

机械磁盘、SSD、FC/SAS协议、HBA卡、存储阵列

文件系统

、分布式文件系统、集群存储

、虚拟化、同步/异步远程复制、Thin Provisioning

、数据迁移配置、Deduplication

N、NAS、iSCSI、FCoE

、磁盘阵列

像、删除、自动分级

CDP持续数据保护、Automati cal Storage Tiering

、存储系统I/O路径、云计算与云存储

VTL虚拟磁带库

、容灾、应用容灾、数据备份

、数据恢复

像、虚拟化、同步/异步远程复制、Thin Provisioning

CDP持续数据保护、Automati cal Storage Tiering

、存储系统I/O路径、云计算与云存储

VTL虚拟磁带库

、容灾、应用容灾、数据备份

、数据恢复

冬瓜头「张冬」著

存储II

大话存储系统架构与
底层原理极限剖析

博学以记对技术
造极

SCSI SAS ATA SATA Over Fibre Channel Over IP iSCSI FCIP FCoE
EVAN VNX DS4200 AXI IPX400 Solaris AS400 OSS390 NetApp IBM HDS HPUX
Cobalt Sied DAS NAS SAN

冬瓜头『张东』著

大话
存储系统
底层原理
极限剖析

清华大学出版社
北京

内 容 简 介

网络存储是一个涉及计算机硬件以及网络协议/技术、操作系统以及专业软件等各方面综合知识的领域。目前国内阐述网络存储的书籍少之又少，大部分是国外作品，对存储系统底层细节的描述不够深入，加之术语太多，初学者很难真正理解网络存储的精髓。

本书以特立独行的行文风格向读者阐述了整个网络存储系统。从硬盘到应用程序，对这条路径上的每个节点，作者都进行了阐述。书中内容涉及：计算机 IO 基本概念，硬盘物理结构、盘片数据结构和工作原理，七种常见 RAID 原理详析以及性能细节对比，虚拟磁盘、卷和文件系统原理，磁盘阵列系统，OSI 模型，FC 协议，众多磁盘阵列架构等。另外，本书囊括了存储领域几乎所有的新兴技术，比如机械磁盘、SSD、FC/SAS 协议、HBA 卡、存储控制器、集群存储系统、FC SAN、NAS、iSCSI、FCoE、快照、镜像、虚拟化、同步/异步远程复制、Thin Provision 自动精简配置、VTL 虚拟磁带库、数据容灾、应用容灾、业务容灾、性能优化、存储系统 IO 路径、云计算与云存储等。

其中每一项技术作者都进行了建模和分析，旨在帮助读者彻底理解每一种技术的原理和本质。本书结尾，作者精心总结和多年来在论坛以及各大媒体发表的帖子内容，超过一百条的问与答，这些内容都是与实际紧密结合的经验总结，颇具参考价值。

本书适合初入存储行业的研发人员、技术工程师、售前工程师和销售人员阅读，同时适合资深存储行业人士用以互相切磋交流提高。另外，网络工程师、网管、服务器软硬件开发与销售人员、Web 开发者、数据库开发者以及相关专业师生等也非常适合阅读本书。

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目(CIP)数据

大话存储 II——存储系统架构与底层原理极限剖析/冬瓜头 著.—北京：清华大学出版社，2011.5
ISBN 978-7-302-24989-4

I. ①大… II. ①冬… III. ①计算机网络—信息存储—研究 IV. ①TP393

中国版本图书馆 CIP 数据核字(2011)第 041652 号

责任编辑：栾大成

版式设计：北京东方人华科技有限公司

责任校对：徐俊伟

责任印制：何 芊

出版发行：清华大学出版社 地 址：北京清华大学学研大厦 A 座

http://www.tup.com.cn 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62795954,jsjjc@tup.tsinghua.edu.cn

质量反馈：010-62772015,zhiliang@tup.tsinghua.edu.cn

印刷者：北京鑫丰华彩印有限公司

装订者：三河市新茂装订有限公司

经 销：全国新华书店

开 本：185×260 **印 张：**57.5 **插 页：**1 **字 数：**1454 千字

版 次：2011 年 5 月第 1 版 **印 次：**2011 年 6 月第 2 次印刷

印 数：5001~10000

定 价：99.00 元

产品编号：040152-01

序 1

我关注张冬这个名字是在《大话存储》一书刚出版的时候。作为一个长期从事信息存储技术研究与教学的大学教师，自认为对于国内外关于网络存储方面的各种书籍和资料比较熟悉，对业界有哪些牛人也算比较了解。但我在书店偶然发现一本名为《大话存储》书的时候，确实感到有点意外和惊喜。好像在熟悉的武林圈子之外，突然出现一位武林高手在那里论道。好奇心驱使我赶紧买了一本书回家研读，结果发现这本书确实与众不同。

