

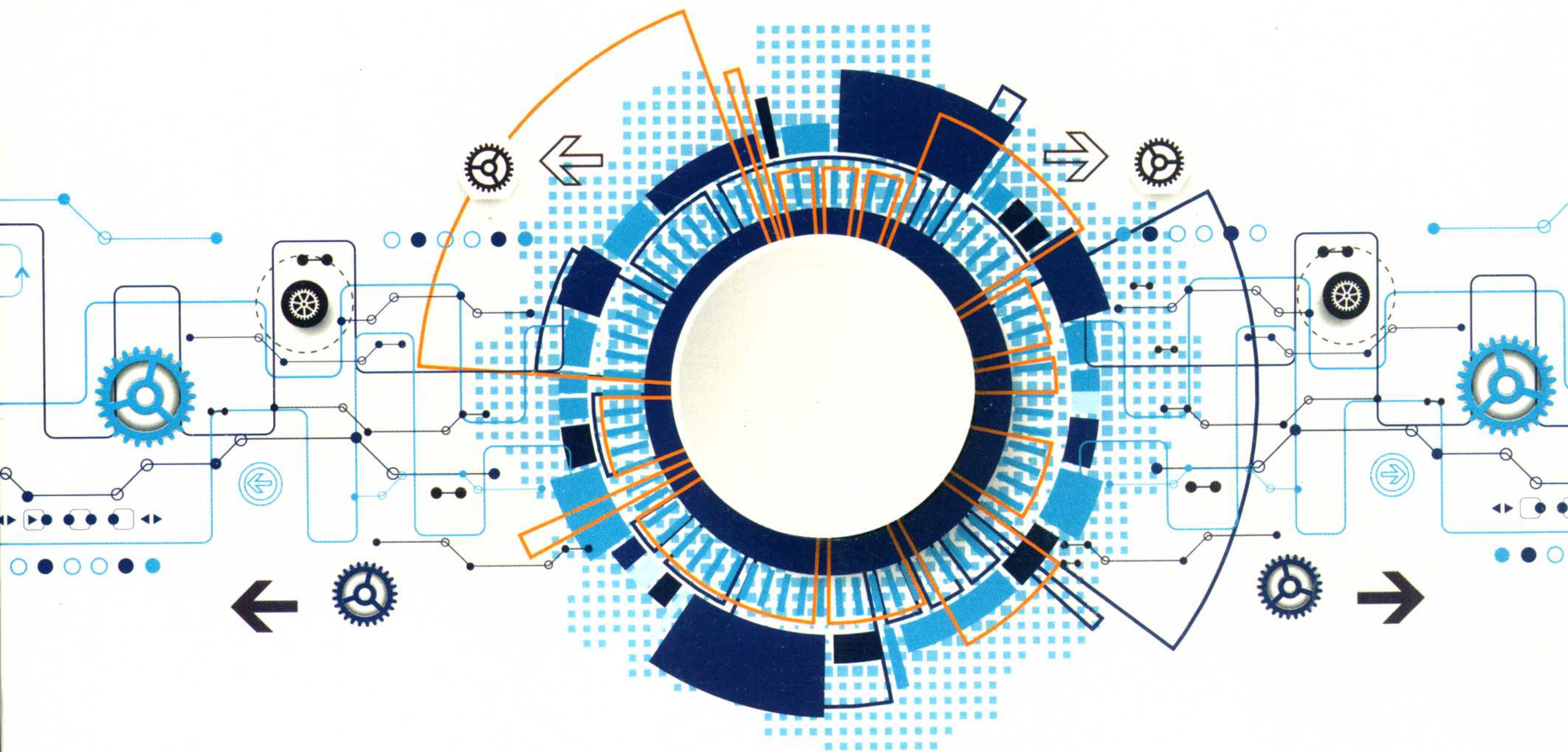


资深大数据专家/一线架构师20000小时实际工作经验总结

从横向视角出发，拉通Hadoop体系技术栈，手把手教你快速构建一个真实可用、安全可靠的企业级大数据平台



技术丛书



The Construction of Big Data Platform for Large Enterprises - Architecture and Implementation

