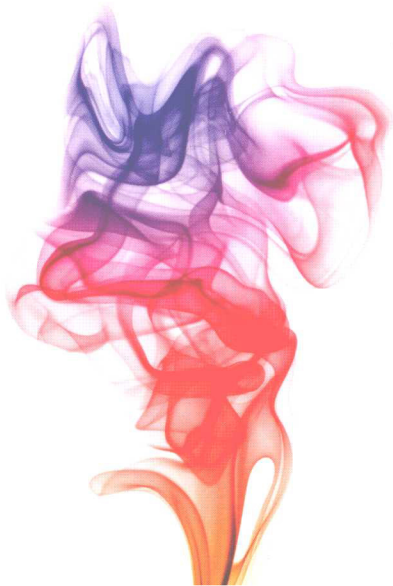


多位Spark的贡献者和专家联袂推荐，详细剖析Spark内核各个模块并辅以相应源码解析的著作。以源码为基础，全面分析Spark内核各个模块的设计思想和实现原理，深入理解其内部运作机制乃至实现细节。帮助Spark领域的从业人员全面掌握Spark核心技术，进而在应用开发中做到游刃有余和性能调优时做到有的放矢。



Spark Internals
Design and Implement Principle of Spark Core

Spark技术内幕

深入解析Spark内核
架构设计与实现原理

张安站◎著



机械工业出版社
China Machine Press



技术丛书

Spark Internals

Design and Implement Principle of Spark Core

Spark技术内幕

**深入解析Spark内核
架构设计与实现原理**

张安站◎著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

Spark 技术内幕: 深入解析 Spark 内核架构设计与实现原理 / 张安站著. —北京: 机械工业出版社, 2015.8
(大数据技术丛书)

ISBN 978-7-111-50964-6

I. S… II. 张… III. 数据处理软件 IV. TP274

中国版本图书馆 CIP 数据核字 (2015) 第 170549 号

Spark 是不断壮大的大数据分析解决方案家族中备受关注的新成员。它不仅为分布式数据集的处理提供了一个有效框架, 而且以高效的方式处理分布式数据集。它支持实时处理、流处理和批处理, 提供了统一的解决方案, 因此极具竞争力。本书以源码为基础, 深入分析 Spark 内核的设计理念和架构实现, 系统讲解各个核心模块的实现, 为性能调优、二次开发和系统运维提供理论支持, 为更好地使用 Spark Streaming、MLlib、Spark SQL 和 GraphX 等奠定基础。

Spark 技术内幕

深入解析 Spark 内核架构设计与实现原理

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: 李 艺

责任校对: 殷 虹

印 刷: 北京瑞德印刷有限公司

版 次: 2015 年 9 月第 1 版第 1 次印刷

开 本: 186mm × 240mm 1/16

印 张: 13.5

书 号: ISBN 978-7-111-50964-6

定 价: 49.00 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

Foreword 序

Apache Spark 项目的高速发展超出了很多人的预期。在 2009 年到 2013 年，Spark 还是 UC Berkeley 大学 AMPLab 的一个研究项目，因为其架构设计得简洁和高效，逐渐吸引了工业界和学术界的广泛关注。我还记得 2013 年 2 月份在 Santa Clara 召开的 Strata Conference 上，虽然是长达一整天的 Spark 技术培训，大厅里还是人满为患，大家都在认真学习这种新的计算框架，并被其高速的性能所折服。尽管在随后的一年半时间内，主流的 Hadoop 厂商并没有接受这个新框架：Cloudera 在忙着开发自己的 Impala 引擎，Hortonworks 经过评估后认为可以改造 Map/Reduce 来实现类似 Spark 的 DAG 机制，也就是后来的 Tez，而 MapR 还在纠结是否要全力投入 Drill 项目。但在 2014 年的夏天，第二次 Spark Summit 召开时，已经在 Spark 上积累大量开发者和用户，从互联网到传统行业，甚至是生物神经学家都用 Spark 来分析脑活动的数据。在这次会议上，大部分的 Hadoop 厂商以及应用开发商开始接受 Spark，并宣布支持 Spark 作为 Hadoop 上的另一个计算引擎。自此以后，Spark 的被接受程度飞速提高。到 2014 年 10 月份，几乎所有的大数据厂商都宣布支持 Spark，Spark 作者们创办的 DataBricks 公司也宣布认证了 50 多个以 Spark 为基础的应用系统。而到了 2015 年，大家在谈论的是 Spark 即将全面替代 Hadoop 中的 Map/Reduce。

