

深入云计算

Hadoop

应用开发实战详解

◎ 万川梅 谢正兰 编著

书中源代码下载地址：
<http://www.tdpress.com/zyzx/tssclfjw>

修订版

全面升级——
重装上市！

精准的内容梯度安排

遵循读者学习习惯和Hadoop技术应用实践，合理安排图书内容；精挑细选经典实例，嵌入完善的代码注释。

精炼的实用经验阐述

作者多年开发经验融入其中，读者在全面掌握Hadoop编程和开发技术的同时更能获得快速分析和解决实际问题的能力。

中国铁道出版社

CHINA RAILWAY PUBLISHING HOUSE

深入云计算

Hadoop

应用开发实战详解

◎ 万川梅 谢正兰 编著



内 容 简 介

本书由浅入深，全面、系统地介绍了 Hadoop 这一高性能处理大量数据集的理想工具。本书内容主要包括 HDFS、MapReduce、Hive、HBase、Mahout、Pig、ZooKeeper、Avro、Chukwa 等与 Hadoop 相关的子项目，各个知识点都配有精心设计的大量经典的小案例，实战性和可操作性很强。

本书旨在帮助云计算初学者迅速掌握 Hadoop 系统，提升读者在云计算实践中的应用和开发能力。同时本书极强的系统性和大量翔实的案例对于有一定基础的中高级用户有非常好的参考价值。

图书在版编目 (CIP) 数据

Hadoop 应用开发实战详解 / 万川梅, 谢正兰编著

.—2 版 (修订本). — 北京 : 中国铁道出版社,

2014. 8

(深入云计算)

ISBN 978-7-113-18625-8

I . ①H… II . ①万… ②谢… III . ①数据处理软件

IV . ①TP274

中国版本图书馆 CIP 数据核字 (2014) 第 115307 号

书 名: 深入云计算: Hadoop 应用开发实战详解 (修订版)

作 者: 万川梅 谢正兰 编著

责任编辑: 荆 波

读者服务热线: 010-63560056

责任印制: 赵星辰

特邀编辑: 赵树刚

封面设计: 多宝格·付 魏

出版发行: 中国铁道出版社 (北京市西城区右安门西街 8 号 邮政编码: 100054)

印 刷: 三河市宏盛印务有限公司

版 次: 2013 年 6 月第 1 版 2014 年 8 月第 2 版 2014 年 8 月第 2 次印刷

开 本: 787mm×1 092mm 1/16 印张: 25.75 字数: 603 千

书 号: ISBN 978-7-113-18625-8

定 价: 59.80 元

版权所有 侵权必究

凡购买铁道版图书, 如有印制质量问题, 请与本社读者服务部联系调换。电话: (010) 51873174

打击盗版举报电话: (010) 51873659

为什么需要云计算？可以通过一个简单的案例说明：一台计算机处理一批数据需要 30 小时，比如处理地震预测、天气预报的数据，这样的计算速度实在是太慢了。提升单台计算机的速度是过去的办法，但 CPU 的速度不可能再大幅度提升了。人们一直希望通过增加计算机数量并行运算提升运算和数据处理速度，例如，希望通过同时在 300 台计算机上处理数据，让处理这批数据的速度变成 10 分钟，并且看起来是在一台计算机上处理的。当然，这是一种理想状态。实际上，人们已经开始设计这样的分布式系统，通过把众多的计算机通过集群方式并行同时运行，以此来提高处理的速度。这就是云计算的源起。

在开源云计算系统中，Hadoop 稳居第一。事实上，Hadoop 非常受欢迎，全球已经安装了数以万计的 Hadoop 系统。在诸多的云计算技术中，Hadoop 具有无与伦比的优势，越来越多的公司和组织选择使用 Hadoop 开源项目作为其解决方案。

