

崭新的视野·扎实的理论基础

存储 网络

技术及应用

鲁士文 编著

构建企业数据中心的向导

▶ 设计数据灾难备份的指南

▶ 实现系统虚拟化存储的助手

清华大学出版社

存储 网络

技术及应用

★ ★
图书馆藏书
鲁士文 编著

中文样本图书

清华大学出版社
北 京

内 容 简 介

本书面向存储网络产品的工程设计人员,特别是针对企事业单位存储网络规划、配置和存储网络系统管理维护人员的需求而编写,全书以通俗易懂的语言和具有实际意义的图表介绍了存储网络的基本原理、体系结构和设计方法,并重点讨论面向企事业单位的存储网络应用技术、解决方案、设备配置和维护管理。

本书最大的特点是理论与技术相结合,注重实用知识的介绍和有据可循的解决方案,可提高读者参加实际网络产品开发和配置使用的能力。

本书可供存储网络技术人员和设备工程师作为技术参考资料,也可用作信息技术相关专业的研究生和大学高年级的教学参考书。

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

存储网络技术及应用 / 鲁士文编著. —北京:清华大学出版社, 2010.2

ISBN 978-7-302-21926-2

I. ①存… II. ①鲁… III. ①计算机网络—信息存贮 IV. ①TP393

中国版本图书馆CIP数据核字(2010)第007704号

责任编辑:夏非彼 郑奎国

责任校对:郑奎国

责任印制:王秀菊

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:010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者:北京市世界知识印刷厂

装 订 者:三河市新茂装订有限公司

经 销:全国新华书店

开 本:190×260 印 张:20.75 插 页:2 字 数:505千字

版 次:2010年2月第1版 印 次:2010年2月第1次印刷

印 数:1~4000

定 价:45.00元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系
调换。联系电话:(010)62770177 转 3103 产品编号:035184-01

前 言

从固定的直接附接的存储向存储网络的过渡对 IT 行业产生了重要的影响。IT 系统也因此正在由以服务器为中心的体系结构转变为以存储为中心的体系结构。第一代光纤通道 SAN 最初只被大的企业采用，现在已向更大范围的中小企业发展。同时在存储网络领域又涌现出如 IP SAN、存储虚拟化和基于 CIM 和 SMI-S 的全面的 SAN 管理等新技术。这些新的 SAN 技术在提供更多的基于 SAN 智能化和存储过程自动化的应用产品的同时，也推动存储网络进入 IT 市场的主流。从当前各个企事业单位对构建数据中心的热情和对 SAN 交换机以及相关存储设备的不断增长的需求就可以清楚地看到这一趋势。

本书为适应存储网络产品的工程设计人员，特别是企事业单位存储网络规划、配置和管理维护人员的需求而编写，介绍了存储网络的基本原理、体系结构和设计方法，并重点讨论面向企事业单位的存储网络应用技术、解决方案、设备配置和维护管理。

全书共有 11 章，包括共享存储的基本概念和 SNIA 模型、SCSI 总线和协议、磁盘子系统、文件系统和网络附接存储、光纤通道和存储区域网、IP 网络存储和 InfiniBand 网络、存储虚拟化、存储网络的应用、网络备份、可移动介质的管理和存储网络管理。每一章都通过通俗易懂的描述和具有实际意义的图表阐明原理、标准、方案和运用技术。

