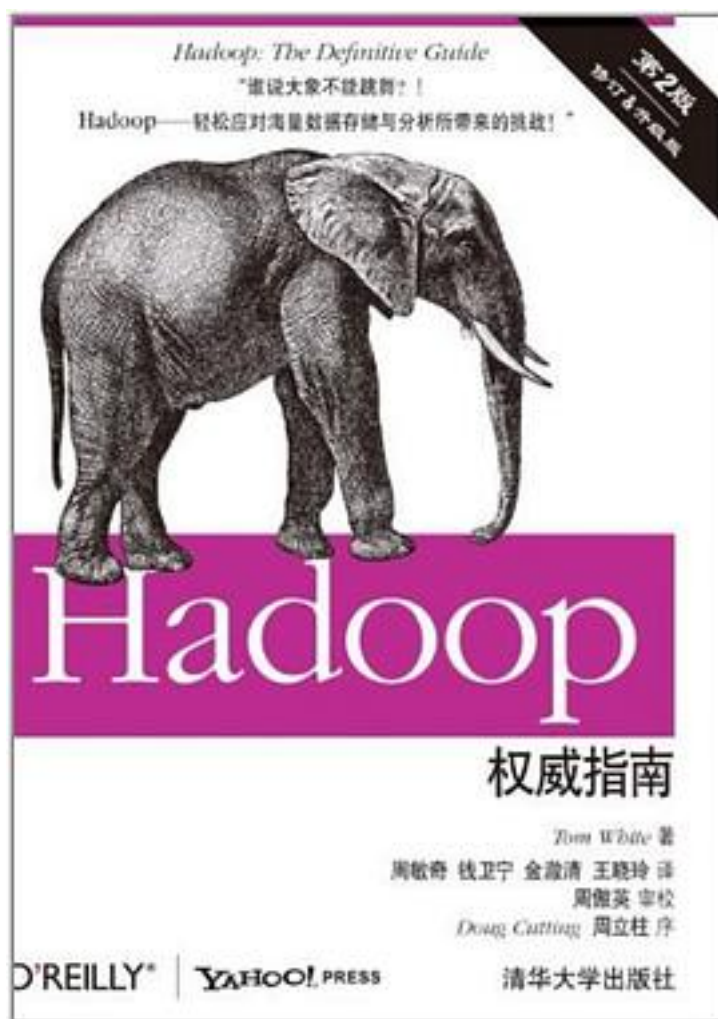


# Hadoop权威指南（第2版）



[Hadoop权威指南（第2版）\\_下载链接1](#)

著者:Tom White

出版者:清华大学出版社

出版时间:2011-7

装帧:平装

isbn:9787302257585

《Hadoop权威指南(第2版)(修订·升级版)》从Hadoop的缘起开始，由浅入深，结合理

论和实践，全方位地介绍Hadoop这一高性能处理海量数据集的理想工具。全书共16章，3个附录，涉及的主题包括：Hadoop简介；MapReduce简介；Hadoop分布式文件系统；Hadoop的I/O、MapReduce应用程序开发；MapReduce的工作机制；MapReduce的类型和格式；MapReduce的特性；如何构建Hadoop集群，如何管理Hadoop；Pig简介；Hbase简介；Hive简介；ZooKeeper简介；开源工具Sqoop，最后还提供了丰富的案例分析。

《Hadoop权威指南(第2版)(修订·升级版)》是Hadoop权威参考，程序员可从中探索如何分析海量数据集，管理员可以从中了解如何安装与运行Hadoop集群。

作者介绍:

Tom White从2007年以来，一直担任Apache Hadoop项目负责人。他是Apache软件基金会的成员之一，同时也是Cloudera的一名工程师。Tom为oreilly.com、java.net和IBM的developerWorks写过大量文章，并经常在很多行业大会上发表演讲。