Hadoop 是 Apache 基金会的开源项目，为开发者提供了一个分布式系统的基础架构，用户可以在不了解分布式系统的底层细节的情况下开发分布式的应用，充分利用集群的强大功能，实现高速运算和存储。Hadoop 项目中包括一个分布式的文件系统 HDFS，一个分布式的并行编程框架 MapReduce，以及包括 Hive、HBase、Mahout、Pig、ZooKeeper、Avro、Chukwa 在内的诸多子项目。

本书的特点

1. 结构合理，内容系统全面；叙述翔实，例程丰富

在内容的安排上，根据读者的学习习惯和内容的梯度合理安排，更加适合读者学习。同时，本书有详细的例程，每个例子都经过精挑细选，有很强的针对性。书中的程序都有完整的代码（读者可在 <http://tdpress.com/zxyz/tsscflwj> 上的下载专区中下载使用），而且代码非常简洁和高效，便于读者学习和调试，读者也可以直接使用这些代码来解决自己的问题。

2. 基础知识和实践并重

本书不仅注重基础知识，而且非常注重实践，让读者快速上手，迅速掌握 Hadoop 知识。

3. 结合实际，实战项目贯穿其中

本书写作时特意给出了大量的实战项目，这些项目的灵活使用，将会让读者事半功倍。同时，为了便于读者高效、直观地学习本书内容，对每章的内容的学习都特意编写了思考与总结。

本书内容体系

全书总共 16 章，分为 3 篇，各篇对应的章节和具体内容如下。

第 1 篇 Hadoop 技术篇，包括 1~10 章，主要介绍了初识 Hadoop、Hadoop 的安装和配置、HDFS 海量存储、初识 MapReduce、分布式开源数据库 HBase、MapReduce 进阶、Hive 数据仓库、Pig 开发应用、Chukwa 数据收集系统、ZooKeeper 开发应用。

第 2 篇 Hadoop 管理和容错篇，包括 11~12 章，主要介绍了 Hadoop 管理和 Hadoop 容错。

第 3 篇 Hadoop 实战篇，包括 13~16 章，主要介绍了综合实战 1：Hadoop 中的数据访问；综合实战 2：一个简单的分布式的 Grep；综合实战 3：打造一个搜索引擎；综合实战 4：移动通信信令监测与查询。

本书读者对象

- Hadoop 初学者。
- 想全面、系统地学习 Hadoop 的人员。
- Hadoop 技术爱好者。
- 利用 Hadoop 进行编程和开发的技术人员。
- 大中专院校的学生和老师。
- 相关培训机构的学员。

由于编者水平有限，加之时间较紧，书中难免存在疏漏和错误之处，敬请读者批评指正，编者也会定期在 www.rzchina.net 进行答疑解惑。

编 者

2014 年 3 月

卖者意见反馈表

亲爱的读者：

感谢您对中国铁道出版社的支持，您的建议是我们不断改进工作的信息来源，您的需求是我们不断开拓创新的基础。为了更好地服务读者，出版更多的精品图书，希望您能在百忙之中抽出时间填写这份意见反馈表发给我们。随书纸制表格请在填好后剪下寄到：北京市西城区右安门西街8号中国铁道出版社综合编辑部 荆波 收（邮编：100054）。或者采用传真（010-63549458）方式发送。此外，读者也可以直接通过电子邮件把意见反馈给我们，E-mail地址是：jb@163.jb18803242@yahoo.com.cn。我们将选出意见中肯的热心读者，赠送本社的其他图书作为奖励。同时，我们将充分考虑您的意见和建议，并尽可能地给您满意的答复。谢谢！

