

 **Hadoop**
核心技术剖析与云端实现

Cloud
Computing

深入浅出

[云计算]



循序渐进讲述云计算的基本概念

基于Hadoop开源云计算平台，讲解如何构建一个基于云计算的应用系统

以Hadoop源代码为对象·深度剖析云计算核心技术

云技术与云端实现，讲述基于Hadoop云计算平台的4个高级应用框架，以及解决实际问题的思路与方法

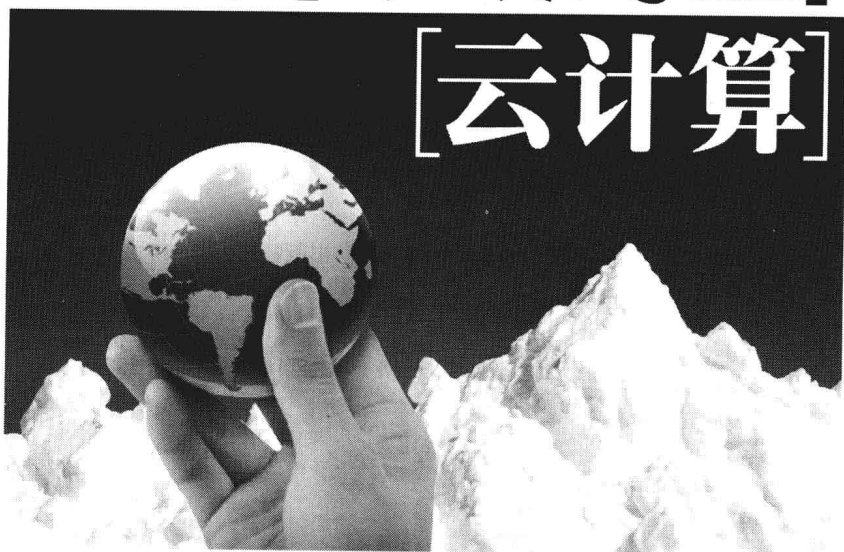
鲍亮 陈荣 编著

清华大学出版社

Cloud
Computing

深入浅出

[云计算]



鲍亮 陈荣 编著

清华大学出版社

北京

内 容 简 介

本书作者以多年实际研发项目为背景,通过项目实战与代码分析,深入浅出地讲述云计算的基本概念,云计算的核心技术细节以及使用云计算平台解决实际问题的思路与方法。

全书共分4篇。第1篇循序渐进地介绍云计算的基本概念,学习云计算需要掌握的基本知识和云计算环境搭建方法;第2篇基于Hadoop开源云计算平台,讲解如何构建一个基于云计算的应用系统,了解云计算应用系统的设计方法;第3篇以开源的Hadoop云计算平台为分析对象,在源代码层次上对分布式文件系统、MapReduce计算模型、NoSQL数据库和集群管理算法与技术等云计算核心技术进行深度剖析;第4篇为云计算应用篇,介绍了基于Hadoop云计算平台的4个高级应用框架,读者可以结合自己的应用需求与场景,使用这些框架解决实际问题。

本书理论联系实际,既有理论深度又有实用价值,可作为高校教材使用,也可作为云计算研发人员以及爱好者的学习和参考手册。

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

深入浅出云计算/鲍亮,陈荣编著. —北京:清华大学出版社,2012.10
ISBN 978-7-302-30238-4

I. ①深… II. ①鲍… ②陈… III. ①计算机网络—基本知识 IV. ①TP393

中国版本图书馆CIP数据核字(2012)第228217号

责任编辑:夏非彼
封面设计:王翔
责任校对:闫秀华
责任印制:何芊

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦A座 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:北京世知印务有限公司

