



普通高等教育“十二五”规划教材

电力系统自动化

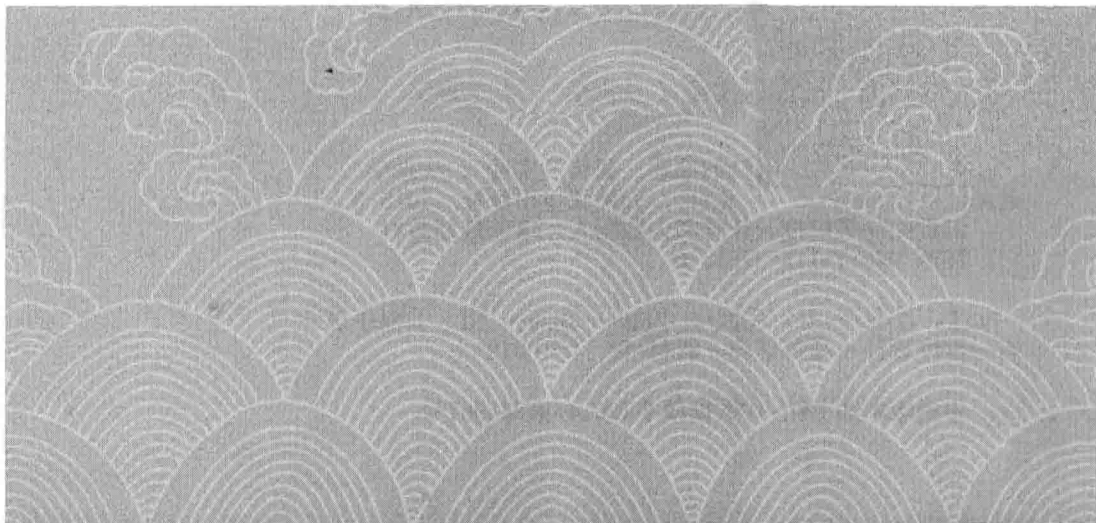
丁坚勇 胡志坚 编



中国电力出版社
CHINA ELECTRIC POWER PRESS



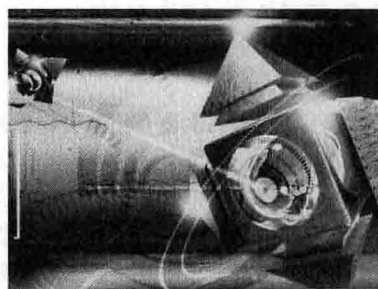
THE CLOUD STORAGE & BIG DATA COMPUTING SYSTEM



21世纪高等院校云计算和大数据人才培养规划教材

云存储系统

——Swift 的原理、架构及实践



The Cloud Storage System

武志学 赵阳 马超英 编著

人民邮电出版社
北京



图书在版编目(CIP)数据

云存储系统：Swift的原理、架构及实践 / 武志学，赵阳，马超英编著. — 北京：人民邮电出版社，2015.2
21世纪高等院校云计算和大数据人才培养规划教材
ISBN 978-7-115-37815-6

I. ①云… II. ①武… ②赵… ③马… III. ①程序语言—程序设计—高等学校—教材 IV. ①TP312

中国版本图书馆CIP数据核字(2014)第281867号

内 容 提 要

本书主要介绍了云存储的起源、概念及特点，文件系统、块存储系统和对象存储系统的原理和使用场景，Swift云存储系统的原理、特性及架构，Swift云存储系统的搭建和维护，Swift云存储系统的各种使用接口；基于Swift的应用开发等方面内容，不仅从理论上介绍了云存储系统的起因、特点、原理、架构和使用场景，更是通过深入浅出地讲解当前国际上最热门的开源云存储系统Swift的原理、架构和使用，使学生在掌握云存储理论知识的同时，能够完全了解、搭建、维护Swift云存储系统，以及开发基于Swift的各类应用。

本书主要面向各级各类院校计算机类专业的学生，对每一个核心概念都进行了严格的定义，并通过各种例题进行详细讲解。学生还可以通过完成每章后面附有的习题和实验，加深对课堂内容的理解和记忆。本书也可供从业人员和计算机爱好者自学参考。

-
- ◆ 编 著 武志学 赵 阳 马超英
责任编辑 王 威
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>
北京隆昌伟业印刷有限公司印刷
 - ◆ 开本：787×1092 1/16
印张：12.25 2015年2月第1版
字数：318千字 2015年2月北京第1次印刷

定价：32.00元

读者服务热线：(010)81055256 印装质量热线：(010)81055316
反盗版热线：(010)81055315

序

数据生产一直贯穿着人类社会的发展历程，利用结绳计数或是在龟甲上刻下书契的人类祖先或许没有料到就是他们这一原始的技术逐步引导并开启了人类数据生产的历史。回顾人类文明早期的几千年，中华民族在数据生产技术上做出了伟大的贡献，印刷术和造纸术的发明为信息记录、复制提供了重要的工具，使人类几千年的文明信息得以保存，这几千年人类度过了信息和数据平稳增长时期。