所购书名：_____

个人资料：

姓名：_____ 性别：_____ 年龄：_____ 文化程度：_____

职业：_____ 电话：_____ E-mail：_____

通信地址：_____ 邮编：_____

您是如何得知本书的：

书店宣传 网络宣传 展会促销 出版社图书目录 老师指定 杂志、报纸等的介绍 别人推荐

其他（请指明）_____

您从何处得到本书的：

书店 邮购 商场、超市等卖场 图书销售的网站 培训学校 其他

影响您购买本书的因素（可多选）：

内容实用 价格合理 装帧设计精美 带多媒体教学光盘 优惠促销 书评广告 出版社知名度

作者名气 工作、生活和学习的需要 其他

您对本书封面设计的满意程度：

很满意 比较满意 一般 不满意 改进建议

您对本书的总体满意程度：

从文字的角度 很满意 比较满意 一般 不满意

从技术的角度 很满意 比较满意 一般 不满意

您希望书中图的比例是多少：

少量的图片辅以大量的文字 图文比例相当 大量的图片辅以少量的文字

您希望本书的定价是多少：

本书最令您满意的是：

1.

2.

您在使用本书时遇到哪些困难：

1.

2.

您希望本书在哪些方面进行改进：

1.

2.

您需要购买哪些方面的图书？对我社现有图书有什么好的建议？

您更喜欢阅读哪些类型和层次的计算机书籍（可多选）？

入门类 精通类 综合类 问答类 图解类 查询手册类 实例教程类

您在学习计算机的过程中有什么困难？

您的其他要求：

第 1 篇 Hadoop 技术篇

第 1 章 初识 Hadoop

| | |
|--|----|
| 1.1 Hadoop 简介 | 2 |
| 1.1.1 Hadoop 的起源 | 2 |
| 1.1.2 什么是 Hadoop | 3 |
| 1.1.3 Hadoop 的核心技术是 Google 核心技术的开源实现 | 4 |
| 1.1.4 Hadoop 的功能与优点 | 5 |
| 1.1.5 Hadoop 的应用现状和发展趋势 | 6 |
| 1.2 Hadoop 的体系结构 | 11 |
| 1.2.1 HDFS 的体系结构 | 12 |
| 1.2.2 MapReduce 的体系结构 | 19 |
| 1.3 Hadoop 与分布式开发 | 21 |
| 1.4 Hadoop 的数据管理 | 23 |
| 1.4.1 HDFS 的数据管理 | 23 |
| 1.4.2 HBase 的数据管理 | 23 |
| 1.4.3 Hive 的数据管理 | 24 |
| 1.5 思考与总结 | 25 |

第 2 章 Hadoop 的安装和配置

| | |
|-----------------------------------|----|
| 2.1 在 Windows 下安装与配置 Hadoop | 27 |
| 2.1.1 JDK 的安装 | 27 |
| 2.1.2 Cygwin 的安装 | 30 |
| 2.1.3 Hadoop 的安装 | 36 |
| 2.2 在 Linux 下安装与配置 Hadoop | 38 |

| | | |
|-------|----------------------------|----|
| 2.2.1 | Ubuntu 的安装 | 38 |
| 2.2.2 | JDK 的安装 | 41 |
| 2.2.3 | Hadoop 的安装 | 41 |
| 2.3 | Hadoop 的执行实例 | 43 |
| 2.3.1 | 运行 Hadoop | 44 |
| 2.3.2 | 运行 wordcount.java 程序 | 44 |
| 2.4 | Hadoop Eclipse 简介和使用 | 45 |
| 2.4.1 | Eclipse 插件介绍 | 45 |
| 2.4.2 | Eclipse 插件开发配置 | 45 |
| 2.4.3 | 在 Eclipse 下运行 WordCount 程序 | 49 |
| 2.5 | Hadoop 的集群和优化 | 56 |
| 2.5.1 | Hadoop 的性能优化 | 57 |
| 2.5.2 | Hadoop 配置机架感知信息 | 58 |
| 2.6 | 思考与总结 | 59 |