星环科技从 2013 年创业的第一天，就开始改造 Spark 引擎来开发批处理和交互式分析引擎。今天在星环的全系列产品中，已经几乎看不到 Map/Reduce 计算框架。星环科技已经证明了在所有 Map/Reduce 擅长的领域，Spark 计算引擎都可以更高效地执行，性能可以提升数倍到数十倍，并且可以 7x24 稳定运行。这也从侧面证明了 Spark 引擎的潜力。

本书详细剖析了 Spark 核心引擎的源代码及其工作原理，内容翔实准确，也是我看到的一本比较全面解析 Spark Core 的不可多得的好书。特别是有志于 Spark 内核开发的研发人员，仔细阅读本书并研读代码，将起到事半功倍的效果。

孙元浩

星环科技创始人兼 CTO

2015 年 8 月上海

Preface 前 言

诞生于 2005 年的 Hadoop 解决了大数据的存储和计算问题，已经成为大数据处理的事实标准。但是，随着数据规模的爆炸式增长和计算场景的丰富细化，使得 Hadoop 越来越难以满足用户的需求。针对不同的计算场景，开源社区和各大互联网公司也推出了各种大数据分析的平台，旨在满足特定应用场景下的计算需求。但是，众多的平台使得用户不得不为平台开发类似的策略，这增加了运维开发成本。

2009 年诞生于 AMPLab 的 Spark，它的设计目标就是为大数据的多种计算场景提供一个通用的计算引擎，同时解决大数据处理的 4V 难题，即 Volume（海量）、Velocity（快速）、Variety（多样）、Value（价值）。正如 Spark 的核心作者之一的 Ion Stoica 所说，“The goal is to build a new generation of data analytics software, to be used across academia and industry.” Hadoop 之父 Doug Cutting 也说过，MapReduce 引擎将被 Spark 替代（Use of MapReduce engine for Big Data projects will decline, replaced by Apache Spark）。可以说，Spark 自诞生之日起就得到了广泛的关注，也是近年来开源社区最活跃的项目之一。Spark 的 1.X 版本的每次发布，都包含了数百位贡献者的上千次提交。最新的版本是发布于 2015 年 6 月 11 日的 1.4.0，是迄今为止 Spark 最大的一次版本发布，涵盖了 210 位开发者的贡献。

Spark 得到了众多大数据公司的支持，这些公司包括 Hortonworks、IBM、Intel、Cloudera、MapR、Pivotal 和星环科技；Spark 也被百度、阿里、腾讯、京东、携程、优酷土豆等互联网公司应用到多种不同的计算场景中，并且在实际的生产环境中获得了很多收益。当前百度的 Spark 已应用于凤巢、大搜索、直达号、百度大数据等业务；

阿里利用 GraphX 构建了大规模的图计算和图挖掘系统，实现了很多生产系统的推荐算法；腾讯 Spark 集群达到 8000 台的规模，是当前已知的世界上最大的 Spark 集群。