本书最突出的特色是理论与技术相结合，注重实用技术的介绍，可提高读者参加实际网络产品开发、配置使用和管理的能力。

本书可供存储网络技术人员和设备工程师用作相关的技术参考资料，也可作为信息技术相关专业的研究生和大学高年级的教学参考书。

编 者

2009 年 10 月 于北京

参考文献

- [1] 郭玉东, 尹青. 基于对象的网络存储. 北京: 电子工业出版社, 2007 年 10 月
- [2] 姜宁康, 时成阁. 网络存储导论. 北京: 清华大学出版社, 2007 年 7 月
- [3] 赵文辉, 徐俊等. 网络存储技术. 北京: 清华大学出版社, 2005 年 3 月
- [4] (美)Tom Clark. 邓劲生, 李宝峰等译. 存储区域网络设计. 北京: 电子工业出版社, 2005 年 1 月
- [5] 李蔚泽. Red Hat Linux 9 网络管理. 北京: 清华大学出版社, 2004 年 1 月
- [6] 魏永明, 郑翔等. 学用 Linux 与 Windows NT. 北京: 电子工业出版社, 1999 年 10 月
- [7] (美)Craig Hunt. 翟炯, 石祥生等译. TCP/IP 网络管理. 北京: 电子工业出版社, 1997 年 8 月
- [8] IBM International Technical Support Organization, "Introduction to Storage Area Networks", Fourth Edition, International Business Machines Corporation, 2006
- [9] Ulf Troppens, Rainer Erkens and Wolfgang, "Storage Networks Explained", John Wiley&Sons, Ltd, 2004
- [10] Cary Orenstein, "IP Storage Networking", Pearson Education, Inc., 2003

目 录

第 1 章 共享存储的基本概念和 SNIA 共享存储模型	1
1.1 以服务器为中心的 IT 体系结构	2
1.2 以存储为中心的 IT 体系结构	3
1.3 SNIA 共享存储模型	5
1.4 共享存储配置方案示例	10
第 2 章 SCSI 总线和协议	11
2.1 I/O 通路	12
2.2 并行 SCSI 总线	14
2.2.1 SCSI 类型	15
2.2.2 SCSI 控制器、设备和电缆	16
2.2.3 终接器	17
2.3 SCSI 协议	17
2.3.1 SCSI 域	18
2.3.2 SCSI 协议模型	18
2.3.3 寻址机制	19
2.3.4 交互方式	21
2.3.5 SCSI 总线信号	22
2.3.6 SCSI 总线的使用阶段	23
2.3.7 异步传输和同步传输	27
2.3.8 SCSI 命令描述块	29
2.3.9 SCSI 的读操作和写操作过程	31
2.4 使用多端口存储设备构建 SCSI 存储网络示例	32
第 3 章 磁盘子系统	33
3.1 硬盘和内部 I/O 通道	36
3.2 JBOD 磁盘阵列	38
3.3 使用 RAID 的存储虚拟化	39
3.4 RAID 等级及配置示例	41
3.4.1 RAID 0: 按块条带	41
3.4.2 RAID 1: 按块镜像	42
3.4.3 RAID 0+1 和 RAID 10: 结合条带与镜像	43
3.4.4 RAID 4 和 RAID 5: 用校验位代替镜像	45
3.4.5 RAID 2 和 RAID 3	49
3.4.6 RAID 等级的比较	50
3.5 使用缓存加速对磁盘的访问	51

3.5.1	在硬盘上的缓存	51
3.5.2	在 RAID 控制器中的写缓存	51
3.5.3	在 RAID 控制器中的读缓存	52
3.6	智能磁盘子系统	52
3.6.1	即时复制	52
3.6.2	远程镜像	53
3.6.3	逻辑设备号掩盖	57
3.7	磁盘子系统的可提供性	59
第 4 章	文件系统和网络附接存储	61
4.1	本地文件系统	62
4.2	网络文件系统和网络附接存储的基本概念	65
4.3	NFS 文件系统的组成结构	68
4.3.1	NFS 通信协议	69
4.3.2	挂载协议	71
4.3.3	文件操作协议	72
4.3.4	NFS 的安全机制	72
4.4	NFS 在 UNIX 系统上的配置示例	73
4.4.1	NFS 守护程序	73
4.4.2	文件系统的输出	74
4.4.3	挂载远程文件系统	75
4.5	CIFS 文件系统	77
4.5.1	SMB 协议	79
4.5.2	CIFS 操作	80
4.5.3	CIFS 的安全机制	85
4.6	CIFS 在 Linux 系统上的配置示例	86
4.6.1	Samba 的主要成分	86
4.6.2	Samba 服务器的配置文件	87
4.6.3	共享 Linux 目录	93
4.6.4	从 Linux 系统中访问 Windows 的共享目录	94
4.6.5	共享 Linux 打印机	96
4.6.6	从 Linux 系统中访问 Windows 的共享打印机	97
4.7	网络文件系统的演变和发展	98
4.7.1	客户标识符和会话	100
4.7.2	服务器名字空间	103
4.7.3	文件句柄	108
4.7.4	并行网络文件系统	109
4.8	网络附接存储的结构及其发展趋势	111
4.9	NAS 和 DAS 的性能比较	113

