

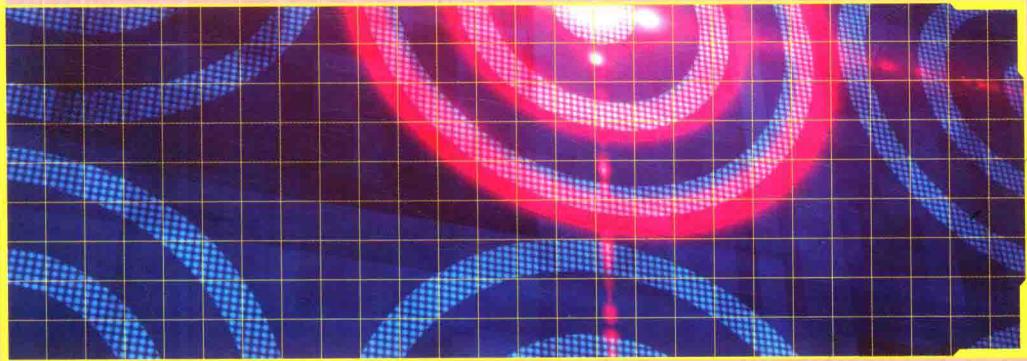


先·进·信·号·处·理·系·列

大数据技术体系 与开源生态



刘驰 罗童 林秋霞 王新科 樊路遥 | 编著



中国工信出版集团



人民邮电出版社
POSTS & TELECOM PRESS

先·进·信·号·处·理·系·列

大数据技术体系 与开源生态

刘驰 罗童 林秋霞 王新科 樊路遥 | 编著

人民邮电出版社
北京

图书在版编目 (C I P) 数据

大数据技术体系与开源生态 / 刘驰等编著. — 北京:
人民邮电出版社, 2018. 8
(先进信号处理系列)
ISBN 978-7-115-49223-4

I. ①大… II. ①刘… III. ①数据处理—研究 IV. ①TP274

中国版本图书馆CIP数据核字(2018)第197403号

内 容 提 要

本书从大数据生命周期的角度阐述大数据技术体系与开源生态的发展。全书分为七篇，包括大数据技术体系与开源生态概述、大数据获取技术、大数据管理技术、大数据处理技术、大数据分析与挖掘技术、大数据可视化与交互技术、大数据安全与治理技术。又分为 15 章，详细介绍大数据的技术概况、发展近况和技术优势、软件架构和应用场景等内容。

本书适合大数据和人工智能业内人员、各大高校相关专业的高年级本科生和研究生，以及对大数据应用中各类框架组件的爱好者阅读。

定价：159.00 元

读者服务热线: (010) 81055488 印装质量热线: (010) 81055316
反盗版热线: (010) 81055315

前 言

大数据是以容量大、类型多、存取速度快、应用价值高为主要特征的数据集合。随着信息技术和人类生产生活的交汇融合，互联网快速普及，全球数据呈现爆发增长、海量集聚的特点，对经济发展、社会治理、国家管理、人民生活都产生了重大影响。大数据技术的应用和发展，为人类提供了认识复杂系统的新思维和新手段，成为推动经济转型升级的新动力，也是提升社会治理能力的新途径，以及提升国家综合能力和保障国家安全的新利器。

美国政府最先对大数据技术革命做出战略反应，从2012年到2016年共实施了4轮政策行动，美国白宫成立了“美国大数据研发高级指导小组”，实施了《大数据研究和开发计划》，加强了在大数据研发和应用方面的布局。欧盟发布了《数据驱动经济战略》，大力推动“数据价值链战略计划”，倡导欧洲各国抢抓大数据发展机遇。英国政府发布了《把握数据带来的机遇：英国数据能力战略》，从提升数据分析技术、加强国家基础设施建设、确保数据安全和共享等方面作出部署。此外，日本、韩国、澳大利亚等国均出台了相关政策，大力推动大数据应用及产业发展。

