



分布式系统 常用技术及案例分析

柳伟卫 编著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

分布式系统 常用技术及案例分析

柳伟卫 编著

电子工业出版社
Publishing House of Electronics Industry
北京•BEIJING

内 容 简 介

本书全面介绍在设计分布式系统时所要考虑的技术方案，内容丰富、案例新颖，相关理论与技术实践较为前瞻。本书不仅仅介绍了分布式系统的原理、基础理论，同时还引入了大量市面上常用的最新分布式系统技术，不仅告诉读者怎么用，同时也分析了为什么这么用，并阐述了这些技术的优缺点。希望本书可以成为读者案头的工具书，供读者随手翻阅。

本书分为三大部分，即分布式系统基础理论、分布式系统常用技术以及经典的分布式系统案例分析。第一部分主要介绍分布式系统基础理论知识，总结一些在设计分布式系统时需要考虑的范式、知识点以及可能会面临的问题，其中包括线程、通信、一致性、容错性、CAP 理论、安全性和并发等相关内容；同时讲述分布式系统的常见架构体系，其中也包括最近比较火的 RESTful 风格架构、微服务、容器技术等。第二部分主要列举了在分布式系统应用中经常用到的一些主流技术，并介绍这些技术的作用和用法；这些技术涵盖了分布式消息服务、分布式计算、分布式存储、分布式监控系统、分布式版本控制、RESTful、微服务、容器等领域的内容。第三部分选取了以淘宝网和 Twitter 为代表的国内外知名互联网企业的大型分布式系统案例，分析其架构设计以及演变过程；这部分相当于是对第二部分零散的技术点做一个“串烧”，让读者可以结合技术的理论，看到实战的效果。

本书主要面向的读者是对分布式系统感兴趣的计算机专业的学生、软件工程师、系统架构师等。

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

分布式系统常用技术及案例分析 / 柳伟卫编著. —北京：电子工业出版社，2017.2
ISBN 978-7-121-30771-3

I. ①分… II. ①柳… III. ①分布式操作系统—研究 IV. ①TP316.4

中国版本图书馆 CIP 数据核字（2016）第 323368 号

责任编辑：陈晓猛

印 刷：北京京科印刷有限公司

装 订：三河市良远印务有限公司

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：43.75 字数：980 千字

版 次：2017 年 2 月第 1 版

印 次：2017 年 5 月第 2 次印刷

定 价：99.00 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）88254888，88258888

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

本书咨询联系方式：010-51260888-819 faq@phei.com.cn。

前 言

写作背景

我一直想写一本关于分布式系统方面的书。一方面是想把个人多年工作中涉及的分布式技术做一下总结，另一方面也想把个人的经验分享给广大的读者朋友。由于我的开发工作大都以 Java 为主，所以一开始的主题设想是“分布式 Java”，书也以开源方式发布在互联网上（网址为 <https://github.com/waylau/distributed-java>）。

后来，陈晓猛编辑看到了这本开源书，以及我关于分布式系统方面的博文，问我是否有兴趣出版分布式相关题材的图书。当然书的内容不仅仅是“分布式 Java”。

对于出书一事，我犹豫良久。首先，本身工作挺忙，实在无暇顾及其他；其次，虽然我之前写过超过一打的书籍（可见 <https://waylau.com/books/>），但多是开源电子书，时间、内容方面自然也就不会有太多约束，几乎是“想写就写，没有时间就不写”，这个跟正式出版还是存在比较大的差异的；最后，这本书涉及面相对较广，需要查阅大量资料，实在是太耗费精力。

但陈晓猛编辑还是鼓励我能够去尝试做这个事情。思索再三，于是我便答应。当然，最后这本书还是在规定时间内完成了。它几乎耗尽了我写作期间所有的业余和休息时间。

“不积跬步，无以至千里；不积小流，无以成江海。”虽然整本书从构思到编写完成的时间不足一年，但书中的大部分知识点，却是我在多年的学习、工作中积累下来的。之所以能够实现快速写作，一方面是做了比较严格的时间管理，另一方面也得益于我多年坚持写博客和开源书的习惯。

