



DPDK
DATA PLANE DEVELOPMENT KIT

NFV的基石

深入浅出DPDK

朱河清 梁存铭 胡雪焜 曹水 等编著



- 挑战网络性能极限，资深开发者伴你在DPDK的世界中成长。
- 展现软件优化的极致魅力，百万兆高性能网络就在身边。
- 剖析数据面软件中的核心算法和关键技术，全面展示DPDK与NFV的深度融合。



机械工业出版社
China Machine Press

深入浅出DPDK

朱河清 梁存铭 胡雪焜 曹水 等编著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

深入浅出 DPDK / 朱河清等编著 . —北京：机械工业出版社，2016.5

ISBN 978-7-111-53783-0

I. 深… II. 朱… III. 应用软件 – 软件包 IV. TP317

中国版本图书馆 CIP 数据核字 (2016) 第 097359 号

近年来，随着半导体和多核计算机体系结构技术的不断创新和市场的发展，越来越多的网络设备基础架构开始向基于通用处理器平台的架构方向融合，期望用更低的成本和更短的产品开发周期来提供多样的网络单元和丰富的功能，如应用处理、控制处理、包处理、信号处理等。为了适应这一新的产业趋势，英特尔公司十年磨一剑，联合第三方软件开发公司及时推出了基于 Intel® x86 的架构 DPDK (Data Plane Development Kit，数据平面开发套件)，实现了高效灵活的包处理解决方案。经过近 3 年的开源与飞速发展，DPDK 已经发展成业界公认的高性能网卡和多通用处理器平台的开源软件工具包，并已成为通用处理器平台上影响力最大的数据平面解决方案。主流的 Linux 发行版都已经将 DPDK 纳入，DPDK 引发了基于 Linux 的高速网络技术的创新热潮，除了在传统的通信网络、安全设施领域应用之外，还被广泛应用于云计算、虚拟交换、存储网络甚至数据库、金融交易系统。

本书汇聚了最资深的 DPDK 技术专家的精辟见解和实战体验，详细介绍了 DPDK 技术的发展趋势、数据包处理、硬件加速技术、虚拟化以及 DPDK 技术在 SDN、NFV、网络存储等领域的实际应用。书中还使用大量的篇幅讲解各种核心软件算法、数据优化思想，并包括大量详尽的实战心得和使用指南。

作为国内第一本全面阐述网络数据面的核心技术的书籍，本书主要面向 IT、网络通信行业的从业人员，以及大专院校的师生，用通俗易懂的文字打开了一扇通向新一代网络处理架构的大门。DPDK 完全依赖软件，对 Linux 的报文处理能力做了重大革新，它的发展历程是一个不可多得的理论联系实际的教科书般的实例。

深入浅出 DPDK

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：姚 蕾

责任校对：殷 虹

印 刷：北京诚信伟业印刷有限公司

版 次：2016 年 5 月第 1 版第 1 次印刷

开 本：186mm×240mm 1/16

印 张：18

书 号：ISBN 978-7-111-53783-0

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

Prologue 序言

2015年的春天，在北京参加DPDK研讨大会时，有幸结识了本书的部分作者和众多DPDK研发的专业人士。这使我对这个专题的感召力深感诧异。DPDK就像一块磁铁，可以把这么多不同行业的专业人士吸引在一起。同时，大家会上也相约来年的春天，国内的同仁们能在DPDK技术进步中展现出自己独到的贡献。

作为运营商研发队伍的一员，我们无时不刻都能感受到NFV这个话题的灼热度。作为网络演进的大趋势，NFV将在未来为运营商实现网络重构扮演重要的角色。然而，大家都知道，NFV技术的发展之路存在各种屏障，性能问题是一道迈不过去的坎。这个问题的复杂性在于，它涉及I/O、操作系统内核、协议栈和虚拟化等多个层面对网络报文的优化处理技术。虽然IT界已发展出多类小众技术来应对，但这些技术对于普通应用技术人员而言比较陌生，即使对于传统网络的开发者而言，全面掌握这些技术也存在巨大的挑战。长久以来，用户更希望在这个领域有系统性的解决方案，能把相关的技术融会贯通，并系统性地组织在一起，同时也需要更为深入的细节技术支持工作。

DPDK的到来正逢其时，它之所以能脱颖而出，并迅速发展为业界在NFV加速领域的一种标杆技术，在于它不仅是上述技术的集大成者，更重要的是它的开放性和持续迭代能力，这些都得益于DPDK背后这支强大的专业研发团队，而本书的专业功力也可见一斑。

作为运营商的网络研发队伍，我们已关注DPDK近3年，尽管学习过DPDK部分源码和大量社区文档，也组织通过大量的DPDK相关NFV测试验证，但我们仍然觉得迫切需要系统性地介绍现代服务器体系架构，以及虚拟化环境下I/O优化的最新技术。令人倍感欣慰的是，本书作者对DPDK的讲解游刃有余，系统全面的同时又不乏敏锐的产业视角。可以说，深入浅出是本书最大的特点。