与我们这些所谓学院派写的中规中矩的书相比，此书风格特立独行，语言形象生动，潇洒洒，颇具武侠之风。书中充满着智慧的思考和有趣的比喻，将各种原本枯燥深奥的技术概念和原理论述得十分透彻明白。不仅如此，该书还收集了大量的实例，使读者在系统获得网络存储知识的同时，还能了解典型实际系统的工作原理和技术细节，具有很好的实用性。我读完之后，对这本书的作者十分好奇。一个 80 后而且还是学化学出身的年轻人，如何就能写出这种行文老到而风格独特的专业技术书籍呢？上网查了一下冬瓜头（张冬的网名）的技术博客和他在各种论坛留下的文字，我得到了答案。这是一个完全由兴趣驱动而对技术极端痴迷的人，也是一位善于思考、富于想象力的人。这种纯粹的、不含任何功利成份的兴趣与痴迷，才是促进科学技术发展的真正源动力。

真正和张冬接触，是因为他来信质疑我们实验室申报的一项专利。收到质疑的来信，我和提出这项专利的博士生经过仔细研究，发现我们提供的图上因为少了一个非门，结果将会因为反相而出错。对如此细致具体的问题，一般人是难以发现的。如果没有打破砂锅问到底的较真精神，哪里会发现如此细节的错误呢？这种质疑的精神，在科学的研究中是极为宝贵的。我们学校被称为“根叔”的李培根校长，在 2010 年的新生开学典礼大会上，就以“质疑”为题作了讲演，激励青年学子发扬质疑精神。有质疑精神的人，不唯上，不唯权威，只认真理，这正是我们这个时代所稀缺的精神。

强烈的兴趣，对技术的痴迷，加上质疑精神，成就了一本存储领域的一本好书。我在研究生新生入学之后，就推荐他们先读一下《大话存储》这本书。一方面此书对研究生而言，确实是一本网络存储技术入门的好书，另一方面我还有一个用意，就是让他们知道，要从事科学研究，强烈的兴趣比什么都重要。

信息存储是信息跨越时间的传递，也是人类传承知识的主要手段。在信息存储技术上，人类有超过万年的发明创造史。早期就地取材，人类利用石刻、泥板、竹简和羊皮来记录信息，后来发明了纸张和活字印刷来保存和传播信息，近代发明了照相、录音和录像技术来存储信息。利用这些发明和创造，人类留下了极为丰富的文字、绘画、图像、语音和视频信息。正是这些信息，记录了人类创造的知识体系，使我们能够传承文明，并在此基础上创造新的文明。

从计算机的发明为开端，人类的信息技术进入了一个以数字化为特征的历史性新阶段。各种形式的信息被转换成数字后，以统一的方式进行处理、传输和存储，然后再转换为各种形式的信息被人们所利用。这种前所未有的方式发明之后，一个以数字化为特征的信息革命浪潮就波澜壮阔地形成。各种信息都被大规模数字化，使数字化的信息呈爆炸性增长。特别是互联网的兴起和普及，大大加快了信息的流通过程，使数字信息加速产生。图灵奖获得者 Jim Gary 观察这种数据急速增长的趋势后，总结出一个规律：人类每 18 个月新增的数据量，

将是历史上所有数据量之和！如此下去，对信息存储的需求将是无止境的，信息存储技术在这种强烈的需求驱动下得到了空前的发展。

为了保存数字化的信息，当代的科学家和工程师在最近的几十年中发明了磁存储、光存储、半导体存储等多种存储技术，其中大容量的硬盘在海量信息存储中扮演了主要的角色。硬盘的密度在短短几十年中增长了一百万倍以上，在近期，硬盘密度每年增长都接近一倍，而且还有不小的增长空间。由硬盘作为基本单元，通过各种总线、网络将硬盘连接成不同层次和不同规模的存储系统，就构成了我们目前的网络存储系统。例如由硬盘组加上冗余纠错技术构成磁盘阵列，再由磁盘阵列通过局部高速网络连接形成存储区域网；又如通过包含硬盘的大规模集群和文件系统形成的海量存储系统成为大型网站和数据中心新的存储架构。人们发明了各种技术来提高存储系统的容量、性能、效率、可用性、安全性和可管理性。存储虚拟化、归档存储、集群存储、云存储、绿色存储等新名词不断涌现，SSD 固态存储、重复数据删除、连续数据保护、数据备份与容灾、数据生命周期管理等新技术层出不穷，令人应接不暇。

在这种情况下，广大的信息领域的从业人员，信息系统的用户，以及学习信息技术的大学生和研究生，迫切需要一本既全面论述网络存储技术原理、又有丰富实例，既反映最新技术进展、又通俗易懂的书来满足他们的需求。冬瓜的《大话存储》就是这样一本恰逢其时的好书。

《大话存储》已在业界产生了很大的影响，对存储技术在我国的普及起到了良好的推动作用。该书还被引进到我国的宝岛台湾，可见其影响深远。张冬再接再厉，以他对技术的痴迷继续钻研，对第一本书作了工作量巨大的改动与增补，并增加了云存储等全新的三章内容，全面反映了他对技术的重新思考和对最新技术的深刻理解。我相信，这些新的内容将给读者带来惊喜。

在技术发展十分迅速的领域，赶时髦的书籍多如牛毛，书店里充满了应景之作，真正经过深入思考、用心写作的书是不多的。而《大话存储 II》却是一位技术高手的呕心沥血之作，书中对每一项技术的介绍都经过深入的思考和反复的推敲，这在当前浮躁的气氛中显得弥足珍贵。在《大话存储 II》即将出版之际，我要向作者表示深深的敬意和衷心的祝贺，并郑重向读者推荐这本学习网络存储技术的好书。