内容介绍

本书分为三大部分，即分布式系统基础理论、分布式系统常用技术以及经典的分布式系统案例分析。第一部分为第 1 章和第 2 章，主要介绍分布式系统基础理论知识，总结一些在设计分布式系统时需要考虑的范式、知识点以及可能会面临的问题。第二部分为第 3 章到第 8 章，主要列举了在分布式系统应用中经常用到的一些主流技术，并介绍这些技术的作用和用法。第

三部分为第 9 章和第 10 章，选取了以淘宝网和 Twitter 为代表的国内外知名互联网企业的大型分布式系统案例，分析其架构设计以及演变过程。

第 1 章介绍分布式系统基础理论知识，总结一些在设计分布式系统时需要考虑的范式、知识点以及可能会面临的问题，其中包括线程、通信、一致性、容错性、CAP 理论、安全性和并发等相关内容。

第 2 章详细介绍分布式系统的架构体系，包括传统的基于对象的体系结构、SOA，也包括最近比较火的 RESTful 风格架构、微服务、容器技术、Serverless 架构等。

第 3 章介绍常用的分布式消息服务框架，包括 Apache ActiveMQ、RabbitMQ、RocketMQ、Apache Kafka 等。

第 4 章介绍分布式计算理论和应用框架方面的内容，包括 MapReduce、Apache Hadoop、Apache Spark、Apache Mesos 等。

第 5 章介绍分布式存储理论和应用框架方面的内容，包括 Bigtable、Apache HBase、Apache Cassandra、Memcached、Redis、MongoDB 等。

第 6 章介绍分布式监控方面常用的技术，包括 Nagios、Zabbix、Consul、ZooKeeper 等。

第 7 章介绍常用的分布式版本控制工具，包括 Bazaar、Mercurial、Git 等。

第 8 章介绍 RESTful API、微服务及容器相关的技术，着重介绍 Jersey、Spring Boot、Docker 等技术的应用。

第 9 章和**第 10 章**分别介绍以淘宝网和 Twitter 为代表的国内外知名互联网企业的大型分布式系统案例，分析其架构设计以及演变过程。

源代码

本书提供源代码下载，下载地址为 <https://github.com/waylau/distributed-systems-technologies-and-cases-analysis>。

勘误和交流

本书如有勘误，会在 <https://github.com/waylau/distributed-systems-technologies-and-cases-analysis> 上进行发布。由于笔者能力有限，时间仓促，难免错漏，欢迎读者批评指正。读者也可以到博文视点官网的本书页面进行交流（www.broadview.com.cn/30771）。

您也可以直接联系我：

博客: <https://waylau.com>

邮箱: waylau521@gmail.com

微博: <http://weibo.com/waylau521>

开源: <https://github.com/waylau>

致谢

首先, 感谢电子工业出版社博文视点公司的陈晓猛编辑, 是您鼓励我将本书付诸成册, 并在我写作过程中审阅了大量稿件, 给予了我很多指导和帮助。感谢工作在幕后的电子工业出版社评审团队对于本书在校对、排版、审核、封面设计、错误改进方面所给予的帮助, 使本书得以顺利出版发行。

其次, 感谢在我十几年求学生涯中教育过我的所有老师。是你们将知识和学习方法传递给了我。感谢我曾经工作过的公司和单位, 感谢和我一起共事过的同事和战友, 你们的优秀一直是我追逐的目标, 你们所给予的压力正是我不断改进自己的动力。

感谢我的父母、妻子 Funny 和两个女儿。由于撰写本书, 牺牲了很多陪伴家人的时间。感谢你们对于我工作的理解和支持。

最后, 特别要感谢这个时代, 互联网让所有人可以公平地享受这个时代的成果。感谢那些为计算机、互联网所做出贡献的先驱, 是你们让我可以站在更高的“肩膀”上! 感谢那些为本书提供灵感的佳作, 包括《分布式系统原理与范型》《Unix Network Programming》《Enterprise SOA》《MapReduce Design Patterns》《Hadoop: The Definitive Guide》《Learning Hbase》《Advanced Analytics with Spark》《Pro Git》《Docker in Action》《淘宝技术这十年》《Hatching Twitter》, 等等, 详细的书单可以参阅本书最后的“参考文献”部分。