第 3 章 HDFS 海量存储

| | | |
|-------|------------------|----|
| 3.1 | 开源的 GFS——HDFS | 60 |
| 3.1.1 | HDFS 简介 | 60 |
| 3.1.2 | HDFS 的体系结构 | 63 |
| 3.1.3 | HDFS 的保障可靠性措施 | 64 |
| 3.2 | HDFS 的常用操作 | 67 |
| 3.2.1 | HDFS 下的文件操作 | 67 |
| 3.2.2 | 管理与更新 | 74 |
| 3.2.3 | HDFS API 详解 | 76 |
| 3.2.4 | HDFS 的读/写数据流 | 88 |
| 3.3 | 用 HDFS 存储海量的视频数据 | 91 |
| 3.3.1 | 场景分析 | 91 |
| 3.3.2 | 设计实现 | 91 |
| 3.4 | 思考与总结 | 93 |

第 4 章 初识 MapReduce

| | | |
|-----|--------------|----|
| 4.1 | MapReduce 简介 | 94 |
|-----|--------------|----|

| | | |
|-------|------------------------------|-----|
| 4.1.1 | MapReduce 要解决什么问题 | 94 |
| 4.1.2 | MapReduce 的理论基础 | 95 |
| 4.1.3 | MapReduce 的编程模式 | 97 |
| 4.2 | MapReduce 的集群行为 | 98 |
| 4.3 | Map/Reduce 框架 | 100 |
| 4.4 | 样例分析：单词计数 | 100 |
| 4.4.1 | WordCount 实例的运行过程 | 100 |
| 4.4.2 | WordCount 的源码分析和程序处理过程 | 103 |
| 4.4.3 | MapReduce 常用类及其接口 | 106 |
| 4.5 | 实例：倒排索引 | 109 |
| 4.5.1 | 倒排索引的分析和设计 | 109 |
| 4.5.2 | 倒排索引完整源码 | 112 |
| 4.5.3 | 运行代码结果 | 116 |
| 4.6 | MapReduce 在日志分析中数据去重案例 | 117 |
| 4.6.1 | 什么是数据去重 | 117 |
| 4.6.2 | 设计思路 | 118 |
| 4.6.3 | 程序代码 | 118 |
| 4.6.4 | 代码运行结果 | 120 |
| 4.7 | 数据排序实例 | 122 |
| 4.7.1 | 实例描述 | 122 |
| 4.7.2 | 设计思路 | 123 |
| 4.7.3 | 程序代码 | 123 |
| 4.8 | 思考与总结 | 126 |

第 5 章 分布式开源数据库 HBase

| | | |
|-------|---------------------|-----|
| 5.1 | HBase 简介 | 127 |
| 5.1.1 | HBase 逻辑视图 | 127 |
| 5.1.2 | HBase 物理存储 | 129 |
| 5.1.3 | 子表 Region 服务器 | 130 |
| 5.1.4 | Hmaster 主服务器 | 132 |
| 5.1.5 | 元数据表 | 132 |
| 5.2 | HBase 的安装配置 | 133 |

| | | |
|-------|---------------------------|-----|
| 5.2.1 | HBase 单机模式 | 133 |
| 5.2.2 | HBase 伪分布模式 | 135 |
| 5.2.3 | HBase 完全分布模式 | 136 |
| 5.3 | 学生成绩表实例 | 140 |
| 5.3.1 | Shell 的基本操作 | 141 |
| 5.3.2 | 代码实现 | 143 |
| 5.3.3 | 关于中文的处理 | 145 |
| 5.3.4 | 常用 HBase 的 Shell 操作 | 149 |
| 5.4 | 思考与总结 | 153 |

