

# Hadoop 2.X

# HDFS源码剖析

徐鹏 著

电子工业出版社

Publishing House of Electronics Industry

北京•BEIJING

## 内 容 简 介

本书以 Hadoop 2.6.0 源码为基础,深入剖析了 HDFS 2.X 中各个模块的实现细节,包括 RPC 框架实现、Namenode 实现、Datanode 实现以及 HDFS 客户端实现等。本书一共有 5 章,其中第 1 章从总体上介绍了 HDFS 的组件、概念以及典型的流程,同时详细介绍了 HDFS 各个组件间 RPC 接口的定义。第 2 章介绍了 Hadoop RPC 框架的实现,Hadoop RPC 是 HDFS 各个组件间通信所依赖的底层框架,可以理解为 HDFS 的神经系统。第 3~5 章分别介绍了 Namenode、Datanode 以及 HDFS 客户端这三个组件的实现细节,同时穿插介绍了 HDFS 2.X 的新特性,例如 Namenode HA、Federation Namenode 等。

阅读本书可以帮助读者从架构设计与源码实现角度了解 HDFS 2.X,同时还能学习 HDFS 2.X 框架中优秀的设计思想、设计模式、Java 语言技巧以及编程规范等。这些对于读者全面提高自己的技术水平有很大的帮助。

未经许可,不得以任何方式复制或抄袭本书之部分或全部内容。  
版权所有,侵权必究。

### 图书在版编目(CIP)数据

Hadoop 2.X HDFS 源码剖析 / 徐鹏著. —北京: 电子工业出版社, 2016.3  
ISBN 978-7-121-28155-6

I. ①H… II. ①徐… III. ①分布式文件系统—研究 IV. ①TP316

中国版本图书馆 CIP 数据核字 (2016) 第 027311 号

策划编辑: 张春雨

责任编辑: 葛 娜

印 刷: 三河市双峰印刷装订有限公司

装 订: 三河市双峰印刷装订有限公司

出版发行: 电子工业出版社

北京市海淀区万寿路 173 信箱 邮编: 100036

开 本: 787×980 1/16 印张: 35.25 字数: 879 千字

版 次: 2016 年 3 月第 1 版

印 次: 2016 年 3 月第 1 次印刷

定 价: 108.00 元

凡所购买电子工业出版社图书有缺损问题, 请向购买书店调换。若书店售缺, 请与本社发行部联系, 联系及邮购电话: (010) 88254888。

质量投诉请发邮件至 [zlts@phei.com.cn](mailto:zlts@phei.com.cn), 盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

服务热线: (010) 88258888。



# 前言

今天 Hadoop 已经成为大数据处理中十分重要的平台，一个以 Hadoop 为基础的活跃的开源生态圈已经逐步形成，Hadoop 的应用也由最初的互联网行业发展到金融行业、电信行业、IT 设备商以及数以万计的中小企业。Hadoop 的 HDFS 组件（Hadoop 分布式文件系统）以及 MapReduce 组件分别为上层框架提供了分布式存储和计算的能力。

HDFS 作为 Hadoop 中解决分布式存储的基础组件，最早是根据 GFS (Google File System) 论文的概念模型来设计实现的。然而，随着 HDFS 上层框架的丰富以及应用场景的扩展，用户对 HDFS 的功能、性能、稳定性、扩展性以及可靠性的要求越来越高，HDFS 2.X 版本也就应运而生。相对于 HDFS 1.X，HDFS 2.X 版本提出了很多振奋人心的新特性，如 Namenode HA、Federation Namenode、集中式缓存、快照等。但令人惋惜的是，至今还没有一本能够深入剖析 HDFS 2.X 内部实现细节，以及介绍 HDFS 2.X 新特性的书籍。本书的出现填补了上述空白，它是国内第一本深入剖析 HDFS 2.X 源码实现的书籍。

本书以 Hadoop 2.6.0 源码为基础，深入剖析了 HDFS 2.X 中各个模块的实现细节，包括 RPC 框架实现、Namenode 实现、Datanode 实现以及 HDFS 客户端实现等。阅读本书可以帮助读者从架构设计与源码实现角度了解 HDFS 2.X，同时还能学习 HDFS 2.X 框架中优秀的设计思想、设计模式、Java 语言技巧以及编程规范等。这些对于读者全面提高自己的技术水平有很大的帮助。