目录: 第1章 初识Hadoop  
数据！数据！  
数据存储与分析  
与其他系统相比  
关系型数据库管理系统  
网格计算  
志愿计算  
1.3.4 Hadoop 发展简史  
Apache Hadoop和Hadoop生态圈  
第2章 关于MapReduce  
一个气象数据集  
数据的格式  
使用Unix工具进行数据分析  
使用Hadoop分析数据  
map阶段和reduce阶段  
横向扩展  
合并函数  
运行一个分布式的MapReduce作业  
Hadoop的Streaming  
Ruby版本  
Python版本  
Hadoop Pipes  
编译运行  
第3章 Hadoop分布式文件系统  
HDFS的设计  
HDFS的概念  
数据块  
namenode和datanode  
命令行接口  
基本文件系统操作  
Hadoop文件系统  
接口  
Java接口  
从Hadoop URL中读取数据  
通过FileSystem API读取数据

写入数据  
目录  
查询文件系统  
删除数据  
数据流  
文件读取剖析  
文件写入剖析  
一致模型  
通过 distcp并行拷贝  
保持 HDFS 集群的均衡  
Hadoop的归档文件  
使用Hadoop归档文件  
不足  
第4章 Hadoop I/O  
数据完整性  
HDFS的数据完整性  
LocalFileSystem  
ChecksumFileSystem  
压缩  
codec  
压缩和输入切分  
在MapReduce中使用压缩  
序列化  
Writable接口  
Writable类  
实现定制的Writable类型  
序列化框架  
Avro  
依据文件的数据结构  
写入SequenceFile  
MapFile  
第5章 MapReduce应用开发  
配置API  
合并多个源文件  
可变的扩展  
配置开发环境  
配置管理  
辅助类GenericOptionsParser， Tool和ToolRunner  
编写单元测试  
mapper  
reducer  
本地运行测试数据  
在本地作业运行器上运行作业  
测试驱动程序  
在集群上运行  
打包  
启动作业  
MapReduce的Web界面  
获取结果  
作业调试  
使用远程调试器  
作业调优  
分析任务  
MapReduce的工作流

将问题分解成MapReduce作业  
运行独立的作业  
第6章 MapReduce的工作机制  
剖析MapReduce作业运行机制  
作业的提交  
作业的初始化  
任务的分配  
任务的执行  
进度和状态的更新  
作业的完成  
失败  
任务失败  
tasktracker失败  
jobtracker失败  
作业的调度  
Fair Scheduler  
Capacity Scheduler  
shuffle和排序  
map端  
reduce端  
配置的调优  
任务的执行  
推测式执行  
重用JVM  
跳过坏记录  
任务执行环境  
第7章 MapReduce的类型与格式  
MapReduce的类型  
默认的MapReduce作业  
输入格式  
输入分片与记录  
文本输入  
二进制输入  
多种输入  
数据库输入（和输出）  
输出格式  
文本输出  
二进制输出  
多个输出  
延迟输出  
数据库输出  
第8章 MapReduce的特性  
计数器  
内置计数器  
用户定义的Java计数器  
用户定义的Streaming计数器  
排序  
准备  
部分排序  
总排序  
二次排序  
联接  
map端联接  
reduce端联接

边数据分布  
利用JobConf来配置作业  
分布式缓存  
MapReduce库类  
第9章 构建Hadoop集群  
集群规范  
网络拓扑  
集群的构建和安装  
安装Java  
创建Hadoop用户  
安装Hadoop  
测试安装  
SSH配置  
Hadoop配置  
配置管理  
环境设置  
Hadoop守护进程的关键属性  
Hadoop守护进程的地址和端口  
Hadoop的其他属性  
创建用户帐号  
安全性  
Kerberos和Hadoop  
委托令牌  
其他安全性改进  
利用基准测试程序测试Hadoop集群  
Hadoop基准测试程序  
用户的作业  
云上的Hadoop  
Amazon EC2上的Hadoop  
第10章 管理Hadoop  
HDFS  
永久性数据结构  
安全模式  
日志审计  
工具  
监控  
日志  
度量  
Java管理扩展（JMX）  
维护  
日常管理过程  
委任节点和解除节点  
升级  
第11章 Pig简介  
安装与运行Pig  
执行类型  
运行Pig程序  
Grunt  
Pig Latin编辑器  
示例  
生成示例  
与数据库比较  
PigLatin  
结构