我国在党的十八届五中全会上将大数据上升为国家战略。2015年8月，国务院印发了《促进大数据发展行动纲要》；2017年1月，工信部印发了《大数据产业发展规划（2016—2020年）》等政策指导文件，为大数据战略的实施指明了具体方向。在党的十九大报告中，习近平总书记明确指出：“推动互联网、大数据、人工智能和实体经济深度融合。”2017年12月8日，在中共中央政治局第二次集体学习时，习近平总书记发表了“大数据发展日新月异，我们应该审时度势、精心谋划、超前布局、力争主动，深入了解大数据发展现状和趋势及其对经济社会发展的影响，分析我国大数据发展取得的成绩和存在的问题，推动实施国家大数据战略，加快完善数字基础设施，推进数据资源整合和开放共享，保障数据安全，加快建设数字中国，更好服务我国经济社会发展和人民生活改善”的讲话。自此，我国大数据发展开启了新的篇章。

在这样的大背景下，全球大数据产业日趋活跃，技术演进和应用创新加速发展，全世界各地的创新人才不断地涌进入到大数据产业中来，进一步反哺大数据技术的发展。再加上开源生态的不断壮大，进一步促进了大数据技术的分享和发展。

大数据技术已成体系，从大数据生命周期的角度，其可以划分为大数据获取、大数据管理、大数据处理、大数据分析与挖掘、大数据可视化、大数据安全与治理 6 个方面，本书也是从这 6 个方面来对大数据技术体系和开源生态的建设进行相应的介绍。

本书共 7 篇 15 章，其中：

第一篇：大数据技术体系与开源生态概述，包括第 1~4 章。第 1 章为大数据技术体系概述。第 2 章描述了开源生态与代码托管平台简介。第 3~4 章着重介绍了大数据技术和云计算技术这两类技术的开源生态建设及发展情况，包括 Apache 软件基金会、Linux 基金会、OpenStack 基金会、Cloud Native Computing Foundation 和开源中国。其中，第 1 章由段雄编写，第 2~4 章由罗童编写。

第二篇：大数据获取技术，包括第 5 章，从当前主流的消息队列相关技术的角度阐述数据如何从异构系统采集并融合在一起，包括 ZeroMQ、RabbitMQ、ActiveMQ、Apache Kafka4 个开源项目。第 5 章由樊路遥编写。

第三篇：大数据管理技术，包括第 6~7 章。第 6 章描述了当前主流数据库技术，包括 MySQL、PostgreSQL、MongoDB、Apache CouchDB、Vertica、Apache HBase、Neo4j、OrientDB、InfiniteGraph、Alluxio、Apache Tajo 11 个开源项目。第 7 章介绍了大数据平台资源管理技术，包括 Apache Zookeeper、Apache Hadoop YARN、Apache Mesos、Apache Mnemonic 4 个开源项目。其中，第 6 章由温琦和胡柏青编写，第 7 章由方久鑫编写。

第四篇：大数据处理技术，包括第 8~9 章。第 8 章介绍了当前主流的大数据批处理平台，包括 Apache Hadoop、Apache Spark、Apache Kylin 这 3 个开源项目。第 9 章介绍了当前主流的实时流处理平台，包括 Apache Storm、Apache Spark Streaming、Apache Flink、Apache Beam 和 Apache Apex 4 个开源项目。其中，第 8 章由王新科编写，第 9 章由林秋霞编写。

第五篇：大数据分析与挖掘技术，包括第 10~11 章。第 10 章介绍了主流大数据分析工具，包括 Apache Mahout、Apache Lens、Apache Spark MLlib、Scikit-Learn4 个大数据开源项目。第 11 章介绍了主流人工智能开源平台，包括 TensorFlow、Caffe、PyTorch、TensorFlow Lite4 个开源项目。其中，第 10 章由陈喆和王新科编写，第 11 章由樊路遥、朴成哲和温琦编写。

