

做个火影般的架构师

架构解密

从分布式到微服务

Microservices, Kubernetes, SOA, Distributed Memory, Elastic Search, Kafka, CAP, NUMA, GlusterFS, Actor, Akka, Replicating, Spring Cloud, ZeroD, etc.

Leader-us 著

架构解密

从分布式到微服务

Leader-us 著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

分布式架构与微服务平台是当今 IT 界的关键技术，也是资深软件工程师和系统架构师必须掌握的核心技术。本书以从传统分布式架构迁移到基于容器技术的微服务架构为主线，全面、透彻地介绍了与分布式架构及微服务相关的知识和技术。本书一开始并没有提及分布式的枯燥理论，而是讲述了一段精彩的 IT 发展史，其中重点讲述了大型机、UNIX 小机器的没落与 X86 平台的崛起，从而巧妙地引出 CPU、内存、网络、存储的分布式演进过程，这恰恰是分布式软件系统赖以运行的“物质基础”。然后简明扼要地介绍了进行系统架构所必需的网络基础，并详细介绍了分布式系统中的经典理论、设计套路及 RPC 通信，对内存、SOA 架构、分布式存储、分布式计算等进行了深度解析，最后详细介绍了全文检索与消息队列中间件，以及微服务架构所涉及的重点内容。

本书是 Leader-us 多年架构经验的倾情分享，主要面向关注分布式架构及微服务，以及有志于成为实力派架构师的 IT 人士。

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

架构解密：从分布式到微服务 / Leader-us 著. —北京：电子工业出版社，2017.7
ISBN 978-7-121-31562-6

I. ①架… II. ①L… III. ①分布式计算机系统—架构 IV. ①TP338.8

中国版本图书馆 CIP 数据核字 (2017) 第 108266 号

策划编辑：张国霞

责任编辑：徐津平

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：18.75 字数：390 千字

版 次：2017 年 7 月第 1 版

印 次：2017 年 7 月第 1 次印刷

印 数：2500 册 定价：79.00 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 zits@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。

作者简介

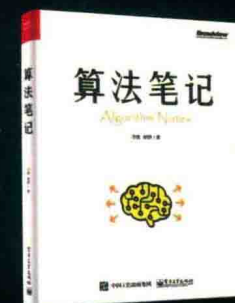
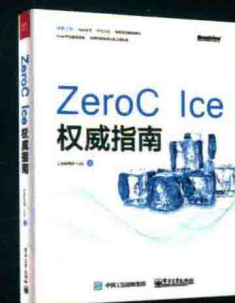
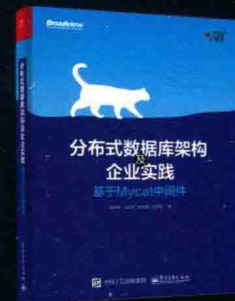
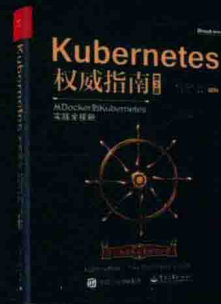


Leader-us

本名吴治辉，惠普资深软件架构师，国内知名开源分布式数据库中间件 Mycat 的发起人，精通 Java 编程，拥有超过 16 年软件研发经验，专注于电信和云计算方面的软件研发，参与过众多分布式与云计算相关的大型项目架构设计和 Coding，是业界少有的具备很强 Coding 能力的 S 级资深架构师；曾经选拔和培养了大批优秀 Java 工程师，他们中的大多数人进入知名软件公司参与核心研发，也有一些人选择创业。

Leader-us 也是《ZeroC Ice 权威指南》《Kubernetes 权威指南：从 Docker 到 Kubernetes 实践全接触》的作者。

好书力荐



作为一名架构师，我们要专业，要能看懂代码，即使光着臂膀去机房，也能独当一面！即使同事搞不定问题，或者撂挑子，你也能给老大一个坚定的眼神：不怕，有我在！还能在会议上滔滔不绝，如若无人，让不懂技术的妹子看你时眼神迷离，就好象落霞与孤鹜齐飞！

—Leader-us—

写给像笔者的你

我们都是 IT 人，所以，我们注定了很像。

我们可能小时候都挺聪明，学习也挺好，也早恋（可能纯洁度不同）。这一切都是有关联的，因为早恋所以你写情书，所以你有了点文采，又所以喜欢读笔者的文字，于是，你成了笔者的第 1 个读者，虽然我们分布在不同的“机房”中。

