTL outwards	146
7.7.2	接收到包含非法 TTL 值的数据包时的处理	147
7.7.3	不能实施 TTL 的匹配	148
7.8	OpenFlow 1.1 中其他的变更	149
7.8.1	支持 MPLS 标签和 VLAN 标签的 Push/Pop	149
7.8.2	OpenFlow 混合交换机	149
7.8.3	支持 SCTP	149
7.8.4	支持 ECN	150
7.8.5	OpenFlow 交换机和控制器之间连接名称的变更	150
7.8.6	紧急事态流缓存的取消	151
7.8.7	Vendor 消息名称的变更	151

第8章

OpenFlow 1.2

153

8.1	OpenFlow 1.2 中的变更点	154
8.2	OpenFlow eXtensible Match (OXM)	155

8.2.1	OXM TLV 的基本结构	155
8.2.2	匹配字段解析规范的取消和 Pre-requisite	157
8.2.3	OXM 匹配字段	158
8.2.4	OXM 中的通配符	159
8.2.5	OXM TLV 示例	160
8.2.6	基于 OXM 的 Set-Field	162
8.2.7	取消 TCP、UDP、SCTP、ICMP 重载使用相同字段	162
8.3	支持基本的 IPv6	163
8.4	支持多台控制器（故障转移和负载均衡）	164
8.4.1	Role	164
8.4.2	Role 变更	165
8.4.3	OpenFlow 控制器之间的协作	165
8.5	OpenFlow 1.2 中的其他变化	167
8.5.1	将虚拟端口分离为逻辑端口和保留端口	167
8.5.2	Flow-Mod 的 MODIFY/MODIFY_STRICT 的规范变更	167
8.5.3	对实验性扩展的支持	167
8.5.4	变更历史记录的添加	168

第 9 章

OpenFlow 1.3

169

9.1	OpenFlow 1.3 中的变更要点	170
9.2	计量表（QoS 支持）	171
9.3	Table-miss 的默认动作改为 Drop	173
9.3.1	Table-miss 流表项	173
9.3.2	流表匹配流程的变更	173
9.4	OpenFlow 1.3 中的其他变更	175
9.4.1	OpenFlow 控制器和 OpenFlow 交换机之间的辅助连接	175
9.4.2	可以通过 UDP、DTLS 等与 OpenFlow 控制器进行通信	175
9.4.3	支持 IPv6 扩展头	176
9.4.4	OXM 匹配字段的添加	176
9.4.5	支持 PBB	176
9.4.6	多框架	177
9.4.7	从握手时的 Features 响应消息中删除端口号	178
9.4.8	流表项构成要素的变更	178
9.5	OpenFlow 1.3.1 和 1.3.2	179
9.5.1	OpenFlow 通道中版本协商的变更	179

9.5.2 建立与 OpenFlow 控制器之间的 OpenFlow 通道	179
---	-----

第 10 章 OpenFlow 的注意事项

181

10.1 Packet-In 消息的处理负载	182
10.1.1 控制面带宽较窄导致的故障	183
10.1.2 Packet-In 导致的消息延迟	183
10.2 匹配和流相关的注意事项	184
10.2.1 未发现 TCP 标志	184
10.2.2 并非数据包及帧的任意字段都可进行匹配	184
10.2.3 匹配字段的依赖关系	185
10.3 取决于实现的事项	187
10.3.1 流表项数量的上限	187
10.3.2 OpenFlow 控制器可同时控制的 OpenFlow 交换机数量的上限	187
10.3.3 通过 buffer_id 表示的数据包未必保存着	187
10.3.4 OpenFlow 通道断开时的重新连接计时器	188
10.4 从下流发送 Flow-Mod	189
10.5 Barrier 消息和错误	190
10.5.1 Flow-Mod 之后的 Packet-Out	190
10.5.2 向不同的 OpenFlow 交换机发送 Flow-Mod 消息和 Packet-Out 消息时	192
10.6 没有检测 Packet-Out 失败的方法	194
10.7 IP 碎片处理	195

第 11 章 OpenFlow 的未来

197

11.1 互通性验证	198
11.1.1 OF-Test	198
11.1.2 PlugFest	198
11.1.3 2012 年进行的第 1 次 PlugFest	198
11.1.4 第 2 次、第 3 次 PlugFest	199
11.1.5 今后的课题	199
11.2 Northbound API	201
11.3 OF-CONFIG	202