第 6 章 MapReduce 进阶

| | | |
|-------|--------------------------------------|-----|
| 6.1 | API 的配置 | 154 |
| 6.1.1 | 一个简单的配置文件 | 155 |
| 6.1.2 | 合并多个源文件 | 156 |
| 6.1.3 | 可变的扩展 | 157 |
| 6.2 | 配置开发环境 | 157 |
| 6.2.1 | 配置文件设置 | 157 |
| 6.2.2 | 设置用户标识 | 159 |
| 6.3 | 复合键值对的使用 | 159 |
| 6.3.1 | 小的键值对如何合并成大的键值对 | 159 |
| 6.3.2 | 巧用复合键让系统完成排序 | 160 |
| 6.4 | 用户定制数据类型 | 164 |
| 6.4.1 | 内置数据类型 | 164 |
| 6.4.2 | 用户自定义数据类型 | 164 |
| 6.5 | 用户定制输入/输出格式 | 166 |
| 6.5.1 | 内置数据的输入格式 | 167 |
| 6.5.2 | 用户定制数据输入格式与 RecordReader | 168 |
| 6.5.3 | Hadoop 内置的数据输出格式 | 172 |
| 6.5.4 | Hadoop 内置的数据输出格式与 RecordWriter | 172 |
| 6.6 | 用户定制 Partitioner 和 Combiner | 173 |
| 6.7 | 组合式的 MapReduce 作业 | 176 |
| 6.7.1 | MapReduce 作业运行机制 | 176 |

| | |
|--------------------------------|-----|
| 6.7.2 组合式 MapReduce 计算作业 | 178 |
| 6.8 DataJoin 连接多数据源 | 183 |
| 6.9 思考与总结..... | 187 |

第 7 章 Hive 数据仓库

| | |
|---------------------------|-----|
| 7.1 Hive 简介..... | 188 |
| 7.2 Hive 安装与配置..... | 189 |
| 7.3 Hive 的服务..... | 191 |
| 7.3.1 Hive shell | 191 |
| 7.3.2 JDBC/ODBC..... | 192 |
| 7.3.3 Thrift 服务 | 192 |
| 7.3.4 Web 接口 | 193 |
| 7.3.5 元数据服务..... | 193 |
| 7.4 HiveQL 查询语言 | 193 |
| 7.5 Hive 实例..... | 202 |
| 7.5.1 UDF 编程实例..... | 202 |
| 7.5.2 UDAF 编程实例 | 204 |
| 7.5.3 Hive 的日志数据统计实战..... | 206 |
| 7.6 思考与总结..... | 211 |

第 8 章 Pig 开发应用

| | |
|--------------------------------|-----|
| 8.1 Pig 简介 | 212 |
| 8.2 Pig 的安装与配置 | 213 |
| 8.3 Pig 的使用 | 215 |
| 8.3.1 Pig 的 MapReduce 模式..... | 215 |
| 8.3.2 Pig 的运行方式..... | 216 |
| 8.4 通过 Grunt 学习 Pig Latin..... | 219 |
| 8.4.1 Pig 的数据模型 | 220 |
| 8.4.2 运算符 | 221 |
| 8.4.3 常用操作 | 222 |
| 8.4.4 各种 SQL 在 Pig 中的实现..... | 229 |

| | |
|--------------------------|-----|
| 8.4.5 Pig Latin 实现 | 233 |
| 8.5 Pig 使用的案例 | 235 |
| 8.6 思考与总结..... | 235 |

第 9 章 Chukwa 数据收集系统

| | |
|-------------------------------|-----|
| 9.1 Chukwa 简介 | 236 |
| 9.1.1 Chukwa 是什么 | 236 |
| 9.1.2 Chukwa 主要解决什么问题 | 240 |
| 9.2 Chukwa 的安装配置 | 240 |
| 9.2.1 Chukwa 的安装..... | 240 |
| 9.2.2 Chukwa 的配置..... | 242 |
| 9.2.3 Chukwa 的启动..... | 245 |
| 9.3 Chukwa 的基本命令 | 248 |
| 9.3.1 Chukwa 端的命令..... | 248 |
| 9.3.2 Agent 端的命令 | 249 |
| 9.4 Chukwa 在数据收集处理方面的运用 | 251 |
| 9.4.1 数据生成 | 251 |
| 9.4.2 数据收集 | 251 |
| 9.4.3 数据处理 | 252 |
| 9.4.4 数据析取 | 252 |
| 9.4.5 数据稀释 | 253 |
| 9.4.6 数据显示 | 253 |
| 9.5 思考与总结..... | 253 |