随着电子信息技术、网络技术的发展，特别是移动互联网技术的发展，数据的生产方式发生了剧烈的变化，数据生产变得无处不在。移动互联网和智能手机的发展使每个人可以随时随地产生数据，数据的生产成为了人们的生活方式，聊微信、刷微博、发照片成为了多数年轻人每天都要做的事。物联网的发展更使数据的生产变得自动化，遍布城市、企业、小区的探头一刻不停地在产生着海量的数据。这一变化使全世界数据量呈现出爆发性的增长，数据的泛滥使技术人员不得不认真应对。吉姆·格雷提出的第四范式标志着“大数据”正式成为了人类需要共同面对的技术难题，未来信息技术的发展“大数据”将成为核心的主题，并将延伸到其他相关学科。

大数据技术的快速发展在产业界已成为共识，如何培养适应新的技术需求的人才是教育界需要迅速面对的问题，人才的缺乏会严重地制约产业发展。在云计算大数据领域，目前图书市场上技术参考书较多，但缺乏适合教学的高品质图书，不少院校苦于没有适合的教材而无法开设相关课程。武志学老师在此时出版的这本《云存储系统——Swift的原理、架构及实践》一书，是对大数据人才培养课程体系研究的一个重要探索。存储问题是大数据需要面对的第一个基础问题，本书将带领大家了解数据存储相关的知识要点，以实践为引导，迅速使学生掌握大数据存储技术的核心理念。

武志学老师毕业于剑桥大学三一学院，拥有丰富的国际云计算企业工作经验，领导创办了国内第一个云计算系，在教学过程中事必亲为，对实验内容要求严格，所有操作及代码均要进行严谨的验证。武老师虽然身兼繁忙的教学和技术领导职务，但一直坚持亲自为本科生进行教学，此书就是武老师在长期教学和工程实践过程中逐步积累起来的，远非一般性的急就篇所能比拟。武老师是一位纯粹的学者，读者在使用本书的过程中可以与武老师进行交流，相信一定会获益匪浅。本书的出版定会为我国云计算大数据人才的培养起到积极的促进作用，也能为其他院校开设相关课程和专业提供有益的参考。

王 鹏 成都信息工程学院并行计算实验室主任
成都信息工程学院并行计算实验室

2014年10月

前 言

随着互联网与移动通信网络的快速发展,数据呈现出爆炸式增长,视频、图片、网站、SaaS应用等都存在无休止的存储需求,这无疑是对数据存储提出新的挑战。2012年国际数据机构(IDC)宣布全球的数字化内容已经超过了2775EB,在最近两年几乎增长了2048EB。

Exabytes(EB)到底有多大?这是一个难以想象的大数量,相当于1000PB、一百万TB或者10亿GB。而一本书的大小一般只有几个MB。世界上最大的图书馆——美国国会图书馆,拥有2300万本书籍,大约23TB的数据。其中,高清照片的大小几倍于2MB。而1EB相当于5亿多张高清照片。

存储数据的高增长空前。IDC预计到2020年,全球的数据量将会超过35000EB,等同于全球每个人都拥有4TB的数据。

面对如此严峻的挑战,2009年,全球三大云计算中心之一的RackSpace开始研发对象存储系统——Swift。该系统旨在解决静态数据的长期存储问题,特别是针对虚拟机镜像、图片存储、邮件存储等。2010年7月,RackSpace将Swift的代码贡献给了OpenStack开源社区,成为了一个开源的云存储系统。

在当前的云计算大数据时代,如何能够安全可靠低成本的存储和使用海量数据是各个企业面临的一个大问题。云存储是通过采用网格、分布式文件系统、服务器虚拟化、集群应用等技术,将网络中海量的异构存储设备构成可弹性扩展、低成本、低能耗的共享存储资源池,并提供数据存储访问、处理功能的一个系统服务。企业和个人都可以通过一个简单的Web服务接口,在任何时间、任何地点存储和检索任意数量的数据,获得高可用、高可靠的数据存储以及稳定廉价的基础存储设施。

Swift具有传统分布式存储所无法比拟的显著优势,为大数据的存储带来了希望,引起了工程师、部署者、管理者们极大的兴趣,但同时,它作为新的存储方式,需要我们从新的角度,以新的思维方式去学习和应用。本书的编写主要是针对高校计算机专业的学生,在全面介绍Swift工作原理和相关技术的同时,还提供了大量的例题和习题,以及实训操作流程。主要内容包括Swift简介、Swift系统架构、Swift工作原理、Swift使用、Swift应用开发、Swift的实现、Swift单机安装、Swift集群安装、Swift集群运维等,从而覆盖了Swift的各个方面,历史、发展趋势、技术术语、整体架构、工作流程、实现方法、调试方法、运维操作等。本书每章配有相关思考题和实训,以巩固学习所用。