第 5 章 光纤通道和存储区域网	115
5.1 光纤通道层次模型	116
5.2 物理结构	117
5.3 FC-0: 物理接口和介质	119
5.4 FC-1: 传输协议	121
5.5 FC-2: 成帧和信令协议	125
5.5.1 交换、序列和帧	125
5.5.2 流控制	127
5.5.3 服务类别	127
5.6 FC-3: 公共服务	131
5.7 链路服务	132
5.7.1 登录	132
5.7.2 编址	134
5.8 交换网服务: 名字服务器和控制器	135
5.9 FC-4: 上层协议映射	137
5.10 存储区域网络	138
5.10.1 SAN 与 DAS 和 NAS 的比较	140
5.10.2 采用光纤通道网络的原因	141
5.10.3 光纤通道 SAN 的组成结构	142
5.11 仲裁环 SAN	146
5.12 交换网 SAN	149
5.12.1 名字服务器	150
5.12.2 交换机的种类和主交换机	151
5.12.3 路由选择	151
5.12.4 存储分区	152
5.13 采用存储区域网的企业信息系统的结构配置示例	153
第 6 章 IP 存储网络和 InfiniBand 网络	155
6.1 因特网 SCSI	156
6.1.1 iSCSI 体系结构	157
6.1.2 目标方发现	158
6.1.3 iSCSI 会话	164
6.1.4 iSCSI 会话协议数据单元	165
6.2 在 TCP/IP 上的光纤通道	170
6.3 因特网光纤通道协议	172
6.4 多协议环境和相关的解决方案	178
6.4.1 主要术语的精确含义	179
6.4.2 需要特别考虑的事项	179
6.4.3 多协议解决方案	182
6.5 IP 存储网络配置示例	184

6.6	InfiniBand 网络	185
第 7 章	存储虚拟化	189
7.1	在 I/O 通路上的虚拟化	191
7.2	块级和文件级的存储虚拟化	195
7.3	在应用服务器中的存储虚拟化	196
7.4	在存储设备中的存储虚拟化	197
7.5	在网络中的存储虚拟化	198
7.5.1	对称的存储虚拟化	199
7.5.2	非对称的存储虚拟化	201
7.6	文件系统和 NAS 虚拟化	204
7.7	存储虚拟化产品及其应用	204
第 8 章	存储网络的应用	211
8.1	存储共享	212
8.1.1	磁盘存储池	212
8.1.2	动态磁带库共享	212
8.1.3	数据共享	214
8.2	数据可提供性	216
8.2.1	预防 I/O 总线故障	216
8.2.2	预防服务器故障	219
8.2.3	预防磁盘子系统故障	220
8.2.4	预防虚拟化故障	223
8.2.5	预防数据中心的崩溃	223
8.3	对 IT 系统的企业适应性和易于扩展性的支持	226
8.3.1	负载分布的集群	227
8.3.2	Web 体系结构	232
8.3.3	Web 应用的实现	234
8.4	园区存储网络配置	237
8.5	因特网提供商的存储网络	239
第 9 章	网络备份	241
9.1	网络备份服务	242
9.2	备份服务器	243
9.1.1	作业调度程序	244
9.1.2	错误处理程序	244
9.1.3	元数据数据库	244
9.1.4	介质管理程序	244
9.3	备份客户	246
9.4	网络备份有助于系统性能的提升	247
9.5	网络备份的性能瓶颈	248