第 10 章 ZooKeeper 开发应用

| | |
|-----------------------------------|-----|
| 10.1 ZooKeeper 简介 | 254 |
| 10.1.1 ZooKeeper 的设计目标 | 254 |
| 10.1.2 ZooKeeper 主要解决什么问题 | 256 |
| 10.1.3 ZooKeeper 的基本概念和工作原理 | 257 |
| 10.2 ZooKeeper 的安装配置 | 260 |
| 10.2.1 单机模式 | 261 |
| 10.2.2 启动并测试 ZooKeeper | 262 |

| | |
|-----------------------------------|-----|
| 10.2.3 集群模式 | 264 |
| 10.3 ZooKeeper 提供的接口 | 267 |
| 10.4 ZooKeeper 事件 | 270 |
| 10.5 ZooKeeper 实例 | 271 |
| 10.5.1 实例 1：一个简单的应用——分布式互斥锁 | 271 |
| 10.5.2 实例 2：进程调度系统 | 276 |
| 10.6 思考与总结 | 283 |

第 2 篇 Hadoop 管理和容错篇

第 11 章 Hadoop 管理

| | |
|-----------------------------|-----|
| 11.1 Hadoop 权限管理 | 286 |
| 11.2 HDFS 文件系统管理 | 292 |
| 11.3 Hadoop 维护与管理 | 298 |
| 11.4 Hadoop 常见问题及解决办法 | 300 |
| 11.5 思考与总结 | 310 |

第 12 章 Hadoop 容错

| | |
|--|-----|
| 12.1 Hadoop 的可靠性 | 311 |
| 12.1.1 HDFS 中的 NameNode 单点失效解决方案 | 311 |
| 12.1.2 HDFS 数据块副本机制 | 313 |
| 12.1.3 HDFS 心跳机制 | 319 |
| 12.1.4 HDFS 负载均衡 | 320 |
| 12.1.5 MapReduce 容错 | 321 |
| 12.2 Hadoop 的 SecondaryNameNode 机制 | 322 |
| 12.2.1 磁盘镜像与日志文件 | 322 |
| 12.2.2 SecondaryNameNode 更新镜像的流程 | 323 |
| 12.3 Avatar 机制 | 325 |
| 12.3.1 Avatar 机制简介 | 325 |
| 12.3.2 Avatars 部署实战 | 326 |
| 12.4 Hadoop_HBase 容错 | 331 |

| | |
|-----------------|-----|
| 12.5 思考与总结..... | 333 |
|-----------------|-----|

第 3 篇 Hadoop 实战篇

第 13 章 综合实战 1：Hadoop 中的数据库访问

| | |
|---|-----|
| 13.1 DBInputFormat 类访问数据库 | 336 |
| 13.1.1 在 DBInputFormat 类中包含的内置类 | 336 |
| 13.1.2 使用 DBInputFormat 读取数据库表中的记录..... | 337 |
| 13.1.3 使用示例 | 337 |
| 13.2 使用 DBOutputFormat 向数据库中写记录..... | 340 |
| 13.3 思考与总结..... | 343 |

第 14 章 综合实战 2：一个简单的分布式的 Grep

| | |
|-----------------|-----|
| 14.1 分析与设计..... | 344 |
| 14.2 实现代码 | 345 |
| 14.3 运行程序 | 346 |
| 14.4 思考与总结..... | 346 |

