



Linux开源存储 全栈详解

从Ceph到容器存储

英特尔亚太研发有限公司 编著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>



Linux开源存储 全栈详解



从Ceph到容器存储

英特尔亚太研发有限公司 编著

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

内 容 简 介

本书致力于帮助读者形成有关 Linux 开源存储世界的细致拓扑。本书从存储硬件、Linux 存储栈、存储加速、存储安全、存储管理、分布式存储等各个角度与层次展开讨论，同时对处于主导地位的、较为流行的开源存储项目进行阐述，包括 SPDK、ISA-L、OpenSDS、Ceph、OpenStack、容器等。本书内容基本不涉及具体源码，主要围绕各个项目的起源与发展、实现原理与框架、要解决的网络问题等展开讨论，以帮助读者对 Linux 开源存储技术的实现与发展形成清晰的认识。本书语言通俗易懂，能够带领读者快速走进 Linux 开源存储的世界并做出自己的贡献。

本书适合希望参与 Linux 开源存储项目开发的读者阅读，也适合互联网应用的开发者、架构师和创业者参考，尤其可作为互联网架构师的开源技术典籍。

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

Linux 开源存储全栈详解：从 Ceph 到容器存储 / 英特尔亚太研发有限公司编著.

北京：电子工业出版社，2019.9

ISBN 978-7-121-36979-7

I . ①L... II . ①英... III . ①Linux 操作系统 IV . ①TP316.85

中国版本图书馆 CIP 数据核字（2019）第 127853 号

责任编辑：孙学瑛

特约编辑：田学清

印 刷：北京季蜂印刷有限公司

装 订：北京季蜂印刷有限公司

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：787×980 1/16 印张：24.75 字数：476 千字

版 次：2019 年 9 月第 1 版

印 次：2019 年 9 月第 1 次印刷

定 价：99.00 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 zlts@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。

作者简介

任桥伟：从事Linux内核、OpenStack、Ceph等开源项目的开发，著有《Linux内核修炼之道》《Linux那些事儿》系列。

李晓燕：活跃于Cinder和Ceph项目，具有多年存储领域经验。

程盈心：Ceph社区的活跃贡献者，专注于分布式系统的分析与优化。

马建朋：Ceph社区的活跃贡献者。

刘春梅：哈尔滨工业大学自动控制专业博士，目前在美国硅谷英特尔工作。从事过网络安全、虚拟化、终端安全、云计算等领域的工作。

尚德浩：Ceph社区的活跃贡献者。

胡伟：从事云计算和边缘计算相关工作。在OpenStack、Ceph和Edge Computing领域支持客户技术方案落地，参与业界多项前沿云计算相关技术验证和评估工作。

杨子夜：从事存储软件开发和优化工作。在虚拟化、存储、云安全等领域拥有诸多专利提交，其中21个专利已经被专利局授予（其中14个在美国，7个在中国）。

曹刚：从事存储软件的开发和优化工作，现为英特尔开发经理。

刘长鹏：从事存储软件和虚拟化研发工作。

刘孝冬：从事存储软件研发及存储相关算法优化的工作。

惠春阳：从事存储软件研发及存储相关算法优化的工作。

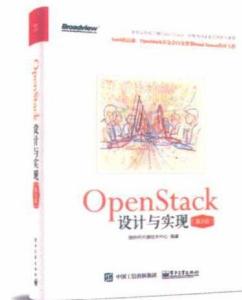
万群：从事测试领域的研究及实践近十年，对测试方法及项目管理有相当丰富的经验。

闫亮：从事存储软件的测试和优化工作。

周雁波：在英特尔实习期间，从事存储软件的开发和优化工作。

徐要昀：在英特尔实习期间，担任存储技术文档工程师，负责SPDK技术文档方面的工作。

好书力荐



拒绝堆砌臃肿，支持纯正原创

投稿邮箱：sxy@hei.com.cn

加入本书读者 QQ 群 465398813, 以书会友, 资源共享!