本书的作者拥有丰富的国际云计算公司的工作经验、国内外高校教学经验,以及实际开发云计算系统的经验。基于这些,本书并不仅仅是将相关技术内容简单地告诉读者,而是结合他们的实践经验将复杂问题简单化,以深入浅出的方式表达了Swift存储系统的方方面面。目的是希望读者通过对本书的阅读能够掌握Swift工作原理,并在此基础上将所学应用于实践,解决工作中遇到的关于大数据存储的实际问题。

在编写过程中,我们得到了电子科技大学成都学院领导和同事的不少帮助和支持,获得了学校教材建设基金给予的经费资助,在这里表示感谢。同时特别感谢刘小珍在认真阅读初稿的基础上进行总结,为每章补充了习题和思考题;感谢汪雪飞对使用Java开发Swift应用章节代码的验证;感谢王娜、宋怡对Swift搭建方法进行实际验证和各种特性的测试。

编者

2014年11月

目 录 CONTENTS

第 1 章 云存储概述 1

1.1 云存储起源	2	1.2 云存储概念	3
1.1.1 云存储技术起源	2	1.3 云存储的特点	4
1.1.2 云存储服务起源	3		

第 2 章 对象存储系统 6

2.1 非结构化数据存储	6	2.2 对象存储系统	9
2.1.1 什么是非结构化数据	6	2.2.1 对象存储的产生	9
2.1.2 非结构化数据的存储要求	7	2.2.2 对象存储的基本概念	10
2.1.3 存储系统的种类	8	2.2.3 对象存储的关键特性与价值	11
2.1.4 传统的共享存储方法的缺点	8	2.2.4 对象存储的主要应用场景	11

第 3 章 Swift 简介 13

3.1 Swift 的开发历史	13	3.3 Swift 应用场景	16
3.2 Swift 的特性	14	3.3.1 常见案例介绍	16
3.2.1 极高的数据持久性	14	3.3.2 存储用于数据分析	16
3.2.2 可扩展性	14	3.3.3 备份、归档和灾难恢复	16
3.2.3 高并发	14	3.3.4 静态网站托管	17
3.2.4 完全对称的系统架构	14	3.4 CAP 理论简介	17
3.2.5 硬件设备要求低	15	3.4.1 CAP 理论	17
3.2.6 开发的友好性	15	3.4.2 一致性种类	17
3.2.7 管理友好性	15	3.4.3 CAP 理论的应用	18

第 4 章 Swift 的工作原理 20

4.1 核心概念	20	4.2 Swift 的总体架构	22
4.1.1 Swift URL	20	4.2.1 代理服务器 (Proxy Server)	22
4.1.2 账号 (Accounts)	21	4.2.2 存储服务器 (Storage Server)	23
4.1.3 容器 (Containers)	21	4.3 Swift 的工作原理	24
4.1.4 对象 (Objects)	21	4.3.1 虚节点	24
4.1.5 Swift API	21	4.3.2 环 (The Ring)	25

4.3.3 一致性服务器 (Consistency Server)	26	4.4 使用场景举例	29
4.3.4 区域 (Zones)	28	4.4.1 上传 (PUT)	29
4.3.5 地区 (Regions)	28	4.4.2 下载 (GET)	30
4.3.6 数据存储点选择算法	29	4.5 总结	31

第 5 章 Swift 的使用 33

5.1 命令行客户端	33	5.3.2 curl 简单使用	60
5.1.1 安装	34	5.3.3 认证	60
5.1.2 认证	34	5.3.4 获取集群存储使用情况	61
5.1.3 访问控制	35	5.3.5 创建容器和获取容器列表	61
5.1.4 访问容器和对象	36	5.3.6 分页返回容器列表	63
5.1.5 swift CLI 命令清单	38	5.3.7 内容格式	64
5.2 存储服务的 HTTP API	39	5.3.8 获取容器的元数据	66
5.2.1 认证	40	5.3.9 删除容器	66
5.2.2 存储账号服务	42	5.3.10 创建对象	67
5.2.3 存储容器服务	47	5.3.11 分页返回对象列表	68
5.2.4 存储对象服务	54	5.3.12 下载、复制和删除对象	69
5.3 利用 curl 使用 Swift 存储服务	59	5.3.13 对象元数据	71
5.3.1 curl 的安装	59	5.4 总结	72

第 6 章 Swift 的高级特性 74

6.1 创建大对象	74	6.3 多版本对象	79
6.1.1 动态大对象	75	6.4 失效对象	82
6.1.2 静态大对象	76	6.5 客户元数据	82
6.1.3 静态和动态大对象的比较	77	6.6 总结	83
6.2 许可和访问控制表	78		

第 7 章 使用 Java 开发 Swift 应用 85

7.1 jclouds 简介	86	7.3 BlobStore API	91
7.1.1 jclouds 的特性	86	7.3.1 连接	91
7.1.2 BlobStore 简介	86	7.3.2 获取 BlobStore 接口	91
7.1.3 BlobStore 的核心概念	86	7.3.3 容器操作命令	91
7.2 jclouds-Swift 的安装	87	7.3.4 blob 操作命令	92
7.2.1 jclouds 简介	87	7.3.5 使用 BlobStore API	93
7.2.2 jclouds 安装	88	7.4 使用 Blob Store API 的高级功能	104