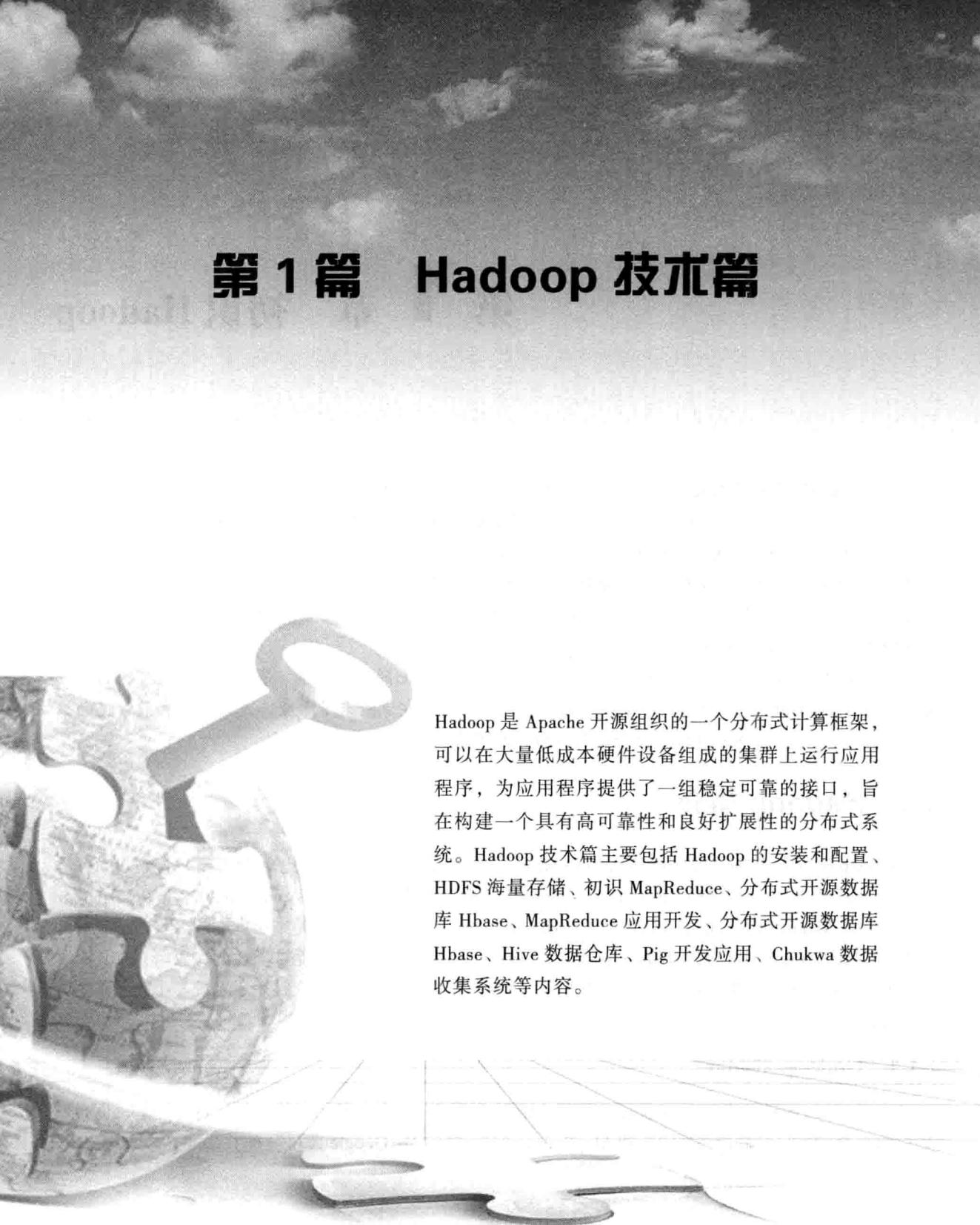
第 15 章 综合实战 3：打造一个搜索引擎

| | |
|---------------------------|-----|
| 15.1 搜索引擎工作原理..... | 348 |
| 15.2 网页搜集与信息提取..... | 350 |
| 15.2.1 设计的主要思想 | 350 |
| 15.2.2 系统设计目标..... | 351 |
| 15.3 网页信息的提取与存储..... | 352 |
| 15.4 MapReduce 的预处理 | 353 |
| 15.4.1 第一步：源数据过滤 | 353 |
| 15.4.2 第二步：生成倒排文件 | 355 |
| 15.4.3 第三步：建立二级索引 | 362 |
| 15.5 建立 Web 信息查询服务 | 365 |
| 15.6 思考与总结..... | 366 |