推荐序一

With an ever-increasing growth of data in both business and consumer uses, the need for storage of the data also continues to grow at an unprecedented pace. Intel technologies such as Intel® Optane™ SSD, and Intel® Optane™ persistent memory serve as examples of disruptive innovation, providing dramatically improved performance and entire new functionality to the storage industry.

Infrastructure software, storage, computing and network functions will serve a broad range of application platform requirements for 5G, Artificial Intelligence, Autonomous Driving, and Big Data Analytics. Open source storage software is critically important in developing and delivering innovative solutions. This book introduces open source storage technologies and techniques for software-defined storage solutions covering both hardware and software. Intel is a leading contributor of open source software and storage hardware IP and products, including contributions for open source storage software development, such as Intelligent Storage Acceleration Library (ISA-L), Storage Performance Development Kit (SPDK), and the Ceph project.

In China, open source and software-defined storage have been well adopted and applied across many different types of application scenarios. A strong ecosystem of storage software developers and solutions providers have been created and are delivering value to

the market through their innovations. We hope the technologies and techniques covered in this book will be useful to the open source communities in China and the rest of the world.

Sandra Rivera

Senior Vice President and General Manager

Network Platforms Group

Intel 5G Executive Sponsor

致谢

感谢所有为本书贡献过智慧和汗水的作者、译者、校对者、编辑者，以及审稿人。他们的辛勤工作使得本书能够顺利出版。特别感谢我的家人，他们在我写作过程中提供了无尽的支持和鼓励。同时，也要感谢我的同事和朋友，他们的帮助和支持让我在写作过程中充满了动力。最后，感谢所有读者，你们的反馈和建议将是我不断进步的动力。

希望本书能为读者提供有价值的参考，帮助大家更好地理解和应用开源存储技术。同时也希望本书能激发读者对开源技术的兴趣，鼓励大家积极参与开源社区，共同推动开源技术的发展。

推荐序二

The evolution of storage has been dramatic. As computing has become an integral part of our daily activities, storage media has had to evolve to keep pace, moving from tape and mechanical hard disk, to solid state hard disk, NAND, and now emerging 3D XPoint technologies.

Demand is only expected to increase. The data boom created by connected devices coupled with growth areas such as Artificial Intelligence and Machine Learning bring higher expectations — greater capacity, scalability, and faster access, not to mention increased availability and reliability. These challenges have resulted in the rise of distributed storage systems and cloud-based services, to help providers meet those demands. Software plays a key role in this transformation. The emergence of software-defined storage enables more scalable and flexible solutions, as well as greater levels of management and orchestration. Offerings such as virtual machine and container persistence storage help meet a growing range of business needs.

Open source provides the foundation for modern storage systems and a platform for rapid innovation. Intel has been active in the open source software community for 20+ years, and our contributions to storage technologies such as Intel® Storage Acceleration Library (ISA-L), Storage Performance Development Kit (SPDK), and the Ceph project

help accelerate performance for storage virtualization platforms.

As a leader in developing storage hardware and software, Intel recognizes the importance of the Chinese ecosystem to the growth of the storage market. It is with great pride we present this book as a resource to the China storage community.

Imad Sousou

Corporate Vice President, Intel Corporation

General Manager, Intel System Software Products

Imad is a world-renowned storage industry thought leader, and has been involved in the storage industry for over 20 years. He currently serves as Corporate Vice President of Intel's System Software Products Group, and General Manager of Intel's Storage Division. He leads the development of Intel's storage products, including the award-winning Intel Pro 1000 SSD, Intel's first solid-state drive, and Intel's first enterprise-class solid-state drive.

Imad is also a member of Intel's Executive Leadership Team, and is responsible for the company's strategy and execution of its storage business. He has extensive experience in the storage industry, having worked at companies such as Seagate, Western Digital, and Hitachi Data Systems. Imad holds a Bachelor's degree in Electrical Engineering from the University of California, Berkeley, and a Master's degree in Computer Science from the University of California, Berkeley. He is a frequent speaker at industry conferences and events, and is a sought-after expert in the field of storage technology.