装 订 者:三河市新茂装订有限公司

经 销:全国新华书店

开 本:190mm×260mm 印 张:28.25 字 数:723千字

版 次:2012年10月第1版 印 次:2012年10月第1次印刷

印 数:1~4000

定 价:59.00元

前 言

“云计算”这一概念最早由 IBM 和 Google 于 2007 年底提出，其本质是一种大规模的、由经济性所驱动的新一代分布式计算技术。采用云计算技术，可以将一组抽象的、虚拟的、动态可伸缩的以及受控的计算能力、存储、平台和服务通过互联网按需发布给海量用户。与传统的分布式计算技术（如集群、网格、P2P 等）相比，云计算的核心特征有其灵活性与经济性：第一，将基础设施架构在大规模的廉价服务器集群之上，通过多个服务器之间的冗余以及应用程序与底层服务协作开发等机制，使得软件具有高可用性，并最大限度地利用资源；第二，采用面向市场的业务模型，为用户提供诸如计算能力、存储和网络传输等基础服务，从中收取使用费用。

基于新型的系统架构，云计算实现了应用系统的可扩展性和高可用性；基于新的业务模型，云计算能够满足企业降低成本、提高工作效率以及简化 IT 管理的需求。目前，云计算技术已经广泛应用于电子商务、电信服务、 workflow 管理与客户关系管理等多个行业中。由于其突出的经济性、性能与可靠性优势，云计算自出现至今短短 5 年的时间，已得到学术界和产业界的广泛关注。

目前市场上与云计算技术相关的书籍较多，但经过作者调研，目前尚无以开源的云计算平台为对象，在源代码层次上对分布式文件系统、MapReduce 计算模型、NoSQL 数据库和集群管理算法与技术等云计算核心技术进行深度剖析的书籍。本书创作团队核心成员自 2005 年起就一直从事分布式计算、SOA 等相关方向的研发工作，具有丰富的项目实践经验。2009 年初，作者以实际研发项目为背景，组织团队分析、理解著名的开源云计算平台 Hadoop 的源代码，经过两年多的努力，其部分成果形成本书最为核心的第 9~12 章。通过项目实战与代码分析，我们积累了大量云计算平台的使用经验，对云计算的核心技术有了较为深刻的理解，认为有必要将自己的经验和认识整理出来，以满足广大读者希望一方面掌握云计算基本概念，能够使用云计算平台解决实际问题；另外一方面希望深入理解云计算核心技术细节的迫切心情，这也正是书名《深入浅出云计算》的由来。

本书适合不同层次的读者阅读。建议读者根据自己的兴趣和目的有选择性地阅读：希望了解云计算的基本概念、本质和发展趋势的读者，可以重点阅读第 1、2 章和附录；对于云计算的初学者，可以重点阅读 1~8 章；已经掌握云计算基本概念，具有一定实践基础，想进一步深入学习、研究云计算核心技术的读者，可以重点阅读 9~12 章；想利用云计算解决各类实际计算问题的读者，可以重点阅读 13~16 章。

感谢书籍创作团队核心成员王玉操、贾世达、王磊、葛军、李朝印、张翔、李江、杨阳和彭恒的辛勤努力。感谢我的导师陈平教授在云计算研究方面对我们的悉心指导。

由于云计算起源于企业界，是企业界推动学术界发展的一项新兴应用技术，相关资料数量不多，加之作者水平有限，时间紧迫，因此书中难免存在错误与不当，敬请读者批评指正。建议和意见请发至作者邮箱 baoliang@mail.xidian.edu.cn。

编 者

目 录

第 1 篇 初始云计算

第 1 章 云计算介绍	3
1.1 云计算相关概念	3
1.1.1 云计算的定义	4
1.1.2 云计算的服务方式	6
1.1.3 云计算的部署模式	8
1.2 云计算的历史	9
1.2.1 虚拟化技术的发展	9
1.2.2 分布式计算技术的发展	10
1.2.3 软件应用模式的发展	13
1.3 云计算的现状	14
1.3.1 产业界现状	14
1.3.2 学术界现状	19
1.3.3 政府机构现状	25
1.4 本章小结	28
第 2 章 云计算技术基础	29
2.1 HDFS 相关技术	29
2.1.1 RPC	29
2.1.2 基于 Socket 的 Java 网络编程	30
2.2 MapReduce 相关技术	31
2.2.1 Java 反射机制	31
2.2.2 序列化和反序列化	33
2.3 HBase 相关技术	34
2.3.1 NoSQL	34
2.3.2 ACID	36
2.3.3 CAP 理论	36
2.3.4 一致性模型	37
2.4 ZooKeeper 相关技术	38
2.4.1 Paxos 算法介绍	38