## 如何阅读本书

由于篇幅原因，本书并没有介绍 HDFS 实现中的一些基础知识，例如 Java NIO、动态代理、protobuf 等。而是直接切入源码分析 HDFS 的设计与实现，同时介绍了一些经典的设计模式、Java 语言技巧在 HDFS 实现中的应用。希望读者在阅读本书之前，先搭建好源码环境，并了解相应的基础知识，这样学习效果会更好。

本书一共有 5 章，相互之间的联系比较紧密，有联系的小节都有注释标注，读者可以根据注释跳跃阅读。

第 1 章是 HDFS 概述，从总体上介绍了 HDFS 的组件、概念以及典型的流程，同时详细介绍了 HDFS 各个组件间 RPC 接口的定义。由于 HDFS 流程大都比较复杂，往往涉及多个组件的配合，读者在阅读后续的具体章节时，可以查阅本章内容获取某个流程的总体描述，以



及 RPC 接口的具体定义。

第 2 章介绍了 Hadoop RPC 框架的实现，Hadoop RPC 是 HDFS 各个组件间通信所依赖的底层框架，可以理解为 HDFS 的神经系统。通过阅读本章，读者可以学习到一个典型的分布式 RPC 框架的实现细节，在本章中会介绍较多的设计模式，编程模型以及语言技巧，请读者注重积累。

第 3~5 章分别介绍了 Namenode、Datanode 以及 HDFS 客户端这三个组件的实现细节，同时穿插介绍了 HDFS 2.X 的新特性，例如 Namenode HA 就放在 Namenode 章介绍，而 Federation Namenode 新特性对 Datanode 的修改比较多，所以就放在 Datanode 章介绍。

读者在阅读本书的过程中，如果发现有任何不当之处，烦请您将意见和建议发往邮箱 [xupeng.bupt@gmail.com](mailto:xupeng.bupt@gmail.com)，不胜感激。

## 本书代码

本书分析的代码版本为 Hadoop 2.6.0，书中部分较长的代码做了省略，完整的代码请从官网 <http://hadoop.apache.org> 下载。

## 致谢

感谢互联网，感谢开源软件，感谢 Hadoop 社区，感谢本书引用文献的所有原作者，是你们为 Hadoop 爱好者打开了一扇大门。

感谢电子工业出版社博文视点的张春雨老师，是您的信任使得这本书的出版成为可能。同时还要感谢许多我不知道名字的编辑为本书最终出版所做的付出和努力。

感谢丁雷在百忙之中抽出时间对本书提出许多建设性意见，同时感谢左谱军、张德阳、闫飞翔以及张涛对本书的审阅。

2015 年写这本书时正是自己很困难的一段时期，我很感恩有许多朋友在生活、工作上给我帮助以及包容。感谢吴佳宁、刘文博，你们一直是第一个伸出援手的哥们。感谢郑晓彤、袁玮、贺子昂、李强、和紫东、刘丹、何一舟为我引荐机会，谢谢你们。

特别感谢远见，书的撰写过程是如此漫长，是你在这段时间里把自信、阳光和快乐传播给我，让我更加积极、勇敢和有信心。没有你，这本书永远无法完成。

最后感谢我的父母和妹妹，谢谢你们默默为我做出的牺牲和付出，你们永远是我前进的动力。



# 目 录

第 1 章	HDFS	1
1.1	HDFS 概述	1
1.1.1	HDFS 体系结构	1
1.1.2	HDFS 基本概念	2
1.2	HDFS 通信协议	4
1.2.1	Hadoop RPC 接口	4
1.2.2	流式接口	20
1.3	HDFS 主要流程	22
1.3.1	HDFS 客户端读流程	22
1.3.2	HDFS 客户端写流程	24
1.3.3	HDFS 客户端追加写流程	25
1.3.4	Datanode 启动、心跳以及执行名字节点指令流程	26
1.3.5	HA 切换流程	27
第 2 章	Hadoop RPC	29
2.1	概述	29
2.1.1	RPC 框架概述	29
2.1.2	Hadoop RPC 框架概述	30
2.2	Hadoop RPC 的使用	36
2.2.1	Hadoop RPC 使用概述	36
2.2.2	定义 RPC 协议	40
2.2.3	客户端获取 Proxy 对象	45
2.2.4	服务器获取 Server 对象	54
2.3	Hadoop RPC 实现	63
2.3.1	RPC 类实现	63
2.3.2	Client 类实现	64
2.3.3	Server 类实现	76