谢长生

华中科技大学计算机学院 教授
信息存储系统教育部重点实验室 主任

序 2

在过去的近二十年中，存储领域在国内外都发生了巨大变化。存储系统已经从早期的服务器附庸形态中脱离出来，而作为独立的产业走向了应用的前台。作为存储领域的一个从业人员，我有幸经历了其中许多变化阶段。而更为幸运的是存储领域的发展势头不是日趋渐微，而是方兴未艾。

存储领域的发展形成了众多的产业环节，也正是这些产业环节支撑着存储领域的持续发展和变化。其中，技术驱动和应用拉动的重要性往往得到人们的青睐，而存储技术及其相关应用技术面向社会的教育和普及则常常受到人们的忽视。而对于专业人员（包括专业应用人员）极度匮乏的今天，这个环节的重要性是不言而喻的。

最初看到《大话存储》，我只是觉得作者想法有新意，没有十分在意。当时国内在存储方面已经有了一些从国外引入的书籍，而且 IT 方面的经典书籍一向出自国外，不知道这本书是否能突出重围成为凤毛麟角的技术经典。与之同时，国内的网络存储研究和应用虽然热火朝天，但占有一席之地的自主存储产品却是少之又少。我所隶属的依托于中科院计算所的蓝鲸存储团队虽然有了较好的技术基础，但市场如何突破也依然处于艰难的摸索阶段。

时至今日，《大话存储》已经在存储技术及其相关应用技术面向社会的教育和普及方面起到了重要的作用，形成了自己的品牌。而我们的蓝鲸存储产品不仅在国内实现了规模化应用，打破了国外产品的垄断，而且已经在国外高端应用环境形成部署。或许是由于同属国产品牌，在与作者的交流中，作者对存储技术和产品的执着态度使我们深受感动，轻松但不乏深度和严谨的表述模式让我们受益匪浅。

转眼之间，《大话存储 II》已经面市。文风深入浅出、举重若轻，内容洋洋洒洒、上下求索。在此之际，我将我们团队实现突围的一些感想与作者和读者一起分享：

- 胸怀理想才能有毅力面对突破过程中内部和外部的挑战；
- 找准定位并持续专注为突破提供了可能性；
- 切中要害、关注细节并保证品质是突破的基础；
- 找到对用户有价值的创新而不盲从是突破的关键。

中科院计算技术研究所 研究员 许鲁

序 3

第一次听说冬瓜头是因为《大话存储》，而见到冬瓜头本人时，已经听他在写《大话存储》第二版了。我们一起有过不长时间的沟通，全是讨论最新的存储技术。有一些内容，我相信他在书中都有写到，在同他见面前，我一直在嘀咕怎么去跟他沟通。但是，看到现实中的他朴实、憨厚、腼腆，但是对技术极其敏感，说起来一套一套的，我就开始被这个山东大汉折服了，还好我是做了准备的，否则非被问倒不可。

他是个靠笔说话和表达的人，在网络论坛里，写的文字常常是洋洋洒洒，时而言辞激烈，时而意气风发，这分明是一位才高八斗的江南才子，又像是一个书场中幽默诙谐的说书人。从此之后，我经常关注他的个人博客。他会经常在博客释放一些思想出来，豪放不羁，甚至还写过长诗以及各种各样的打油诗，有些还写得非常棒，配上那个流着鼻涕的冬瓜头漫画形象，真是绝配了。

这么一个内秀的北方大汉，用豪放的气势描述一个个生涩、枯燥的技术领域，他的思想遍布他的文字，通过笔端流放在读者面前，并且还一直这么坚持。听他说，大话 2 出来后，还会继续写大话 3，毕竟技术日新月异，尤其是在 IT 领域。正如他所言，昨天的中国同仁在存储技术还是个初学者，今天已经开始从蹒跚学步到自主创新了，而明天，有什么理由不能期盼他们引领潮流的身影呢。我想，这也是冬瓜头要写大话 3 的动力所在吧。

《大话存储 II》的初稿篇幅已经超过了 1500 页。在浏览了全部章节之后，发现这 1500 页中真的是字字珠玑！看得出来，是冬瓜头一个字一个字写出来的。更加可贵的是，全书字里行间透着他那独特的思想，对技术、对世界的理解以及他做人的态度。能够将这些世界观的东西融入一本技术书籍，这在以前是绝无仅有的！比如书中多次提到“轮回”、“阴阳”等，最后还有一节是用中医的思想来“诊治”系统性能瓶颈，看后真是令我等感叹至极！世间万物都是相互联系的，都可以找到类比和轮回，这也是冬瓜头所描述的世界观的一种。

在和冬瓜头的交谈中获知，他大学学习的专业是化学，因为高中时他化学成绩最好，所以就报了化学专业，而且期间还自学过分子生物学领域的内容，我更加惊讶了！按照他的话来说，就是“兴趣是第一驱动力”。是的，好奇和探索正是人类不断发展的第一动力。说到这里我对《大话存储 II》中关于冬瓜头所设想的“机器如何认知自身”这段内容产生了强烈共鸣，人可以认识自身认识世界，那么机器为何不能呢？冬瓜头说他学习和接触存储业不过三四年的时间，在这短短的几年内，竟然能有如此造诣，这使我感觉到，强烈的好奇心可以创造奇迹！可以让机器开口，可以让机器进化！这也正是冬瓜头所表述的世界观的一种！