9.6	提高网络备份系统性能的技术途径	249
9.6.1	针对以服务器为中心的 IT 系统结构的措施	249
9.6.2	面向新一代网络备份的措施	252
9.7	文件系统备份	259
9.7.1	文件服务器备份	259
9.7.2	文件系统备份	260
9.7.3	NAS 服务器备份	261
9.7.4	网络数据管理协议 NDMP	262
9.8	数据库备份	267
9.8.1	数据库的操作方法	268
9.8.2	经典的数据库备份	269
9.8.3	下一代数据库备份	271
第 10 章	可移动介质的管理	273
10.1	可移动介质	274
10.2	常用术语	275
10.3	介质库和驱动器	276
10.3.1	驱动器	277
10.3.2	介质更换器	277
10.3.3	对介质更换器的控制	277
10.4	可移动介质管理的问题和需求	279
10.4.1	对可提供资源的有效使用	281
10.4.2	访问控制	282
10.4.3	访问同步	284
10.4.4	访问优先级和安装请求队列	284
10.4.5	介质跟踪	285
10.4.6	组合	286
10.4.7	监视	288
10.4.8	报告	289
10.4.9	寿命周期管理	290
10.4.10	储藏库管理	292
10.5	IEEE 1244 介质管理系统标准简介	292
10.6	介质管理系统结构	293
10.6.1	介质管理模块	294
10.6.2	库和驱动器管理模块	295
10.6.3	特权和非特权客户	295
10.7	介质管理系统数据模型	296
10.8	IEEE 1244 通信协议	298
第 11 章	存储网络管理	301
11.1	对管理系统的要求	302

11.2	管理接口	303
11.3	标准的和专用的机制	305
11.3.1	标准机制	305
11.3.2	专用机制	305
11.4	带内管理	306
11.4.1	管理服务	308
11.4.2	发现	308
11.4.3	监视	309
11.4.4	消息	309
11.4.5	辖区问题	309
11.5	带外管理	310
11.5.1	使用 SNMP	311
11.5.2	基于 Web 的企业管理标准 WBEM	315
11.5.3	存储管理倡议规范	319
11.6	对存储网络管理配置的选择	320
	参考文献	322

第 1 章 共享存储的基本概念和 SNIA 共享存储模型

学习要点

- 以服务器为中心的 IT 体系结构
- 以存储为中心的 IT 体系结构
- SNIA 共享存储模型
- 共享存储配置方案示例

存储网络的设计需要对依赖存储资源的高层应用以及满足应用需求的存储系统的体系结构有深入的了解。本章先描述传统的以服务器为中心的 IT (Information Technology, 信息技术) 体系结构, 指出其局限性; 然后介绍以存储为中心的另一类 IT 体系结构, 说明其优越性; 接着讨论 SNIA (Storage Networking Industry Association, 存储网络行业协会) 提出的共享存储模型; 最后给出一个共享存储配置方案示例。

1.1 以服务器为中心的 IT 体系结构

在传统的 IT 体系结构中, 存储设备通常只连接到单个服务器 (参见图 1-1)。为了增加容错能力, 有时也把存储设备连接到两个服务器, 但在任一时刻仅一个服务器能够实际地直接使用存储设备。在上述两种情况下, 存储设备都隶属于它所连接的服务器。其他的服务器不能直接访问记录在该存储设备上的数据, 它们只能通过连接该存储设备的服务器访问。因此, 这种传统的 IT 体系结构被称作是以服务器为中心的。在这种方法中, 服务器和存储设备一般是通过 SCSI (The Small Computer System Interface, 小型计算机系统接口) 电缆连接在一起的。

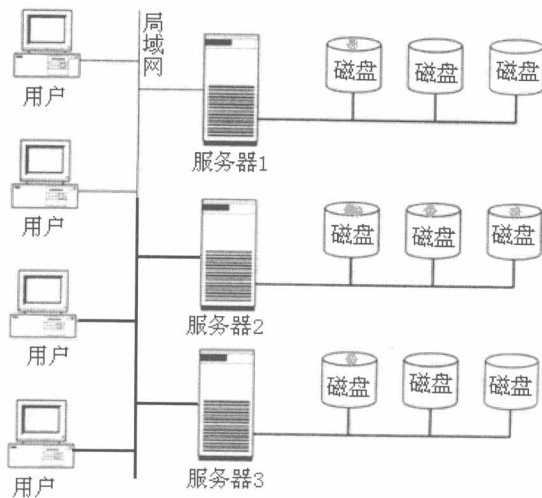


图 1-1 存储设备被静态地分配给它所连接的计算机

在传统的以服务器为中心的 IT 体系结构中, 由于存储设备隶属于它们连接的一个或两个服务器, 所以如果所连接的服务器都失效了, 那么在存储设备上的数据将不能够被访问。这对于大多数企业都是不可接受的, 例如, 医院病历文件和银行账户的数据需要随时都可以提供访问。