企业级大数据平台构建 架构与实现

朱凯◎著



机械工业出版社
China Machine Press

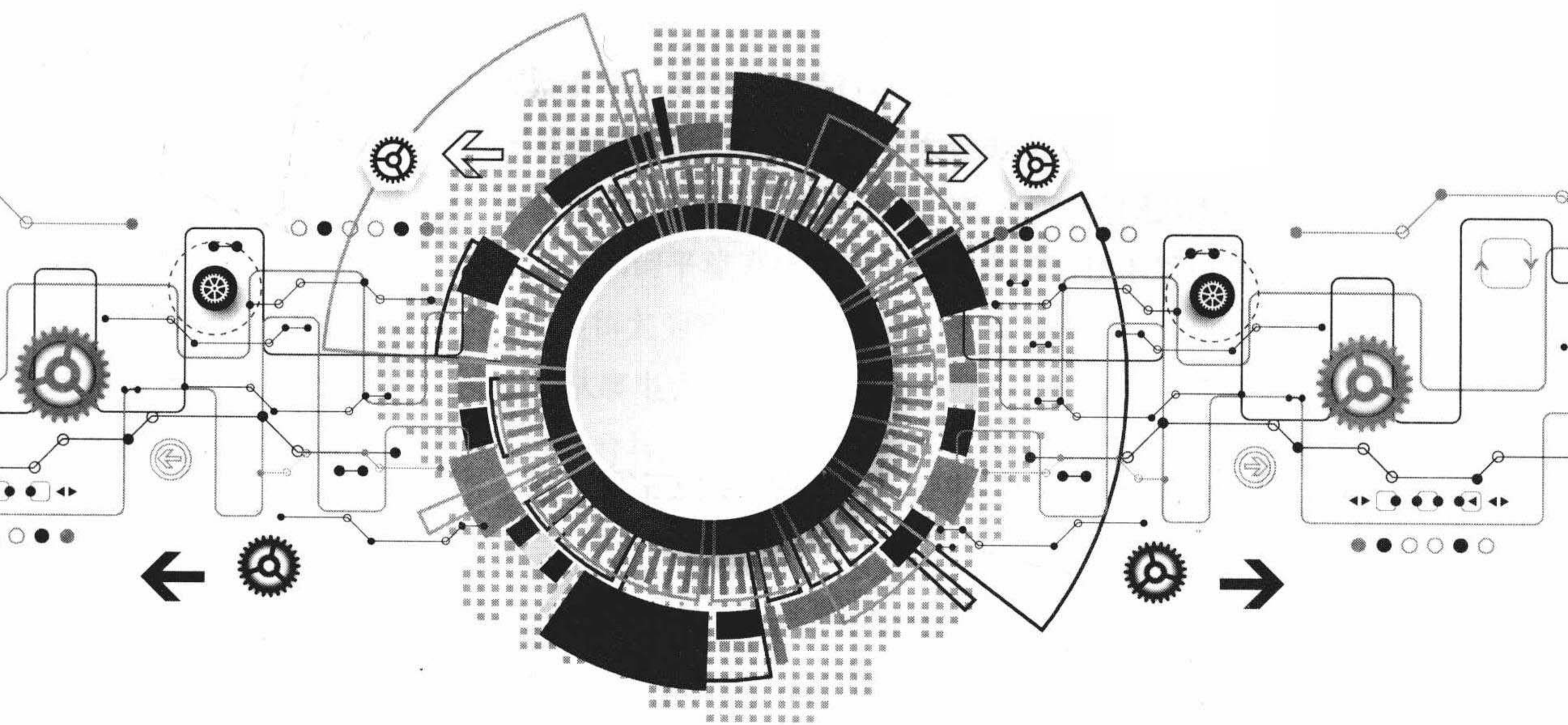


技术丛书

The Construction of Big Data Platform for Large Enterprises - Architecture and Implementation

企业级大数据平台构建 架构与实现

朱凯◎著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

企业级大数据平台构建：架构与实现 / 朱凯著. —北京：机械工业出版社，2018.3
(大数据技术丛书)