形而上者谓之道，形而下者谓之器。书中一方面透彻地讲解了现代处理器体系架构、网络I/O、内核优化和I/O虚拟化的原理与技术发展史，在这个“道用”的基础上，另一方面也

清晰地介绍了 DPDK 细节性的“器用”知识，包括并行处理、队列调度、I/O 调优、VNF 加速等大量方法与应用，两方面相得益彰。结合 DPDK 社区的开源代码和动手实践，相信读者仔细学习完本书，必能加快对 NFV 性能关键技术的领悟。本书的受益对象首先是那些立志跨界转型的 NFV 研发工程师，也面向高等院校计算机专业希望在体系架构方面有更深发展的在校生，更包括像我们这样关注 DPDK 应用场景、NFVI 集成和测试技术的最终用户。我们衷心感谢作者为业界带来的全新技术指引。这本书就像一粒种子，其中蕴含的知识未来定会在 NFV 这片沃土上枝繁叶茂，开花结果。

严格地讲，我们的团队只是 DPDK 用户的用户，我们研究 DPDK 的目的并非针对 DPDK 本身，而是为 NFV 的集成和开发提供一个准确的、可供评估的 NFVI 性能基准，减少各类网络功能组件在私有的优化过程中存在的不稳定风险。从对 DPDK 的初步评测来看，结果令人满意甚至超出预期，但我们仍应清醒地认识到，DPDK 作为 NFV 加速技术架构仍有很长的路要走，打造成熟、规范和完善的产业链是近期要解决的重要课题。我们呼吁也乐见有更多的朋友加入 DPDK 应用推广的行列，众志成城，汇聚成一股 SDN 时代的创新洪流。

欧亮博士 中国电信广州研究院

动机

2015 年 4 月，第一届 DPDK 中国峰会在北京成功召开。来自中国移动、中国电信、阿里巴巴、IBM、Intel、华为以及中兴的专家朋友登台演讲，一起分享了以 DPDK 为中心的技术主题。表 1 列出了 2015 DPDK 中国峰会的主题及演讲者。

表 1 2015 DPDK 中国峰会主题及演讲者

| 主 题 | 演 讲 者 | 公 司 |
|---------------------------|-------|-------|
| 利用 DPDK 加速 NFV | 邓辉 | 中国移动 |
| 利用 DPDK 优化云基础设施 | 孙成浩 | 阿里巴巴 |
| 构建 core 以及高能效应用的最佳实践 | 梁存铭 | Intel |
| 基于英特尔 ONP 构建虚拟化的 IP 接入方案 | 欧亮 | 中国电信 |
| DPDK 加速无线数据核心网络 | 陈东华 | 中兴 |
| 电信业务场景下的数据面挑战 | 刘郡 | 华为 |
| 运行于 Power 架构下的 DPDK 和数据转发 | 祝超 | IBM |

这次会议吸引了来自各行业、科研单位与高校的 200 多名开发人员、专家和企业代表参会。会上问答交流非常热烈，会后我们就想，也许是时间写一本介绍 DPDK、探讨 NFV 数据面的技术书籍。现在，很多公司在招聘网络和系统软件人才时，甚至会将 DPDK 作为一项技能罗列在招聘要求中。DPDK 从一个最初的小众技术，经过 10 年的孕育，慢慢走来，直至今日已经逐渐被越来越多的通信、云基础架构厂商接受。同时，互联网上也出现不少介绍 DPDK 基础理论的文章和博客，从不同的角度对 DPDK 技术进行剖析和应用，其中很多观点非常新颖。作为 DPDK 的中国开发团队人员，我们意识到如果能够提供一本 DPDK 的书籍，进行一些系统性的梳理，将核心的原理进行深入分析，可以更好地加速 DPDK 技术的普及，触发更多的软件创新，促进行业的新技术发展。于是，就萌发了写这本书的初衷。当然，我

们心里既有创作的激动骄傲，也有些犹豫忐忑，写书不是一件简单的事情，但经过讨论和考量，我们逐渐变得坚定，这是一本集结团队智慧的尝试。我们希望能够把 DPDK 的技术深入浅出地解释清楚，让更多的从业人员和高校师生了解并使用 DPDK，促进 DPDK 发展日新月异，兴起百家争鸣的局面，这是我们最大的愿景。