2.4.2	Java NIO 库.....	38
2.5	本章小结.....	39
第 3 章	云计算开发环境搭建.....	40
3.1	集群环境介绍.....	40
3.2	Hadoop 环境搭建.....	41
3.2.1	Hadoop 简介.....	41
3.2.2	安装前准备.....	41
3.2.3	安装环境搭建.....	42
3.2.4	详细安装步骤.....	42
3.3	Hadoop 集群配置.....	49
3.3.1	配置 Hadoop 守护进程的运行环境.....	49
3.3.2	配置 Hadoop 守护进程的运行参数.....	50
3.4	HBase 环境搭建.....	51
3.4.1	HBase 简介.....	51
3.4.2	HBase 的数据模型.....	51
3.4.3	HBase 安装前的准备.....	52
3.4.4	HBase 的安装配置.....	52
3.4.5	HBase 的运行.....	55
3.5	ZooKeeper 环境搭建.....	57
3.5.1	ZooKeeper 简介.....	57
3.5.2	安装前的准备.....	58
3.5.3	独立服务器的安装与配置.....	59
3.5.4	集群服务器的安装与配置.....	60
3.6	本章小结.....	62

第 2 篇 浅出云计算

第 4 章	应用实例：图像百科系统.....	65
4.1	应用背景.....	65
4.2	需求分析.....	66
4.2.1	功能需求.....	66
4.2.2	非功能需求.....	69
4.3	核心业务处理流程.....	70
4.3.1	查询百科条目处理流程.....	70
4.3.2	编辑百科条目处理流程.....	72
4.3.3	更新百科条目处理流程.....	74

4.4	总体设计.....	75
4.5	本章小结.....	78
第 5 章	使用 HDFS 存储海量图像数据.....	79
5.1	HDFS 介绍.....	79
5.1.1	HDFS 架构.....	79
5.1.2	HDFS 的特点.....	80
5.1.3	HDFS 存取机制简介.....	81
5.2	HDFS 接口介绍.....	83
5.3	图像百科系统中的图像存储.....	87
5.3.1	图像存储基本思想.....	87
5.3.2	图像存储设计目标.....	88
5.3.3	图像存储体系结构.....	88
5.3.4	图像百科系统的功能结构.....	89
5.4	系统实现.....	90
5.4.1	存储模块类交互图.....	90
5.4.2	核心类详细介绍.....	92
5.4.3	HDFS 存储小文件.....	97
5.5	本章小结.....	98
第 6 章	使用 MapReduce 处理图像.....	99
6.1	分布式数据处理 MapReduce.....	99
6.1.1	MapReduce 简介.....	99
6.1.2	编程模型.....	100
6.1.3	执行概括.....	101
6.2	使用 MapReduce 编程模型.....	102
6.2.1	MapReduce 程序模板.....	102
6.2.2	MapReduce 编程思想.....	107
6.3	更新图像百科条目的 MapReduce 设计.....	107
6.3.1	设计目标.....	107
6.3.2	更新条目的体系结构.....	109
6.3.3	更新条目的逻辑流程.....	110
6.4	MapReduce 对更新条目的实现.....	112
6.4.1	更新条目的核心类.....	112
6.4.2	MapReduce 核心类实现.....	113
6.4.3	编译运行.....	118
6.5	本章小结.....	121