ISBN 978-7-111-59595-3

I. 企… II. 朱… III. 数据处理 IV. TP274

中国版本图书馆 CIP 数据核字 (2018) 第 058184 号

企业级大数据平台构建：架构与实现

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：孙海亮

责任校对：李秋荣

印刷：三河市宏图印务有限公司

版次：2018 年 4 月第 1 版第 1 次印刷

开本：186mm × 240mm 1/16

印张：16.5

书号：ISBN 978-7-111-59595-3

定价：69.00 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 推荐序

思者常新，厚积薄发

相比以 BAT 为引领的互联网公司的天生“数质”（业务高度数字化，技术更具创新性和开放性），大多数传统集团企业经过之前 ERP 时代积累了海量的业务数据。但是由于业务的复杂性与新老技术升级变革的压力，面对汹涌而来的大数据浪潮，这些企业却依旧停留在探索技术如何稳步更替升级、分散在不同部门的数据如何更有效地集中统一、数据本身以及数据技术如何有效形成企业级治理体系等一系列“知易行难”的问题当中。

相比两年前全民热捧大数据概念的疯狂，数据本身和大数据主流技术显然已经稳步度过了“过高期望的峰值期”和“泡沫化的低谷期”，正式进入“稳步爬升的光明期”。正因为这样，在这个黎明前的时期，传统企业如何平稳完成数字化变革带来的技术架构变迁，找到去伪存真的企业级大数据实战指南就显得尤为重要了。

最近十年，我一直在远光软件从事电力企业信息化相关的工作，组织、带领了包括企业大数据平台、企业新一代敏捷 BA 平台与能源 BDaaS 平台在内的研发团队。电力企业在国内正是信息化水平较高，但业务运营和技术管控模式最为复杂的一类企业。正是意识到“生于互联网的大数据技术对于集团企业的大数据应用支撑不足”这一事实，我们很早就开始孵化相关的团队、探索相关的应用。直到我们在公司正式组建第一支“企业大数据产品商业化团队”的时候，本书作者通过“普通社招”成为第一批加入的开发工程师之一。在短短的半年时间内，他就在如何快速学习新技术、实践新架构方面展现出高于常人、高于前辈的能力和素养。

四年的时间，我们的大数据产品 EDT、创新数字化平台产品 Realinsight 相继诞生，我们和用户一起完成了一个又一个企业大数据解决方案实战。用户数屡创新高、嘉奖年年不断，获得了市场和行业的肯定。当年那支由 20 人组成的产品团队，一年一个台阶发展为如今公

司的一级产品事业部，当年那个“普通开发工程师”也当仁不让地成为我们整个大数据产品线中最为核心的系统架构师和技术布道师。

在此过程中，本书作者和各个技术同仁、产品经理、业务部门同事紧密合作，而这本著作就融入了他在这些实战项目中所积累的丰富经验。所以，本书最大的闪光点在于，它的内容不局限于技术本身，而是考虑到了在不同企业应用场景下，这些技术如何得到更合理地应用。除此之外，作者文艺青年的背景让这本书读起来极其顺畅，他的钻研精神又让这本书在理论上更具深度。因此，本书除了适合集团企业的技术管理人员通读外，也非常适合从事大数据产品相关工作的设计者、产品经理或者架构师阅读。我想，对于希望利用大数据技术解决业务痛点的读者而言，本书更是不可或缺的良好益友。

当得知朱凯有出书的打算时，我们都很兴奋。谁会比他更能胜任这件事呢？毫无疑问，这会是国内企业大数据技术领域的一本不可多得的图书。“思者常新，厚积薄发”正是我对本书作者这几年状态的一个真实表述，但同时也是对于正走向真正落地的企业大数据时代的共勉。数字经济时代已经到来，作为这个时代积极的参与者，我们渴望和更多的思考者共同分享、一起创造，实现企业大数据技术应用的爆发。

远光软件大数据事业部总经理 解来甲

为什么要写这本书

近年来，大数据这个概念越来越火爆，特别是在国家层面，大数据被提升到了国家战略的高度。在这样的背景下，很多传统企业开始涉足大数据领域并研发自己的大数据技术平台。在这股技术升级与转型的浪潮中，传统领域的程序员纷纷转型投向大数据的怀抱。目前大数据技术开源领域以 Hadoop 生态构建的技术体系为主。现在市面上有很多与 Hadoop 体系相关的技术书籍，Hadoop、Spark 这类火爆的技术已经有大量优秀的专业书籍进行讲解。但我发现这类书籍多是以纵向的视角去讲解某一类具体的技术，而大数据领域涉及的知识繁多，在构建大数据平台的过程中我们不仅需要精通单个技术组件的知识，还需要拥有横向整合拉通 Hadoop 体系技术栈的能力。而这类横向拉通 Hadoop 体系技术栈的书籍并不多见。所以我将自己在构建大数据平台上的一些经验和实践进行了整理，分享给各位读者。希望本书能够为各位读者构建大数据平台或解决方案提供一定的帮助。