多核

2005 年的夏天，刚加入 Intel 的我们畅想着 CPU 多核时代的到来给软件业带来的挑战与机会。如果要充分利用多核处理器，需要软件针对并行化做大量改进，传统软件的并行化程度不高，在多核以前，软件依靠 CPU 频率提升自动获得更高性能。并行化改进不是一件简单的工作，许多软件需要重新设计，基本很难在短期实现，整个计算机行业都对此纠结了很久。2005 年以前，整个 CPU 的发展历史，是不断提升芯片运算频率核心的做法，软件性能会随着处理器的频率升高，即使软件不做改动，性能也会跟着上一个台阶。但这样的逻辑进入多核时代已无法实现。首先我们来看看表 2 所示的 Intel® 多核处理器演进。

表 2 Intel® 多核处理器演进的历史图表

| Xeon 处理器代码 | 制造工艺 | 最大核心数量 | 发布时间 | 超线程 | 双路服务器可使用核心数量 |
|----------------|------|--------|------------|-----|--------------|
| WoodCrest | 65nm | 2 | 2006 年 6 月 | 否 | 4 |
| Nehalem-EP | 45nm | 4 | 2009 年 7 月 | 是 | 16 |
| Westmere-EP | 32nm | 6 | 2010 年 2 月 | 是 | 24 |
| SandyBridge-EP | 32nm | 8 | 2012 年 3 月 | 是 | 32 |
| IvyBridge-EP | 22nm | 12 | 2013 年 9 月 | 是 | 48 |
| Haswell-EP | 22nm | 18 | 2014 年 9 月 | 是 | 72 |

在过去 10 年里，服务器平台的处理器核心数目扩展了很多。表 2 参考了英特尔至强系列的处理器的核心技术演进历史，这个系列的处理器主要面向双通道（双路）服务器和相应的硬件平台。与此同时，基于 MIPS、Power、ARM 架构的处理器也经历着类似或者更加激进的并行化计算的路线图。在处理器飞速发展的同时，服务器平台在硬件技术上提供了支撑。基于 PCI Express 的高速 IO 设备、内存访问与带宽的上升相辅相成。此外，价格和经济性优势越发突出，今天一台双路服务器的价格可能和 10 年前一台高端笔记本电脑的价格类似，但计算能力达到甚至超越了当年的超级计算机。强大的硬件平台为软件优化技术创新蕴蓄了温床。

以太网接口技术也经历了飞速发展。从早期主流的 10Mbit/s 与 100Mbit/s，发展到千兆网（1Gbit/s）。到如今，万兆（10Gbit/s）网卡技术成为数据中心服务器的主流接口技术，近年来，Intel 等公司还推出了 40Gbit/s、100Gbit/s 的超高速网络接口技术。而 CPU 的运行频率基本停

留在 10 年前的水平，为了迎接超高速网络技术的挑战，软件也需要大幅度创新。

结合硬件技术的发展，DPDK (Data Plane Development Kit)，一个以软件优化为主的数据面技术应时而生，它为今天 NFV 技术的发展提供了绝佳的平台可行性。

IXP

提到硬件平台和数据面技术，网络处理器是无法绕过的话题。电信行业通常使用网络处理器或类似芯片技术作为数据面开发平台首选。Intel 此前也曾专注此领域，2002 年收购了 DEC 下属的研究部门，在美国马萨诸塞州哈德逊开发了这一系列芯片，诞生了行业闻名的 Intel Exchange Architecture Network Processor (IXP4xx、IXP12xx、IXP24xx、IXP28xx) 产品线，曾取得行业市场占有率第一的成绩。即使今日，相信很多通信业的朋友，还对这些处理器芯片有些熟悉或者非常了解。IXP 内部拥有大量的微引擎 (MicroEngine)，同时结合了 XSCALE 作为控制面处理器，众所周知，XSCALE 是以 ARM 芯片为核心技术的一种扩展。

2006 年，AMD 向 Intel 发起了一场大战，时至今日结局已然明了，Intel 依赖麾下的以色列团队，打出了新一代 Core 架构，迅速在能效比上完成超车。公司高层同时确立了 Tick-Tock 的研发节奏，每隔两年推出新一代体系结构，每隔两年推出基于新一代制造工艺的芯片。这一战略基本保证了每年都会推出新产品。当时 AMD 的处理器技术一度具有领先地位，并触发了 Intel 在内部研发架构城门失火的状况下不得不进行重组，就在那时 Intel 的网络处理器业务被进行重估，由于 IXP 芯片系列的市场容量不够大，Intel 的架构师也开始预测，通用处理器多核路线有取代 IXP 专用处理芯片的潜力。自此，IXP 的研发体系开始调整，逐步转向使用 Intel CPU 多核的硬件平台，客观上讲，这一转型为 DPDK 的产生创造了机会。时至今日，Intel 还保留并发展了基于硬件加速的 QuickAssist 技术，这和当日的 IXP 息息相关。由此看来，DPDK 算是生于乱世。