附录

203

附录1 各版本的行动一览	204
附录 1.1 OpenFlow 1.0 的行动	204
附录 1.2 OpenFlow 1.1 的行动	205
附录 1.3 OpenFlow 1.2 的行动	206
附录 1.4 OpenFlow 1.3 的行动	207
附录2 各版本的消息一览	209
附录 2.1 OpenFlow 1.0 的消息	209
附录 2.2 OpenFlow 1.1 的消息	210
附录 2.3 OpenFlow 1.2 的消息	211
附录 2.4 OpenFlow 1.3 的消息	212
附录3 OpenFlow 从 1.0 到 1.3.2 的变更之处	213
附录 3.1 行动集	213
附录 3.2 指令	213
附录 3.3 行动列表	213
附录 3.4 组表	214
附录 3.5 计量表	214
附录 3.6 行动	214
附录4 参考文献及 URL	215

第 1 章

OpenFlow 概要

本章我们将要介绍 OpenFlow 及与 OpenFlow 相关的内容。这里我们仅对 OpenFlow 规范进行总结性的说明，更加具体的信息请参阅第 2 章之后的章节。

1.1

OpenFlow 的发展历程

OpenFlow 协议由斯坦福大学提出，最初的出发点是为了更加轻松地构建用于研究的网络。下面简单介绍 OpenFlow 的发展历程。

1.1.1 OpenFlow 的历史

构建实验网络是在进行网络相关研究时需要完成的一项重要工作。以往构建网络时，普遍采用的方法是将一定台数的计算机进行物理连接。但随着 VMware 等虚拟软件的上市和普及，通过组合使用虚拟机创建虚拟平台，构建实验环境的方法逐渐成为了主流。像这样构建的虚拟实验环境中，比较著名的有提供全球平台、支持广域分布式环境的 PlanetLab[1]；也有运用校园网，不通过整个互联网，而只是在特定的范围内虚拟地创建出来的网络。

以美国为中心，从零开始重新创建网络的新一代网络技术研究已成为了近几年来的发展趋势。虽然技术人员已经提出了各种新一代网络技术的方案，但在进行这些新技术的实验时，如果有某种架构能够实现比以往互联网通信设备更加精确的控制，则非常方便。

在此背景下，为了更轻松地构建各种实验网络，斯坦福大学开发了 OpenFlow。实际上，在由美国 NSF (National Science Foundation，美国国家科学基金会) 提供支持的大规模虚拟实验网络环境 GENI (Global Environment for Network Innovations，全球网络创新环境) [2] 项目中就使用了 OpenFlow[3]。

OpenFlow 不仅用于网络的实验及研究，还设立了广受关注、由多家网络设备供应商参加的 OpenFlow 交换机论坛。2008 年，该论坛制定了 OpenFlow 的基本规范 OpenFlow Switching Specification 0.2.0。从 OpenFlow 交换机论坛设立时开始，本书 3 名作者中 2 人就职的 NEC 便与斯坦福大学、Deutsche Telekom、Hewlett-Packard、Nicira 等一起，共同参与了相关研究▼。

▼ Nicira 公司于 2012 年被 VMware 公司收购。

1.1.2 ONF

在本书写作时，制定 OpenFlow 规范的并非是 OpenFlow 交换机论坛，而是 Open Networking Foundation (ONF，开放网络基金会)。

随着 OpenFlow 的不断升温，希望参加 OpenFlow 交换机论坛的组织越来越多，希望参与 OpenFlow 标准化工作的呼声也越来越高。因此，以继承 OpenFlow 交换机论坛活动的形式，Open Networking Foundation 成立了。

在 OpenFlow 规范中，OpenFlow 1.1 之前的规范由 OpenFlow 交换机论坛

制定，OpenFlow 1.2 以后的规范由 ONF 制定（截至本书写作时的最新版本为 OpenFlow 1.3.2^①）。但是，OpenFlow 1.0 和 1.1 并非与 ONF 毫无关系。ONF 曾对 OpenFlow 1.0 和 1.1 进行审查，对版本中模糊不清的部分进行了修订，然后才作为正式版本加以认定。因此，OpenFlow 1.0 和 1.1 也符合 ONF 标准。

^① 截至本书出版时，最新版本为 OpenFlow 1.4。——编者注