读者对象

- **想了解大数据技术，想进入大数据领域的工程师：**作为一个想进入大数据领域的“新人”，你可以通过本书从宏观的视角迅速对大数据的基础设施和技术栈有一个全面的认识 and 了解。本书可以作为你的入门指南和技术栈索引目录。
- **大数据领域的中高级工程师：**作为一个大数据领域的中高级工程师，对 Hadoop 生态体系的技术应该早已运用自如。通过本书的学习，相信你对大数据领域多种技术栈的整合会有一个更深刻的认识。同时本书中的一些平台级方案也会帮助你提升在平台架构方面的造诣。
- **平台架构师：**作为一个平台架构师，本书中的一些解决方案和设计思路可以作为你进

行系统架构的参考资料。

本书主要内容

本书从企业的实际需求出发，完整地介绍了构建一个真实可用、安全可靠的企业级大数据平台所需要运用的知识体系，并详细地描述了构建企业级大数据平台的设计方案和实施步骤。

本书逻辑上可分为3大部分，共8章，每个章节循序渐进：

- 第一部分（第1、2章）描述了企业级大数据平台的需求和能力。
- 第二部分（第3～5章）着重讲述了如何去搭建并配置一个大数据平台，以及如何构建非常重要的平台安全方案。
- 第三部分（第6～8章）以实战的形式讲解如何以Java编码的方式实现平台的基础管理功能，以提升其易用性与可用性。

具体各章内容如下：

第1章 阐述企业级大数据平台的重要性，并解释了为什么需要构建一个统一的企业级大数据平台。接着介绍作为一个企业级大数据平台应当具备的能力，并解释其原因。

第2章 介绍通过Hadoop生态体系去构建一个企业级大数据平台可以使用的技术栈，如HDFS、HBase、Spark等，并一一介绍了它们的核心概念。

第3章 介绍集群管理工具Ambari，并站在集群服务器的角度分类解释如何去设计一个Hadoop集群，详细描述了如何使用Ambari来安装、管理和监控一个Hadoop集群。

第4章 介绍企业级大数据平台中非常重要的安全部分。首先阐述了企业级大数据平台面临的一些安全隐患，接着展示了一套初级解决方案并介绍了如何使用Knox和Ranger解决访问控制和数据授权与管理的问题。

第5章 着重介绍Hadoop服务的安全方案，并说明如何通过Kerberos协议等一系列措施来保障Hadoop集群的安全。

第6章 阐述大数据平台在易用性上的一些遗留问题，接着介绍如何通过CAS实现平台的单点登录功能，最后描述如何使用Java程序实现统一的用户管理服务。

第7章 简单阐述服务化的重要性以及如何将大数据平台管理端的功能封装成RESTful服务。首先介绍了如何使用Spring-Boot快速搭建一套RESTful服务的程序框架，接着详细描述如何实现Kerberos用户查询、Hive数据仓库查询和元数据查询等一系列RESTful服务。

第8章 介绍如何使用Java程序实现Spark的任务提交与任务调度功能。首先着重介绍使用Java程序实现Spark任务提交到YARN的三种方式，接着描述如何通过Quartz实现任

务调度功能。

如何阅读本书

本书内容会涉及大数据领域相关的技术知识，所以假定读者已具有一定的编程经验，了解分布式、多线程、集群等概念。本书部分内容涉及集群服务的实战安装示例，所以需要准备至少两台用于搭建测试环境的 Linux 服务器或虚拟机。

勘误和支持

由于水平有限，编写时间仓促，书中难免会出现一些错误或者不准确的地方，恳请读者批评指正。为此，我特意创建了一个提供在线支持与应急方案的站点 <https://github.com/nauu/bigdatabook>。你可以将书中的错误发布在 Bug 勘误表页面中，如果你遇到任何问题，也可以访问 Q&A 页面，我将尽量在线上为读者提供最满意的解答。如果你有更多的宝贵意见，也欢迎发送邮件至邮箱 yawface@gmail.com 或者访问新浪微博 <http://weibo.com/boness>，期待能够得到你们的真挚反馈。