DPDK 的历史

网络处理器能够迅速将数据报文收入系统，比如将 64 字节的报文以 10Gbit/s 的线速也就是 14.88Mp/s (百万报文每秒) 收入系统，并且交由 CPU 处理，这在早期 Linux 和服务器平台上无法实现。以 Venky Venkataren、Walter Gilmore、Mike Lynch 为核心的 Intel 团队开始了可行性研究，并希望借助软件技术来实现，很快他们取得了一定的技术突破，设计了运行在 Linux 用户态的网卡程序架构。传统上，网卡驱动程序运行在 Linux 的内核态，以中断方式来唤醒系统处理，这和历史形成有关。早期 CPU 运行速度远高于外设访问，所以中断处理方式十分有效，但随着芯片技术与高速网络接口技术的一日千里式发展，报文吞吐需要高达

10Gbit/s 的端口处理能力，市面上已经出现大量的 25Gbit/s、40Gbit/s 甚至 100Gbit/s 高速端口，主流处理器的主频仍停留在 3GHz 以下。高端游戏玩家可以将 CPU 超频到 5GHz，但网络和通信节点的设计基于能效比经济性的考量，网络设备需要日以继夜地运行，运行成本（包含耗电量）在总成本中需要重点考量，系统选型时大多选取 2.5GHz 以下的芯片，保证合适的性价比。I/O 超越 CPU 的运行速率，是横在行业面前的技术挑战。用轮询来处理高速端口开始成为必然，这构成了 DPDK 运行的基础。

在理论框架和核心技术取得一定突破后，Intel 与 6wind 进行了合作，交由在法国的软件公司进行部分软件开发和测试，6wind 向 Intel 交付了早期的 DPDK 软件开发包。2011 年开始，6wind、Windriver、Tieto、Radisys 先后宣布了对 Intel DPDK 的商业服务支持。Intel 起初只是将 DPDK 以源代码方式分享给少量客户，作为评估 IA 平台和硬件性能的软件服务模块，随着时间推移与行业的大幅度接受，2013 年 Intel 将 DPDK 这一软件以 BSD 开源方式分享在 Intel 的网站上，供开发者免费下载。2013 年 4 月，6wind 联合其他开发者成立 www.dpdk.org 的开源社区，DPDK 开始走上开源的大道。

开源

DPDK 在代码开源后，任何开发者都被允许通过 www.dpdk.org 提交代码。随着开发者社区进一步扩大，Intel 持续加大了在开源社区的投入，同时在 NFV 浪潮下，越来越多的公司和个人开发者加入这一社区，比如 Brocade、Cisco、RedHat、VMware、IBM，他们不再只是 DPDK 的消费者，角色向生产者转变，开始提供代码，对 DPDK 的代码进行优化和整理。起初 DPDK 完全专注于 Intel 的服务器平台技术，专注于利用处理器与芯片组高级特性，支持 Intel 的网卡产品线系列。

DPDK 2.1 版本在 2015 年 8 月发布，几乎所有行业主流的网卡设备商都已经加入 DPDK 社区，提供源代码级别支持。另外，除了支持通用网卡之外，能否将 DPDK 应用在特别的加速芯片上是一个有趣的话题，有很多工作在进行中，Intel 最新提交了用于 Crypto 设备的接口设计，可以利用类似 Intel 的 QuickAssit 的硬件加速单元，实现一个针对数据包加解密与压缩处理的软件接口。

在多架构支持方面，DPDK 社区也取得了很大的进展，IBM 中国研究院的祝超博士启动了将 DPDK 移植到 Power 体系架构的工作，Freescale 的中国开发者也参与修改，Tilera 与 Ezchip 的工程师也花了不少精力将 DPDK 运行在 Tile 架构下。很快，DPDK 从单一的基于 Intel 平台的软件，逐步演变成一个相对完整的生态系统，覆盖了多个处理器、以太网和硬件加速技术。

在 Linux 社区融合方面，DPDK 也开始和一些主流的 Linux 社区合作，并得到了越来越多的响应。作为 Linux 社区最主要的贡献者之一的 RedHat 尝试在 Fedora Linux 集成 DPDK；接着 RedHat Enterprise Linux 在安装库里也加入 DPDK 支持，用户可以自动下载安装 DPDK 扩展库。RedHat 工程师还尝试将 DPDK 与 Container 集成测试，并公开发布了运行结果。传统虚拟化的领导者 VMware 的工程师也加入 DPDK 社区，负责 VMXNET3-PMD 模块的维护。Canonical 在 Ubuntu 15 中加入了 DPDK 的支持。