第六篇：大数据可视化与交互技术，包括第 12~13 章。第 12 章介绍了两个主流大数据：可视化与交互技术，Tableau 和 Apache Zeppelin。第 13 章介绍了其他大数据可视化与交互技术，包括 BIRT、KNIME、Jaspersoft Community 3 个开源

项目。其中，第 12 章由张晶和罗童编写，第 13 章由陈喆、胡柏青和林秋霞编写。

第七篇：大数据安全与治理技术，包括第 14~15 章。在第 14 章，通过介绍 Apache Falcon 和 Apache Atlas 两个开源项目阐述了大数据治理技术。第 14 章介绍了 Apache Kerberos、Apache Ranger、Apache Sentry 和 Apache Metron 4 个大数据安全开源项目。其中，第 14 章由陈喆编写，第 15 章由陈喆和张晶编写。

由于编者及撰写者的认识和水平所限，本书内容仅从大数据技术体系及其开源生态发展的一个视角写作，既不全面，也难免偏颇。但能够为读者提供一定的参考，则本书目的已达。欢迎广大读者对本书进行批评指正。

作者

2018 年 4 月

目 录

第一篇 大数据技术体系与开源生态概述

第1章 大数据技术体系概述	3
1.1 大数据技术的主要内容	3
1.2 大数据开源框架	4
1.2.1 大数据获取技术	4
1.2.2 大数据管理技术	5
1.2.3 大数据处理技术	5
1.2.4 大数据安全与治理技术	5
1.2.5 大数据分析与挖掘技术	6
1.2.6 大数据可视化技术	6
1.3 本章小结	7
第2章 开源生态与代码托管平台简介	8
2.1 开源和开源软件的简介	8
2.1.1 开源的简介	8
2.1.2 开源软件的简介	8
2.2 开源代码托管平台——GitHub	9
2.3 本章小结	10
第3章 大数据开源生态的介绍	11
3.1 Apache 软件基金会	11
3.1.1 发展历史	11
3.1.2 主要参与者	12

3.1.3 开源项目	13
3.2 Linux 基金会	14
3.2.1 发展历史	15
3.2.2 主要参与者	15
3.2.3 开源项目	17
3.3 开源中国	18
3.3.1 发展历史	18
3.3.2 主要参与者	19
3.4 本章小结	19
第 4 章 云计算开源生态的介绍	20
4.1 OpenStack 基金会	20
4.1.1 发展历史	21
4.1.2 主要参与者	21
4.1.3 开源项目	22
4.2 Cloud Native Computing Foundation	23
4.2.1 发展历史	23
4.2.2 主要参与者	23
4.2.3 开源项目	25
4.3 本章小结	25

第二篇 大数据获取技术

第 5 章 消息队列相关技术	29
5.1 ZeroMQ	29
5.1.1 技术概况	29
5.1.2 发展近况和技术优势	30
5.1.3 软件架构	31
5.1.4 应用场景	33
5.2 RabbitMQ	34
5.2.1 技术概况	34
5.2.2 发展近况和技术优势	35
5.2.3 软件架构	36
5.2.4 应用场景	38

5.3 Active MQ	40
5.3.1 技术概况	40
5.3.2 发展近况和技术优势	40
5.3.3 软件架构	42
5.3.4 应用场景	43
5.4 Apache Kafka	44
5.4.1 技术概况	44
5.4.2 发展近况和技术优势	45
5.4.3 软件架构	46
5.4.4 应用场景	47
5.5 本章小结	50

第三篇 大数据管理技术

第6章 数据库相关技术	53
6.1 传统关系型数据库	53
6.1.1 MySQL	53
6.1.2 PostgreSQL	60
6.2 文档型数据库	65
6.2.1 MongoDB	65
6.2.2 Apache CouchDB	69
6.3 列存储数据库	73
6.3.1 Vertica	73
6.3.2 Apache HBase	76
6.4 键/值对型数据库	80
6.4.1 Redis	80
6.4.2 Riak	82
6.5 图形数据库	85
6.5.1 Neo4j	85
6.5.2 OrientDB	90
6.5.3 InfiniteGraph	93
6.6 基于内存的分布式文件系统之 Alluxio	95
6.6.1 技术概况	95
6.6.2 发展近况和技术优势	96