我们因为都受过严格、系统的全面教育，所以骨子里是温顺的，性格上是温柔的。我们因为在智商上高于情商的概率是 99%，所以多年独占风云榜之状元称谓——呆。我们一起努力的结果，是验证了那句话——科学无国界。在《生活大爆炸》《IT 狂人》等热播美剧中，我们终于找到祖国之外的同类，于是我们开始自恋地打广告：我很呆但我很幽默。

🙄不知道女友要什么东西啊，分手了才知道人家暗示了那么多次自己都不懂啊！

我们都是高学历的概率是 99%，我们都近视的概率是 99%，我们未富先胖的概率是 99%，我们未老先白头的概率是 99%，我们目前在北上广或者未来在北上广的概率是 99%，我们背井离乡的概率是 99%，我们的计算机内存超过 4GB 的概率是 99%，我们喜欢新计算机胜过于喜欢新女友的概率是 99%。

我们曾经是众人眼中的宠儿，但不知从何时起，沦落为新一代的“农民工”。在《死神来了之中国特供版》里，每年都有几个 IT 精英注定被永远地带走，我们在默默悲伤的同时，心里也在默默祈祷：O My God，让我活到 82 岁吧，就算没有 28 岁的小娇妻。

我们听过最多的公司是微软，我们最离不开的品牌是 Windows，虽然对于它们的评价，我们无法达成一致；同样，对于马云及乔布斯，我们也有着不同的评价，虽然他们缔造的“帝国”对我们的生活都产生了重要影响。

最后也是最重要的一点，我们都生活在一个有意思的时代，这个时代无法用任何哲学理论来左右我们的思想和行为。金钱向左，理想朝右，我们始终不放弃一个宏伟梦想：寻找最优秀的算法，收获金钱，实现理想。

吐槽归吐槽，言归正传，笔者假设你跟笔者一样是个有为青年，目标是成为 IT 精英，目前烦透了信息系统、Web 及低水平的重复编码工作，打算进阶架构师队伍，并下定决心潜心修行一年半载，脱掉程序猿的旧外套，换上“土豪金”的 IT 新人套装，那么，请你准备如下软硬件，开始和笔者一起，探秘分布式架构的奥义，走向“云端”。

- 有 8GB 内存的计算机一台，4GB 勉强过关。
- 计算机保持联网，遇到问题能随时“谷歌”。
- 笔者的 QQ 号码，该号码在本书某个 DEMO 的代码中。
- Eclipse 或你熟悉的 Java 开发工具。

除此之外，更重要的是以下几点。

- 不求快，但求坚持到底，系统学习比局部掌握更重要。
- 不怕错，就怕蒙混过关，尝试和出错是学编程的王者之道。
- 不怕动手，就怕只动眼，原理与实践都重要，技术都是实践和总结出来的。

Leader-us

2017 年 5 月 31 日

目 录

第 1 章 大话分布式系统

1

1.1 IT 争霸战	1
1.1.1 划时代的第一台计算机	1
1.1.2 IT 界的恐龙时代	4
1.1.3 贵族的没落与平民的胜利	6
1.1.4 ARM 新贵的爆发	10
1.1.5 超级计算机的绝地反击	11
1.2 分布式系统的开国元勋	13
1.3 分布式系统的基石：TCP/IP	17
1.4 从无奈到崛起的 CDN 网	19
1.5 这是一个最好的时代	21

第 2 章 “知识木桶”中的短板——网络基础

23

2.1 即使高手也不大懂的网络	23
2.2 NIO，一本难念的经	30
2.2.1 难懂的 ByteBuffer	30
2.2.2 晦涩的“非阻塞”	39
2.2.3 复杂的 Reactor 模型	41

