



HZ BOOKS

华章 IT

Apress®

世界级大数据专家亲笔撰写，深入剖析Hadoop系统实现原理及其工程实践应用，包含大量真实案例和数据

从Hadoop使用者角度深入浅出地讲解Hadoop各组件的运行机理，深入到系统实现源码，帮助你从架构、开发、应用和运维等多方面全面掌握Hadoop系统



# Pro Apache Hadoop

Second Edition

# 深入理解Hadoop

(原书第2版)

萨米尔·瓦德卡 (Sameer Wadkar)

[美] 马杜·西德林埃 (Madhu Siddalingaiah) 著  
杰森·文纳 (Jason Venner)

于博 冯傲风 译



机械工业出版社  
China Machine Press

01 010011000

01 010011000

010 010011000

010



技术丛书

# Pro Apache Hadoop

Second Edition

# 深入理解Hadoop

(原书第2版)

萨米尔·瓦德卡 (Sameer Wadkar)

[美] 马杜·西德林埃 (Madhu Siddalingaiah) 著

杰森·文纳 (Jason Venner)

于博 冯傲风 译



机械工业出版社  
China Machine Press

## 图书在版编目 (CIP) 数据

深入理解 Hadoop (原书第 2 版) / (美) 瓦德卡 (Wadkar, S.), (美) 西德林埃 (Siddalingaiah, M.), (美) 文纳 (Venner, J.) 著; 于博, 冯傲风译. —北京: 机械工业出版社, 2015.10  
(大数据技术丛书)

书名原文: Pro Apache Hadoop, Second Edition

ISBN 978-7-111-51565-4

I. 深… II. ①瓦… ②西… ③文… ④于… ⑤冯… III. 数据处理软件 IV. TP274

中国版本图书馆 CIP 数据核字 (2015) 第 227389 号

---

本书版权登记号: 图字: 01-2015-4154

Sameer Wadkar, Madhu Siddalingaiah and Jason Venner: Pro Apache Hadoop, Second Edition (ISBN : 978-1-4302-4863-7).

Original English language edition published by Apress L.P., 2560 Ninth Street, Suite 219, Berkeley, CA 94710 USA. Copyright © 2014 by Sameer Wadkar and Madhu Siddalingaiah. Simplified Chinese-language edition copyright © 2016 by China Machine Press. All rights reserved.

This edition is licensed for distribution and sale in the People's Republic of China only, excluding Hong Kong, Taiwan and Macao and may not be distributed and sold elsewhere.

本书原版由 Apress 出版社出版。

本书简体字中文版由 Apress 出版社授权机械工业出版社独家出版。未经出版者预先书面许可，不得以任何方式复制或抄袭本书的任何部分。

此版本仅限在中华人民共和国境内（不包括中国香港、台湾、澳门地区）销售发行，未经授权的本书出口将被视为违反版权法的行为。

## 深入理解 Hadoop (原书第 2 版)

---

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 吴 怡

责任校对: 殷 虹

印 刷: 北京诚信伟业印刷有限公司

版 次: 2016 年 1 月第 1 版第 1 次印刷

开 本: 186mm×240mm 1/16

印 张: 25

书 号: ISBN 978-7-111-51565-4

定 价: 79.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88379426 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

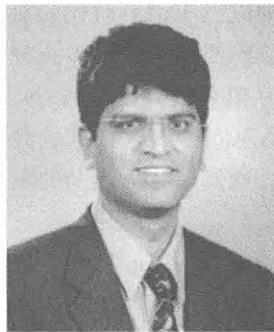
读者信箱: hzit@hzbook.com

版权所有 · 侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

## 作者简介 *About the Author*



萨米尔·瓦德卡 (Sameer Wadkar) 在软件架构与开发方面有多达 16 年的工作经验。他在孟买大学取得了电子工程方向的硕士学位，及金融专业的 MBA 学位，在 National Center of Software Technology (现为 Center for Development of Advanced Computing) 取得了软件工程专业的硕士学位，在约翰霍普金斯大学获得了应用与计算数学专业的博士学位。他还实现了分布式系统及高访问量的 Web 网站，从政府机构到银行投资公司多种领域都有客户使用这个 Web 网站。他从 2011 年就积极地参与到了 Hadoop/HBase 的开发实现当中，也是开源软件的贡献者。他的 GitHub 的主页是：<https://github.com/sameerwadkar>。