《大话存储 II》对各项存储技术的细节描述已经可以说是达到了研发级别，有很多部分甚至可以指导我们的研发！但是他却并没有用代码来表述，而是用通俗的语言和详实的图示，将原本通过阅读代码才可以理解透彻的原理，就这么轻而易举地表述了出来，这是目前我所看到的任何存储书籍或者文章都没有做到的。就这一点我曾经问过冬瓜头，问他如何做到的。他每次的回答很简单，一针见血，实实在在，他说：“因为我就是一个从不懂钻到懂的草根，我深知一个根本不懂存储的人最想了解的东西和切入角度，并且愿意毫无保留地帮助其他草根生长！”是啊，只有亲历过那悬梁刺股的学习之路的不易，才能产出精华！

在与冬瓜头的交谈中，他还常提到一句口号：“振兴民族科技。”从他说话的眼神和口

气看得出来，振兴民族科技已经成为他的信仰。他也说到，他现在所做的一切都围绕着这个信仰，他愿意为中国存储事业鞠躬尽瘁死而后已，出版《大话存储》只是他要做的第一个环节而已，今后他还会有一系列的动作来兑现他的诺言。信仰可以改变一个人的心态与行为，我们目前太缺乏信仰，我想如果我们所有人都有这种信仰，那么“振兴民族科技”这句口号早就实现了。

信息存储已经成为了一个时刻影响人们生产、生活的新兴产业，它的发展也代表着世界未来的发展，让我们再来看看国产存储信息产业的发展正在经历着怎样的变革和转变。

我国即将进入“十二五”时期，“十二五”期间我国将要实现三大转变目标：从国强到民富、从外需到内需、从高碳到低碳。这也意味着国家的发展需要依靠科技，需要大力發展新技术，尤其以信息化技术为主轴，信息化技术的发展带动重点工程的进行，势必对国产产品催生更大的需求，我相信存储业也会有更多的民族产业佼佼者诞生。

IT 环境日益复杂，数据量快速膨胀，存储业也进入了一个技术更新极为活跃的黄金发展时期，产业发展迅速，技术活跃度高，这对国内厂家来说，无疑是一个脱颖而出的良好契机。那么，我们如何在这个时代背景下产生代表着民族存储业的国产佼佼者？

在这个时代，我们应该遵守什么？我们应该坚持什么？商业道德、创新精神、客户意识，我想只有将这些融入到企业性格中才能为企业注入新的活力。作为一家有理想的企业，需要具备一定的时代精神，而在存储技术日新月异的今天，企业打造独有的技术张性，终究才会超越历史，才会产生新时代的民族企业。我相信现在越来越多的企业正朝这个方向发展。

《大话存储》以通俗易懂的语言、风趣的行文手法向读者阐述枯燥难懂的技术精髓，致力于存储信息技术发展的民族企业也同样可以在深刻理解本土文化精髓的前提下为中国写下辉煌的历史篇章。我想这个世界没有什么不可能的，只要有这份热情、专注和执著，又有什么是不可能实现的呢？

这样的一本特立独行的书，它就是时代的产物，它就是时代的精髓。

爱数软件股份有限公司
李基亮

序 4

存储是个大市场，有意向在数据和信息系统上做投资规划的企业逐年增加，这标志着越来越多的企业意识到自身的数据安全问题。

在我十几年前刚刚踏入存储圈子之时，数据安全问题只被金融、电信等少数行业所考虑，而如今，几乎各个行业都存在数据保护与信息安全的需求。随着用户需求的急速增长，无论是硬件设备还是软件产品都是生机一片。但是，多年来我国的这个领域一直被国外产品所垄断，究其原因，是我国存储领域技术相对滞后。

我们在经营企业的过程中，花费了大量的精力进行人才的培养。在国内，计算机行业的传统教育大多集中于软件应用与网络维护上，对于专业存储的技术培训几乎为零，而存储行业又在飞速地发展着，因此，存储市场的需求与人才滞后的落差越拉越大，我们急切渴望拥有存储专业的人才去发展存储领域。“人才为本，教育当先”，人才的培养离不开教育。多年以来，存储领域的教材乃至书籍几乎是一片空白，有的也只是太过于教条以及模式化的书籍，当看到张冬先生的《大话存储》后，我深刻地体会到我国存储领域开始有了专业的教科书，我国的存储业生机盎然。

之所以赋予《大话存储》如此高的评价，是因为它的语言通俗而不失专业，幽默而不失严谨。张冬先生用读者极易接受的语言道出了存储领域的精髓。对于初学者来说，能使存储领域不再陌生，而又充满吸引。我曾了解到，《大话存储》已经成为某院校计算机专业的教材，这不仅是存储业的幸事，同时也是现代教育的幸事。坦率地讲，我们做企业，时刻关心教育的发展，我们需要新鲜的血液来继承和发展我们的事业。《大话存储》作为能够真正做到学以致用的教材之一，使我们倍感欣慰。我为我们选择的存储道路之前景充满信心，为振兴我们的民族工业充满信心，同时，为张冬这样的后继人才而倍感骄傲。