第 16 章 综合实战 4：移动通信信令监测与查询

| | |
|----------------------------------|------------|
| 16.1 分析与设计..... | 367 |
| 16.1.1 CDR 数据文件的检测与索引创建任务调度..... | 369 |
| 16.1.2 从 HDFS 读取数据并创建索引 | 370 |
| 16.1.3 查询 CDR 信息..... | 371 |
| 16.2 代码实现 | 371 |
| 16.2.1 CDR 文件检测和索引创建任务程序 | 371 |
| 16.2.2 读取 CDR 数据和索引创建处理 | 375 |
| 16.2.3 CDR 查询..... | 383 |
| 16.3 思考与总结..... | 384 |
| 附录 A Hadoop 命令大全 | 385 |
| 附录 B HDFS 命令大全 | 392 |

第1篇 Hadoop 技术篇



Hadoop 是 Apache 开源组织的一个分布式计算框架，可以在大量低成本硬件设备组成的集群上运行应用程序，为应用程序提供了一组稳定可靠的接口，旨在构建一个具有高可靠性和良好扩展性的分布式系统。Hadoop 技术篇主要包括 Hadoop 的安装和配置、HDFS 海量存储、初识 MapReduce、分布式开源数据库 Hbase、MapReduce 应用开发、分布式开源数据库 Hbase、Hive 数据仓库、Pig 开发应用、Chukwa 数据收集系统等内容。

第 1 章 初识 Hadoop

云计算是 IT 界的第三次浪潮。在这次浪潮中，各大厂商面临着极大的挑战——他们需要从 TB 乃至 PB 级数据中挖掘出有用的数据，并对这些海量的数据进行更快捷、高效率的处理。于是 IT 厂商推出了自己的云计算平台。Google 的 MapReduce、GFS、BigTable 成为互联网的领头羊，然而它的技术是保密的，Google 公司并没有开源 MapReduce 的实现细节。Amazon 的 AWS、微软的 Azure 和 IBM 的蓝云等也是云计算的典型代表，但它们都是商业性平台，对想要继续研究和发展云计算技术的人员或科研团体来说，无法获得更多的了解。

Hadoop 作为 Apache 基金会资助的开源项目，模仿 Google 的核心技术，是一个分布式系统的基础架构，由 Doug Cutting 带领的团队进行开发。它的出现给研究者带来了希望，是最典型和最常见的云计算平台。未来，Hadoop 将作为一个幕后英雄，应用于越来越多的行业。

1.1 Hadoop 简介

Hadoop 起源于从 2002 年开始的 Apache Nutch，它是 Apache Lucene 的子项目之一。直到 2006 年，Hadoop 才逐渐成为一套完整而独立的软件，并被正式命名，最大支持者是 Yahoo!。2008 年初，Hadoop 开始应用到 Yahoo! 以外的很多互联网公司。Hadoop 并不是一个缩写，而是一个虚构的名字，该项目的创建者 Doug Cutting 这样解释 Hadoop 的得名：“这个名字是我的孩子给一个棕黄色的大象填充玩具命名的。我的命名标准就是简短，容易发音和拼写，没有太多的意义，并且不会被用于别处。小孩子是这方面的高手。”

1.1.1 Hadoop 的起源

MapReduce 编程思想是由 Google 工程师 Jeffrey Dean 于 2004 年提出来的，与此同时，Google 也发表了 GFS、BigTable 等底层系统以应用 MapReduce 模型。2007 年，Google 公司发布了 *Google's MapReduce*