他在开源领域的贡献包括 Latent Dirichlet Allocation 文本挖掘算法，该算法能够在单台机器上处理数百万级别的文档。他还是一位活跃的国际象棋玩家，经常在网站 [www.chessclub.com](http://www.chessclub.com) 上下国际象棋。



马杜·西德林埃 (Madhu Siddalingaiah) 是一名拥有 25 年工作经验的技术顾问，其工作经验涉及多个领域，包括航空航天、医疗保险、金融、能源、科学研究等。多年以来，他专门从事电子工程、互联网技术、大数据等领域。最近，他发布了几个著名的大数据系统及其解决方案。

他在马里兰大学取得了物理学的硕士学位。专业之外，Madhu 还是一名私人直升飞机驾驶员，他喜欢旅行，不断参与越来越多的爱好活动。

杰森·文纳 (Jason Venner) 有 20 多年的软件开发经验，涉及软件工程、数据挖掘、架构设计等领域，曾在 IT 企业中担任过副总裁、主管、顾问，近些年来他的钻研方向是 Java、Hadoop、云计算等领域。

## *The Translator's Words* 译 者 序

我们正处在一个数据量爆炸增长的时代，也是一个越来越依靠数据进行决策的时代，对海量数据的挖掘分析也从抽样数据分析变成了全量数据分析。为此，业界涌现出了许多新技术，Apache Hadoop 是其中当之无愧的王者。当前 Hadoop 已然发展为拥有 HDFS、YARN、MapReduce、Pig、HCatalog、HBase 等系统的较为完整的大数据系统生态圈，为大数据的开发利用提供了比较全面的解决方案。毫无疑问，Hadoop 已成为大数据行业发展背后强劲的驱动力。

本书详细地讲述了 Hadoop 生态圈中最为重要的几个组件。不仅介绍了 Hadoop 涉及的分布式理论基础知识，还着重讲解 Hadoop 系统的工程实践应用。为了深入浅出地讲述 Hadoop 各个组件的运行机理，作者使用了贴切的实战用例贯穿全书，同时，又深入到系统实现源码，使读者做到知其然，更知其所以然。通过对本书的学习，读者可以从架构、开发、应用和运维等多方面全面地掌握 Hadoop 系统，从而成为名副其实的 Hadoop 专家。

我曾经拜读过本书英文版的第 1 版，机缘巧合，华章给了我翻译该书第 2 版的机会，感谢我的小伙伴百度公司的冯傲风同学与我并肩奋战，完成了本书的翻译；感谢吴怡编辑对本书认真负责的审校和中肯的修改建议，使得本书能够顺利地出版，她的认真和敬业令人钦佩！

最后还要感谢家人，只有他们的理解和支持，我才能每天抽出那么多时间来完成本书的翻译工作。

书中概念和术语颇多，有一些目前尚无固定中文译法，加上译者水平有限，译文中的不当之处在所难免，我们真诚地希望同行和读者们不吝赐教。

于博

2015 年 9 月于北京

## 技术评审人简介



威米什·奥姆·米塔尔 (Vimlesh Om Mittal) 具有多达 14 年的技术实现经验，是 Apache Hadoop CDH4 的 Cloudera 认证开发者 (Cloudera Certified Developer)。他具有非常广博的技术背景，他的经验主要包括关于现代商业处理的相关项目、客户应用开发、ETL 开发、数字应用开发、商业智能、数据库设计和数据转换策略。

他拥有构建健壮的信息管理系统、解决方案架构设计和基于多种开发框架进行软件开发的综合经验。这些开发框架包括 Hadoop 生态系统中的大数据框架、数据获取、ETL、报告 / 分析、.NET、J2EE 和不同技术供应商提供的脚本语言，这些技术供应商包括 Cloudera、SAS、IBM、Oracle、Microsfot、iOS 以及 BEA。