第 3 章	Namenode (名字节点)	88
3.1	文件系统目录树	88
3.1.1	INode 相关类	89
3.1.2	Feature 相关类	102
3.1.3	FSEditLog 类	117
3.1.4	FSImage 类	138
3.1.5	FSDirectory 类	158
3.2	数据块管理	162
3.2.1	Block、Replica、BlocksMap	162
3.2.2	数据块副本状态	167
3.2.3	BlockManager 类 (done)	177
3.3	数据节点管理	211
3.3.1	DatanodeDescriptor	212
3.3.2	DatanodeStorageInfo	214
3.3.3	DatanodeManager	217
3.4	租约管理	233
3.4.1	LeaseManager.Lease	233
3.4.2	LeaseManager	234
3.5	缓存管理	246
3.5.1	缓存概念	247
3.5.2	缓存管理命令	247
3.5.3	HDFS 集中式缓存架构	247
3.5.4	CacheManager 类实现	248
3.5.5	CacheReplicationMonitor	250
3.6	ClientProtocol 实现	251
3.6.1	创建文件	251
3.6.2	追加写文件	254
3.6.3	创建新的数据块	257
3.6.4	放弃数据块	265
3.6.5	关闭文件	266
3.7	Namenode 的启动和停止	268
3.7.1	安全模式	268
3.7.2	HDFS High Availability	276
3.7.3	名字节点的启动	301
3.7.4	名字节点的停止	306



第 4 章 Datanode (数据节点)	307
4.1 Datanode 逻辑结构	307
4.1.1 HDFS 1.X 架构	307
4.1.2 HDFS Federation	308
4.1.3 Datanode 逻辑结构	310
4.2 Datanode 存储	312
4.2.1 Datanode 升级机制	312
4.2.2 Datanode 磁盘存储结构	315
4.2.3 DataStorage 实现	317
4.3 文件系统数据集	334
4.3.1 Datanode 上数据块副本的状态	335
4.3.2 BlockPoolSlice 实现	335
4.3.3 FsVolumeImpl 实现	342
4.3.4 FsVolumeList 实现	345
4.3.5 FsDatasetImpl 实现	348
4.4 BlockPoolManager	375
4.4.1 BPServiceActor 实现	376
4.4.2 BPOfferService 实现	389
4.4.3 BlockPoolManager 实现	396
4.5 流式接口	398
4.5.1 DataTransferProtocol 定义	398
4.5.2 Sender 和 Receiver	399
4.5.3 DataXceiverServer	403
4.5.4 DataXceiver	406
4.5.5 读数据	408
4.5.6 写数据 (done)	423
4.5.7 数据块替换、数据块拷贝和读数据块校验	437
4.5.8 短路读操作	437
4.6 数据块扫描器	437
4.6.1 DataBlockScanner 实现	438
4.6.2 BlockPoolSliceScanner 实现	439
4.7 DirectoryScanner	442
4.8 DataNode 类的实现	443
4.8.1 DataNode 的启动	444
4.8.2 DataNode 的关闭	446

- 第 5 章 HDFS 客户端 ..... 447
  - 5.1 DFSCClient 实现 ..... 447
    - 5.1.1 构造方法 ..... 448
    - 5.1.2 关闭方法 ..... 449
    - 5.1.3 文件系统管理与配置方法 ..... 450
    - 5.1.4 HDFS 文件与目录操作方法 ..... 451
    - 5.1.5 HDFS 文件读写方法 ..... 452
  - 5.2 文件读操作与输入流 ..... 452
    - 5.2.1 打开文件 ..... 452
    - 5.2.2 读操作——DFSInputStream 实现 ..... 461
  - 5.3 文件短路读操作 ..... 481
    - 5.3.1 短路读共享内存 ..... 482
    - 5.3.2 DataTransferProtocol ..... 484
    - 5.3.3 DFSCClient 短路读操作流程 ..... 488
    - 5.3.4 Datanode 短路读操作流程 ..... 509
  - 5.4 文件写操作与输出流 ..... 512
    - 5.4.1 创建文件 ..... 512
    - 5.4.2 写操作——DFSOutputStream 实现 ..... 516
    - 5.4.3 追加写操作 ..... 543
    - 5.4.4 租约相关 ..... 546
    - 5.4.5 关闭输出流 ..... 548
  - 5.5 HDFS 常用工具 ..... 549
    - 5.5.1 FsShell 实现 ..... 550
    - 5.5.2 DFSAdmin 实现 ..... 552