- 语句
- 表达式
- 1.4.4 类型
- 模式
- 函数
- 用户自定义函数
- 过滤UDF
- 计算UDF
- 加载UDF
- 数据处理操作
- 加载和存储数据
- 过滤数据
- 分组与连接数据
- 对数据进行排序
- 组合和分割数据
- Pig实战
- 并行处理
- 参数代换
- 第12章 Hive
- 1.1 安装Hive
- 1.1.1 Hive外壳环境
- 1.2 示例
- 1.3 运行Hive
- 1.3.1 配置Hive
- 1.3.2 Hive服务
- 1.3.3 Metastore
- 1.4 和传统数据库进行比较
- 1.4.1 读时模式（Schema on Read）vs.写时模式（Schema on Write）
- 1.4.2 更新、事务和索引
- 1.5 HiveQL
- 1.5.1 数据类型
- 1.5.2 操作和函数
- 1.6 表
- 1.6.1 托管表（Managed Tables）和外部表（External Tables）
- 1.6.2 分区（Partitions）和桶（Buckets）
- 1.6.3 存储格式
- 1.6.4 导入数据
- 1.6.5 表的修改
- 1.6.6 表的丢弃
- 1.7 查询数据
- 1.7.1 排序（Sorting）和聚集（Aggregating）
- 1.7.2 MapReduce脚本
- 1.7.3 连接
- 1.7.4 子查询
- 1.7.5 视图（view）
- 1.8 用户定义函数（User-Defined Functions）
- 1.8.1 编写UDF
- 1.8.2 编写UDAF
- 第13章 HBase
- 2.1 HBasics
- 2.1.1 背景
- 2.2 概念
- 2.2.1 数据模型的“旋风之旅”
- 2.2.2 实现

- 2.3 安装
  - 2.3.1 测试驱动
- 2.4 客户机
  - 2.4.1 Java
  - 2.4.2 Avro, REST, 以及Thrift
- 2.5 示例
  - 2.5.1 模式
  - 2.5.2 加载数据
  - 2.5.3 Web查询
- 2.6 HBase和RDBMS的比较
  - 2.6.1 成功的服务
  - 2.6.2 HBase
  - 2.6.3 实例: HBase在Streamy.com的使用
- 2.7 Praxis
  - 2.7.1 版本
  - 2.7.2 HDFS
  - 2.7.3 用户接口 (UI)
  - 2.7.4 度量 (metrics)
  - 2.7.5 模式设计
  - 2.7.6 计数器
  - 2.7.7 批量加载 (bulkloading)
- 第14章 ZooKeeper
  - 安装和运行ZooKeeper
  - 示例
  - ZooKeeper中的组成员关系
  - 创建组
  - 加入组
  - 列出组成员
  - ZooKeeper服务
  - 数据模型
  - 操作
  - 实现
  - 一致性
  - 会话
  - 状态
  - 使用ZooKeeper来构建应用
  - 配置服务
  - 具有可恢复性的ZooKeeper应用
  - 锁服务
  - 生产环境中的ZooKeeper
  - 可恢复性和性能
  - 配置
- 第15章 开源工具Sqoop
  - 获取Sqoop
  - 一个导入的例子
  - 生成代码
  - 其他序列化系统
  - 深入了解数据库导入
  - 导入控制
  - 导入和一致性
  - 直接模式导入
  - 使用导入的数据
  - 导入的数据与Hive
  - 导入大对象