Vimlesh 还专注于金融服务数据的解决方案、移动应用设计及开发。他对软件开发的整个生命周期中的各个方面都十分熟悉，在项目管理科学方面有宝贵经验。

## 前　　言 *preface*

Hadoop 已经进入 Apache 社区发展五年多了，使用 Hadoop 系统进行开发的工作仍然富于挑战但收获丰厚。本书第 1 版在若干年前就已经出版了，在这期间，Hadoop 系统已经被越来越多的企业使用，自身也得到了飞速发展。Hadoop2.0 基于 YARN 框架做了全新升级，重写了 Hadoop 系统的底层平台。本书从 Hadoop 使用者的角度出发讲解 Hadoop 的实现原理，浓缩了 Hadoop 软件系统的精华。作为作者，我们希望可以深入到源代码级别来理解 Hadoop 的运行原理及其背后的设计目标，渴望与你分享 Hadoop 带给我们的启迪。通过本书，你不仅能够深入学习 Hadoop 系统，还可以对 Java 语言在大数据处理过程中的运用有更深入的认识。

本书总体上介绍了大数据的范畴概念，然后对 Hadoop 进行了详细讲解。因为如果对大数据没有一个总体上的认知，是不可能真正理解 Hadoop 系统的。本书的读者对象具备中级 Java 开发能力的 Hadoop 开发人员。本书可以让你成为企业中的 Hadoop 专家。通过本书的学习，你可以获得大量实用技巧，这些技巧都是我们在从事基于 Hadoop 的分布式数据处理系统实践过程中总结出来的。

本书循序渐进地讲解了大量的实例，会帮助你从 Hadoop 集群初级使用者成长为能胜任运行任何复杂应用程序的专家。下面是本书的一个内容纲要：

- 第 1 章，讲解大数据系统的作用，介绍各种各样的大数据系统。
- 第 2 章，高屋建瓴地介绍 Hadoop2.0 和 YARN，讲解 Hadoop 平台的关键概念。
- 第 3 章，开始 Hadoop 学习，构建你的第一个 MapReduce 程序。
- 第 4 章，讲解 Hadoop 平台管理的关键概念要素。
- 第 5 章、第 6 章和第 7 章是本书的重点，深入分析讲解了 MapReduce 框架。你将会学习 MapReduce 框架所有知识点。我们以当前广泛使用的语言 SQL 作为对照，来学习理解 MapReduce 框架的使用。利用 MapReduce 框架来实现 SQL 语言中 SELECT、WHERE、GROUP BY 和 JOIN 这些关键字的功能。Hadoop 系统一个常用用途就是来做数据 ETL。这几章会让你掌握如何使用 MapReduce 框架支持通用的数据处理功能。我们不仅仅学习 MapReduce 框架的 API，还会学习 MapReduce 框架中更复杂的概念。

及其设计理念。

- 第 8 章，介绍一种测试框架，这种测试框架支持 MapReduce 程序的单元测试和集成测试。
- 第 9 章，讲解 Hadoop 框架的监控和日志。
- 第 10 章，讲解 Hive 框架，这是一个基于 MapReduce 的数据仓库框架。
- 第 11 章，讲解 Pig 和 Crunch 框架。这些框架可以帮助使用者在 Hadoop 平台上构建数据处理管道。
- 第 12 章，讲解 HCatalog 框架。该框架帮助企业用户可以像访问数据库中的数据表一样，来访问存放在 Hadoop 文件系统上的海量数据。
- 第 13 章，讲解如何使用 Hadoop 来分析处理日志流。
- 第 14 章，介绍 HBase，这是一个基于 Hadoop 系统的 NoSQL 数据库。你会学习许多 HBase 用法范例。
- 第 15 章，对数据科学做简要介绍。介绍了 MapReduce 框架在数学科学上的局限和不足之处。同时引入介绍新的框架，比如 Spark、Hama，这些框架突破了 MapReduce 框架中的某些局限。
- 第 16 章，简要介绍了 Hadoop 系统在云计算中的应用。讲解如何在一个真实的生产环境中的 Hadoop 集群上工作。
- 第 17 章，粗略地介绍 Hadoop2.0 的一个关键附加功能：在 Hadoop 上开发一个自己的类似于 MapReduce 框架的分布式数据处理框架。讲解了如何基于 Hadoop2.0 开发一个简单的分布式下载服务。