Imad is a highly regarded industry leader, and is known for his deep technical expertise and his ability to translate complex storage concepts into practical, real-world solutions. He is a true pioneer in the field of storage technology, and is a valuable resource for anyone looking to learn more about the latest developments in the industry.

前言

自 1991 年 Linux 诞生，时间已经走过了近三十年。即将而立之年的 Linux 早已没有了初生时的稚气，它正在各个领域展示自己成熟的魅力。

以 Linux 为基础，各种开源生态，如网络、存储都出现了。而生态离不开形形色色的开源项目，在人人谈开源的今天，一个又一个知名的开源项目正在全球快速生长。当然，本书的主题仅限于 Linux 开源存储生态，面对其中一个又一个扑面而来且快速更迭的新项目、新名词，我们会有一种紧迫感，想去了解它们背后的故事，也会有一定的动力想要踏上 Linux 开源存储世界的旅程。而无论是否强迫，面对这样的一段旅程，我们心底浮现的最为愉悦的开场白或许应该是：“说实话，我学习的热情从来都没有低落过。Just for Fun。”正如 Linus 在自己的自传 *Just for Fun* 中所希望的那样。

面对 Linux 开源存储这么一个庞大而又杂乱的世界，让人最为惴惴不安的问题或许便是：我该如何更快、更好地适应这个全新的世界？人工智能与机器学习领域里研究的一个很重要的问题是，“为什么我们小时候有人牵一匹马告诉我们那是马，于是之后我们看到其他的马就知道那是马了？”针对这个问题的一个结论是：我们在头脑里形成了一种生物关系的拓扑结构，我们所认知的各种生物都会放进这个拓扑结构里，而生物不断成长的过程就是形成并完善各种各样的或树形、或环形等拓扑结构的过程，并以此来认知我们所面对的各种新事物。

由此可见，或许我们认知 Linux 开源存储世界最快也最为自然的方式就是努力在脑海里形成它的拓扑结构，并不断细化，比如这个生态包括什么样的层次，每个层次里又有什么样的项目去实现，各个项目又实现了哪些服务及功能，这些功能又是以什么样的方式实

现的，等等。对于感兴趣的项目，我们可以更为细致地去勾勒其中的脉络，就好似我们头脑里形成的有关一个城市的地图，它有哪些区，区里又有哪些标志性建筑及街道，对于熟悉的地方，我们甚至可以将它的周围放大并细化到一个微不足道的角落。

本书的组织形式

本书正是为帮助读者形成有关 Linux 开源存储世界的细致拓扑而组织的。

第 1 章主要对 Linux 开源存储的生态进行整体描述，包括开源存储领域研究的热点方向、相关的开源基金会等。

第 2 章从存储硬件的角度介绍了存储技术的发展历史，包括存储介质的进化、存储协议的更新等。

第 3 章作为整个 Linux 开源存储世界的基础，描述了 Linux 存储栈（Linux Storage Stack），对 I/O 在 Linux 内核里的处理流程及所涉及的主要模块进行介绍。

第 4~9 章的内容分别从存储加速、存储安全、存储管理与软件定义存储、分布式存储与 Ceph、OpenStack 存储、容器存储等角度与层次对处于主导地位的、较为流行的项目进行介绍，以帮助读者对相应项目形成比较细致的拓扑。