柳伟卫

2016年11月13日于杭州

目 录

第1章 分布式系统基础知识.....	1
1.1 概述	2
1.1.1 什么是分布式系统	2
1.1.2 集中式系统 VS. 分布式系统	3
1.1.3 如何设计分布式系统	4
1.1.4 分布式系统所面临的挑战	5
1.2 线程	6
1.2.1 什么是线程	6
1.2.2 进程和线程	7
1.2.3 编程语言中的线程对象	8
1.2.4 SimpleThreads 示例	11
1.3 通信	14
1.3.1 网络基础知识	14
1.3.2 网络 I/O 模型的演进	19
1.3.3 远程过程调用 (RPC)	33
1.3.4 面向消息的通信	41
1.4 一致性	43
1.4.1 以数据为中心的一致性模型	44
1.4.2 以客户为中心的一致性	45
1.5 容错性	46
1.5.1 基本概念	46
1.5.2 故障分类	47
1.5.3 使用冗余来掩盖故障	48
1.5.4 分布式提交	48
1.6 CAP 理论.....	52
1.6.1 什么是 CAP 理论.....	52

1.6.2 为什么说 CAP 只能三选二.....	53
1.6.3 CAP 常见模型.....	55
1.6.4 CAP 的意义	56
1.6.5 CAP 最新发展.....	56
1.7 安全性	57
1.7.1 基本概念	58
1.7.2 加密算法	60
1.7.3 安全通道	63
1.7.4 访问控制	72
1.8 并发	74
1.8.1 线程与并发	75
1.8.2 并发与并行	75
1.8.3 并发带来的风险	76
1.8.4 同步 (Synchronization)	78
1.8.5 原子访问 (Atomic Access)	83
第 2 章 分布式系统架构体系.....	85
2.1 基于对象的体系结构	86
2.1.1 分布式对象	86
2.1.2 微软 DCOM (COM+)	87
2.1.3 CORBA	88
2.1.4 Java RMI.....	90
2.2 面向服务的架构 (SOA)	93
2.2.1 架构 VS. 标准	94
2.2.2 SOA 的基本概念	95
2.2.3 基于 Web Services 的 SOA.....	97
2.2.4 SOA 的演变	112
2.3 REST 风格的架构.....	112
2.3.1 什么是 REST.....	112
2.3.2 REST 有哪些特征.....	113
2.3.3 Java 实现 REST 的例子	115
2.3.4 REST API 最佳实践	125
2.4 微服务架构 (MSA)	128
2.4.1 什么是 MSA.....	128

2.4.2 MSA VS. SOA	130
2.4.3 何时采用 MSA	134
2.4.4 如何构建微服务	135
2.5 容器技术	139
2.5.1 虚拟化技术	139
2.5.2 容器 VS. 虚拟机	139
2.5.3 基于容器的持续部署	142
2.6 Serverless 架构	149
2.6.1 什么是 Serverless 架构	150
2.6.2 Serverless 典型的应用场景	151
2.6.3 常见的 Serverless 框架	153
2.6.4 Serverless 架构原则	155
2.6.5 例子：使用 Serverless 实现游戏全球同服	157
第3章 分布式消息服务	164
3.1 Apache ActiveMQ	165
3.1.1 Apache ActiveMQ 简介	165
3.1.2 Apache ActiveMQ 安装配置	166
3.1.3 例子：producer-consumer	173
3.1.4 例子：使用 JMX 来监控 ActiveMQ	174
3.1.5 例子：使用 Java 实现 producer-consumer	176
3.2 RabbitMQ	180
3.2.1 RabbitMQ 简介	180
3.2.2 RabbitMQ 安装配置	181
3.2.3 例子：Work Queues	185
3.2.4 例子：Publish/Subscribe	191
3.2.5 例子：Routing	195
3.2.6 例子：Topics	200
3.2.7 例子：RPC	204
3.3 RocketMQ	210
3.3.1 RocketMQ 简介	210
3.3.2 RocketMQ 安装配置	213
3.3.3 例子：使用 Java 实现 producer-consumer	214
3.3.4 RocketMQ 最佳实践	219