致谢

感谢我的家人，如果没有你们的悉心照顾和鼓励，我不可能完成本书。

感谢我的公司远光软件，为我提供了学习和成长的环境，本书中的很多知识都来自工作中的实践。

感谢我的挚友李根，为本书提出了许多宝贵的建议。

感谢我的同事兼伙伴们——解来甲、张琛、杨柯、潘登、胡艺、李国威、陈世宾、陈泽华，以及名单之外的更多朋友，感谢你们在工作中的照顾和支持，十分荣幸能够与你们在这个富有激情和活力的团队共事。

感谢机械工业出版社华章公司的编辑杨福川老师、孙海亮老师，在这一年多的时间中始终支持我的写作，你们的鼓励和帮助引导我顺利地全部书稿。

朱 凯

目 录 *Contents*

推荐序 思者常新，厚积薄发

前 言

第 1 章 浅谈企业级大数据平台的重要性

1.1 缺乏统一大数据平台的问题	2
1.1.1 资源浪费	2
1.1.2 数据孤岛	2
1.1.3 服务孤岛	3
1.1.4 安全存疑	3
1.1.5 缺乏可维护性和可扩展性	3
1.1.6 缺乏可复制性	4
1.2 构建统一大数据平台的优势	4
1.3 企业级大数据平台需要具备的基本能力	6
1.3.1 集群管理与监控	7
1.3.2 数据接入	7
1.3.3 数据存储与查询	7
1.3.4 数据计算	8
1.3.5 平台安全与管理	10
1.4 平台辅助工具	12
1.5 本章小结	13

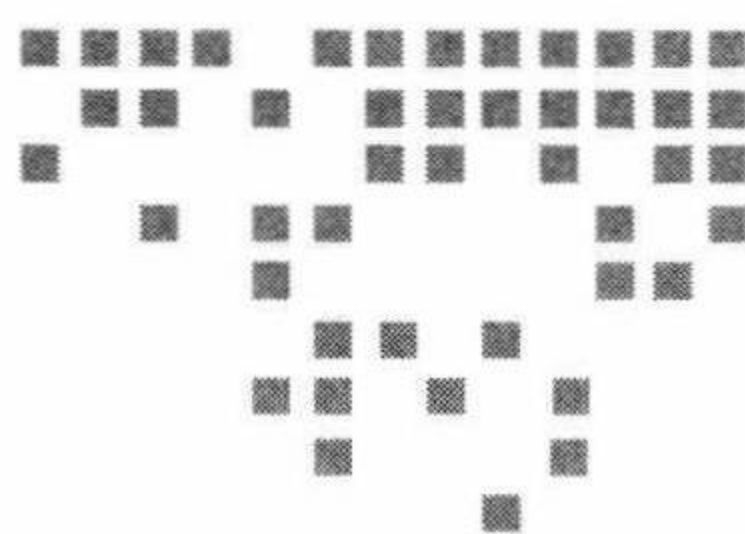
第 2 章 企业级大数据平台技术栈

介绍	15
2.1 HDFS	16
2.1.1 概述	16
2.1.2 RAID 技术	17
2.1.3 核心设计目标	18
2.1.4 命名空间	19
2.1.5 数据模型	20
2.1.6 Namenode 和 Datanode	20
2.1.7 使用场景	21
2.2 Zookeeper	22
2.2.1 概述	22
2.2.2 核心特性	23
2.2.3 命名空间	24
2.2.4 数据模型	24
2.2.5 节点状态监听	25
2.2.6 原子消息广播协议	25
2.2.7 使用场景	32
2.3 HBase	33
2.3.1 概述	33
2.3.2 数据模型	34
2.3.3 Regions	34