延伸

由于 DPDK 主体运行在用户态，这种设计理念给 Linux 或者 FreeBSD 这类操作系统带来很多创新思路，也在 Linux 社区引发一些讨论。

DPDK 的出现使人们开始思考，Linux 的用户态和内核态，谁更适合进行高速网络数据报文处理。从简单数据对比来看，在 Intel 的通用服务器上，使用单核处理小包收发，纯粹的报文收发，理想模型下能达到大约 57Mp/s (每秒百万包)。尽管在真实应用中，不会只收发报文不处理，但这样的性能相对 Linux 的普通网卡驱动来说已经是遥不可及的高性能。OpenVSwitch 是一个很好的例子，作为主流的虚拟交换开源软件，也尝试用 DPDK 来构建和加速虚拟交换技术，DPDK 的支持在 OVS2.4 中被发布，开辟了在内核态数据通道之外一条新的用户态数据通道。目前，经过 20 多年的发展，Linux 已经累积大量的开源软件，具备丰富的协议和应用支持，无所不能，而数据报文进出 Linux 系统，基本都是在 Linux 内核态来完成处理。因为 Linux 系统丰富强大的功能，相当多的生产系统（现有软件）运行在 Linux 内核态，这样的好处是大量软件可以重用，研发成本低。但也正因为内核功能强大丰富，其处理效率和性能就必然要做出一些牺牲。

使用

在专业的通信网络系统中，高速数据进出速率是衡量系统性能的关键指标之一。大多通信系统是基于 Linux 的定制系统，在保证实时性的嵌入式开发环境中开发出用户态下的程序完成系统功能。利用 DPDK 的高速报文吞吐优势，对接运行在 Linux 用户态的程序，对成本降低和硬件通用化有很大的好处，使得以软件为主体的网络设备成为可能。对 Intel® x86 通用处理器而言，这是一个巨大的市场机会。

对于通信设备厂商，通用平台和软件驱动的开发方式具有易采购、易升级、稳定性、节约成本的优点。

- 易采购：通用服务器作为主流的基础硬件，拥有丰富的采购渠道和供应商，供货量巨大。

- **易升级：**软件开发模式简单，工具丰富，最大程度上避免系统升级中对硬件的依赖和更新，实现低成本的及时升级。
- **稳定性：**通用服务器平台已经通过大量功能的验证，产品稳定性毋庸置疑。而且，对于专用的设计平台，系统稳定需要时间累积和大量测试，尤其是当采用新一代平台设计时可能需要硬件更新，这就会带来稳定性的风险。
- **节约研发成本和降低复杂性：**传统的网络设备因为功能复杂和高可靠性需求，系统切分为多个子系统，每个子系统需要单独设计和选型，独立开发，甚至选用单独的芯片。这样的系统需要构建复杂的开发团队、完善的系统规划、有效的项目管理和组织协调，来确保系统开发进度。而且，由于开发的范围大，各项目之间会产生路径依赖。而基于通用服务器搭建的网络设备可以很好地避免这些问题。

版权

DPDK 全称是 Data Plane Development Kit，从字面解释上看，这是专注于数据面软件开发的套件。本质上，它由一些底层的软件库组成。目前，DPDK 使用 BSD license，绝大多数软件代码都运行在用户态。少量代码运行在内核态，涉及 UIO、VFIO 以及 XenDom0，KNI 这类内核模块只能以 GPL 发布。BSD 给了 DPDK 的开发者和消费者很大的自由，大家可以自由地修改源代码，并且广泛应用于商业场景。这和 GPL 对商业应用的限制有很大区别。作为开发者，向 DPDK 社区提交贡献代码时，需要特别注意 license 的定义，开发者需要明确 license 并且申明来源的合法性。

社区

参与 DPDK 社区 [Ref1-1] (www.dpdk.org) 就需要理解它的运行机制。目前，DPDK 的发布节奏大体上每年发布 3 次软件版本 (announce@dpdk.org)，发布计划与具体时间会提前公布在社区里。DPDK 的开发特性也会在路标中公布，一般是通过电子邮件列表讨论 dev@dpdk.org，任何参与者都可以自由提交新特性或者错误修正，具体规则可以参见 www.dpdk.org/dev，本书不做详细解读。对于使用 DPDK 的技术问题，可以参与 user@dpdk.org 进入讨论。

子模块的维护者名单也会发布到开源社区，便于查阅。在提交代码时，源代码以 patch 方式发送给 dev@dpdk.org。通常情况下，代码维护者会对提交的代码进行仔细检查，包括代码规范、兼容性等，并提供反馈。这个过程全部通过电子邮件组的方式来完成，由于邮件量可能巨大，如果作者没有得到及时回复，请尝试主动联系，提醒代码维护人员关注，这是参