7.4.1 上传大型数据	104	7.5 Swift Client 接口	107
7.4.2 大型列表	106	7.5.1 SwiftClient 接口简介	107
7.4.3 目录标识	106	7.5.2 SwiftClient 接口使用	108
7.4.4 Content Disposition	106		

第 8 章 Swift 的实现原理 114

8.1 环 (Ring) 的实现原理	114	8.3.2 账号 (accounts) 目录	127
8.1.1 普通 Hash 算法与场景分析	115	8.3.3 容器 (containers) 目录	132
8.1.2 一致性哈希算法	116	8.3.4 临时 (tmp) 目录	134
8.1.3 虚节点 (Partition)	118	8.3.5 async_pending 目录	134
8.1.4 副本 (Replica)	120	8.3.6 隔离 (quarantined) 目录	137
8.1.5 分区 (Zone)	122	8.3.7 小结	138
8.1.6 权重 (Weight)	122	8.4 容器间同步的实现	138
8.1.7 小结	123	8.4.1 简介	138
8.2 环的数据结构	123	8.4.2 设置容器同步	138
8.3 存储节点的实现	124	8.4.3 容器同步的实现	140
8.3.1 对象 (objects) 目录	125	8.5 总结	142

第 9 章 Swift 的单机搭建 144

9.1 安装说明	145	9.3.7 生成相关 ring 以及 builder 文件	149
9.1.1 安装环境	145	9.4 安装存储节点	151
9.1.2 单机版 Swift 结构	145	9.4.1 安装存储服务相关包	151
9.2 环境准备	146	9.4.2 配置各个存储节点	152
9.2.1 系统要求	146	9.4.3 更改 rsyncd.conf 文件	158
9.2.2 更新配置操作系统	146	9.4.4 设置 rsyncd 文件	160
9.3 安装代理 (Proxy) 节点	148	9.4.5 建立存储点	160
9.3.1 创建 Swift 目录	148	9.5 安装成功验证	161
9.3.2 创建 swift.conf 文件	148	9.5.1 检测 Swift 运行状态	161
9.3.3 创建 Swift 服务	148	9.5.2 上传和列出文件	161
9.3.4 创建 SSL 自签名证书	148	9.5.3 下载文件	162
9.3.5 更改 memcached 监听地址	148	9.6 常见问题说明	162
9.3.6 创建代理节点配置文件	149		

第 10 章 Swift 的多机搭建 163

10.1 基本结构和术语	163	10.2.1 操作系统配置	165
10.2 安装环境准备	165	10.2.2 添加下载源	165

10.2.3	创建 Swift 用户	166	10.4.1	安装存储服务相关包	171
10.2.4	创建 Swift 的工作目录	166	10.4.2	存储点的设置	171
10.3	安装代理节点	167	10.4.3	创建 Swift 工作目录	172
10.3.1	安装代理节点 Proxy	167	10.4.4	复制配置文件	172
10.3.2	创建工作目录	167	10.4.5	创建/etc/rsyncd.conf	173
10.3.3	配置 memcached 监听默认端口	167	10.4.6	修改/etc/default/rsync	173
10.3.4	创建 swift.conf 文件	168	10.4.7	创建配置文件	173
10.3.5	创建 SSL 自签名证书	168	10.4.8	开启存储节点服务	175
10.3.6	创建代理节点配置文件	168	10.5	安装成功验证	176
10.3.7	构建创建 ring 的 builder 文件	169	10.5.1	检测 Swift 运行状态	176
10.3.8	添加 Zone 的命令	170	10.5.2	上传和列出文件	176
10.3.9	启动代理服务	171	10.5.3	下载文件	177
10.4	安装存储节点	171	10.6	常见问题说明	177

第 11 章 运行维护 Swift 集群 178

11.1	增加存储容量	179	11.3	处理硬件故障	184
11.1.1	Swift 安置数据的方法	179	11.3.1	处理有故障的磁盘驱动器	185
11.1.2	添加新磁盘的方法	179	11.3.2	处理写满的磁盘驱动器	185
11.1.3	平滑添加存储容量的方法	180	11.3.3	处理磁盘区域失效故障	185
11.1.4	添加新的存储节点	181	11.3.4	处理失去联系的节点故障	186
11.2	移出存储设备	182	11.3.5	处理故障节点	186
11.2.1	移出存储节点	182	11.4	观察和优化集群性能	187
11.2.2	移出存储磁盘	183	11.5	总结	187

主要内容：



- 云存储起源
- 云存储概念
- 云存储特点

本章目标：



- 了解什么是云存储
- 了解云存储的技术起源以及服务起源
- 了解云存储的典型特征

近年来，被业界称为第三次 IT 革命的“云计算”技术的快速发展，掀起了全球信息技术的变革浪潮。Google CEO 埃里克·施密特对“云计算”概念的解释是“云计算把数据分布在大量的分布式计算机上，从而使得存储获得很强大的可扩展能力。”在百度文库查询云计算的定义，可以看到这样一段文字“云计算就是这样一种变革——由谷歌、IBM 这样的专业网络公司来搭建计算机存储、运算中心，用户通过一根网线借助浏览器就可以很方便地访问，把‘云’作为资料存储以及应用服务的中心。”云计算的粗略概念如图 1.1 所示。