## 致谢

Hadoop 在过去的 10 年发展过程中，有许多厂商和独立开发者为之贡献代码。Hadoop 系统中也使用了大量的开源函数库。我们要感谢为 Hadoop 系统贡献源码的所有开源作者所付出的辛勤劳动。我们也要感谢那些在论坛上回答了 Hadoop 系统相关问题的博客作者及专家学者。他们的积极参与，使得 Hadoop 这样复杂的系统变得易于使用并成为应用主流。

我们还想感谢 Apress 的全体员工，他们的努力付出使得本书内容清晰易读。

我们最后要感谢我们的家庭、朋友及同事，本书在编写过程中得到他们的一贯支持。

# 目 录 *Contents*

译者序		
作者简介		
前言		
<b>第1章 为什么会有大数据</b>	<b>1</b>	
1.1 什么是大数据	1	
1.2 大数据技术背后的核心思想	2	
1.2.1 把数据分发到多个节点	2	
1.2.2 把计算逻辑移动到数据 附近	3	
1.2.3 计算节点进行本地数据处理	3	
1.2.4 优选顺序读，次之随机读	4	
1.2.5 一个例子	4	
1.3 大数据的编程模型	5	
1.3.1 大规模并行处理数据库系统	5	
1.3.2 内存数据库系统	6	
1.3.3 MapReduce 系统	6	
1.3.4 整体同步并行系统	8	
1.4 大数据和事务性系统	8	
1.5 我们能处理多大的数据量	9	
1.5.1 一个计算密集型的例子	10	
1.5.2 Amdhal 定律	10	
1.6 大数据商业用例	11	
1.7 本章小结	12	
<b>第2章 Hadoop 中的概念</b>	<b>13</b>	
2.1 Hadoop 简介	13	
2.2 MapReduce 编程模型简介	15	
2.3 Hadoop 系统的组成	19	
2.3.1 Hadoop 分布式文件系统	20	
2.3.2 辅助名称节点	25	
2.3.3 任务跟踪器	26	
2.3.4 作业跟踪器	26	
2.4 Hadoop 2.0	27	
2.4.1 容器	29	
2.4.2 节点管理器	29	
2.4.3 资源管理器	30	
2.4.4 应用程序管理器	30	
2.4.5 分步详解 YARN 请求	31	
2.5 HDFS 的高可用性	33	
2.6 本章小结	33	
<b>第3章 初识 Hadoop 框架</b>	<b>34</b>	
3.1 安装类型	34	
3.1.1 单机模式	35	
3.1.2 伪分布式集群模式	35	
3.1.3 多节点集群安装模式	35	
3.1.4 基于 Amazon EMR 预安装 模式	35	