6.6.3 软件架构	97
6.6.4 应用场景	98
6.7 数据仓库系统之 ApacheTajo	99
6.7.1 技术概况	99
6.7.2 发展近况和技术优势	100
6.7.3 软件架构	101
6.7.4 应用场景	103
6.8 本章小结	105
第7章 大数据平台资源管理技术	106
7.1 Apache ZooKeeper	106
7.1.1 技术概况	106
7.1.2 发展近况和技术优势	107
7.1.3 软件架构	108
7.1.4 应用场景	110
7.2 Apache Hadoop YARN	111
7.2.1 技术概况	111
7.2.2 发展近况和技术优势	112
7.2.3 软件架构	113
7.2.4 应用场景	116
7.3 Apache Mesos	119
7.3.1 技术概况	119
7.3.2 发展近况和技术优势	120
7.3.3 软件架构	120
7.3.4 应用场景	122
7.4 Apache Mnemonic	123
7.4.1 技术概况	123
7.4.2 发展近况和技术优势	124
7.5 本章小结	125

第四篇 大数据处理技术

第8章 开源批处理平台	129
8.1 Apache Hadoop	129

8.1.1 技术概况	129
8.1.2 发展近况和技术优势	130
8.1.3 软件架构	131
8.1.4 应用场景	136
8.2 Apache Spark	142
8.2.1 技术概况	142
8.2.2 发展近况和技术优势	142
8.2.3 软件架构	144
8.2.4 应用场景	146
8.3 Apache Kylin	150
8.3.1 技术概况	150
8.3.2 发展近况和技术优势	150
8.3.3 软件架构	152
8.3.4 应用场景	153
8.4 本章小结	159
第9章 开源实时处理平台	160
9.1 Apache Storm	160
9.1.1 技术概况	160
9.1.2 发展近况和技术优势	161
9.1.3 软件架构	162
9.1.4 应用场景	163
9.2 Apache Spark Streaming	169
9.2.1 技术概况	169
9.2.2 发展近况和技术优势	170
9.2.3 软件架构	170
9.2.4 应用场景	171
9.3 Apache Flink	173
9.3.1 技术概况	173
9.3.2 发展近况和技术优势	174
9.3.3 软件架构	175
9.3.4 应用场景	176
9.4 Apache Beam	179
9.4.1 技术概况	179
9.4.2 发展近况和技术优势	180

9.4.3 软件架构	181
9.4.4 应用场景	182
9.5 Apache Apex	186
9.5.1 技术概况	186
9.5.2 发展近况和技术优势	187
9.5.3 软件架构	188
9.5.4 应用场景	191
9.6 本章小结	194

第五篇 大数据分析与挖掘技术

第 10 章 开源数据分析平台	199
-----------------------	-----

10.1 Apache Mahout	199
10.1.1 技术概况	199
10.1.2 发展近况和技术优势	200
10.1.3 应用场景	202
10.2 Apache Spark MLlib	204
10.2.1 技术概况	204
10.2.2 发展近况和技术优势	204
10.2.3 软件架构	205
10.2.4 应用场景	207
10.3 Apache Lens	208
10.3.1 技术概况	208
10.3.2 发展近况及技术优势	209
10.3.3 软件架构	213
10.3.4 应用场景	214
10.4 Scikit-Learn	217
10.4.1 技术概况	217
10.4.2 发展近况与技术优势	217
10.4.3 软件架构	218
10.4.4 应用场景	220
10.5 本章小结	223