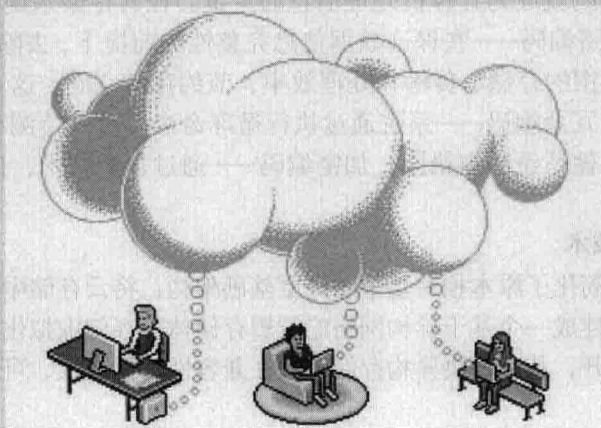


图 1.1 云计算

由“云计算”的定义中，我们看到另一个新的概念——云存储。云存储是云计算技术的重要组成部分，是云计算的重要应用之一。在云计算技术发展过程中，伴随着数据存储技术的云化发展历程。云计算的起源包含了技术模式和服务模式两个方面，云存储的起源同样也包含了技术和服务两个方面。

1.1 云存储起源

1.1.1 云存储技术起源

任何一项新技术的出现与发展，都有着与其密不可分的、推动其向前的背景技术。云存储技术的发展，同样源于集群技术、网格技术、分布式存储技术、虚拟化存储技术的发展。

1. 宽带网络的发展

随着互联网技术的不断提升、宽带网络建设速度的加快、大容量数据传输技术的实现和普及，传统的基于 PC 的存储技术将逐渐被云存储技术所替代。

2. Web 2.0 技术的出现

Web2.0 的出现改变了网络信息的传播方式。在 Web1.0 的时代，用户对网络信息的使用强调的是“获取”和“下载”；而在 Web2.0 时代，Web2.0 更多强调的是“分享”与“互动”。移动终端的广泛使用，令网络数据的应用方式更加灵活多样。

3. 分布式文件系统

传统的网络存储系统采用集中的存储服务器存放所有的数据信息，面对庞大的数据存储需求，存储服务器自身成为系统性能的瓶颈，也是可靠性和安全性的焦点问题。与其不同，分布式文件系统不是将数据信息放在一块磁盘上，由上层操作系统来管理，而是存放在一个由主控服务器（Master/NameNode）、数据服务器（ChunkServer/DataNode）和客户服务器组成的服务器集群上，文件的目录结构独立存储在一个主控服务器上，而具体的文件数据拆分成若干块，冗余地存放在不同的数据服务器上。集群中的服务器共同协作，提供整个文件系统的服务。

4. 数据编码技术

数据编码技术就是将需要加工处理的数据转化成代码或编码字符形式，以便于数据的可靠传输和迅速调用。对数据进行编码，可以方便地进行信息分类、检索、校对、计算等操作。因此，数据编码就成为计算机处理信息的关键。在云存储系统中，一般采用以下 3 种编码方式：数据压缩编码——在保证数据信息完整性的前提下，去除信息中的数据冗余，减少数据量，提高数据的存储、传输和处理效率，节约存储空间，这对大数据处理技术的应用有着重要作用；冗余编码——系统通过执行循环命令，可以检测和纠正数据在传输中发生的错误，提高存储系统的容错性；加密编码——通过加密技术，保证所存储数据的保密性和完整性。

5. 存储虚拟化技术

存储虚拟化技术简化了原本相对复杂的底层基础架构，将云存储中数量庞大、分布地域广泛的服务器集群构建成一个基于异构网络的逻辑存储体。存储虚拟化的思想是将资源的逻辑映像与物理存储分开，从而解决异构存储系统在兼容性、扩展性、可靠性、容错容灾等方面存在的问题。

1.1.2 云存储服务起源

云存储服务萌发于互联网早期的 E-mail 系统, 最早由 Hotmail 推出。随着 Web2.0 和宽带网络的发展, 各种云存储服务呈现爆炸式增长。从 Google 的 GFS 到 Amazon 的 S3, 从微软的 Azure Blob 到苹果的 iCloud, 从 Facebook 到 Twitter, 云存储服务无处不在。可以说, 不是每个互联网用户都知道云存储的存在, 但是几乎所有的互联网用户都已经站立在“云”端。