但是，当前并没有一本系统介绍 Spark 内核实现原理的书，而 Spark 内核是 Spark SQL、Spark Streaming、MLlib、GraphX 等多个模块的基础，这些模块最终的计算执行都是由内核模块完成的。为了在应用开发中做到游刃有余，在性能调优时做到有的放矢，需要了解内核模块的实现原理。笔者从 Spark 发布版本 0.8.1 时开始关注 Spark，并深入学习内核模块的架构实现原理。Spark 在 1.0 发布后，内核模块趋于稳定，虽然内核模块依旧会有不断地改进和完善，但是整体的设计思想和实现方法是不会变的，因此笔者决定为 Spark 社区的用户和关注者写一本书，详细介绍 Spark 内核模块的实现原理。最终，笔者基于 Spark 1.2.0 版本完成了本书。

写作是一件严肃的事情，同样是一份苦差事，尤其是在工作比较忙的时候。本书在半年前就完成了基本的框架，但是随后又对本书进行了多次修改和完善。笔者认为，对一本架构分析的书，一个最基本的要求就是基于源码如实描述系统的实现，能做到这点就是一本及格的书；如果能做到分析这个架构的好坏，指出架构改进的方案，那么就是一本质量比较好的书；如果能高屋建瓴地进行再次抽象，指出类似架构不同实现的优劣，抽象出一些理论，那么这就是一本质量上乘，可以当作教科书的书。我深知自己的能力水平，希望这本书最起码是一本及格的书，即能基于源码如实描述系统的实现，对那些希望深入学习 Spark 架构实现的同仁有所帮助。

目标读者

本书适合大数据领域的架构师、运维人员，尤其是 Spark 领域的从业人员阅读，也适合作为研究生和高年级的本科生大数据领域分布式架构具体实现原理的参考资料。

内容概述

第 1 章介绍了 Spark 的技术背景和特点，给出了架构的整体概述，并简单介绍了 Spark 的生态圈。

第 2 章介绍了 Spark 源码如何获取和学习环境如何搭建。

第 3 章是 RDD 的详细介绍，介绍了 RDD 的定义和 Spark 对于 DAG 的实现，最后通过 RDD 计算的详细介绍，讲解了 Spark 对于计算的实现原理。

第 4 章详细介绍任务调度的实现，包括如何通过 DAG 来生成计算任务，最后通过“Word Count”来加深对这个实现过程的理解。

第 5 章介绍了 Spark 的运行模式，尤其是 Standalone 模式。Standalone 是 Spark 自身实现的资源管理和调度的模块，这里会详细介绍它的实现原理。

第 6 章是 Executor 模块的详细讲解。Executor 是最终执行计算任务的单元，这章将详细介绍 Executor 的实现原理，包括 Executor 的分配、Task 在 Executor 的详细执行过程。

第 7 章详细介绍了 Spark 对于 Shuffle 的实现原理，包括基于 Hash 和基于排序的实现。除了详细阐述基于 Hash 和排序的 Shuffle 写和 Shuffle 读之外，还介绍了 Shuffle Pluggable 框架，为需要实现特定 Shuffle 逻辑的读者介绍其实现原理。

第 8 章详细介绍了 Spark 的 Storage 模块，在详细介绍了模块的架构后详细解析了不同存储级别的实现细节。

第 9 章介绍了 Spark 在百度、腾讯和阿里等国内互联网领域的应用现状。

致谢

本书在写作的过程中，得到了很多朋友、同仁的帮助和支持，在此表示衷心感谢！

感谢七牛云的技术总监陈超先生、蘑菇街的资深架构师刘旭晖先生、百度公司的高级架构师柴华先生、Databricks 软件工程师连城先生在百忙之中为本书审稿，并提出了很多宝贵的修改意见。尤其感谢星环科技的创始人兼 CTO 孙元浩先生对本书的完成给予了很大的支持，还在百忙之中为本书作序。感谢华为诺亚方舟实验室的董振华博士，在 Spark 上做机器学习方面给了我很多指导。

感谢百度上海研发中心网页搜索部的同事们。在一年多的工作中，笔者得到了很多同事的指导、支持和帮助，尤其感谢曲晶莹、吴永巍、汪韬、葛俊、刘桐仁、段雪涛、周波涛、马鑫云、李战平、杨大毛、朱晓阳、赵鹏程等。