3.4 Apache Kafka	223
3.4.1 Apache Kafka 简介	224
3.4.2 Apache Kafka 的核心概念	225
3.4.3 Apache Kafka 的使用场景	228
3.4.4 Apache Kafka 的安装、配置、使用.....	229
第 4 章 分布式计算	236
4.1 MapReduce	237
4.1.1 MapReduce 简介	237
4.1.2 MapReduce 的编程模型	238
4.1.3 MapReduce 的实现.....	243
4.1.4 MapReduce 的使用技巧	248
4.2 Apache Hadoop.....	251
4.2.1 Apache Hadoop 简介.....	252
4.2.2 Apache Hadoop 核心组件.....	253
4.2.3 Apache Hadoop 单节点上的安装配置	254
4.2.4 Apache Hadoop 集群上的安装配置	258
4.2.5 例子：词频统计 WordCount 程序	267
4.3 Apache Spark	272
4.3.1 Apache Spark 简介	272
4.3.2 Apache Spark 与 Apache Hadoop 的关系	274
4.3.3 Apache Spark 2.0 的新特性	275
4.3.4 Apache Spark 的安装和使用	279
4.3.5 Apache Spark 集群模式	280
4.4 Apache Mesos	282
4.4.1 Apache Mesos 简介	283
4.4.2 Apache Mesos 的安装、使用	285
4.4.3 设计高可用的 Mesos framework.....	289
第 5 章 分布式存储	296
5.1 Bigtable	297
5.1.1 Bigtable 简介	297
5.1.2 Bigtable 的数据模型	298
5.1.3 Bigtable 的实现.....	300

5.1.4 Bigtable 的性能优化	304
5.2 Apache HBase	308
5.2.1 Apache HBase 简介	308
5.2.2 Apache HBase 基本概念	310
5.2.3 Apache HBase 架构	318
5.2.4 Apache HBase 的安装、配置、使用	332
5.2.5 本地伪分布式	336
5.2.6 完全分布式	338
5.3 Apache Cassandra	342
5.3.1 Apache Cassandra 简介	342
5.3.2 Apache Cassandra 的应用场景	345
5.3.3 Apache Cassandra 的架构、数据模型	346
5.3.4 用于配置 Apache Cassandra 的核心组件	347
5.3.5 Apache Cassandra 的安装、配置、使用	349
5.4 Memcached	352
5.4.1 Memcached 简介	352
5.4.2 Memcached 的架构	353
5.4.3 Memcached 的安装、使用	355
5.4.4 Memcached 客户端	358
5.5 Redis	370
5.5.1 Redis 简介	370
5.5.2 Redis 的下载、安装、使用	372
5.5.3 Redis 的数据类型及抽象	372
5.6 MongoDB	392
5.6.1 MongoDB 简介	392
5.6.2 MongoDB 的安装、配置、运行	394
5.6.3 MongoDB 核心概念	401
5.6.4 MongoDB 的数据模型	406
5.6.5 示例：Java 连接 MongoDB	420
第6章 分布式监控	422
6.1 Nagios	423
6.1.1 Nagios 简介	423
6.1.2 Nagios 的安装、使用	424

6.1.3 Nagios 监控	428
6.1.4 Nagios 插件	446
6.2 Zabbix	448
6.2.1 Zabbix 简介	449
6.2.2 安装 Zabbix	451
6.2.3 Zabbix 对于容器的支持	460
6.2.4 Zabbix 基本概念	463
6.3 Consul	474
6.3.1 Consul 简介	475
6.3.2 Consul 架构	476
6.3.3 Consul 的安装和使用	478
6.3.4 Consul agent	492
6.4 ZooKeeper	501
6.4.1 ZooKeeper 简介	501
6.4.2 ZooKeeper 的安装和使用	505
6.4.3 ZooKeeper 内部工作原理	509
6.4.4 例子：ZooKeeper 实现 barrier 和 producer-consumer queue	514
第 7 章 分布式版本控制系统	522
7.1 Bazaar	523
7.1.1 Bazaar 简介	523
7.1.2 Bazaar 的核心概念	525
7.1.3 Bazaar 的安装	526
7.1.4 Bazaar 的使用	528
7.2 Mercurial	533
7.2.1 Mercurial 简介	533
7.2.2 Mercurial 的核心概念	533
7.2.3 Mercurial 的安装	537
7.2.4 Mercurial 的使用	538
7.3 Git	545
7.3.1 Git 简介	545
7.3.2 Git 的安装	546
7.3.3 Git 的基础概念	548
7.3.4 Git 的使用	551