《大话存储》能够成为教材是张冬对于存储领域不懈努力的成果，《大话存储 II》的出版，更是他不断追求与探索的结果。《大话存储 II》在《大话存储》的基础上更加深入地剖析了存储技术，以及存储在如今市场的广泛应用。书中不乏一些当今企业的存储实例，也包含了国内外软硬件厂家的存储技术应用，加入了更多实际范例，使读者更易理解，同时具有很强的应用性。

我相信《大话存储 II》会给广大读者很大的帮助，同时也希望此书能够带领更多的有识青年进入存储领域，为我国民族产业的振兴而奋斗。

火星高科 总经理
龚平

序 5

认识张冬，是因他的《大话存储》，我曾在去年拜读此书，感觉一个80后的小伙子能用如此通俗的语言诠释存储技术，实属存储行业的一大喜事。这本书，可以让不了解存储的人认识存储，能够了解到存储并不是高深莫测的，即使一个存储行业以外的人去阅读《大话存储》，也一定能够读懂。用什么样的语言和叙述方式不重要，重要的是把要说的说明白。

张冬本人就像他的书一样，饱含着严谨的作风和真诚的态度，而又不乏幽默的风格。看过他的BLOG，人气一直很旺，这个致力于为国产存储业做出贡献的年轻人更是让我对他刮目相看。他在博客中写到：“我所能够做的，只有让中国人，让所有中国存储行业的人，以及中国存储行业本身，有一个扎实的基础。如果能够促进国产存储软件硬件的发展，那鄙人就是鞠躬尽瘁，死而后已，死而无憾！”一个80后年轻人有这样的雄心壮志，我们有什么理由不去努力不去发展国产存储业呢？

记得十几年前，我刚刚进入存储领域，那时候相关的书籍非常少，完全要靠自己进行反复的试验。那时（IT行业根本不成形，姑且称作计算机行业）计算机业的从业者都是抱着掌握20世纪末最具科技含量的技术的心态进行工作，从根本上说，对存储技术充满了崇拜，甚至有一丝恐惧。在探索期间，也走了不少弯路，耽误了很多时间。如果那个时候有这样一本关于存储的书籍，那简直是一大幸事！书中并没有把存储看做是多么高深的技术，而是任何一个普通人都能掌握的技术。我和张冬开玩笑说，如果你早生10年，你就可以带领我们走向一条存储道路的捷径。

看到张冬最新写作的《大话存储II》时，我就感觉到这又是一本好书。不仅延续了《大话存储》中通俗易懂的语言及“武侠”式的章节回目，在技术深度上，也有很深的挖掘。书中不仅囊括了时下最先进的“云”技术以及持续数据保护（CDP）技术，还涉及到了很多非常底层的架构。在《大话存储I》的基础上，有了更为深刻的剖析。值得一提的是，张冬在最后还加入了Q&A的内容，把几年来读者以及网友提出的问题一一列出，并作出详细的解答，能够体会张冬在这一年多的时间里，对于存储技术的探索花了很多的心思。最可贵的是，这个年轻人不以如此成就为骄傲，继续孜孜不倦地探求。

《大话存储II》是一本好书，作者那严谨而真诚的态度以及致力于发展本国存储业的信心注定能够成就这样一部优秀的作品。我完全有理由相信此书能够给从业者乃至热爱存储的读者带来帮助。从中，你会受益匪浅，并乐意向你的朋友推荐此书。

火星高科 技术总监
黄疆

关于冬瓜头和《大话存储》

问：你平时怎么都不说话，一点八卦都不聊？

答：我从来不浪费精力在那些八卦上，省点精力就能多思考和多写一些东西出来。

问：那你关心周围发生的事情么？

答：我回去就打开收音机，发生的国家大事我都知道。

问：平时有什么爱好？

答：思考周围的事情并写出来，感叹一下社会，玩一玩电玩，弹一弹吉他，看一看存储界的新闻和技术。

问：你咋不懂生活呢？来北京不去看看天安门么？

答：这就是我的生活，你敢说你比我懂生活？我的理解是，好生活就是感到幸福。我现在很幸福，因为我把我的知识共享给了所有人，别人看了提高很快，这难道不幸福么？看个天安门能让你感到幸福？我不理解，也不想去理解，如果你觉得看天安门你就幸福，那我祝福你，多去看一看。

问：有人说你书里的那些武侠情节，还有那些诗，很烂，你怎么看？

答：说实话，我现在看看觉得有些地方也确实挺烂的，但是当初我写这些东西是有原因的，写书得有个引子，带领你去不断地写，很少有人能够只写而不联想不类比不思考，就这么写下去写出一本书来的。而这些武侠情节，就是我当时自己给自己的一个引子。

问：《大话存储 II》为何没有去掉这些引子？

答：保持原汁原味吧，大家将就将就吧，留着这个引子，也是对当初写书时的心境的一种保留和回忆。

问：听接触过你的某些人说你好像脑袋缺根筋，不知道你对这个说法怎么看？

答：我不是缺根筋，我是缺很多根筋，而且少了很多心眼。正因如此，才能够让我有足够的精力来完成这样一本大部头、高容量的书籍，要那么多筋作甚？与人斗一斗，或者偷个懒什么的，有意思么？无聊透顶！为了这本书我可以说是废寝忘食，整个人的状态和精神病患者接近，不过现在好了，终于完成了，死而无憾了。