第7章 使用 HBase 存储百科数据	122
7.1 HBase 的基本特征	122
7.1.1 RDBMS 与 HBase	122
7.1.2 面向列的 NoSQL 数据库	123
7.1.3 HBase 数据库架构	126
7.1.4 HBase 的特点	128
7.2 使用 HBase 编程	129
7.2.1 HBase 的 Java API	129
7.2.2 HBase 客户端编程	130
7.2.3 HBase 编程示例	150
7.3 Fotospedia 系统的数据库设计	153
7.3.1 数据库模块总体设计	154
7.3.2 数据库模块详细设计	154
7.3.3 数据库模块交互设计	158
7.4 Fotospedia 系统的数据库实现	160
7.4.1 数据库模块类交互图	160
7.4.2 数据库模块核心类实现	161
7.5 本章小结	167
第8章 使用 ZooKeeper 管理集群	168
8.1 ZooKeeper 详细介绍	168
8.2 ZooKeeper 的使用方法及 API 介绍	172
8.2.1 ZooKeeper 的使用方法	172
8.2.2 基本类和接口	173
8.2.3 常用类与方法的实例介绍	173
8.3 图像百科系统集群管理详细设计	179
8.3.1 集群管理	179
8.3.2 配置管理	181
8.4 图像百科系统集群管理实现	182
8.4.1 集群管理实现	182
8.4.2 配置管理实现	188
8.4.3 测试	194
8.5 本章小结	197

第 3 篇 深入云计算

第 9 章 深入分析 HDFS	201
9.1 HDFS 核心设计机制.....	201
9.1.1 Namenode 和 Datanode	201
9.1.2 数据副本策略.....	201
9.1.3 数据组织.....	204
9.1.4 健壮性.....	204
9.1.5 存储空间回收.....	205
9.2 HDFS 源码总体介绍.....	206
9.3 核心代码分析.....	208
9.3.1 HDFS 的通信协议.....	208
9.3.2 HDFS 读文件源码分析.....	214
9.3.3 HDFS 写文件源码分析.....	219
9.4 Hadoop 支持的其他文件系统	222
9.4.1 KFS 文件系统体系架构	223
9.4.2 KFS 各模块关键技术	224
9.4.3 HDFS 与 KFS 写数据的区别	225
9.5 本章小结.....	227
第 10 章 深入分析 MapReduce	228
10.1 MapReduce 框架结构	228
10.1.1 MapReduce 中的角色	228
10.1.2 MapReduce 流程	230
10.2 代码静态分析.....	233
10.2.1 创建 Job 的相关类.....	233
10.2.2 初始化 Job 的相关类	234
10.2.3 作业调度相关类.....	234
10.2.4 执行 MapTask 的相关类.....	235
10.3 代码详细分析.....	236
10.3.1 JobClient 提交 Job.....	236
10.3.2 JobTracker 初始化作业.....	237
10.3.3 TaskTracker 启动.....	240
10.3.4 JobTracker 调度作业.....	242
10.3.5 TaskTracker 加载 Task.....	245
10.3.6 子进程执行 MapTask.....	247

10.3.7	子进程执行 ReduceTask.....	251
10.4	本章小结.....	254
第 11 章	深入分析 HBase	255
11.1	HBase 体系与原理.....	255
11.1.1	HBase 的集群架构.....	255
11.1.2	HBase 的系统架构.....	258
11.1.3	HBase 的存储架构.....	259
11.2	HBase 总体结构.....	264
11.2.1	总体包图.....	265
11.2.2	常用类分析.....	266
11.3	HBase 关键剖析.....	269
11.3.1	集群启动与关闭.....	269
11.3.2	HBase 配置过程.....	279
11.3.3	读取图像百科数据.....	282
11.3.4	写入图像百科数据.....	289
11.4	本章小结.....	294
第 12 章	深入分析 ZooKeeper	295
12.1	概述.....	295
12.1.1	ZooKeeper 角色.....	295
12.1.2	ZooKeeper 工作原理.....	295
12.2	代码静态分析.....	297
12.2.1	包概述.....	297
12.2.2	核心类浅析.....	298
12.3	代码情景分析.....	302
12.3.1	服务器的启动.....	302
12.3.2	Leader 服务器.....	311
12.3.3	Follower 服务器.....	318
12.3.4	客户端服务请求.....	320
12.4	本章小结.....	325