虽然由于技术的进步, 硬盘和磁带的存储密度一直在增加, 但是对于安装的存储容量的需求增加得更快, 因此, 需要把更多的存储设备连接到计算机。然而现实的情

况是每台计算机只能安装有限数目的 I/O 卡（例如 SCSI 卡），更糟的是，SCSI 电缆的长度最大只允许 25m。这就意味着，使用常规技术连接到一台计算机的存储容量是有限的。因此常规技术已经不能满足对存储容量日益增长的需求。

在以服务器为中心的 IT 环境中，存储设备被静态地分配给它所连接的计算机。一般情况下，一台计算机不能够访问连接到另一台计算机的存储设备。这样做一方面是出于数据安全的考虑，另一方面也是为了避免每台计算机都要执行授权和权限管理的复杂性。其结果是，如图 1-1 所示，在服务器 2 已经用完了它连接的磁盘空间的情况下，尽管服务器 1 和服务器 3 仍然有剩余磁盘空间，但它却不能够利用。另外，把存储设备分散在建筑物的各个房间内，既不易保证空调、防尘和防潮等机房条件，也不利于对非授权访问的防范。

1.2 以存储为中心的 IT 体系结构

存储网络可以解决前述以服务器为中心的 IT 体系结构的问题。它还为我们开辟了一种数据管理的新方式。如图 1-2 所示，在存储网络背后的思想是用一个网络代替 SCSI 电缆，并且该网络主要用于在计算机和存储设备之间的数据交换。

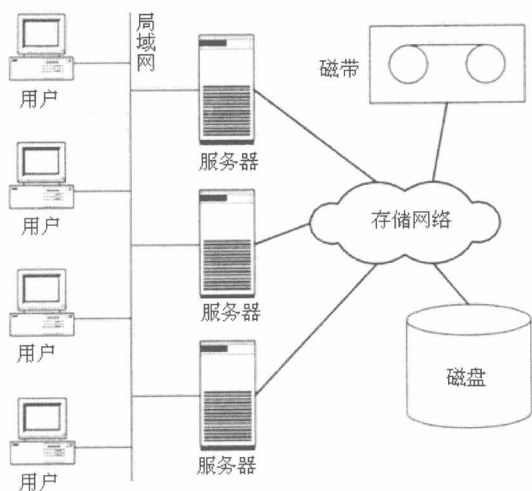


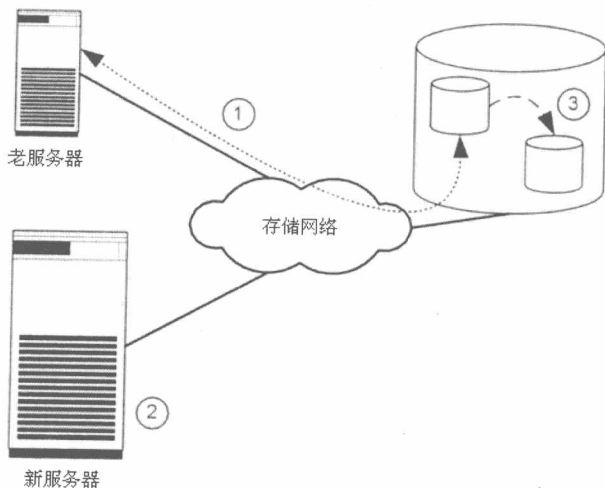
图 1-2 把 SCSI 电缆用一个网络代替

和以服务器为中心的 IT 体系结构不同，在存储网络中的存储设备完全独立于任何计算机，多个服务器可以直接在存储网络上访问同一个存储设备，而不必通过另一个服务器，因此存储设备被放到了 IT 体系结构的中心位置；而在另一方面，服务器却成了存储设备的附属品，它只是处理数据。也正是因为如此，使用存储网络的 IT 体系结构被称作以存储为中心的 IT 体系结构。