问：你经常提到科学家，是不是有科学家情结？

答：是的，你说对了。我从小就有较强的好奇心，记得小时候去了哪见到抽屉就想翻一翻。印象最深的就是在祖母和外祖母家，老人一般不会那么严格管教小孩，所以每次去都要翻个底朝天看看有什么好东西玩。现在没有这个癖好了，呵呵。不过我看到现在的八九岁的小孩好像也很爱翻抽屉，我想这是人成长中必要的阶段。可惜，现在的父母缺乏思想境界，每次总是强烈禁止这种好奇行为，殊不知无意中就扼杀了孩子的好奇心。而现在遇到了不懂的技术问题，也是一定要搞清楚，通过各种手段，不弄清楚我就睡不好觉。晚上躺下之后我就一直在想白天没有弄明白的问题，当然，往往是没想出个什么来就呼呼大睡了。可惜，我没能成为科学家。

问：你性格很倔强，这样能容到人堆里么？

答：我在人堆里就一傻子，有时候是真傻，有时候是装傻。傻一点好，事少，能有更多

的精力来钻研技术，要那么精作甚？我比较讨厌人精，在一起感觉特别累！所以一般情况下我都愿意一个人呆着，一堆人在一起那种娱乐对我来说是一种最难受的负担。一个月甚至一年不说话也没问题。

问：你这种性格和追求估计在国外会更加舒服，不考虑出国么？

答：出国干什么，中国人就好好在中国呆着。振兴民族科技，这是我当前阶段的信仰，基于这个信仰，我可以做出一些之前不想去做或者做不到的事情。

问：还有人说你所谓的支持国产存储属于狭隘的民族主义，你怎么看？

答：我是个俗人，不懂什么是狭隘什么是宽大，我只知道有了信仰就去做。不管什么狭隘还是宽大，我都不懂这些是什么意思。我把知识共享出来，让国人提高，我支持国产存储，就说我狭隘，那他们不狭隘，他们宽大。道不同不相为谋。

问：好么，你整一个老学究啊！

答：我是小学究，跟得上潮流，不是老古董，上得天堂下得地狱，光脚不怕穿鞋，我就知道想好了就做。只是性格上很学究而已，总喜欢穷根究底，问烦过很多人，也问爽过很多人，有些人喜欢我这种性格和做事方法，而有些人则不喜欢。喜不喜欢的，我就这样，生来不是让人喜欢的。不过我希望大家都喜欢我的书，不穷根究底也就出不来这本书，如果你喜欢这书，你就必须知道这书当初是怎么出来的，书如其人，不喜欢我的性格的，不要看这书。

问：你现在最大的愿望是什么？

答：我希望中国做存储的人都能够打牢基础，届时不愁做不出超越西方的产品，还希望到大学去讲课，让大学生学习到真正有用的东西，把他们从那些垃圾文娱的侵害中解救出来！我还有一个希望，就是如果能够靠这本书养活下半生就好了，这样我就可以没有后顾之忧，拿出全部精力来做存储教育事业了，当然这个愿望恐怕是难以实现的了。

前　言

各位读者好，很高兴再次为大家“大话”存储。记得上一次是在3年前，当《大话存储》一书在2008年出版面世之后，我当时就许下承诺，要写《大话存储II》。当时之所以敢夸下海口要继续写第二本，是因为《大话存储》介绍了存储领域最基本的概念和架构，但并没有涉及与深入存储领域最新的技术，比如重复数据删除、Thin Provision、动态分级存储、CDP连续数据保护、SSD固态硬盘、FCoE、SAS、云计算和云存储等。

当年的《大话存储》确实满足了广大读者的需求，出版之后也获得了诸多好评和官方的民间的很多奖项。这些成果逐渐让我感觉到更大的责任和压力。正因如此，所以我深知绝对不能就此停歇，学习是永无止境的，技术是不断发展的，所以我先向大家做了承诺，这样就可以无时无刻的激励我继续学习研究下去了。

写作过程是极其困难的，尤其是当一字一句都需要精雕细琢，并且时刻以通俗表达且让所有人都能看懂的原则和基准去写的时候，其所耗费的精力和脑力是巨大的。记得在一年前撰写本书主体的时候，基本上每天都是早晨七八点钟起来，从床上直接到书桌前开始写，直到中午吃饭，吃饭过程中依然在脑海中构思着，就这样一直到晚上，最晚的一次记得是做一个实验，通宵达旦，直到第二天天亮，实在体力不支，去床上躺到中午，然后继续写。每次睡觉之前，都会带着一个疑问入睡，躺下之后就在脑海中构思、建模，一旦想到某些重要的东西，就用笔记下几个关键词，否则第二天准忘。大部分情况一般都是没想到什么思路就已经呼呼大睡了。这种状态持续了半年之久，当完成了主体稿件之后，真的有一种如释重负的感觉。可惜，好景不长，随着不断的学习和深入，逐渐发现已经写完的内容当中有大量需要补充完善、修饰的部分，在修饰完善的过程中，继续思考，结果发现又引申出更多的东西，有些甚至推翻了以前的结论。这种状态又持续了半年，最终定稿交给编辑之后，依然发现还有零碎的东西需要完善甚至推翻，结果一再将更新的内容同步给编辑，导致出版日期一推再推，出版社相关编辑、校对叫苦不迭，还好咱的老战友大成编辑一如既往的支持，我们都顶住了压力，直到最后一个月时间内没有再发现需要完善的内容，达到了最终收敛。后面这个过程感觉更加耗费精力，因为当你重新审视之前内容的时候，一旦发现不完善甚至错误，就会感觉到一种挫败感和愧疚感，使你的激情和斗志有所丧失。