2.3	AIO, 大道至简的设计与苦涩的现实	45
2.4	网络传输中的对象序列化问题	50
第 3 章	分布式系统的经典基础理论	55
3.1	从分布式系统的设计理念说起	55
3.2	分布式系统的一致性原理	58
3.3	分布式系统的基石之 ZooKeeper	61
3.3.1	ZooKeeper 的原理与功能	61
3.3.2	ZooKeeper 的场景案例分析	65
3.4	经典的 CA 理论	69
3.5	BASE 准则, 一个影响深远的指导思想	72
3.6	重新认识分布式事务	73
3.6.1	数据库单机事务的实现原理	73
3.6.2	经典的 X/OpenDTP 事务模型	75
3.6.3	互联网中的分布式事务解决方案	78
第 4 章	聊聊 RPC	83
4.1	从 IPC 通信说起	83
4.2	古老又有生命力的 RPC	85
4.3	从 RPC 到服务治理框架	91
4.4	基于 ZeroC Ice 的微服务架构指南	94
4.4.1	微服务架构概述	95
4.4.2	ZeroC Ice 微服务架构指南	100

第 5 章 深入浅析内存

107

5.1 你所不知道的内存知识	107
5.1.1 复杂的 CPU 与单纯的内存	107
5.1.2 多核 CPU 与内存共享的问题	110
5.1.3 著名的 Cache 伪共享问题	113
5.1.4 深入理解不一致性内存	115
5.2 内存计算技术的前世今生	118
5.3 内存缓存技术分析	123
5.3.1 缓存概述	123
5.3.2 缓存实现的几种方式	125
5.3.3 学习 Memcache 的内存管理技术	127
5.3.4 Redis 的独特之处	129
5.4 内存计算产品分析	131
5.4.1 SAP HANA	131
5.4.2 Hazelcast	133
5.4.3 VoltDB	135

第 6 章 深入解析分布式存储

138

6.1 数据存储进化史	138
6.2 经典的网络文件系统 NFS	145
6.3 高性能计算领域的分布式文件系统	148
6.4 企业级分布式文件系统 GlusterFS	150
6.5 创新的 Linux 分布式存储系统——Ceph	153
6.6 软件定义存储	160

第 7 章	聊聊分布式计算	166
7.1	不得不说的 Actor 模型	166
7.2	Actor 原理与实践	170
7.3	初识 Akka	177
7.4	适用面很广的 Storm	185
7.5	MapReduce 及其引发的新世界	194
第 8 章	全文检索与消息队列中间件	201
8.1	全文检索	201
8.1.1	什么是全文检索	201
8.1.2	起于 Lucene	202
8.1.3	Solr	206
8.1.4	ElasticSearch	209
8.2	消息队列	217
8.2.1	消息队列概述	217
8.2.2	JEE 专属的 JMS	221
8.2.3	生生不息的 ActiveMQ	226
8.2.4	RabbitMQ	231
8.2.5	Kafka	238
第 9 章	微服务架构	244
9.1	微服务架构概述	244
9.1.1	微服务架构兴起的原因	244
9.1.2	不得不提的容器技术	246
9.1.3	如何全面理解微服务架构	249

9.2	几种常见的微服务架构方案	253
9.2.1	ZeroC IceGrid 微服务架构	253
9.2.2	Spring Cloud 微服务架构	256
9.2.3	基于消息队列的微服务架构	259
9.2.4	Docker Swarm 微服务架构	261
9.3	深入 Kubernetes 微服务平台	263
9.3.1	Kubernetes 的概念与功能	263
9.3.2	Kubernetes 的组成与原理	268
9.3.3	基于 Kubernetes 的 PaaS 平台	272

第 1 章

大话分布式系统

分布式世界是一个很复杂的世界，任何技术都不是孤立的存在，任何技术都无法适应所有场景。作为一名分布式系统架构师或资深研发人员，你必须尽可能多地学习与之相关的各种知识，掌握各种技术的演进路线，从一名编程狂人逐渐升级成为一名博学的 IT 专家，实践与理论并行、代码与页码齐飞，唯有如此，你才能更好地成就未来。

1.1 IT 争霸战

1.1.1 划时代的第一台计算机

在两亿多年前的侏罗纪时代，地球上生活着数量庞大的恐龙家族，它们统治着海洋、陆地和天空，并幸福地生活了上亿年，却突然灭亡。在那个时代称霸的恐龙体型巨大，它们所在的时代是如此令人着迷，以至于直到现在，好莱坞导演们也不忘和我们 IT 界的同行携手，用先进的 IT 技术制造出一个个在视觉上令人震撼的史前怪兽。