第 4 章讲解了存储领域的加速技术，包括 FPGA、QAT、NVDIMM 等硬件加速技术，以及 ISA-L、SPDK 等开源的软件加速方案。

第 5 章从可用性、可靠性、数据完整性、访问控制、加密与解密等方面讨论了存储安全问题。

第 6 章介绍存储管理与软件定义存储方面的主要开源项目，包括 OpenSDS、Libvirt 等。

第 7 章讨论分布式存储并详细介绍了目前流行的开源分布式存储项目 Ceph 的设计与实现。

第 8 章与第 9 章分别对 OpenStack 与 Kubernetes 两种主要云平台中的存储支持进行讨论。

感谢

作为英特尔的开源技术中心，参与各个 Linux 开源存储项目的开发与推广是再自然不过的事情了。除了为各个开源项目的完善与稳定贡献更多的思考和代码，我们还希望通过

这本书让更多的人更快地融入 Linux 开源存储世界的大家庭。

如果没有 Sandra Rivera（英特尔高级副总裁兼网络平台事业部总经理）、Imad Sousou（英特尔公司副总裁兼系统软件产品部总经理）、Mark Skarpness（英特尔系统软件产品部副总裁兼数据中心系统软件总经理）、Timmy Labatte（英特尔网络平台事业部副总裁兼软件工程总经理）、练丽萍（英特尔系统软件产品部网络与存储研发总监）、冯晓焰（英特尔系统软件产品部安卓系统工程研发总监）、周林（英特尔网络平台事业部中国区软件开发总监）、梁冰（英特尔系统软件产品部市场总监）、王庆（英特尔系统软件产品部网络与存储研发经理）的支持，这本书不可能完成，谨在此感谢他们在本书编写过程中给予的关怀与帮助。

感谢本书编辑孙学瑛老师，从选题到最后的定稿，在整个过程中，都给予我们无私的帮助和指导。

感谢参与各章内容编写的各位同事，他们是李晓燕、程盈心、马建朋、尚德浩、胡伟、刘春梅、任桥伟、杨子夜、曹刚、刘长鹏、刘孝冬、惠春阳、万群、闫亮、周雁波、徐雯昀。为了本书的顺利完成，他们付出了很多努力。

感谢所有对 Linux 开源存储技术抱有兴趣或从事各个 Linux 开源存储项目工作的人，没有你们提供的源码与大量技术资料，本书便会成为无源之水。

目 录

第1章 Linux 开源存储	1
1.1 Linux 和开源存储	1
1.1.1 为什么需要开源存储	3
1.1.2 Linux 开源存储技术原理和解决方案	6
1.2 Linux 开源存储系统方案介绍	8
1.2.1 Linux 单节点存储方案	8
1.2.2 存储服务的分类	11
1.2.3 数据压缩	13
1.2.4 重复数据删除	16
1.2.5 开源云计算数据存储平台	27
1.2.6 存储管理和软件定义存储	29
1.2.7 开源分布式存储和大数据解决方案	33
1.2.8 开源文档管理系统	37
1.2.9 网络功能虚拟化存储	39
1.2.10 虚拟机/容器存储	40
1.2.11 数据保护	43
1.3 三大顶级基金会	44
第2章 存储硬件与协议	47
2.1 存储设备的历史轨迹	47
2.2 存储介质的进化	53
2.2.1 3D NAND	53

2.2.2 3D XPoint	55
2.2.3 Intel Optane	58
2.3 存储接口协议的演变	59
2.4 网络存储技术	62
第3章 Linux 存储栈	67
3.1 Linux 存储系统概述	67
3.2 系统调用	69
3.3 文件系统	72
3.3.1 文件系统概述	73
3.3.2 Btrfs	75
3.4 Page Cache	80
3.5 Direct I/O	82
3.6 块层 (Block Layer)	83
3.6.1 bio 与 request	84
3.6.2 I/O 调度	86
3.6.3 I/O 合并	88
3.7 LVM	90
3.8 bcache	93
3.9 DRBD	96
第4章 存储加速	99
4.1 基于 CPU 处理器的加速和优化方案	100
4.2 基于协处理器或其他硬件的加速方案	103
4.2.1 FPGA 加速	103
4.2.2 智能网卡加速	105
4.2.3 Intel QAT	107
4.2.4 NVDIMM 为存储加速	110
4.3 智能存储加速库 (ISA-L)	111
4.3.1 数据保护: 纠删码与磁盘阵列	112
4.3.2 数据安全: 哈希	113
4.3.3 数据完整性: 循环冗余校验码	115
4.3.4 数据压缩: IGZIP	116
4.3.5 数据加密	117
4.4 存储性能软件加速库 (SPDK)	117
4.4.1 SPDK NVMe 驱动	119
4.4.2 SPDK 应用框架	133