与社区非常有效的方式。开发者也可以在第一次提交代码时明确抄送相关的活跃成员和专家，以得到更加及时的反馈和关注。

目前，开源社区的大量工作由很多自愿开发者共同完成，因此需要耐心等待其他开发者来及时参与问答。通过提前研究社区运行方式，可以事半功倍。这对在校学生来说更是一个很好的锻炼机会。及早参与开源社区的软件开发，会在未来选择工作时使你具有更敏锐的产业视角和技术深度。作为长期浸润在通信行业的专业人士，我们强烈推荐那些对软件有强烈兴趣的同学积极参与开源社区的软件开发。

贡献

本书由目前 DPDK 社区中一些比较资深的开发者共同编写，很多作者是第一次将 DPDK 的核心思想付诸于书面文字。尽管大家已尽最大的努力，但由于水平和能力所限，难免存在一些瑕疵和不足，也由于时间和版面的要求，很多好的想法无法在本书中详细描述。但我们希望通过本书能够帮助读者理解 DPDK 的核心思想，引领大家进入一个丰富多彩的开源软件的世界。在本书编纂过程中我们得到很多朋友的帮助，必须感谢上海交通大学的金耀辉老师、中国科学技术大学的华蓓和张凯老师、清华大学的陈渝教授、中国电信的欧亮博士，他们给了我们很多中肯的建议；另外还要感谢李训和刘勇，他们提供了大量的资料和素材，帮助我们验证了大量的 DPDK 实例；还要感谢我们的同事杨涛、喻德、陈志辉、谢华伟、戴启华、常存银、刘长鹏、Jim St Leger、MJay、Patrick Lu，他们热心地帮助勘定稿件；最后还要特别感谢英特尔公司网络产品事业部的领导周晓梅和周林给整个写作团队提供的极大支持。正是这些热心朋友、同事和领导的支持，坚定了我们的信心，并帮助我们顺利完成此书。最后我们衷心希望本书读者能够有所收获。

作者介绍 (按姓名排序) *About the Authors*

曹水: 黑龙江省佳木斯人, 2001 年毕业于复旦大学计算机系, 硕士。现为英特尔软件经理, 从事嵌入式开发和软件行业超过 15 年, 现主要负责 DPDK 软件测试工作。

陈静: 湖北省沙市人, 2006 年毕业于华中科技大学, 硕士。现为英特尔软件开发工程师, 主要从事 DPDK 网卡驱动的开发和性能调优工作。

何少鹏: 江西省萍乡人, 毕业于上海交通大学, 硕士。现为英特尔 DPDK 软件工程师, 开发网络设备相关软件超过十年, 也有数年从事互联网应用和 SDN 硬件设计工作。

胡雪焜: 江西省南昌人, 毕业于中国科学技术大学计算机系, 硕士。现为英特尔网络通信平台部门应用工程师, 主要研究底层虚拟化技术和基于 IA 架构的数据面性能优化, 以及对网络演进的影响, 具有丰富的 SDN/NFV 商业实践经验。

梁存铭: 英特尔资深软件工程师, 在计算机网络领域具有丰富的实践开发经验, 提交过多项美国专利。作为 DPDK 早期贡献者之一, 在 PCIe 高性能加速、I/O 虚拟化、IA 指令优化、改善闲时效率、协议栈优化等方面有较深入的研究。

刘继江: 黑龙江省七台河人, 毕业于青岛海洋大学自动化系, 现主要从事 DPDK 网卡驱动程序和虚拟化研发, 和 overlay 网络的性能优化工作。

陆文卓: 安徽省淮南人, 2004 年毕业于南京大学计算机系, 硕士。现为英特尔中国研发中心软件工程师。在无线通信、有线网络方面均有超过十年的从业经验。

欧阳长春: 2006 年毕业于华中科技大学计算机系, 硕士。目前在阿里云任开发专家, 从事网络虚拟化开发及优化, 在数据报文加速、深度报文检测、网络虚拟化方面具有丰富开发经验。

仇大玉: 江苏省南京人, 2012 年毕业于东南大学, 硕士。现为英特尔亚太研发有限公司软件工程师, 主要从事 DPDK 软件开发和测试工作。

陶喆: 上海交通大学学士, 上海大学硕士。先后在思科和英特尔从事网络相关的设备、

协议栈和虚拟化的开发工作。曾获 CCIE R&S 认证。

万群：江西省南昌人，毕业于西安交通大学计算机系，硕士。现为英特尔上海研发中心研发工程师。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。

王志宏：四川省绵阳人，2011 年毕业于华东师范大学，硕士。现为英特尔亚太研发中心高级软件工程师，主要工作方向为 DPDK 虚拟化中的性能分析与优化。

吴菁菁：江苏省扬州人，2007 年毕业于西安交通大学电信系，硕士。现为英特尔软件工程师，主要从事 DPDK 软件开发工作。

许茜：浙江省杭州市人，毕业于浙江大学信电系，硕士，现为英特尔网络处理事业部软件测试人员，主要负责 DPDK 相关的虚拟化测试和性能测试。

杨子夜：2009 年毕业于复旦大学软件学院，硕士。现为英特尔高级软件工程师，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 5 个相关专利以及 20 项申请。