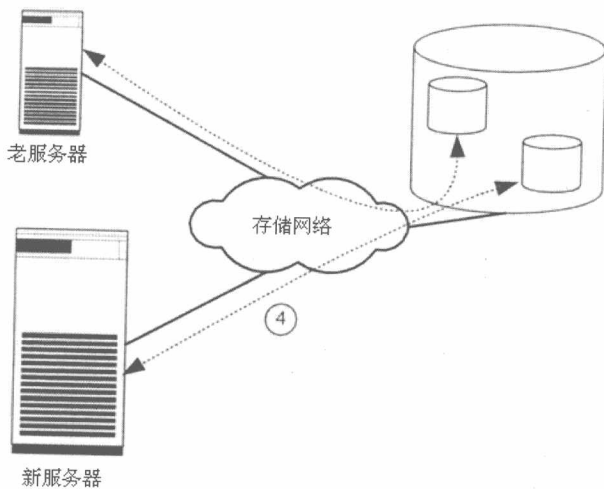
通常在引入存储网络之后，存储设备也被强化。这涉及到把附接到若干个计算机

的许多小的硬盘用一个大的磁盘子系统代替。现在的磁盘子系统的最大容量可以有数十个太字节。存储网络允许连接到它的所有计算机都有可能访问这个磁盘子系统，也就是说，磁盘子系统被共享了。这样，空闲存储容量就可以被灵活地分配给需要它的计算机。同样地，我们也可以把许多小的磁带库用一个大的磁带库代替。

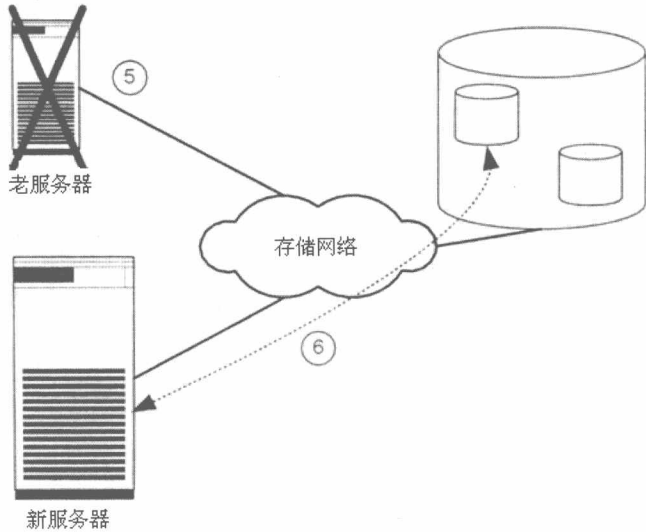
下面我们将通过考察一个实例来说明以存储为中心的 IT 体系结构所具有的优点。随着公司业务的发展，在一个生产环境中的一个应用服务器在被使用了多年之后需要被淘汰，决定用一个具有更高性能的计算机替换。在传统的以服务器为中心的 IT 体系结构中，这将是一个非常复杂的过程；而在存储网络中，这个过程就变得很容易了，只要通过执行下列几个步骤，就能妥善地完成（参见图 1-3（a）、(b)和(c)）。