有趣的是，计算机领域也呈现出与侏罗纪时代的恐龙同样的发展轨迹：从早期个体的强大逐渐发展为群体的强大。据记载，世界上第一台电子数字式计算机于 1946 年情人节（2 月 14

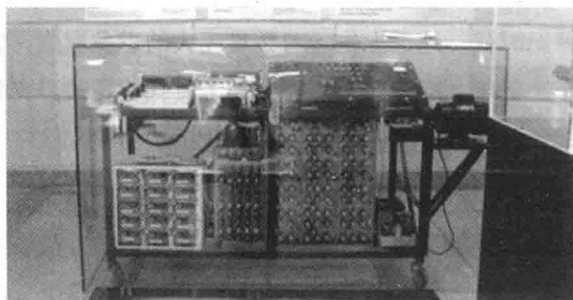
日) 诞生并在美国宾夕法尼亚大学正式投入运行, 它的名字是 ENIAC, 其主要设计制造者毛克利申请获得了美国专利。ENIAC 有 17468 个真空电子管, 并使用电容器进行数值存储, 以电量表示数值, 数据输入时采用打孔读卡, 并采用二进位制计算, 耗电 174 千瓦, 占地 170 平方米, 重达 30 吨, 每秒钟可进行 5000 次加法运算。



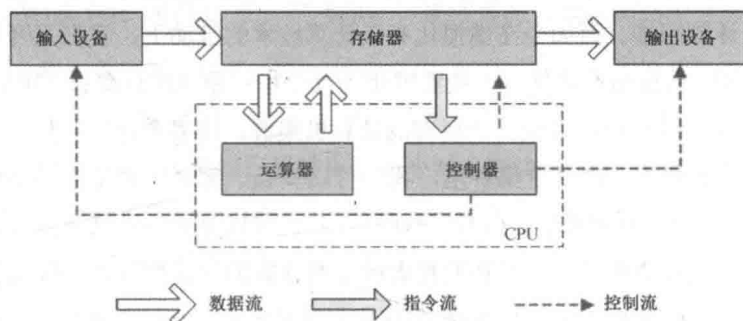
恐龙的王者——霸王龙的体重最大为 14 吨, 第一台巨无霸电子计算机重 30 吨, 该计算机仅从重量上已经完胜霸王龙。虽然每秒只能进行 5000 次加法运算, 现在任意一台 PC 的计算能力都超过它千倍, 但看看下面这段记录, 作为 IT 人, 你应该为这个鼻祖的诞生感到骄傲:

中国的古代科学家祖冲之利用算筹, 耗费 15 年心血, 才把圆周率计算到小数点后 7 位数。一千多年后, 英国人香克斯以毕生精力计算圆周率, 才计算到小数点后 707 位。而使用 ENIAC 进行计算, 仅用了 40 秒就达到了这个记录, 还发现香克斯的计算中, 第 528 位是错误的。

水落石出的真相: 依俄华州立大学物理系约翰·文森特·阿坦那索夫和研究生克利福德·贝里在 1939 年制造出一台完整的具备现代计算机 4 个核心要素(二进制、内存、I/O、计算单元)的样机 ABC (Atanasoff Berry Computer), 该样机后来被销毁而没有流传下来。1973 年美国明尼苏达地区法院正式宣判, 吊销毛克利的专利, 并肯定了阿坦那索夫才是真正的现代计算机的发明人。如下所示是 ABC 的复制品, 即使该复制品是人们靠回忆还原的, 它也是世界上第一台现代计算机, 让我们默默地瞻仰 1 分钟……



但凡新的学科出现，总有天才人物披荆斩棘地做开路先锋，他们有着常人所不具备的敏锐洞察力和想象力。让我们再看看奠定了现代计算机系统结构的经典理论——冯·诺依曼体系：计算机硬件由运算器、控制器、存储器、输入设备和输出设备五大部分组成。直到今天，计算机仍没有跳出该体系的范畴。冯·诺依曼洋洋洒洒的 101 页关于计算机系统结构的技术报告，奠定了他在计算机领域的地位，但他却亲手把“计算机之父”的头衔戴在了同时代的天才阿兰·图灵的头，可见图灵对计算机发展所做出的贡献。