张合林：湖南省湘潭人，2004 年毕业于东华大学，工学硕士。现主要从事 DPDK 网卡驱动程序研发及性能优化工作。

张帆：湖南省长沙人，爱尔兰利莫里克大学计算机网络信息学博士。现为英特尔公司爱尔兰分部网络软件工程师，湖南省湘潭大学兼职教授。近年专著有《Comparative Performance and Energy Consumption Analysis of Different AES Implementations on a Wireless Sensor Network Node》等。发表 SCI/EI 检索国际期刊及会议论文 3 篇。目前主要从事英特尔 DPDK 在 SDN 应用方面的扩展研究工作。

朱河清：江苏省靖江人，毕业于电子科技大学数据通信与计算机网络专业，硕士，现为英特尔 DPDK 与 Hyperscan 软件经理，在英特尔、阿尔卡特、华为、朗讯有 15 年通信网络设备研发与开源软件开发经验。

Venky Venkatesan：毕业于印度孟买大学，现为英特尔网络产品集团高级主任工程师（Sr PE），DPDK 初始架构师，在美国 Oregon 负责报文处理与加速的系统架构与软件创新工作。

目 录 *Contents*

| | | |
|-----------------------------|-------------------------|----|
| 序 言 | 1.5.3 DPDK 加速存储节点 | 15 |
| 引 言 | 1.5.4 DPDK 的方法论 | 16 |
| 作者介绍 | 1.6 从融合的角度看 DPDK | 16 |
| 第一部分 DPDK 基础篇 | 1.7 实例 | 17 |
| 第 1 章 认识 DPDK | 1.7.1 HelloWorld | 17 |
| 1.1 主流包处理硬件平台 | 1.7.2 Skeleton | 19 |
| 1.1.1 硬件加速器 | 1.7.3 L3fwd | 22 |
| 1.1.2 网络处理器 | 1.8 小结 | 25 |
| 1.1.3 多核处理器 | | |
| 1.2 初识 DPDK | | |
| 1.2.1 IA 不适合进行数据包处理吗 | | |
| 1.2.2 DPDK 最佳实践 | | |
| 1.2.3 DPDK 框架简介 | | |
| 1.2.4 寻找性能优化的天花板 | | |
| 1.3 解读数据包处理能力 | | |
| 1.4 探索 IA 处理器上最艰巨的任务 | | |
| 1.5 软件包处理的潜力——再识 DPDK | | |
| 1.5.1 DPDK 加速网络节点 | | |
| 1.5.2 DPDK 加速计算节点 | | |
| 第 2 章 Cache 和内存 | | 26 |
| 2.1 存储系统简介 | | 26 |
| 2.1.1 系统架构的演进 | | 26 |
| 2.1.2 内存子系统 | | 28 |
| 2.2 Cache 系统简介 | | 29 |
| 2.2.1 Cache 的种类 | | 29 |
| 2.2.2 TLB Cache | | 30 |
| 2.3 Cache 地址映射和变换 | | 31 |
| 2.3.1 全关联型 Cache | | 32 |
| 2.3.2 直接关联型 Cache | | 32 |
| 2.3.3 组关联型 Cache | | 33 |
| 2.4 Cache 的写策略 | | 34 |
| 2.5 Cache 预取 | | 35 |