3.2 使用 Cloudera 虚拟机搭建开发环境	36
3.3 一个 MapReduce 程序的组成	37
3.4 第一个 Hadoop 程序	38
3.4.1 以本地模式运行程序的必要条件	39
3.4.2 使用旧 API 编写的单词计数程序	39
3.4.3 构建程序	42
3.4.4 在集群模式下运行单词计数程序	42
3.4.5 使用新 API 编写的单词计数程序	43
3.4.6 构建程序	44
3.4.7 在集群模式下运行单词计数程序	45
3.5 Hadoop 作业中的第三方函数库	45
3.6 本章小结	50
<b>第 4 章 Hadoop 系统管理</b>	<b>51</b>
4.1 Hadoop 的配置文件	51
4.2 配置 Hadoop 守护进程	52
4.3 Hadoop 配置文件的优先级	53
4.4 深入探究 Hadoop 配置文件	54
4.4.1 core-site.xml	54
4.4.2 hdfs-*.xml	55
4.4.3 mapred-site.xml	56
4.4.4 yarn-site.xml	58
4.4.5 YARN 中的内存分配	60
4.5 调度器	61
4.5.1 计算能力调度器	62
4.5.2 公平调度器	65
4.5.3 公平调度器配置	65
4.5.4 yarn-site.xml 配置	66
4.5.5 策略文件的格式和配置	67
4.5.6 按照 drf 策略来确定优势资源的分配	68
4.6 从属文件	69
4.7 机架感知	69
4.8 集群管理工具	71
4.8.1 检查 HDFS	71
4.8.2 HDFS 管理命令行	73
4.8.3 均衡 HDFS 上的数据分布	75
4.8.4 从 HDFS 中复制海量数据	76
4.9 本章小结	76
<b>第 5 章 MapReduce 开发基础</b>	<b>78</b>
5.1 Hadoop 和数据处理	78
5.2 航空公司数据集介绍	79
5.2.1 准备开发环境	80
5.2.2 准备 Hadoop 系统	81
5.3 MapReduce 编程模式	81
5.3.1 只有 Map 阶段的作业（SELECT 和 WHERE 查询）	82
5.3.2 问题定义——SELECT 子句	82
5.3.3 问题定义——WHERE 子句	90
5.3.4 Map 和 Reduce 作业（聚合查询）	93
5.3.5 问题定义——GROUP BY 和 SUM 子句	93
5.3.6 应用 Combiner 提高 Aggregation 性能	99
5.3.7 问题定义——优化后的 Aggregators	99
5.3.8 Partitioner 的作用	104
5.3.9 问题定义——按月分离航空数据	105
5.4 综合分析	108
5.5 本章小结	110

<b>第6章 MapReduce 开发进阶</b>	111	作业	147
6.1 MapReduce 编程模式	111	6.6.3 总结探讨 Map-Only 连接时的 Hadoop 关键特性	149
6.2 Hadoop I/O 介绍	111		
6.3 问题定义——排序	114	6.7 在 MR 作业中保存结果到多 输出文件	149
6.3.1 主要挑战：全排序	115		
6.3.2 在 Cluster 中运行 Sorting 作业	125	6.8 使用计数器收集统计数据	151
6.3.3 仅根据 Writable 键排序	125		
6.3.4 根据排序回顾 Hadoop 的关键 特性	128	6.9 本章小结	153
6.4 问题定义——分析连续的记录	128		
6.4.1 支持二次排序的重要组件	129		
6.4.2 在没有 Grouping Comparator 的 情况下实现 Secondary Sort	136		
6.4.3 在 Cluster 中运行 Secondary Sort 作业	137		
6.4.4 利用 Secondary Sort 回顾 Hadoop 的关键特性	137		
6.5 问题定义——使用 MapReducer 进行连接	138		
6.5.1 处理多输入：Multiple- Inputs 类	138		
6.5.2 具备多个输入的 Mapper 类	139		
6.5.3 自定义 Partitioner: Carrier- CodeBasedPartitioner	141		
6.5.4 在 Reducer 中实现连接	141		
6.5.5 在集群中运行 MapReduce 连接 作业	143		
6.5.6 探讨与 MapReduce 相关的 Hadoop 主要特性	144		
6.6 问题定义——使用 Map-Only 作业 进行连接	144		
6.6.1 基于 DistributeCache 的解决 方案	145		
6.6.2 在集群中运行 Map-Only 的连接			
<b>第7章 Hadoop 输入 / 输出</b>	155		
7.1 压缩方式	155		
7.1.1 压缩内容的选择	156		
7.1.2 各种压缩方式	157		
7.1.3 配置压缩方式	158		
7.2 Hadoop 的 I/O 处理过程内部	159		
7.2.1 InputFormat	159		
7.2.2 OutputFormat	161		
7.2.3 自定义 OutputFormat: 将文本 转换成 XML	161		
7.2.4 自定义 InputFormat: 使用 自定义的 XML 文件	165		
7.3 Hadoop 文件	173		
7.3.1 SequenceFile	173		
7.3.2 MapFiles	178		
7.3.3 Avro Files	180		
7.4 本章小结	185		
<b>第8章 测试 Hadoop 程序</b>	186		
8.1 回顾一下单词统计的程序	186		
8.2 MRUnit 概述	188		
8.2.1 安装 MRUnit	188		
8.2.2 MRUnit 核心类	188		
8.2.3 编写一个 MRUnit 测试用例	189		
8.2.4 测试计数器	191		
8.2.5 MRUnit 的特性	194		
8.2.6 MRUnit 的局限性	194		