第 4 篇 应用云计算

第 13 章	应用 Pig 实现并行数据处理	329
13.1	Apache Pig 简介.....	329
13.2	Pig 的安装与配置.....	330
13.2.1	Pig 安装准备.....	330

13.2.2	安装配置过程.....	331
13.2.3	运行模式.....	331
13.3	深入分析 Pig	335
13.3.1	Pig 数据模型	335
13.3.2	Pig 常用命令和数据读写操作	337
13.3.3	Pig 诊断操作	338
13.3.4	Pig 关系操作	339
13.3.5	Pig 表达式和函数	340
13.3.6	Pig 用户自定义函数 (UDF)	342
13.3.7	探索逻辑执行计划.....	346
13.4	Pig 实例分析.....	347
13.4.1	Pig Latin 示例	347
13.4.2	简单实例解析.....	348
13.4.3	深入使用 Pig	352
13.5	Pig 与 SQL 比较	355
13.6	本章小结.....	356
第 14 章	应用 Hive 构建数据处理平台	357
14.1	Hive 简介	357
14.1.1	Hive 架构	357
14.1.2	Hive 和 Hadoop 关系.....	358
14.1.3	Hive 和传统数据库进行比较	359
14.1.4	Hive 的数据存储	360
14.1.5	Hive 元数据 Metastore	361
14.2	Hive 安装配置	363
14.2.1	安装前准备.....	363
14.2.2	安装 Hive.....	363
14.2.3	安装 MySQL 与 Hive 配置	365
14.3	Hive 使用与操作	369
14.3.1	Hive 基本操作	369
14.3.2	查询数据 Hive Select	375
14.3.3	Hive 函数	380
14.4	实例介绍.....	384
14.5	本章小结.....	389
第 15 章	应用 Mahout 实现机器学习算法	390
15.1	Mahout 概述	390
15.1.1	Mahout 简介	390

15.1.2	机器学习简介.....	390
15.2	Mahout 安装配置.....	392
15.2.1	安装前准备.....	392
15.2.2	Mahout 安装.....	393
15.3	Mahout 使用简介.....	394
15.3.1	使用 Mahout 实现集群.....	395
15.3.2	使用 Mahout 实现分类.....	406
15.3.3	使用 Mahout 实现决策树.....	408
15.3.4	使用 Mahout 实现推荐挖掘.....	409
15.4	本章小结.....	411
第 16 章	应用 HAMA 实现分布式计算.....	412
16.1	HAMA 简介.....	412
16.1.1	HAMA 系统架构.....	412
16.1.2	BSPMaster.....	413
16.1.3	GroomServer.....	414
16.1.4	ZooKeeper.....	414
16.2	HAMA BSP 介绍.....	415
16.2.1	BSP 并行计算.....	416
16.2.2	创建自定义的 BSP.....	417
16.2.3	用户接口.....	418
16.3	HAMA 安装配置.....	421
16.3.1	安装前准备.....	421
16.3.2	安装和环境配置.....	421
16.3.3	HAMA 运行模式.....	423
16.3.4	运行 HAMA.....	425
16.3.5	HAMA Web 接口.....	426
16.3.6	在 Eclipse 中创建 HAMA 工程.....	427
16.4	实例介绍.....	429
16.4.1	打印“Hello BSP”.....	429
16.4.2	估算 PI 值.....	430
16.5	本章小结.....	432
附 录	433

第 1 篇

初始云计算

本篇是入门篇，主要介绍云计算的基本概念、学习云计算需要掌握的基本知识和云计算环境搭建方法 3 个部分的内容。通过本篇的学习，读者可以了解到云计算的相关概念，云计算的历史以及云计算在产业界、学术界和政府机构等领域的发展现状，掌握 Java 网络编程、序列化与反序列化、数据库基础理论和分布式算法等基础知识，并能够搭建一个小型的云计算环境。