感谢机械工业出版社的姚蕾编辑，她不但积极策划和推动本书的出版，还容忍我蜗牛般的写作速度；感谢她在这一年多的时间中给予的理解与支持。感谢机械工业出版社的李艺编辑为本书做了非常辛苦和专业的编辑工作。

还要感谢我的家人一直以来对我的支持和宽容。感谢父亲、母亲和三个姐姐，你们是我漫漫求学路的最强大支柱和后盾；感谢我的妻子王珊珊，不但在家庭方面付出很多，也为本书的顺利出版做出了很重要的贡献；感谢我的女儿，你的微笑是爸爸消除疲惫的良药。

联系方式

由于本书的写作都是在业余时间完成，加上笔者自身水平有限，错误在所难免，敬请读者谅解，如果有任何问题，可以通过下列方式与 Spark 的关注者和使用者进行交流沟通：

Spark 架构实现原理交流 QQ 群：473842835

Spark 用户使用交流 QQ 群：473853269

也可以直接与笔者联系：

邮箱：anzhsoft@gmail.com

新浪微博：@anzhsoft

个人博客：<http://blog.csdn.net/anzhsoft>

Contents 目 录

序

前言

第 1 章 Spark 简介	1
1.1 Spark 的技术背景.....	1
1.2 Spark 的优点.....	2
1.3 Spark 架构综述.....	4
1.4 Spark 核心组件概述.....	5
1.4.1 Spark Streaming.....	5
1.4.2 MLlib.....	6
1.4.3 Spark SQL.....	7
1.4.4 GraphX.....	8
1.5 Spark 的整体代码结构规模.....	8
第 2 章 Spark 学习环境的搭建	9
2.1 源码的获取与编译.....	9
2.1.1 源码获取.....	9
2.1.2 源码编译.....	10
2.2 构建 Spark 的源码阅读环境.....	11
2.3 小结.....	15

第3章 RDD 实现详解	16
3.1 概述.....	16
3.2 什么是 RDD.....	17
3.2.1 RDD 的创建.....	19
3.2.2 RDD 的转换.....	20
3.2.3 RDD 的动作.....	22
3.2.4 RDD 的缓存.....	23
3.2.5 RDD 的检查点.....	24
3.3 RDD 的转换和 DAG 的生成.....	25
3.3.1 RDD 的依赖关系.....	26
3.3.2 DAG 的生成.....	30
3.3.3 Word Count 的 RDD 转换和 DAG 划分的逻辑视图.....	30
3.4 RDD 的计算.....	33
3.4.1 Task 简介.....	33
3.4.2 Task 的执行起点.....	33
3.4.3 缓存的处理.....	35
3.4.4 checkpoint 的处理.....	37
3.4.5 RDD 的计算逻辑.....	39
3.5 RDD 的容错机制.....	39
3.6 小结.....	40
第4章 Scheduler 模块详解	41
4.1 模块概述.....	41
4.1.1 整体架构.....	41
4.1.2 Scheduler 的实现概述.....	43
4.2 DAGScheduler 实现详解.....	45
4.2.1 DAGScheduler 的创建.....	46
4.2.2 Job 的提交.....	48
4.2.3 Stage 的划分.....	49
4.2.4 任务的生成.....	54

4.3	任务调度实现详解	57
4.3.1	TaskScheduler 的创建	57
4.3.2	Task 的提交概述	58
4.3.3	任务调度具体实现	61
4.3.4	Task 运算结果的处理	65
4.4	Word Count 调度计算过程详解	72
4.5	小结	74
第 5 章	Deploy 模块详解	76
5.1	Spark 运行模式概述	76
5.1.1	local	77
5.1.2	Mesos	78
5.1.3	YARN	82
5.2	模块整体架构	86
5.3	消息传递机制详解	87
5.3.1	Master 和 Worker	87
5.3.2	Master 和 Client	89
5.3.3	Client 和 Executor	91
5.4	集群的启动	92
5.4.1	Master 的启动	92
5.4.2	Worker 的启动	96
5.5	集群容错处理	98
5.5.1	Master 异常退出	98
5.5.2	Worker 异常退出	99
5.5.3	Executor 异常退出	101
5.6	Master HA 实现详解	102
5.6.1	Master 启动的选举和数据恢复策略	103
5.6.2	集群启动参数的配置	105
5.6.3	Curator Framework 简介	106
5.6.4	ZooKeeperLeaderElectionAgent 的实现	109
5.7	小结	110