第 8 章 RESTful API、微服务及容器技术	578
8.1 Jersey	579
8.1.1 Jersey 简介	579
8.1.2 Jersey 的模块和依赖	580
8.1.3 JAX-RS 核心概念	583
8.1.4 例子：用 SSE 构建实时 Web 应用	595
8.2 Spring Boot	603
8.2.1 Spring Boot 简介	603
8.2.2 Spring Boot 的安装	604
8.2.3 Spring Boot 的使用	610
8.2.4 Spring Boot 的属性与配置	615
8.3 Docker	620
8.3.1 Docker 简介	621
8.3.2 Docker 的特性	621
8.3.3 Docker 的概念和原理	622
8.3.4 Docker Engine 的安装	628
8.3.5 Docker 的使用	633
第 9 章 淘宝网：“双 11”神话的缔造者	636
9.1 从 LAMP 到 Java 平台的转变	637
9.1.1 淘宝网的诞生与发展	637
9.1.2 “平民英雄” LAMP 架构	638
9.1.3 数据库更改为 Oracle	639
9.1.4 向 Java 平台转变	642
9.2 坚定不移地走“去 IOE”的道路	643
9.2.1 使用小型机、EMC 存储	644
9.2.2 考虑“去 IOE”	644
9.2.3 如何去“I”	646
9.2.4 如何去“O”	649
9.2.5 如何去“E”	650
9.3 打造云计算，决战“双 11”	653
9.3.1 “大淘宝”战略简介	653
9.3.2 成立阿里云，专注云计算	656

9.3.3 利用大数据优化物流	660
9.3.4 技术是决胜“双11”的关键.....	661
第10章 Twitter：实时信息传递的王者	664
10.1 缓存，让响应更快	665
10.1.1 Twitter 的诞生	665
10.1.2 RoR 的蛮荒时代	666
10.1.3 使用 Memcached.....	667
10.2 服务拆分与治理	668
10.2.1 关系数据库不是万灵药	668
10.2.2 系统拆分，平台转换	670
10.2.3 Finagle	670
10.3 抗击流量的洪流	672
10.3.1 业务的重新设计	673
10.3.2 Storm 处理实时的大数据.....	675
10.3.3 从 Storm 到 Heron.....	676
参考文献	680



第 1 章

分布式系统基础知识

1.1 概述

自从有了计算机，人类的生活就发生了巨大的变化，原来人力需要数年才能解决的计算问题，计算机“分分钟”就搞定了。回顾历史，在19世纪50年代，计算机在诞生初期，主要用于军事方面，而天才的数学家图灵建造了计算机来破译德军的“英格玛(Enigma)”密码系统。那时候的计算机如一个庞大笨重的巨大人，足足占据了好几间房子的空间。在19世纪80年代，计算机的体积越来越小，开始被广泛应用于科研计算，但那时候计算机的价格仍然是极其昂贵的，只有少数几个科研单位或者大学才能买得起。今天，个人计算机或者以手机为代表的智能设备已经走进寻常百姓家了。每个人几乎都拥有手机，手机不仅仅是通信工具，还能发语音、看视频、玩游戏，让人与人之间的联系变得更加紧密。智能手环随时监控你的身体状况，并根据你每天的运动量、身体指标来给你提供合理的饮食运动建议。出门逛街甚至不需要带钱包了，吃饭购物搭车时使用手机就可以支付费用，多么方便快捷。智能家居系统更是你生活上的“管家”，什么时候该睡觉了，智能家居系统就自动拉上窗帘，关灯；早上起床了，智能家居系统会自动拉开窗帘，并播放动人的音乐，让你可以愉快地享受新的一天的来临；你再也不用担心家里的安全情况，智能家居系统会帮你监控一切，有异常情况时会及时发送通知到你的手机，让你第一时间掌握家里的状况。