(a) 步骤①~③



(b) 步骤④



(c) 步骤⑤~⑥

图 1-3 使用存储网络的服务器替换过程

(1) 假定在替换之前，老的服务器是通过一个存储网络连接到一个存储设备的，而且存储空间仅被使用了一部分。

(2) 首先，在新的服务器上安装必要的应用软件。由于使用了存储网络，新服务器可以安装在跟存储系统和老服务器不同的物理位置上。

(3) 接着，在磁盘子系统内复制生产数据，建立测试数据。现代存储系统可以在几秒钟内复制数太 (10^{12}) 字节的数据。

(4) 然后把复制的数据分配给新服务器，并对新服务器进行强化的测试性运行。

(5) 在成功地进行了测试之后，把两台服务器都关机，再把生产数据分配给新服务器。分配过程也只需花几秒钟的时间。

(6) 最后使用生产数据启动新服务器，替换过程即告结束。

1.3 SNIA 共享存储模型

在传统的意义上，每个计算机都可以配置直接附接的存储设备 (DAS: Direct Attached Storage)。为了在多个服务器和工作站之间共享存储资源，则需要有一个连接目标机器和源机器的对等网络。网络的构成及其上传输的存储数据类型因结构而异。一般来说，共享存储体系结构主要划分为 SAN (Storage Area Network, 存储区域网络) 和 NAS (Network Attached Storage, 网络附接存储) 两大类。对于 SAN，网络基础设施可以是光纤通道，也可以是千兆位或万兆位以太网，其上传输的数据类型是 SCSI 块数据。对于 NAS，典型的网络基础设施是以太网 (快速以太网，千兆位或万兆位以太

网)，其上传输的存储数据类型是基于文件的。在最抽象的层次上，SAN 和 NAS 的共同特征是二者都允许存储资源在多个计算机上的多个用户之间共享，而不论数据是基于块还是基于文件的。

对于 IT 工程师和负责存储设备购置的决策人员来说，一个重要的问题是如何理解 DAS、SAN 和 NAS 方案的不同角色，以及如何将它们应用于统一的 IT 存储战略中。SNIA 的共享存储模型（SSM: Shared Storage Model）提供了一个有用的框架，它一方面有助于统一术语和描述模型，另一方面也有助于理解高层应用、存储网络和存储设备之间的关系。

如图 1-4 所示，SNIA 共享存储模型把一个共享的存储环境划分为 4 个基本部分，即应用、文件/记录层、块层和服务子系统。

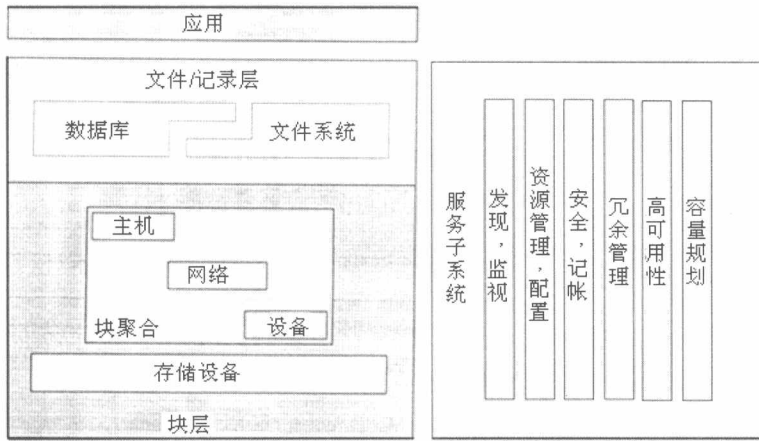


图 1-4 SNIA 共享存储模型的主要成分

SNIA 共享存储模型对应用没有进一步描述，它被看成是该模型的用户。一般来说，应用支持诸如在线事务处理、数据挖掘或 Web 服务等用户行为。

文件/记录层包含数据库和文件系统。

块层包括存储设备和块聚合。SNIA 共享存储模型使用术语“聚合”代替通常使用的术语“存储虚拟化”。

服务子系统定义了管理其他部件的功能，包括存储特定的应用，比如管理、安全、备份、可用性维护以及容量规划等。这样，该模型就把处于高层的终端用户或商业应用与监测、支持底层存储设施的辅助应用区分开了。

SNIA 共享存储模型为运行在计算机上的用户应用与底层的存储设施建立了普遍的联系。具体地讲，它定义了下列成分。

1. 互连网络

互连网络表示存储网络，即把共享存储环境的各个成分互相连接的基础设施，网络必须能够提供高的性能和易于扩展的连接。在实际的安装中，当前主要使用光纤通