在传统存储模式下, 当我们使用某一个独立的存储设备时, 我们必须了解这个存储设备是什么型号, 什么接口和传输协议, 必须清楚地知道存储系统中有多少块磁盘, 分别是什么型号, 多大容量, 存储设备和服务器之间采用什么样的连接线缆。为了保证数据安全和业务的连续性, 我们还需要建立相应的数据备份系统和容灾系统。除此之外, 对存储设备进行定期的状态监控、维护、软硬件更新和升级也是必需的。而在云存储系统中的所有设备对使用者来讲都是完全透明的, 任何地方的任何一个经过授权的使用者都可以通过一根接入线缆与云存储连接, 对云存储进行数据访问。

如同云计算一样, 云存储对使用者来讲, 不是指某一个具体的设备, 而是指一个由许许多多存储设备和服务器所构成的集合体。使用者使用云存储, 并不是使用某一个存储设备, 而是使用整个云存储系统带来的一种数据访问服务。所以严格来讲, 云存储不是存储, 而是一种服务。

云存储服务模式的出现, 改变了长期以来人们对数据信息存取的方式。种类繁多的云存储服务更是让人们充分体验了现代信息技术带来的全新感受。从最初简单的邮件服务发展到现在, 以国际著名 5 大存储服务(苹果 iCloud、Google、亚马逊 Cloud Drive、Windows Live SkyDrive 和 Dropbox)为代表的云存储服务商们为用户提供了空间服务、类型文件服务、搜索引擎服务、音乐服务、离线支持服务等。云存储的核心是应用软件与存储设备相结合, 通过应用软件来实现存储设备向存储服务的转变。

云存储服务与传统存储相比较, 最大的区别在于人们摆脱了对本地设备的依赖。对存储空间、数据信息的获取和使用, 由原来的购买、下载保存变为通过一个简单的 Web 服务接口, 便可以在任何时间、任何地点存储和检索任意数量的数据, 获得高可用、高可靠的数据存储以及稳定廉价的基础存储设施。

1.2 云存储概念

通过上一节的学习, 我们不可避免地要提出一个问题: 什么是云存储?

云存储是伴随云计算衍生出来的新概念, 尚没有一个标准的定义。但是, 业界对于云存储已达成基本共识, 即云存储不仅是数据信息存储的新技术、新设备模型, 也是一种服务的创新模型。云存储是通过采用网格技术、分布式文件系统、服务器虚拟化、集群应用等技术, 将网络中海量的异构存储设备同构成可弹性扩展、低成本、低能耗的共享存储资源池, 并提供数据存储访问、处理功能的一个系统服务。

当云计算系统核心应用海量数据存储、访问管理时, 就需要配置大量的云存储设备, 云计算系统就转换成了一个云存储系统。所以, 从另外一个角度来讲, 云存储实际上也是一个以数据存储和管理为核心的云计算系统。

从技术上的角度看, 云存储基于网络, 利用分布式协同工作软件, 将用户数据分散存储

于若干通用存储服务器上，并通过副本或编码方法实现容错，向用户提供可靠的统一的逻辑存储空间；从业务模型角度看，云存储是指将用户数据通过网络存储在共享存储空间里，方便用户使用各种终端访问和共享。

云存储在服务架构方面，包含了云计算三层服务架构的技术体系。云存储服务在 IaaS 层为用户提供了数据存储、归档、备份的服务，在 PaaS 层为用户提供各种不同的类型文件及数据库服务。作为云存储在 SaaS 层的使用，涉及的内容就丰富和广泛得多了，包括我们熟悉的云网盘、照片的保存与共享、在线音乐、网络影院、在线备份、文档笔记的保存、在线游戏等。

1.3 云存储的特点

1. 低成本

传统的存储系统的架构主要是针对具体应用领域而采用专门、特定的硬件组件（服务器、磁盘阵列、控制器、系统接口）构成的架构，提供的服务类型比较单一，并且一般来讲是通过硬件来实现系统的可靠性和性能的。

云存储通常是通过大量的普通廉价主机构建成集群，甚至是跨地域的多个数据中心，可靠性和性能多是采用软件架构的方式来获取的。容灾机制开始就包含在架构体系设计和每一个开发环节中。与传统存储系统中的故障恢复机制不同，云存储系统的快速更换单位通常是一个存储主机，而不是单个 CPU、内存等内部硬件部件。当某个节点出现硬件故障时，管理人员只需将此节点替换为新的节点，数据就能自动得到恢复。所以，云存储可以大大降低企业级存储的成本，包括硬件设备购置成本、运维存储服务的成本、修复存储的成本以及管理存储的成本。

2. 服务模式

按需使用、按量付费是云存储的一大亮点。云存储实际上不仅仅是一个采用集群式的分布式架构，而且是通过硬件以及软件虚拟化而提供了一种存储服务。企业和个人不是通过购买和部署硬件设备来完成数据存储，而是通过购买服务把数据存储到云数据中心。

3. 可动态伸缩性

存储系统的动态伸缩性主要包含读/写性能和存储容量的扩展和缩减。随着业务量的增加，存储系统需要提高其读/写性能和存储容量来满足新的需求。有时候因为季节因素或者市场变化，为了节约成本，存储系统可以根据实际情况缩减其性能和容量。