执行导出  
深入了解导出  
导出与事务  
导出和SequenceFile  
第16章 实例分析  
Hadoop 在Last.fm的应用  
Last.fm：社会音乐史上的革命  
Hadoop a Last.fm  
用Hadoop产生图表  
Track Statistics程序  
总结  
Hadoop和Hive在Facebook的应用  
概要介绍  
Hadoop a Facebook  
假想的使用情况案例  
Hive  
问题与未来工作计划  
Nutch 搜索引擎  
背景介绍  
数据结构  
Nutch系统利用Hadoop进行数据处理的精选实例  
总结  
Rackspace的日志处理  
简史  
选择Hadoop  
收集和存储  
日志的MapReduce模型  
关于Cascading  
字段、元组和管道  
操作  
Tap类，Scheme对象和Flow对象  
Cascading实战  
灵活性  
Hadoop和Cascading在ShareThis的应用  
总结  
在Apache Hadoop上的TB字节数量级排序  
使用Pig和Wukong来探索10亿数量级边的 网络图  
测量社区  
每个人都在和我说话：Twitter回复关系图  
(度) degree  
对称链接  
社区提取  
附录A 安装Apache Hadoop  
附录B Cloudera's Distribution for Hadoop  
附录C 准备NCDC天气数据  
索引  
• • • • • ([收起](#))

[Hadoop权威指南（第2版）](#) [下载链接1](#)



## 标签

hadoop

分布式

MapReduce

云计算

大数据

计算机

O'Reilly

编程

## 评论

hadoop真牛逼

---

很快地读完了，没特别了解具体的语法，对于一个宏观上的掌握以及技术选型来说足够了~没有学过Hadoop的人值得一读。

---

如果一本书是由好几个人合伙翻译的，那么译者通常会漏掉Google，请务必做好精分的准备。

---

好书烂翻译

-----  
感觉这本书字体比较小，印刷的间距比较大，内容上比第1版没什么大的更新——那个SQL导入到Hadoop的工具倒是很有意思

-----  
简单读过

-----  
中规中矩吧，理论科普和实操手册，仅此

-----  
暂时不深究了。

-----  
科普，并非深入

-----  
翻译质量实在是不敢恭维…

-----  
太晦涩了，需要看多遍。。

-----  
只有看完之后以后才可能知道有用的部分，第一遍不要求看得多深入，只求知道有用的特性，在以后需要时能回想起来，这也许是阅读这样的工具书的有效方法；除去中文版的翻译问题，整本书还不错；鉴于内容过于啰嗦的问题，建议新手先看《Hadoop实战》，等有了一定的经验之后再看这本书，这样可以保证不陷入繁琐的细节又能增加涉猎。

-----  
知道翻译的不好，也知道可能会看起来很痛苦，还是买了，看英文实在太慢。没看完。

-----

比第一版增加的HIVE和Sqoop

-----  
工具书，粗读一遍，不求甚解

-----  
暂时还用不上，细看了前10章以及第14章的Zookeeper。

-----  
当手册用还OK，不适合入门

-----  
书很权威，这本翻译也很好！译本的页数与原本的一直，这对于对照阅读很有好处！囫圇吞枣的看了一遍！大概算是了解了hadoop以及相关项目是干嘛的吧！完全分布式模式还没搭建起来！还要继续努力啊！

-----  
书的内容不错，学习Hadoop的经典，但是翻译的太粗糙了！

-----  
O'REILLY的书还是很不错的

-----  
[Hadoop权威指南（第2版）\\_下载链接1\\_](#)

## 书评

买了第一版，时间太紧，没来得及看，后来出了个号称修订升级的第二版，毫不犹豫又买了，后来听说第二版比第一版翻译得好，心中窃喜，再后来看了第二版，我震惊了，我TM就是一傻子，放着好好的英文版不看，赶什么时髦买中文版呢。在这个神奇的国度，牛奶里放的是三聚氰胺，火腿...

-----  
其实也不算全部读完了，读它主要是为了技术选型，考虑升级持久层架构、提高系统可

扩展性，仔细研读了前几章，对Hadoop、MapReduce、HDFS的模型、机制、使用场景有了一定了解。后面几章及其生态圈内的其他项目抱着了解的心态简单浏览了一下。整体感觉还行，至少从我看过的章节来...

-----  
中文版412页：

所以理论上，任何东西都可以表示成二进制形式，然后转化成为长整型的字符串或直接对数据结构进行序列化，来作为键值。原文460页： ..., so theoretically anything can serve as row key, from strings to binary representations of long or even serialized ...