8.3 用 LocalJobRunner 测试.....	195	10.1.4 HiveQL 编译基础.....	217
8.3.1 setUp() 方法 .....	196	10.1.5 Hive 使用的概念 .....	218
8.3.2 LocalJobRunner 的局限性 .....	197	10.1.6 HiveQL 编译细节.....	222
8.4 用 MiniMRCluster 测试.....	198	10.1.7 数据定义语言 .....	226
8.4.1 配置开发环境 .....	198	10.1.8 数据操作语言 .....	226
8.4.2 MiniMRCluster 例子.....	199	10.1.9 扩展接口 .....	227
8.4.3 MiniMRCluster 的局限性 .....	201	10.1.10 Hive 脚本 .....	229
8.5 对访问网络资源的 MR 作业进行 测试.....	202	10.1.11 性能表现 .....	229
8.6 本章小结.....	202	10.1.12 整合 MapReduce.....	230
<b>第 9 章 Hadoop 的监控.....</b>	<b>203</b>	10.1.13 创建分区 .....	230
9.1 在 Hadoop MapReduce Jobs 中写 日志消息.....	203	10.1.14 用户定义函数 .....	232
9.2 在 Hadoop MapReduce Jobs 中查看 日志消息.....	206	<b>10.2 Impala.....</b>	<b>234</b>
9.3 在 Hadoop 2.x 中使用日志管理.....	208	10.2.1 Impala 架构 .....	234
9.3.1 Hadoop 2.x 中的日志存储.....	208	10.2.2 Impala 特性 .....	235
9.3.2 日志管理提升 .....	210	10.2.3 Impala 的局限 .....	235
9.3.3 使用基于 Web 的界面查看 日志 .....	210	<b>10.3 Shark.....</b>	<b>235</b>
9.3.4 命令行界面 .....	211	10.4 本章小结 .....	237
9.3.5 日志的保存 .....	211	<b>第 11 章 使用 Pig 进行数据处理 .....</b>	<b>238</b>
9.4 Hadoop 集群性能监控.....	211	11.1 Pig 简介.....	238
9.5 使用 YARN REST API .....	212	11.2 运行 Pig .....	240
9.6 使用供应商工具管理 Hadoop 集群.....	213	11.2.1 在 Grunt Shell 中执行 .....	241
9.7 本章小结.....	214	11.2.2 执行 Pig 脚本 .....	241
<b>第 10 章 使用 Hadoop 构建数据 仓库.....</b>	<b>215</b>	11.2.3 嵌入式 Java 程序 .....	242
10.1 Apache Hive.....	215	11.3 Pig Latin .....	243
10.1.1 安装 Hive .....	216	11.3.1 Pig 脚本中的注释 .....	243
10.1.2 Hive 的架构 .....	217	11.3.2 Pig 语句的执行 .....	243
10.1.3 元数据存储.....	217	11.3.3 Pig 命令 .....	244
11.4 UDF .....	249	11.4.1 Mapper 中的 Eval 函数 调用 .....	249
11.4.2 Reducer 中的 Eval 函数 调用 .....	250	11.4.2 Reducer 中的 Eval 函数 调用 .....	250
11.4.3 编写并使用自定义 Filter- Func .....	256	11.4.3 编写并使用自定义 Filter- Func .....	256