传统的存储系统一般按照其型号有规定的硬件配置以及性能和容量确定的扩展功能，但是当业务需要超出系统的支持范围，就需要更新整套硬件设备来满足需求。

可动态伸缩性是云存储与传统存储系统相比的最大亮点之一。一个设计良好的云存储系统可以在系统运行过程中简单地通过添加或移除节点来自由扩展和缩减，并且这些操作对用户来说都是透明的。

4. 超大容量

云存储具备海量存储的特点，可以支持数十 PB 级的存储容量，高效地管理上百亿个文件，并且具有很好的线性可扩展性。

5. 高可靠性

传统的存储系统一般通过冗余磁盘阵列（RAID）来提供数据冗余技术。这种方法是通过

在一台高性能的主机上挂载多块磁盘形成一个阵列，然后通过数据镜像、数据分条和奇偶校验等技术，将一个文件或其数据切片存放到多块磁盘上形成冗余。

云存储系统通常是通过大量的普通廉价主机构建成集群，甚至是跨地域的多个数据中心，来提供并行读/写和冗余存储，从而达到高吞吐量和高可靠性。云存储系统从实际失效数据分析和建立统计模型着手，寻找软硬件失效规律，根据不间断的服务需求设计多种冗余编码模式，并据此在系统中构建具有不同容错能力、存取和重构性能等特性的功能区。通过负载、数据集和设备在功能区之间自动匹配和流动，实现系统内数据的最优化布局，并在站点之间提供全局精简配置和公用网络数据及带宽复用等高效容灾机制，从而提高系统的整体运行效率，满足高可靠性要求。

6. 高可用性

云存储服务可以为在不同时区的用户提供服务并保证 7×24 小时服务。云存储方案中包括多路径、控制器、不同光纤网、端到端的架构控制/监控和成熟的变更管理过程，从而大大提高了云存储的可用性。另外，云存储服务按照 CAP 理论在不影响应用使用正确性的前提下，通过适当放松对数据一致性的要求来提高数据的可用性。

7. 安全性

与传统的存储系统相比，云存储对于一个企业和个人来讲已经没有一个物理边界。所有云存储服务间传输以及保存的数据都有潜在被截取或篡改的隐患。云存储服务都需要在数据传输过程中，以及在云服务中心保存时采用加密技术来限制对数据的访问。另外，云存储系统还采用数据分片混淆存储作为实现用户数据私密性的一种方案。

8. 规范化

2010年4月 SNIA 公布了云存储标准——CDMI 规范，其提供了数据中心利用云存储的方式。尽管 SNIA 号称可以使大多数非云存储产品访问方式演进成云存储访问，但是 CDMI 并没有提供可靠性和质量来衡量云存储服务提供商质量的方式。并且，业界并没有大量采用 CDMI 规范。市场现有的云存储服务平台，包含 Amazon S3、Google Drive、Microsoft Azure 都是采用了自己的私有接口规范。因此，云存储数据管理的规范化工作还需进一步努力。

习题

- 1.1 什么是分布式文件系统？
- 1.2 什么是云存储？
- 1.3 云存储与传统存储方式的主要差异有哪些？
- 1.4 描述云存储的主要特点。

主要内容：



- 非结构化数据存储
- 对象存储系统

本章目标：



- 了解非结构化数据的存储要求
- 了解对象存储产生原因
- 了解对象存储系统的特点及价值
- 了解对象存储的应用场景

在本章，我们将首先介绍“结构化数据”和“非结构化数据”的概念及其特点，然后介绍 3 种不同类型的存储系统：块存储系统、文件存储系统和对象存储系统，最后描述什么是对象存储系统。

2.1 非结构化数据存储

2.1.1 什么是非结构化数据

如图 2.1 所示，在信息社会，信息可以划分为 3 大类。一类信息能够用数据或统一的结构加以表示，我们称之为结构化数据，如数字、符号；另一类信息无法用数字或统一的结构表示，如文本、图像、声音、网页等，我们称之为非结构化数据。结构化数据属于非结构化数据，是非结构化数据的特例；而所谓半结构化数据，就是介于完全结构化数据和完全无结构的数据之间的数据，HTML 文档就属于半结构化数据。它一般是自描述的，数据的结构和内容混在一起，没有明显的区分。



图 2.1 结构化数据与非结构化数据对比

结构化数据就是存储在传统数据库里的数据，也就是说可以用二维表结构来逻辑表达的数据。结合到典型场景中更容易理解，比如企业 ERP、财务系统、医疗数据库、教育一卡通、政府行政审批、其他核心数据库等。这些应用对存储方案的需求包括高速存储应用需求、数据备份需求、数据共享需求以及数据容灾需求。

相对于结构化数据而言，不方便用数据库二维逻辑表来表现的数据即称为非结构化数据，包括所有格式的办公文档、文本、图片、XML、HTML、各类报表、图像和音频/视频信息等。也就是说，数据没有数据模型，以文件形式存储，而不是存放在数据库系统。具体到典型案例中，比如医疗影像系统、教育视频点播、视频监控、国土 GIS、设计系统、文件服务器（PDM/FTP）、媒体资源管理等具体应用，这些行业对于存储的需求包括数据存储、数据备份以及数据共享等。