-----  
-- china-pub 赠书活动 -- <http://www.douban.com/group/topic/20965935/>  
一直比较忙，整本书还没读完，只是粗略翻了个大概，其中有两三章细读了一遍。先做个大体评价吧，有时间全部细读后再评论。  
从书的内容上来讲，大致上与网上该书的内容介绍一致。简单点概括：这本书对...

-----  
参加豆瓣China-pub抽奖，比较幸运的得到这本Hadoop权威指南中文第二版，拿来与第一版相比，发现新加入了Hive和Sqoop章节，译文质量也提高了不少，并且保留了英文索引。  
这本书对Hadoop的介绍还算全面，有实践冲动的朋友基本可以拿着书、配合Google百度马上实现梦想。个人感觉 “...

-----  
看了几章中文版的，各种错误，太低级，实在是看不下去了。建议还是看原版吧。译者们的脸皮可真厚，英文译不明白也就罢了，中文都组织的不通顺，好意思吗！！什么叫 “但是，.....，但是” 啊，“但是体” 啊。

-----  
很好的Hadoop教程，比Apache和Yahoo  
!网页版guide详细很多，很多想不明白的Hadoop实现细节都可以在这本书里找到。

-----  
详见：<http://www.cnblogs.com/aprilrain/archive/2013/03/07/2947664.html>  
-----

很多地方翻译的不行，需要对照英文看才能明白。。。不过对于快速学习，仍然是不错的选择。建议译者看看每部分内容的重要性，不重要的瞎翻翻就算了，重要的部分还是好好花点功夫，不要本末倒置了。比如第三章的数据流部分，这么经典的地方居然被翻译烂的一塌糊涂。不知道译者会...

书中没有透露太多实现架构方面的细节，更多的是从使用者的角度上介绍了Hadoop的各种知识，包括MapReduce, HDFS, Hive, Pig, HBase, ZooKeeper。几乎涉及了Hadoop的所有关于使用方面的知识，包括安装和使用。你甚至可以直接在自己的电脑上装上一个Hadoop，对着书中的例子实际演...

是我遇到过的翻译最烂的一本书，在译者的“妙语连珠”里折腾了半个钟头就再也没兴趣了。略举几例如下： P.6 任然 -> 仍然 P.21 输入键（为什么不像后面那样有个“的”？），输入的值，输出的键…… P.27 “计数器” (Counter)，译文附原文； "Context Object"(上下文对象)，原...

专门登录来评论的，翻译也太烂了吧，真的真的建议强烈英语阅读能力好的人去读原版书，不要花冤枉钱在这上面，除了文字错误外，里边的图居然也有错，就比如260页的图最后两个年份应该是1901结果这里竟然是1900，我是真滴服了，一本神书被翻译成这样，作者得气死。zsbd zsbd zsbd...

你的履历添了一笔<hadoop权威指南>译者,但是你不配 这是我见过的最不用心的翻译, 字里行间行文不通顺, 请别勉强自己,map reduce shuffle机制都没翻译的好 虽然原作者写作功底也实在是一般 第 1 2 5 6 7 这几章 翻译的实在是太烂了 请不要呐Google翻译糊弄人阿 误人子弟 ...

首先，翻译太差，很多句子就是瞎翻，根本不通顺，很多时候你要停下来断句，慢慢去理解。然后，这本书是很多人去翻译的，很多人连代码都不懂，曾经一段代码看到我蒙圈，去看了一下源代码，好家伙，四行有五个错误。另外，从代码瞎缩进也可以看出这是群没写过代码的人翻的，而且...

-----

-----

Cobub Razor APP数据统计分析工具官网上有篇文章是讲Hadoop Yarn调度器的选择和使用的，我觉得写的挺好的，推荐<http://www.cobub.com/the-selection-and-use-of-hadoop-yarn-scheduler/>

-----

[Hadoop权威指南（第2版）\\_下载链接1\\_](#)