2.3.4	HBase Master	35	3.3	Ambari 的安装、配置与启动	55
2.3.5	Region Server	36	3.3.1	安装前的准备	55
2.3.6	MemStore 与 HFile	37	3.3.2	安装 Ambari-Server	62
2.3.7	使用场景	37	3.3.3	Ambari-Server 目录结构	64
2.4	YARN	38	3.3.4	配置 Ambari-Server	65
2.4.1	概述	38	3.3.5	启动 Ambari-Server	66
2.4.2	资源模型和 Container	40	3.4	新建集群	67
2.4.3	ResourceManager	40	3.4.1	设置集群名称并配置 HDP 安装包	67
2.4.4	ApplicationMaster	40	3.4.2	配置集群	69
2.4.5	NodeManager	41	3.5	Ambari 控制台功能简介	77
2.4.6	单一集群架构	41	3.5.1	集群服务管理	78
2.4.7	工作流程	41	3.5.2	集群服务配置	80
2.4.8	使用场景	43	3.5.3	辅助工具	82
2.5	Spark	43	3.6	本章小结	86
2.5.1	概述	43	第 4 章 构建企业级平台安全方案		87
2.5.2	数据模型	45	4.1	浅谈企业级大数据平台面临的安全 隐患	88
2.5.3	编程模型和作业调度	45	4.1.1	缺乏统一的访问控制机制	88
2.5.4	依赖	46	4.1.2	缺乏统一的资源授权策略	88
2.5.5	容错	47	4.1.3	缺乏 Hadoop 服务安全保障	89
2.5.6	集群模式	47	4.2	初级安全方案	89
2.5.7	使用场景	48	4.2.1	访问控制	89
2.6	本章小结	49	4.2.2	数据授权与管理	97
第 3 章 使用 Ambari 安装 Hadoop 集群		50	4.3	本章小结	110
3.1	概述	50	第 5 章 Hadoop 服务安全方案		111
3.2	集群设计	52	5.1	Kerberos 协议简介	111
3.2.1	主控节点	52	5.2	使用 FreeIPA 安装 Kerberos 和 LDAP	113
3.2.2	存储与计算节点	53			
3.2.3	安全认证与管理节点	54			
3.2.4	协同管理与其他节点	54			

5.2.1	安装 FreeIPA	115	7.2.4	实现服务层	181
5.2.2	IPA-Server 管理控制台功能 介绍	119	7.2.5	实现 RESTful 服务	181
5.2.3	IPA CLI 功能介绍	122	7.2.6	整合用户管理	183
5.3	开启 Ambari 的 Kerberos 安全 选项	127	7.3	RESTful 服务安全认证	184
5.3.1	集成前的准备	127	7.3.1	用户登录服务	185
5.3.2	集成 IPA	129	7.3.2	使用 JWT 认证	185
5.3.3	测试 Kerberos 认证	133	7.3.3	创建用户登录 RESTful 服务	188
5.4	本章小结	136	7.3.4	认证过滤器	194
			7.3.5	测试服务安全认证	198
第 6 章	单点登录与用户管理	137	7.4	数据仓库数据查询	200
6.1	集成单点登录	139	7.4.1	创建 JDBC 连接	200
6.1.1	CAS 简介	140	7.4.2	Kerberos 登录	202
6.1.2	安装 CAS-Server	141	7.4.3	使用 JDBC 协议查询	202
6.1.3	集成 Knox 网关与 CAS- Server	148	7.4.4	实现服务层与 RESTful 服务	206
6.1.4	集成 Ranger 与 CAS-Server	151	7.4.5	测试查询	207
6.1.5	集成 Ambari 与 CAS-Server	152	7.5	数据仓库元数据查询	208
6.2	实现统一的用户管理系统	155	7.5.1	使用 query 服务查询数仓元 数据	208
6.3	使用 Java 程序调用脚本	161	7.5.2	引入 JdbcTemplate 模块	209
6.4	创建 Ranger 扩展用户	166	7.5.3	增加 Hive 元数据库配置	210
6.5	本章小结	169	7.5.4	实现元数据持久层	211
			7.5.5	实现元数据服务层与 RESTful 服务	216
第 7 章	搭建平台管理端 RESTful 服务	170	7.5.6	测试元数据查询	218
7.1	搭建 RESTful 服务框架	170	7.6	本章小结	219
7.2	用户查询	174			
7.2.1	引入 LDAP 模块	174	第 8 章	Spark 任务与调度服务	220
7.2.2	配置 LDAP	174	8.1	提交 Spark 任务的 3 种方式	220
7.2.3	实现持久层	177	8.1.1	使用 Spark-Submit 脚本 提交	220