第 6 章 Executor 模块详解	112
6.1 Standalone 模式的 Executor 分配详解.....	113
6.1.1 SchedulerBackend 创建 AppClient.....	114
6.1.2 AppClient 向 Master 注册 Application.....	116
6.1.3 Master 根据 AppClient 的提交选择 Worker.....	119
6.1.4 Worker 根据 Master 的资源分配结果创建 Executor.....	121
6.2 Task 的执行.....	122
6.2.1 依赖环境的创建和分发.....	123
6.2.2 任务执行.....	125
6.2.3 任务结果的处理.....	128
6.2.4 Driver 端的处理.....	130
6.3 参数设置.....	131
6.3.1 spark.executor.memory.....	131
6.3.2 日志相关.....	132
6.3.3 spark.executor.heartbeatInterval.....	132
6.4 小结.....	133
第 7 章 Shuffle 模块详解	134
7.1 Hash Based Shuffle Write.....	135
7.1.1 Basic Shuffle Writer 实现解析.....	136
7.1.2 存在的问题.....	138
7.1.3 Shuffle Consolidate Writer.....	139
7.1.4 小结.....	140
7.2 Shuffle Pluggable 框架.....	141
7.2.1 org.apache.spark.shuffle.ShuffleManager.....	141
7.2.2 org.apache.spark.shuffle.ShuffleWriter.....	143
7.2.3 org.apache.spark.shuffle.ShuffleBlockManager.....	143
7.2.4 org.apache.spark.shuffle.ShuffleReader.....	144
7.2.5 如何开发自己的 Shuffle 机制.....	144
7.3 Sort Based Write.....	144

7.4	Shuffle Map Task 运算结果的处理	148
7.4.1	Executor 端的处理	148
7.4.2	Driver 端的处理	150
7.5	Shuffle Read	152
7.5.1	整体流程	152
7.5.2	数据读取策略的划分	155
7.5.3	本地读取	156
7.5.4	远程读取	158
7.6	性能调优	160
7.6.1	spark.shuffle.manager	160
7.6.2	spark.shuffle.spill	162
7.6.3	spark.shuffle.memoryFraction 和 spark.shuffle.safetyFraction	162
7.6.4	spark.shuffle.sort.bypassMergeThreshold	163
7.6.5	spark.shuffle.blockTransferService	163
7.6.6	spark.shuffle consolidateFiles	163
7.6.7	spark.shuffle.compress 和 spark.shuffle.spill.compress	164
7.6.8	spark.reducer.maxMbInFlight	165
7.7	小结	165
第 8 章	Storage 模块详解	167
8.1	模块整体架构	167
8.1.1	整体架构	167
8.1.2	源码组织结构	170
8.1.3	Master 和 Slave 的消息传递详解	173
8.2	存储实现详解	181
8.2.1	存储级别	181
8.2.2	模块类图	184
8.2.3	org.apache.spark.storage.DiskStore 实现详解	186
8.2.4	org.apache.spark.storage.MemoryStore 实现详解	188
8.2.5	org.apache.spark.storage.TachyonStore 实现详解	189
8.2.6	Block 存储的实现	190