毫无疑问，计算机改变了人类的工作和生活方式，而计算机系统也正在进行一场变革。没错，任何一个手机应用，或者智能App，都离不开背后那个神秘的巨人——分布式系统。正是那些看不见的分布式系统，每天处理着数以亿计的计算，提供可靠而稳定的服务。

本章就揭开分布式系统的神秘面纱。

1.1.1 什么是分布式系统

《分布式系统原理与范型》一书中是这样定义分布式的：

“分布式系统是若干独立计算机的集合，这些计算机对于用户来说就像单个相关系统”。

这里面包含了2个含义：

- 硬件独立；
- 软件统一。

什么是硬件独立？所谓硬件独立，是指计算机机器本身是独立的。一个大型的分布式系统，会由若干台计算机来组成系统的基础设施。而软件统一，是指对于用户来说，用户就像是跟单个系统打交道。就好比我们每天上网看视频，视频网站对我们来说就是一个系统软件，它们背后是如何运作的，部署了几台服务器，每台服务器是干什么的，这些对用户来说是透明的，不

可见的。用户不关心背后的这些服务器，用户所关心的是，今天访问的这个网站能提供什么样的节目，视频运行是否流畅，卡不卡帧、清晰度如何等。

而软件统一的另外一个方面是指，分布式系统的扩展和升级都比较容易。分布式系统中的某些节点发生故障，不会影响整体系统的可用性。用户和应用程序交互时，不会察觉哪些部分正在替换或者维修，也不会感知到新加入的部分。

1.1.2 集中式系统 VS. 分布式系统

集中式系统主要部署在 HP、IBM、SUN 等小型机以上档次的服务器中，把所有的功能都集成到主服务器上，这样对服务器的要求很高，性能也极其苛刻。它们的主要特色在于年宕机时间只有几小时，所以又统称为 z 系列（zero，零）。AS/400 主要应用在银行和制造业，还用于 Domino，主要的技术在于 TIMI（技术独立机器界面）、单级存储，有了 TIMI 技术可以做到硬件与软件相互独立。RS/6000 比较常见，用于科学计算和事务处理等。这类主机的单机性能一般都很强，带多个终端。终端没有数据处理能力，运算全部在主机上进行。现在的银行系统大部分都是这种集中式的系统，此外，在大型企业、科研单位、军队、政府等也有应用。

集中式系统主要流行于 20 世纪。现在还在使用集中式系统的，很大一部分原因是为了沿用原来的软件，而这些软件往往很昂贵。优点是便于维护，操作简单。但这样的系统也有缺陷，就是不出问题还好，一出问题，就会造成单点故障，所有功能就都不能正常工作了。由于集中式的系统相关的技术只被少数厂商所掌握，个人要对这些系统进行扩展和升级往往也比较麻烦，一般的企业级应用很少会用到集中式系统。图 1-1 是一个典型的集中式系统的示意图。

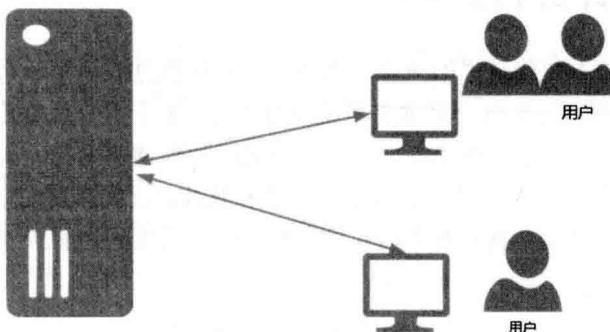


图 1-1 集中式系统示意图

而分布式系统恰恰相反。分布式系统是通过中间软件来对现有计算机的硬件能力和相应的软件功能进行重新配置和整合。它是一种多处理器的计算机系统，各处理器通过互连网络构成统一的系统。系统采用分布式计算结构，即把原来系统内中央处理器处理的任务分散给相应的处理器，实现不同功能的各个处理器相互协调，共享系统的外设与软件。这样就加快了系统的