第 1 章 云计算介绍

自从 2007 年 10 月份云计算诞生至今，这一技术在短短的 4 年时间里对整个 IT 行业产生了巨大的影响。学术界、产业界和政府都对云计算产生了浓厚的兴趣：全球范围内讨论云计算技术的学术活动如火如荼；谷歌、亚马逊、IBM、微软等 IT 巨头大力推动云计算技术的宣传和产品的普及；各国政府纷纷斥巨资打造大规模的数据中心与计算中心。云计算技术目前已经得到了业界的高度认同，逐渐走向成熟。那么，云计算该如何定义？云计算的发展历史如何？云计算目前的现状如何？本章将重点解答这些问题，以帮助读者初步理解云计算。

1.1 云计算相关概念

云计算（Cloud Computing）这一概念于 2007 年 10 月 8 日正式出现，其标志性事件是谷歌和 IBM 宣布联合加入云计算的研究工作，并给出云计算的定义。同年 11 月 15 日，IBM 上海和阿莫科（Armonk, NY）同时发布了 Blue Cloud，Blue Cloud 是一系列的云计算产品，使得共同的数据中心像互联网一样运作。

2007 年 10 月 8 日，纽约时报报道了谷歌和 IBM 联合加入云计算研究的新闻，如图 1.1 所示。



图 1.1 谷歌和 IBM 联合加入云计算研究的新闻

1.1.1 云计算的定义

自云计算这一概念诞生至今，尚未形成业界广泛认可的统一定义。本书将列举 4 种有代表性的云计算定义，并对每种定义方式进行解读。

1. IBM 的云计算定义

2007 年 10 月，IBM 的 Greg Boss 等人以技术白皮书^①的形式给出了云计算的定义：“云计算是同时描述一个系统平台或者一类应用程序的术语。云计算平台按需进行动态部署（Provision）、配置（Configuration）、重新配置（Reconfigure）以及取消服务（Deprovision）等。在云计算平台中的服务器可以是物理或虚拟的服务器。高级的云计算通常包含一些其他的计算资源，如存储区域网络（SANs）、网络设备、防火墙以及其他安全设备等。

在应用方面，云计算描述了一类可以通过互联网进行访问的可扩展应用程序。这类云应用基于大规模数据中心及高性能服务器来运行网络应用程序与 Web 服务。用户可以通过合适的互联网接入设备及标准的浏览器访问云计算应用程序。”

IBM 技术白皮书中的云计算定义原文：

Cloud computing is a term used to describe both a platform and type of application. A cloud computing platform dynamically provisions, configures, reconfigures, and deprovisions servers as needed. Servers in the cloud can be physical machines or virtual machines. Advanced clouds typically include other computing resources such as storage area networks (SANs), network equipment, firewall and other security devices.

Cloud computing also describes applications that are extended to be accessible through the Internet. These cloud applications use large data centers and powerful servers that host Web applications and Web services. Anyone with a suitable Internet connection and a standard browser can access a cloud application.

IBM 的定义明确指出云计算概念的内涵包含两个方面：平台和应用。平台即基础设施，其地位相当于 PC 机上的操作系统，云计算应用程序需要构建在平台之上；云计算应用所需的计算与存储通常在“云端”完成，客户端需要通过互联网访问计算与存储能力。

2. 加州大学伯克利分校的云计算定义

2009 年 2 月 10 日，加州大学伯克利分校电子工程和计算机学院的 Michael Armbrust 等人发布技术报告《Above the Clouds: A Berkeley View of Cloud Computing》^②，介绍了对云计算

^① IBM 云计算技术白皮书.http://download.boulder.ibm.com/ibmdl/pub/software/dw/wes/hipods/Cloud_computing_wp_final_8Oct.pdf

^② 加州大学伯克利分校的云计算定义. <http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.pdf>