8.1.2 使用 Spark Client 提交	226	8.3.4 改进空间	241
8.1.3 使用 YARN RESTful API 提交	229	8.4 本章小结	241
8.2 查询 Spark 日志	234	附录 A Hadoop 简史	242
8.3 任务调度	236	附录 B Hadoop 生态其他常用组件 一览	245
8.3.1 引入 Quartz 模块	237	附录 C 常用组件配置说明	248
8.3.2 增加 Quartz 配置	237		
8.3.3 编写调度任务	240		



浅谈企业级大数据平台的重要性

不论你愿不愿意承认，大数据时代已经来临了。大数据潮流引领的技术变革正在悄无声息地改变着各行各业。虽说“大数据”是近些年才火热起来的词汇，但可以说“大数据”其实一直存在，只是由于技术的局限性使得人们在很长的一段时间里没有办法能够使用全量数据。但是随着技术的发展与革新，现在人们可以使用大数据技术来处理海量的数据了，这使得很多之前只能停留在理论研究层面的算法和思想现在能够付诸行动，比如现在很火爆的深度学习。与此同时，大数据技术这一新兴的工具也让人们拥有了一种新的思维模式，即大数据思维。

大数据思维注重全量样本数据而不是局部数据，注重相关性而不是因果关系。通过分析和挖掘数据将其转化为知识，再由知识提炼成智慧以获取洞察。大数据思维在很多行业都有用武之地，比如在银行业，基于大数据的风险控制体系就是一个很好的例子。通过大数据技术重构的机器学习算法不仅可以在全量样本数据上进行训练，还能引入更多的维度参与学习，从而构建一个比传统技术更高效、更准确的信用征信评分体系。同样，在电商行业也有很多大数据应用的例子。比如电商企业通过对手中大量的用户行为数据进行分析挖掘，可以得知用户的喜好并绘制出完善的用户画像。这使得电商企业能够更加了解自己的客户，从而对他们进行精准营销和相关商品推荐。

类似的例子数不胜数，这些案例的背后大数据技术功不可没。作为这个时代的参与者，我们的企业理应做好充足向大数据领域转型的技术准备，以免在这个时代落伍。

在这个转型的过程中最为重要的环节之一便是技术平台的建设。

1.1 缺乏统一大数据平台的问题

大数据思维需要依托大数据技术的支撑才能得以实现，所以隐藏在背后的支撑平台非常重要。正所谓下层基础决定上层建筑，没有一个牢固的地基是建不成摩天大楼的。我们不妨设想一下作为一个投身于大数据领域的企业，如果没有一个统一的大数据平台会出现什么问题。

1.1.1 资源浪费

通常在一个企业的内部会有多个不同的技术团队和业务团队。如果每个团队都搭建一套自己的大数据集群，那么宝贵的服务器资源就这样被随意地分割成了若干个小块，没有办法使出合力，服务器资源的整体利用率也无法得到保证。这种做法无疑是对企业资源的一种浪费。

其次大数据集群涉及的技术繁杂，其搭建和运维也是需要学习和运营成本的。这种重复的建设费时费力且没有意义，只会造成无谓的资源浪费。

1.1.2 数据孤岛

如果企业内部存在多个分散的小集群，那么首先各种业务数据从物理上便会被孤立地存储于各自的小集群之中，我们就没有办法对数据进行全量的整合使用，数据便失去了关联的能力，大数据技术使用全量数据进行分析的优势也丧失了。

其次，在这种情况下也很难实现对业务数据进行统一的模型定义与存储，一些相同的数据被不同的部门赋予了不同的含义，同一份数据就这样以不同的模型定义重复

地存储到了多个集群之中，不仅造成了不必要的存储资源浪费，还造成不同部门之间沟通成本的增长。

1.1.3 服务孤岛

企业内部各自为政的小集群的首要任务是支撑团队或项目组自身的业务场景来满足自身的需求，所以在实现功能的时候不会以面向服务的思维来抽象提炼服务，很可能都没有可以暴露出来供小集群外部使用的服务。退一步讲就算这些小集群有提供出来的服务，那么它们也缺乏统一的顶层设计，在做服务设计的时候没有统一的规则，导致提供的服务参差不齐，其访问入口也很有可能不统一。同时这些服务被分散在不同的集群之中，应用程序不能跨越多个集群使用所有的服务。