写书不但是给他人共享知识的过程，它更是一个总结自身知识体系、提高自身修养以及让自己学习更多知识的途径。比如，我在写书过程中，不但通过各方面渠道纠正了之前对某项技术的一些错误认识，而且还学习了更多的知识，并且将这些知识进行深度理解分析，之后通俗的表达出来。当你发现其他人通过你的知识快速提高之后，这种感觉是最充实的。只有在奉献之后才会感到充实，而不是一味的去索取，这样只能更加空虚。

感谢家人对我的支持！长达半年的无业状态，没有家人支持就没有这本书。

感谢那些曾经帮助过我的不计其数的网友和同事！没有鼓励也不会有这本书。

感谢清华大学出版社的工作人员为本书所付出的工作！没有信任更不会有这本书。

另外感谢爱数软件、火星高科、中科蓝鲸与华为赛门铁克这四家国内存储厂商对本书的大力支持！希望国产存储越做越强！没有他们的技术支持和指导，本书专业性会大打折扣。

感谢华中科技大学武汉光电国家重点实验室博士生导师谢长生教授以及中科院计算技术研究所研究员许鲁对本书的大力支持！

感谢本书的广大读者，你们的支持给了我持续前进的动力！

作者联系方式：

QQ: 122567712

Email: 122567712@qq.com; myprotein@sina.com

MSN: myprotein0007@hotmail.com

Blog: <http://space.doit.com.cn/35700>

目 录

第 1 章 混沌初开——存储系统的前世今生	1
1.1 存储历史	2
1.2 信息、数据和数据存储	4
1.2.1 信息	4
1.2.2 什么是数据	6
1.2.3 数据存储	6
1.3 用计算机来处理信息、保存数据	7
第 2 章 IO 大法——走进计算机 IO 世界	9
2.1 IO 的通路——总线	10
2.2 计算机内部通信	11
2.2.1 IO 总线是否可以看作网络	12
2.2.2 CPU、内存和磁盘之间通过网络来通信	13
2.3 网中之网	14
第 3 章 磁盘大挪移——磁盘原理与技术详解	15
3.1 硬盘结构	16
3.1.1 盘片上的数据组织	17
3.1.2 硬盘控制电路简介	22
3.1.3 磁盘的 IO 单位	23
3.2 磁盘的通俗演绎	25
3.3 磁盘相关高层技术	27
3.3.1 磁盘中的队列技术	27
3.3.2 无序传输技术	27
3.3.3 几种可控磁头扫描方式概论	28
3.3.4 关于磁盘缓存	29
3.3.5 影响磁盘性能的因素	31
3.4 硬盘接口技术	31

3.4.1 IDE 硬盘接口	32
3.4.2 SATA 硬盘接口	34
3.5 SCSI 硬盘接口	37
3.6 磁盘控制器、驱动器控制电路和磁盘控制器驱动程序	43
3.6.1 磁盘控制器	43
3.6.2 驱动器控制电路	44
3.6.3 磁盘控制器驱动程序	44
3.7 内部传输速率和外部传输速率	45
3.7.1 内部传输速率	45
3.7.2 外部传输速率	46
3.8 并行传输和串行传输	46
3.8.1 并行传输	46
3.8.2 串行传输	48
3.9 磁盘的 IOPS 和传输带宽 (吞吐量)	48
3.9.1 IOPS	48
3.9.2 传输带宽	49
3.10 固态存储介质和固态硬盘	50
3.10.1 SSD 固态硬盘的硬件组成	50
3.10.2 从 Flash 芯片读取数据的过程	53
3.10.3 向 Flash 芯片中写入数据的过程	54
3.10.4 Flash 芯片的通病	55
3.10.5 SSD 给自己开的五剂良药，药到是否病除	57
3.10.6 SSD 如何处理 Cell 损坏	59
3.10.7 SSD 的前景	60
3.11 小结：网中有网，网中之网	61
第 4 章 七星北斗——大话/详解七种 RAID	63
4.1 大话七种 RAID 武器	64