## 第 1 章

*Chapter 1*

# 为什么会有大数据

随着近 20 多年来计算技术的不断革新，企业积累了大量数据。数字传感器的进步使得通信系统越来越广泛的使用，尤其是移动平台和移动终端的飞速增长；系统运行产生的大量日志以及越来越多的企业采用无纸化办公的工作方式，这些情况都使得企业积攒起了海量数据资源。并且随着人们对现代科技越来越多的依赖，数据将会以更快的速度增长下去。

摩尔定律告诉我们，大约每隔两年，计算机的性能将历史性地提升一倍。最初，计算资源的提升速度满足了飞速增长的数据的处理需求。好景不长，到 2005 年前后，数据处理需求增长的速度已经快于计算资源处理能力的提升速度。

计算机工业界想到了另外一个更加经济有效的解决办法，就是数据的并行处理。既然单台计算机无法满足大数据量的计算需求，那就用多台计算机来并行处理这些海量数据。Hadoop 就是利用互相联网的多台计算机使用 MapReduce（一种改进的单指令多数据流 [SIMD] 计算技术）来并行地处理计算大量数据。

像亚马逊、谷歌、微软这样的公司为我们提供了较为廉价的云计算服务，使得分布式并行计算的概念更加深入人心，相比于购买计算资源，我们只要花少量的费用就可以租用相应的计算资源。

本书是使用 Hadoop 平台来开发和运行软件的实用指南。Hadoop 项目起初由 Apache Software Foundation 来负责管理，目前像 Cloudera、MapR 和 Hortonworks 这样的公司也加入到维护发展它的阵营中。本章会从整体上介绍大数据的由来以及 Hadoop 项目。

## 1.1 什么是大数据

就本书而言，大数据（Big Data）姑且定义为无法被符合服务等级协议（service level agreement，

SLA) 的单台计算机处理或(有些时候)存储的任何数据集。后续内容会有更精确的解释。理论上讲,单台计算机可以处理任意规模的数据。对于超过单台计算机存储量的海量数据,可以存放到类似网络附属存储(network attached storage, NAS)这样的共享存储设备中,然后输入到单台计算机去计算处理。但是,这样处理数据所花费的时间往往会大大超过允许的数据处理时间。

举个简单的例子。某个业务单元每次平均需要处理 200GB 的数据,数据的读取速率是每秒钟 50MB。如果按照每秒钟 50MB 的数据读取速率来计算,顺序读取硬盘上的 100MB 数据,需要两秒钟,读取全部的 200GB 数据就大约需要一个小时了。然而这 200GB 的数据要求在 5 分钟之内处理完毕。如果这个业务处理的 200GB 数据可以被平均分布到 100 个计算节点,每个节点处理分配的数据(举一个简单的用例,假设用 SALES\_YEAR>2001 这样一条简单的准则来划分数据集),去掉这 100 个计算节点的 CPU 处理时间和数据结果整合的时间,整个数据处理过程会在 1 分钟之内完成。

这个简单的例子说明,大数据是一个相对的概念,是相对于我们的业务需求而言的。

---

 **注意** Jeff Dean 博士有一篇关于并行计算的论文,详见 <http://www.cs.cornell.edu/projects/ladis2009/talks/dean-keynote-ladis2009.pdf>。从本地磁盘顺序读取 1MB 数据需要 2 千万纳秒。从速度为 1 Gbps 的网络中读取 1MB 的数据需要 2 亿 5 千万纳秒(假定传输 2KB 的数据需要 25 万纳秒,还要返回传输确认,这样每传输 2KB 数据需要网络往返共 50 万纳秒的时间)。虽然文中的描述稍显过时,数据的传输读取绝对速度已经有了提高,但是这不妨碍我们在本章中拿来做比照,因为这些数据之间的相对值变化不大。

---

## 1.2 大数据技术背后的核心思想

上文中的例子我们作了诸多假设,要表明的核心问题是虽然我们可以很快地处理数据,但是从持久性的存储设备中读取的速度受到限制,这是整个数据处理流程上的关键瓶颈所在。相对于读写本地节点存储设备上的数据,通过网络来传输数据会更慢。