半结构化数据，包括邮件、HTML、报表、资源库等，典型场景如邮件系统、WEB 集群、教学资源库、数据挖掘系统、档案系统等。这些应用有数据存储、数据备份、数据共享以及数据归档等基本存储需求。

2.1.2 非结构化数据的存储要求

随着网络技术的发展，特别是 Internet 和 Intranet 技术的飞快发展，非结构化数据的数量日趋增大。这些非结构化数据的绝大部分是来自不断扩散的照片、录像、电子信、文档、IM 等，用户产生和使用的数据也比任何时候都多。据统计，2012 年底 Facebook 的用户每天上传 3.5 亿张新照片，一个月将产生 7 个 PB 数据。IDC 的一项调查报告指出：企业中 80% 的数据都是非结构化数据，这些数据每年都按指数增长 60%。报道指出：平均只有 1%~5% 的数据是结构化的数据。如今，这种迅猛增长的从不使用的数据在企业里消耗着复杂而昂贵的一级存储的存储容量。如何更好地保留那些在全球范围内具有潜在价值的不同类型的文件，而不是因为处理它们却干扰日常的工作？这时，主要用于管理结构化数据的关系数据库的局限性暴露得越来越明显。

非结构化数据存储需要保证持续性、可访问性、低成本以及可管理性。

持续性：尽管你可能再也不会去看你过去度假的照片，但是为用户提供照片存储服务的

公司，比如，Flickr，却需要永久存放它们。实际上，用户期望或者法律规则使得大多数非结构化数据需要永久性存储。

可访问性：非结构化数据还需要能够通过各种设备主要是移动手机和浏览器，实现即时访问。尽管有些数据可以存档，但是用户还是期待他们的大部分数据能够立即使用。

低成本：非结构化数据需要低成本的存储。如果有足够的钱，任何存储问题都可以解决，但现实生活并非如此，有限的预算需要低成本的存储。

可管理性：超大型在线数据存储系统的可管理性是非常关键的。为了使得数据中心的数据管理变得简单，我们需要将数据控制和数据存储分离，从而大量减少系统管理工作。

2.1.3 存储系统的种类

不同类型的数据具有不同的访问模式，需要使用不同类型的存储系统。总体来讲有 3 大类存储系统：块存储系统、文件存储系统和对象存储系统。

块存储系统直接访问原始的未格式化的磁盘。这种存储的特点是速度快，空间使用率高。块存储多用于数据库系统，它可以使用未格式化的磁盘对结构化数据进行高效读写。而数据库最适合存放的是结构化数据。

文件存储是最常用的存储系统。文件存储系统使用格式化的硬盘，为用户提供了文件系统的使用界面。当你在计算机上打开和关闭文档的时候，所看到的就是文件系统。尽管文件系统在磁盘上提供了一层有用的抽象，但是它并不适合于管理大量的数据，或者超量使用文件中的部分数据。

对象存储可能是大家最不熟悉的存储系统。对象存储不提供对未格式化数据模块的访问，也不提供基于文件系统的访问，而是提供对整个对象（Object）的访问。一般来讲，通过特定的 API 对其进行访问。对象存储的优势在于它可以存放无限增长的内容，最适合用来存储包含备份、存档、静态 web 页面、视频、照片等非结构化或半结构化的数据。除此之外，对象存储还具有低成本、高可靠的优点。

2.1.4 传统的共享存储方法的缺点

首先我们来简单分析一下现有的共享存储系统。

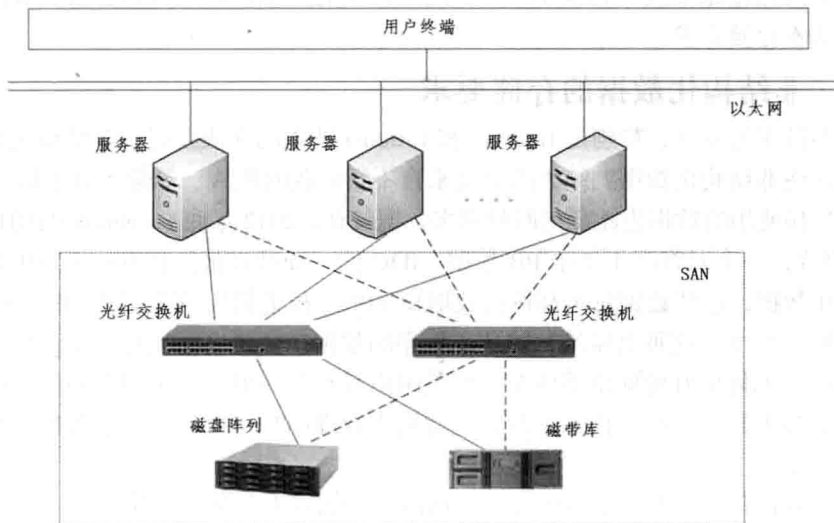


图 2.2 SAN 结构