第 11 章 开源深度学习平台.....	225
11.1 TensorFlow	225
11.1.1 技术概况	225
11.1.2 发展近况和技术优势	226
11.1.3 软件架构	226
11.1.4 应用场景	230
11.2 Tensorflow Lite.....	233
11.2.1 技术概况	233
11.2.2 发展近况和技术优势	233
11.2.3 软件架构	234
11.3 Caffe	237
11.3.1 技术概述	237
11.3.2 发展近况和技术优势	237
11.3.3 软件架构	239
11.3.4 应用场景	241
11.4 PyTorch.....	243
11.4.1 技术概况	243
11.4.2 发展近况和技术优势	243
11.4.3 软件架构	245
11.4.4 应用场景	247
11.5 本章小结	248

第六篇 大数据可视化与交互技术

第 12 章 主流大数据可视化与交互工具	251
12.1 Tableau	251
12.1.1 技术概况	251
12.1.2 发展近况和技术优势	252
12.1.3 软件架构	255
12.1.4 应用场景	256
12.2 Apache Zeppelin.....	260
12.2.1 技术概况	260
12.2.2 发展近况和技术优势	261
12.2.3 软件架构	262

12.2.4 应用场景	263
12.3 本章小结	266
第 13 章 其他大数据可视化与交互工具	267
13.1 Jaspersoft Community	267
13.1.1 技术概况	267
13.1.2 发展近况和技术优势	268
13.1.3 软件架构	270
13.1.4 应用场景	271
13.2 BIRT	274
13.2.1 技术概况	274
13.2.2 发展近况和技术优势	275
13.2.3 软件架构	278
13.2.4 应用场景	280
13.3 KNIME	281
13.3.1 技术概况	281
13.3.2 发展近况和技术优势	281
13.3.3 软件架构	283
13.3.4 应用场景	285
13.4 本章小结	285

第七篇 大数据安全与治理技术

第 14 章 大数据治理技术	289
14.1 Apache Falcon	289
14.1.1 技术概况	290
14.1.2 发展近况和技术优势	290
14.1.3 软件架构	292
14.1.4 应用场景	294
14.2 Apache Atlas	297
14.2.1 技术概况	298
14.2.2 发展近况和技术优势	301
14.2.3 软件架构	306
14.3 本章小结	314

第 15 章 大数据安全技术	316
15.1 Apache Ranger	316
15.1.1 技术概况	316
15.1.2 发展近况和技术优势	318
15.1.3 软件架构	321
15.1.4 应用场景	322
15.2 Apache Sentry	324
15.2.1 技术概况	324
15.2.2 发展近况和技术优势	326
15.2.3 软件架构	332
15.3 Apache Kerberos	334
15.3.1 技术概况	335
15.3.2 发展近况和技术优势	336
15.3.3 软件架构	337
15.4 Apache Metron	339
15.4.1 技术概况	339
15.4.2 发展近况及技术优势	340
15.4.3 软件架构	344
15.5 Hyperledger	346
15.5.1 技术概况	346
15.5.2 发展近况和技术优势	347
15.5.3 软件架构	348
15.5.4 应用场景	349
15.6 本章小结	351
结束语	353
名词索引	355

第一篇

大数据技术体系与开源生态概述

随着互联网的飞速发展，我们的生活中到处都充斥着数据，并且数据的规模正在呈指数级增长，显然大数据时代已经到来。对于如此庞大的数据量，我们需要对其进行相应地获取、存储、管理、处理、分析和可视化等操作。大数据技术也就成了解决这些问题的核心。同时随着开源软件的不断发展，大数据技术也得到不断的发展。现如今，大数据技术的发展已经相对成熟，开源生态系统也不断地壮大起来。第一篇，主要对大数据技术体系进行相应介绍，介绍了现在比较流行的一些大数据技术，然后对开源生态中的大数据开源组织和云计算开源组织进行相应介绍。本篇的目的是使读者对大数据技术体系和开源生态有一个理论上的认识，为后面学习具体的大数据技术打下基础。