下面列出了所有大数据处理方法中的一些共同特征:

- 数据分布在多个节点(网络 I/O 速度 << 本地磁盘 I/O 速度)。
- 计算程序离数据更近(集群上的节点),而不是相反。
- 数据的处理尽量在本地完成(网络 I/O 速度 << 本地磁盘 I/O 速度)。
- 使用可顺序读取磁盘 I/O 代替随机读取磁盘 I/O(数据交换速度 << 数据寻道时间)。

所有大数据计算处理模式都有一个目的,就是使输入 / 输出(I/O) 并行化,从而提高数据处理性能。

### 1.2.1 把数据分发到多个节点

根据定义,大数据是无法仅靠单台计算机资源来处理的。大数据的一个特点就是使用商

用服务器。一台典型的商用服务器拥有 2TB 至 4TB 的磁盘。因为大数据是指远远超出这个容量的数据集，所以还是要将数据分布在多个节点上。

需要注意的是，我们把要处理的数据分布存放在多个计算节点，并不是因为这些数据的数据量有数十 T 之巨。你会发现大数据处理系统会有条不紊地在各个计算节点上计算处理所分配的数据。把数据分布到各个服务器节点的最终目的是为了让大量计算节点同时参与到数据的处理计算过程中来。哪怕仅有 500GB 的数据集，这么大的数据量，集群中的单台计算节点完全可以存储，也会把数据分发到多台计算节点。这样的数据分发有两点好处：

- 每个数据块会在多个节点上有多份拷贝（Hadoop 默认是一个数据块有 3 份拷贝），这使得系统具备容错性，当一个节点发生故障，其他节点还备份有故障节点上的数据。
- 为了达到数据并行处理的目的，多个节点可以同时参与数据处理过程。比如，50GB 的数据被分配到 10 台计算节点，每个计算节点仅处理分配给自己的数据集，可以达到数据处理性能 5 ~ 10 倍的提升。读者可能会问，为什么不把这些数据存放到网络文件系统（NFS）中，这样每个节点可以去读取它要处理的部分。答案是本地存储磁盘的数据读取速度要远远高于通过网络来读取数据的速度。当数据分析任务启动的时候，分析任务程序运行所需的函数库被拷贝到各个计算节点，大数据处理系统会尽量在本地处理计算数据。我们会在后面章节中详谈。

## 1.2.2 把计算逻辑移动到数据附近

对于我们这些精通 J2EE 编程的人来说，三层架构思想深植脑海。在三层编程模型中，所有的数据会通过网络集中到一起，交由应用层来处理。我们由此形成了固有的观念，就是数据应该是分散的，而程序应该是集中的。

大数据系统无法处理网络过载的问题。传输动辄数 T 的数据量给应用层，使得网络带宽耗尽，网络拥挤不堪，导致传输效率大幅下降，甚至有可能导致系统故障。从大数据的观念来看，应该把数据分布存放到各个计算节点，程序也要移动到数据附近。要做到这一点，是一件很不容易的事情。除了程序要移动到存放数据的节点，程序运行所依赖的函数库也要移动到数据处理节点才行。如果大数据系统的集群拥有数百个计算节点，显然那将是程序维护 / 部署人员的噩梦。所以，大数据系统可以让我们集中式地部署程序代码，大数据系统后台会在计算任务启动之前把这些程序移动到各个数据处理节点。

## 1.2.3 计算节点进行本地数据处理

前文提到的大数据系统的两个特点，决定了分配到计算节点的数据要在计算节点本地处理。所有的大数据编程模型都是基于分布式和并行处理的。网络 I/O 比本地磁盘 I/O 慢了好几个数量级。数据被分发到各个计算节点，程序运行依赖库也移动到了数据所在节点，计算节点就地计算处理数据的条件完备了。

虽然典型的大数据处理系统都希望把数据处理过程放在拥有数据的节点本地完成，但并