| | | | |
|------------------------------------|----|---------------------------------------|-----|
| 2.5.1 Cache 的预取原理 | 35 | 3.2.2 单指令多数据 | 68 |
| 2.5.2 NetBurst 架构处理器上的 预取 | 36 | 3.3 小结 | 70 |
| 2.5.3 两个执行效率迥异的程序 | 37 | 第 4 章 同步互斥机制 | 71 |
| 2.5.4 软件预取 | 38 | 4.1 原子操作 | 71 |
| 2.6 Cache 一致性 | 41 | 4.1.1 处理器上的原子操作 | 71 |
| 2.6.1 Cache Line 对齐 | 41 | 4.1.2 Linux 内核原子操作 | 72 |
| 2.6.2 Cache 一致性问题的由来 | 42 | 4.1.3 DPDK 原子操作实现和应用 | 74 |
| 2.6.3 一致性协议 | 43 | 4.2 读写锁 | 76 |
| 2.6.4 MESI 协议 | 44 | 4.2.1 Linux 读写锁主要 API | 77 |
| 2.6.5 DPDK 如何保证 Cache 一致性 | 45 | 4.2.2 DPDK 读写锁实现和应用 | 78 |
| 2.7 TLB 和大页 | 47 | 4.3 自旋锁 | 79 |
| 2.7.1 逻辑地址到物理地址的转换 | 47 | 4.3.1 自旋锁的缺点 | 79 |
| 2.7.2 TLB | 48 | 4.3.2 Linux 自旋锁 API | 79 |
| 2.7.3 使用大页 | 49 | 4.3.3 DPDK 自旋锁实现和应用 | 80 |
| 2.7.4 如何激活大页 | 49 | 4.4 无锁机制 | 81 |
| 2.8 DDIO | 50 | 4.4.1 Linux 内核无锁环形缓冲 | 81 |
| 2.8.1 时代背景 | 50 | 4.4.2 DPDK 无锁环形缓冲 | 82 |
| 2.8.2 网卡的读数据操作 | 51 | 4.5 小结 | 89 |
| 2.8.3 网卡的写数据操作 | 53 | 第 5 章 报文转发 | 90 |
| 2.9 NUMA 系统 | 54 | 5.1 网络处理模块划分 | 90 |
| 第 3 章 并行计算 | 57 | 5.2 转发框架介绍 | 91 |
| 3.1 多核性能和可扩展性 | 57 | 5.2.1 DPDK run to completion 模型 | 94 |
| 3.1.1 追求性能水平扩展 | 57 | 5.2.2 DPDK pipeline 模型 | 95 |
| 3.1.2 多核处理器 | 58 | 5.3 转发算法 | 97 |
| 3.1.3 亲和性 | 61 | 5.3.1 精确匹配算法 | 97 |
| 3.1.4 DPDK 的多线程 | 63 | 5.3.2 最长前缀匹配算法 | 100 |
| 3.2 指令并发与数据并行 | 66 | 5.3.3 ACL 算法 | 102 |
| 3.2.1 指令并发 | 67 | 5.3.4 报文分发 | 103 |
| | | 5.4 小结 | 104 |

| | |
|--------------------------------------|-----|
| 第6章 PCIe与包处理I/O | 105 |
| 6.1 从PCIe事务的角度看包处理 | 105 |
| 6.1.1 PCIe概览 | 105 |
| 6.1.2 PCIe事务传输 | 105 |
| 6.1.3 PCIe带宽 | 107 |
| 6.2 PCIe上的数据传输能力 | 108 |
| 6.3 网卡DMA描述符环形队列 | 109 |
| 6.4 数据包收发——CPU和I/O的协奏 | 111 |
| 6.4.1 全景分析 | 111 |
| 6.4.2 优化的考虑 | 113 |
| 6.5 PCIe的净荷转发带宽 | 113 |
| 6.6 Mbuf与Mempool | 114 |
| 6.6.1 Mbuf | 114 |
| 6.6.2 Mempool | 117 |
| 6.7 小结 | 117 |
| 第7章 网卡性能优化 | 118 |
| 7.1 DPDK的轮询模式 | 118 |
| 7.1.1 异步中断模式 | 118 |
| 7.1.2 轮询模式 | 119 |
| 7.1.3 混合中断轮询模式 | 120 |
| 7.2 网卡I/O性能优化 | 121 |
| 7.2.1 Burst收发包的优点 | 121 |
| 7.2.2 批处理和时延隐藏 | 124 |
| 7.2.3 利用Intel SIMD指令进一步并行化包收发 | 127 |
| 7.3 平台优化及其配置调优 | 128 |
| 7.3.1 硬件平台对包处理性能的影响 | 129 |
| 7.3.2 软件平台对包处理性能的影响 | 133 |
| 7.4 队列长度及各种阈值的设置 | 136 |
| 7.4.1 收包队列长度 | 136 |
| 7.4.2 发包队列长度 | 137 |
| 7.4.3 收包队列可释放描述符数量阈值(rx_free_thresh) | 137 |
| 7.4.4 发包队列发送结果报告阈值(tx_rs_thresh) | 137 |
| 7.4.5 发包描述符释放阈值(tx_free_thresh) | 138 |
| 7.5 小结 | 138 |
| 第8章 流分类与多队列 | 139 |
| 8.1 多队列 | 139 |
| 8.1.1 网卡多队列的由来 | 139 |
| 8.1.2 Linux内核对多队列的支持 | 140 |
| 8.1.3 DPDK与多队列 | 142 |
| 8.1.4 队列分配 | 144 |
| 8.2 流分类 | 144 |
| 8.2.1 包的类型 | 144 |
| 8.2.2 RSS | 145 |
| 8.2.3 Flow Director | 146 |
| 8.2.4 服务质量 | 148 |
| 8.2.5 虚拟化流分类方式 | 150 |
| 8.2.6 流过滤 | 150 |
| 8.3 流分类技术的使用 | 151 |
| 8.3.1 DPDK结合网卡Flow Director功能 | 152 |
| 8.3.2 DPDK结合网卡虚拟化及Cloud Filter功能 | 155 |