4.1.1 RAID 0 阵式.....	64
4.1.2 RAID 1 阵式.....	66
4.1.3 RAID 2 阵式.....	68
4.1.4 RAID 3 阵式.....	70
4.1.5 RAID 4 阵式.....	74
4.1.6 RAID 5 阵式.....	75
4.1.7 RAID 6 阵式.....	79
4.2 七种 RAID 技术详解.....	81
4.2.1 RAID 0 技术详析.....	83
4.2.2 RAID 1 技术详析.....	85
4.2.3 RAID 2 技术详析.....	86
4.2.4 RAID 3 技术详析.....	88
4.2.5 RAID 4 技术详析.....	90
4.2.6 RAID 5 技术详析.....	93
4.2.7 RAID 6 技术详析.....	96
第 5 章 降龙传说——RAID、虚拟磁盘、卷和文件系统实战	99
5.1 操作系统中 RAID 的实现和配置.....	100
5.1.1 Windows Server 2003 高级磁盘管理	100
5.1.2 Linux 下软 RAID 配置示例	105
5.2 RAID 卡	107
5.3 磁盘阵列	118
5.3.1 RAID 50.....	119
5.3.2 RAID 10 和 RAID 01	119
5.4 虚拟磁盘	120
5.4.1 RAID 组的再划分	121
5.4.2 同一通道存在多种类型的 RAID 组	121
5.4.3 操作系统如何看待逻辑磁盘	121
5.4.4 RAID 控制器如何管理逻辑磁盘	121
5.5 卷管理层	123
5.5.1 有了逻辑盘就万事大吉...	123
5.5.2 深入卷管理层	124
5.5.3 Linux 下配置 LVM 实例	125
5.5.4 卷管理软件的实现	127
5.5.5 低级 VM 和高级 VM.....	129
5.5.6 VxVM 卷管理软件配置简介	130
5.6 大话文件系统	133
5.6.1 成何体统——没有规矩的仓库	133
5.6.2 慧眼识人——交给下一代去设计	134
5.6.3 无孔不入——不浪费一点空间	134
5.6.4 一箭双雕——一张图解决两个难题	135
5.6.5 宽容似海——设计也要像心胸一样宽	137
5.6.6 老将出马——权威发布...	137
5.6.7 一统江湖——所有操作系统都在用	138
5.7 文件系统中的 IO 方式	138
第 6 章 阵列之行——大话磁盘阵列	141
6.1 初露端倪——外置磁盘柜应用探索	142
6.2 精益求精——结合 RAID 卡实现外置磁盘阵列	143
6.3 独立宣言——独立的外部磁盘阵列	144
6.4 双龙戏珠——双控制器的高安全性磁盘阵列	146
6.5 龙头凤尾——连接多个扩展柜..	148
6.6 锦上添花——完整功能的模块化磁盘阵列	149
6.7 一脉相承——主机和磁盘阵列本是一家	150
6.8 天罗地网——SAN	151

第 7 章 熟读宝典——系统与系统之间的语言 OSI	153
7.1 人类模型与计算机模型的对比剖析	154
7.1.1 人类模型	154
7.1.2 计算机模型	155
7.1.3 个体间交流是群体进化的动力	156
7.2 系统与系统之间的语言——OSI 初步	156
7.3 OSI 模型的七个层次	157
7.3.1 应用层	157
7.3.2 表示层	158
7.3.3 会话层	158
7.3.4 传输层	158
7.3.5 网络层	159
7.3.6 数据链路层	160
7.3.7 物理层	162
7.4 OSI 与网络	163
第 8 章 勇破难关——Fibre Channel 协议详解	167
8.1 FC 网络——极佳的候选角色	168
8.1.1 物理层	168
8.1.2 链路层	168
8.1.3 网络层	170
8.1.4 传输层	175
8.1.5 上三层	176
8.1.6 小结	176
8.2 FC 协议中的七种端口类型	177
8.2.1 N 端口和 F 端口	177
8.2.2 L 端口	177
8.2.3 NL 端口和 FL 端口	178
8.2.4 E 端口	180
8.2.5 G 端口	180
8.3 FC 适配器	181
8.4 改造盘阵前端通路——SCSI 迁移到 FC	182
8.5 引入 FC 之后	183
8.6 多路径访问目标	186
第 9 章 天翻地覆——FC 协议的巨大力量	191
9.1 FC 交换网络替代并行 SCSI 总线的必然性	192
9.1.1 面向连接与面向无连接	192
9.1.2 串行和并行	193
9.2 不甘示弱——后端也升级换代为 FC	193
9.3 FC 革命——完整的盘阵解决方案	195
9.3.1 FC 磁盘接口结构	195
9.3.2 一个磁盘同时连入两个控制器的 Loop 中	196
9.3.3 共享环路还是交换——SBOD 芯片级详解	196
9.4 SAS 大革命	206
9.4.1 SAS 物理层	206
9.4.2 SAS 链路层	208
9.4.3 SAS 网络层	209
9.4.4 SAS 传输层和应用层	211
9.4.5 SAS 的应用设计和实际应用示例	213
9.4.6 SAS 目前的优势和面临的挑战	214
9.5 中高端磁盘阵列整体架构简析	215
9.5.1 IBM DS4800 和 DS5000 控制器架构简析	217
9.5.2 NetApp FAS 系列磁盘阵列控制器简析	223
9.5.3 IBM DS8000 简介	225
9.5.4 富士通 ETERNUS DX8000 磁盘阵列控制器结构简析	225
9.5.5 EMC 公司 Clariion CX/CX3 及 DMX 系列盘阵介绍	228
9.5.6 HDS 公司 AMS2000 和 USP 系列盘阵介绍	232