图灵是罕见的天才数学家和计算机科学家，天生悟性过人，16岁就能弄懂爱因斯坦的相对论，并能运用那深奥的理论，独立推导出力学定律。1935年，图灵年仅23岁且刚刚大学毕业，就被剑桥大学国王学院甄选为研究员，成为剑桥大学有史以来最年轻的研究员。1936年，图灵在伦敦权威的数学杂志上发表了一篇划时代的重要论文《可计算数字及其在判断性问题中的应用》，在该论文中首次提出了奠定现代计算机的理论基础的“图灵机”理论。在曼彻斯特大学，图灵度过了其短暂生命的最后几年，“人工智能”是他发出的最后的生命之光，他是这一领域开天辟地的大师。从计算机理论到实践，图灵贡献无人能出其右。1951年，图灵39岁，被英国皇家学会选为会员，成为其家族中第4位皇家学会会员。曼彻斯特大学也因为图灵的存在，被英国皇家学会认定为国家计算机科学的最高学术机构。

天才或许真的只是为了拯救我们这些碌碌无为的凡人而落入凡间的灵魂吧，他们在人间绽放了无与伦比的耀眼光芒后，就匆忙离去，只留下我们无限怅惘。

——Leader-us

1954年6月8日，42岁的图灵吃了一小口含有氰化钾的苹果，绝世而去。1998年6月22日，世界各地的计算机大师齐聚伦敦纪念他们的“创业领袖”，英国下议院向科学家们道歉，承认在44年前对图灵做出了不公正的审判，并且当即修改法律，同性恋不再非法。至于那个被咬

了一小口的苹果，则又被另一个 IT 奇才——乔布斯发扬光大，国外媒体通过苹果自 iPhone 上市后的每一个季度财报的相关统计得知，这款革命性的移动设备目前的全球销量已经超过 5 亿台。

第一台电子计算机诞生以后，一个日新月异的 IT 时代到来了。一方面单台计算机的性能每年都在提升：从最早的 8 位 CPU 到现在的 64 位 CPU；从最早的 MB 级内存到现在的 GB 级内存；从慢速的机械硬盘存储到现在的固态 SSD 硬盘存储。另一方面，分布式架构技术让我们把单台计算机的计算能力、内存、I/O 等传统部件“分布”到联网的各个单独的计算机节点上，最终组成一个超级计算网格。而如今在虚拟化和云计算技术的推动下，我们又开始实现另一个极致的新技术：一台计算机通过软件方式被虚拟化为几个相互独立的计算机（虚机），从而使得一台计算机变成了 N 台计算机，形成一个局部的计算机集群，许多个这样的集群互联进而形成一个规模更大的计算机集群，在这个集群里，我们可以安装、部署多种不同的操作系统，彼此相互独立地完成各种任务，而当某个虚机发生故障后，我们可以立即自动转移到其他机器上重建，也可以根据系统的负载情况动态创建和消耗虚机。通过虚拟化管理软件，仅仅需要几十秒，一个完全虚拟的 Linux 机器便可以立即使用，也仅仅只需几十秒，一个安装好了 MySQL、Tomcat、JDK 的虚拟机像被启动，你可以部署你的 Java 程序，供别人使用，这是怎样一个神奇的世界呢？本书将为你揭露上述神奇表象背后的密码，使你也拥有实现上述目标的魔力。

1.1.2 IT 界的恐龙时代

ENIAC 之后，电子计算机便进入了 IBM 主导的大型机时代，IBM 大型机之父吉恩·阿姆达尔被认为是有史以来最伟大的计算机设计师之一。1964 年 4 月 7 日，在阿姆达尔的带领下，历时三年，耗费 50 亿美元，第一台 IBM 大型机 SYSTEM/360（简称 S/360）诞生。这项 50 亿美元的投资甚至超过了原子弹的研究资费，但最终被证实这是一次改变了商业运作的历史性变革，这使得 IBM 在 20 世纪 50~60 年代统治整个大型计算机工业，奠定了 IBM 计算机帝国的江山。IBM 的大型机过去曾支撑美国航天登月计划，而近 50 年以来，IBM 主机一直服务于金融等核心行业的关键业务领域，IBM 更是集中精力研发大型机和小型机。由于高可靠性和超强的计算能力，即便在目前 X86 和云计算飞速发展的情况下，IBM 的大型机仍然牢牢占据着一定的高端市场份额。