1.1.4 安全存疑

企业内部各项目组或团队自身维护的小集群通常都是只为支撑自身业务而实现的，不会同时面对多个用户。企业通过一些行政管理手段可以在一定程度上保障集群的安全。但是当团队人员扩充、集群规模扩大或是大数据集群的服务同时面向多个技术团队和业务部门的时候，很多问题就会显露出来。首当其冲的便是需要面对多用户的问题，集群不再只有一个用户，而是需要面对多个不同的用户。这就自然而然地引出一系列需要切实面对和解决的问题，比如用户的管理、用户的访问控制、服务的安全控制和数据的授权等。小集群通常都处于“裸奔状态”，基本没有什么安全防护的能力。集群安全涉及方方面面，是一个非常复杂的系统工程，不是轻易能够实现的。

1.1.5 缺乏可维护性和可扩展性

大数据领域的技术发展日新月异，其本身正处于一个高速的发展期，我们的集群服务会时不时需要进行更新以获得新的能力，或是需要安装补丁以修复 Bug。在这种情况下对多个小集群进行维护就会变得非常麻烦。同时当某个小集群性能达到瓶颈的时候也没有办法很容易地做到横向扩容。

1.1.6 缺乏可复制性

各自为政的小集群缺乏统一的技术路线，导致大数据集群的运维工作会缺乏可复制性。因为一个部门或者团队与其他部门使用的技术组件可能完全不一样，这样一个集群的安装、维护和调试等经验就没有办法快速复制和推广到其他团队或部门。同时在大数据应用研发方面也会存在同样的问题，正常来讲当我们做过的项目越多，从项目中获得的经验也就越多，我们能从这个过程中提炼、抽象和总结一些经验、规则或是开发框架来帮助我们加速今后的应用研发。但是技术路线的不统一很可能导致这些先验经验丧失后续的指导意义。

1.2 构建统一大数据平台的优势

如果我们能够化零为整，在企业内部从宏观、整体的角度设计和实现一个统一的大数据平台，引入单一集群、单一存储、统一服务和统一安全的架构思想就能较好地解决上述的种种问题。

1. 资源共享

使用单一集群架构，可以实现通过一个大集群来整合所有可用的服务器资源，通过一个大集群对外提供所有的能力。这样将所有服务器资源进行统一整合之后，能够更加合理地规划和使用整个集群的资源，并且能够实现细粒度的资源调度机制，从而使其整体的资源利用率更加高效。同时集群的存储能力和计算能力也能够突破小集群的极限。

不仅如此，因为只使用了一个大集群，所以我们现在只需要部署和维护一个集群，不需要重复投入人力资源进行集群的学习和维护。

2. 数据共享

使用单一存储架构，可以实现将企业内部的所有数据集中存储在一个集群之内，

方便进行各种业务数据的整合使用。这样我们便能够结合业务实际场景对数据进行关联使用，从而充分利用大数据技术全量数据分析的优势。同时，在这种单一存储架构之下，各种业务数据也可以进行统一的定义和存储，自然的也就不会存在数据重复存储和沟通成本增长的问题了。

3. 服务共享

通过统一服务架构，我们可以站在宏观服务化设计的角度来考虑问题，可将一套统一服务设计规则应用到所有的服务实现之上，同时也能统一服务的访问入口与访问规则。

除此之外，因为所有的服务是由一个统一的大集群提供的，这便意味着这些服务不存在孤岛问题，可以进行整合使用。

4. 安全保障

通过统一安全架构，可以从平台层面出发，设计并实现一套整体的安全保证方案。在单一集群架构的基础之上，可以实现细粒度的资源隔离；在单一存储架构的基础之上可以实现细粒度的数据授权；在单一服务架构的基础之上可以实现细粒度的访问控制；如此等等。

5. 统一规则

由于统一的大集群实现了技术路线的统一，这使得我们在后续应用开发的过程中有很多施展拳脚的空间。比如我们可以通过在大数据应用的开发过程中得到的一些经验总结，将这些经验整理为方法论和模型，再基于这些理论和模型实现一套大数据平台开发的 SDK。最终通过这套 SDK，可以很方便地将这些经验快速复制推广到整个企业内部。

6. 易于使用

在开发一款大数据产品或者业务的时候，我们应当将主要的精力放在业务的梳理