8.3 性能调优	194
8.3.1 spark.local.dir	194
8.3.2 spark.executor.memory	194
8.3.3 spark.storage.memoryFraction	194
8.3.4 spark.streaming.blockInterval	195
8.4 小结	195
第9章 企业应用概述	197
9.1 Spark 在百度	197
9.1.1 现状	197
9.1.2 百度开放云 BMR 的 Spark	198
9.1.3 在 Spark 中使用 Tachyon	199
9.2 Spark 在阿里	200
9.3 Spark 在腾讯	200
9.4 小结	201



Spark 简介

Apache Spark 是一种快速、通用、可扩展的大数据分析引擎。它是不断壮大的大数据分析解决方案家族中备受关注的明星成员，为分布式数据集的处理提供了一个有效框架，并以高效的方式处理分布式数据集。Spark 集批处理、实时流处理、交互式查询与图计算于一体，避免了多种运算场景下需要部署不同集群带来的资源浪费。Spark 在过去的一年中获得了极大关注，并得到广泛应用，Spark 社区也成为大数据领域和 Apache 软件基金会最活跃的项目之一，其活跃度甚至远超曾经只能望其项背的 Hadoop。

1.1 Spark 的技术背景

无论是工业界还是学术界，都已经广泛使用高级集群编程模型来处理日益增长的数据，如 MapReduce 和 Dryad。这些系统将分布式编程简化为自动提供位置感知性调度、容错以及负载均衡，使得大量用户能够在商用集群上分析超大数据集。大多数现有的集群计算系统都是基于非循环的数据流模型。即从稳定的物理存储（如分布式文件系统）中加载记录，记录被传入由一组确定性操作构成的 DAG（Directed Acyclic Graph，有向无环图），然后写回稳定存储。DAG 数据流图能够在运行时自动实现任务

调度和故障恢复。

尽管非循环数据流是一种很强大的抽象方法，但仍然有些应用无法使用这种方式描述。这类应用包括：①机器学习和图应用中常用的迭代算法（每一步对数据执行相似的函数）；②交互式数据挖掘工具（用户反复查询一个数据子集）。基于数据流的框架并不明确支持工作集，所以需要将数据输出到磁盘，然后在每次查询时重新加载，这会带来较大的开销。针对上述问题，Spark 实现了一种分布式的内存抽象，称为**弹性分布式数据集**（Resilient Distributed Dataset，RDD）。它支持基于工作集的应用，同时具有数据流模型的特点：自动容错、位置感知性调度和可伸缩性。RDD 允许用户在执行多个查询时显式地将工作集缓存在内存中，后续的查询能够重用工作集，这极大地提升了查询速度。

RDD 提供了一种高度受限的共享内存模型，即 RDD 是只读记录分区的集合，只能通过在其他 RDD 执行确定的转换操作（如 map、join 和 groupBy）而创建，然而这些限制使得实现容错的开销很低。与分布式共享内存系统需要付出高昂代价的检查点和回滚机制不同，RDD 通过 Lineage 来重建丢失的分区：一个 RDD 中包含了如何从其他 RDD 衍生所必需的相关信息，从而不需要检查点操作就可以重建丢失的数据分区。尽管 RDD 不是一个通用的共享内存抽象，但它具备了良好的描述能力、可伸缩性和可靠性，能够广泛适用于数据并行类应用。

1.2 Spark 的优点

Spark 凭借什么在众多的大数据分析处理平台中脱颖而出呢？这得益于 Spark 的几个主要优点，具体如下所示。

第一个优点是速度。与 Hadoop 的 MapReduce 相比，Spark 基于内存的运算要快 100 倍以上（如图 1-1 所示）；而基于硬盘的运算也要快 10 倍以上。Spark 实现了高效的 DAG 执行引擎，可以通过基于内存来高效处理数据流。

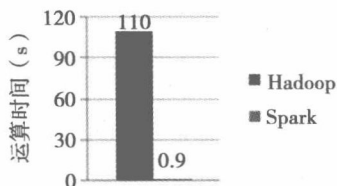


图 1-1 逻辑回归在 Hadoop、Spark 的运算速度对比