4.4.3 SPDK 用户态块设备层	136
4.4.4 SPDK vhost target	150
4.4.5 SPDK iSCSI Target	156
4.4.6 SPDK NVMe-oF Target	163
4.4.7 SPDK RPC	165
4.4.8 SPDK 生态工具介绍	172
第 5 章 存储安全	181
5.1 可用性	181
5.1.1 SLA	181
5.1.2 MTTR、MTTF 和 MTBF	182
5.1.3 高可用方案	183
5.2 可靠性	185
5.2.1 磁盘阵列	186
5.2.2 纠删码	187
5.3 数据完整性	188
5.4 访问控制	189
5.5 加密与解密	191
第 6 章 存储管理与软件定义存储	194
6.1 OpenSDS	194
6.1.1 OpenSDS 社区	195
6.1.2 OpenSDS 架构	195
6.1.3 OpenSDS 应用场景	198
6.1.4 与 Kubernetes 集成	200
6.1.5 与 OpenStack 集成	200
6.2 Libvirt 存储管理	201
6.2.1 Libvirt 介绍	201
6.2.2 Libvirt 存储池和存储卷	205
第 7 章 分布式存储与 Ceph	206
7.1 Ceph 体系结构	209
7.1.1 对象存储	211
7.1.2 RADOS	212
7.1.3 OSD	212
7.1.4 数据寻址	214
7.1.5 存储池	219
7.1.6 Monitor	220

7.1.7	数据操作流程	227
7.1.8	Cache Tiering.....	228
7.1.9	块存储	230
7.1.10	Ceph FS	232
7.2	后端存储 ObjectStore	235
7.2.1	FileStore	236
7.2.2	BlueStore	240
7.2.3	SeaStore	243
7.3	CRUSH 算法	244
7.3.1	CRUSH 算法的基本特性	244
7.3.2	CRUSH 算法中的设备位置及状态	246
7.3.3	CRUSH 中的规则与算法细节	249
7.3.4	CRUSH 算法实践	254
7.3.5	CRUSH 算法在 Ceph 中的应用	261
7.4	Ceph 可靠性	262
7.4.1	OSD 多副本	263
7.4.2	OSD 纠删码	264
7.4.3	RBD mirror	265
7.4.4	RBD Snapshot	267
7.4.5	Ceph 数据恢复	271
7.4.6	Ceph 一致性	274
7.4.7	Ceph Scrub 机制	278
7.5	Ceph 中的缓存	279
7.5.1	RBDCache 具体实现	285
7.5.2	固态硬盘用作缓存	287
7.6	Ceph 加密和压缩	289
7.6.1	加密	289
7.6.2	压缩	291
7.6.3	加密和压缩的加速	294
7.7	QoS	294
7.7.1	前端 QoS	294
7.7.2	后端 QoS	295
7.7.3	dmClock 客户端	297
7.8	Ceph 性能测试与分析	298
7.8.1	集群性能测试	299
7.8.2	集群性能数据	304

7.8.3 综合测试分析工具	307
7.8.4 高级话题	311
7.9 Ceph 与 OpenStack	315
第 8 章 OpenStack 存储	318
8.1 Swift	321
8.1.1 Swift 体系结构	321
8.1.2 环	327
8.1.3 Swift API	330
8.1.4 认证	331
8.1.5 对象管理与操作	333
8.1.6 数据一致性	337
8.2 Cinder	338
8.2.1 Cinder 体系结构	338
8.2.2 Cinder API	341
8.2.3 cinder-scheduler	342
8.2.4 cinder-volume	343
8.2.5 cinder-backup	347
第 9 章 容器存储	348
9.1 容器	348
9.1.1 容器技术框架	350
9.1.2 Docker	353
9.1.3 容器与镜像	355
9.2 Docker 存储	356
9.2.1 临时存储	357
9.2.2 持久化存储	366
9.3 Kubernetes 存储	369
9.3.1 Kubernetes 核心概念	370
9.3.2 Kubernetes 数据卷管理	376
9.3.3 Kubernetes CSI	380