





## 图书在版编目(CIP)数据

分布式系统设计实践 / 李庆旭编著. -- 北京: 人民邮电出版社, 2019. 11  
ISBN 978-7-115-51945-0

I. ①分… II. ①李… III. ①分布式操作系统—系统设计 IV. ①TP316.4

中国版本图书馆CIP数据核字(2019)第190630号

## 内 容 提 要

本书对近年来涌现出的各种主流分布式技术做了简要介绍和全面梳理。本书将分布式系统中涉及的技术分为前端构造技术、分布式中间件技术和分布式存储技术三大类,对每类技术都详细介绍了其原理、主流实现的设计思想和架构,以及相关应用场景。此外,本书还总结了分布式系统的构建思想,并分别针对业界几个非常成功的大型分布式系统(谷歌搜索系统、淘宝网电商平台、阿里云公有云平台、领英社交平台)进行了案例研究。

本书适合业界的架构师、工程师、项目管理人员,也适合大中专院校的高年级本科生和研究生参考和阅读。

- 
- ◆ 编 著 李庆旭  
责任编辑 杨海玲  
责任印制 马振武
  - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号  
邮编 100164 电子邮件 315@ptpress.com.cn  
网址 <http://www.ptpress.com.cn>  
北京鑫正大印刷有限公司印刷
  - ◆ 开本: 800×1000 1/16  
印张: 14  
字数: 292千字 2019年11月第1版  
印数: 1-2500册 2019年11月北京第1次印刷
- 

定价: 59.00元

读者服务热线: (010)81055410 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147号

# 序

我很喜欢看书。从小到大，我看过很多书。

每看一本好书，我都会对书的作者充满感激和敬意，因为他愿意花那么多精力，写一本那么好的书给别人看，实在是非常值得尊敬。我知道，写书是赚不了多少钱的，一个人不为名、不为利，却愿意花很多精力写一本好书给别人去读、去欣赏，实在是难能可贵！

因此，我一直有一个想法，就是有一天，我也要尽力去写一本“好点儿”的书，一则对自己的所学做一个系统的整理，二则也为读者做些贡献。

我不知道是否有人会喜欢我的这本书，但无论如何，我已经尽力把它写好，因为我想和那些自己曾经读过的好书的作者一样，奉献给读者一本值得读的书。

我从1991年开始接触计算机，绝对算是一个“老码农”了。屈指一算，从那时到现在，计算机技术经历了几次重要的更新换代。

第一次发生在20世纪90年代初，是从DOS到Windows的换代。这一次换代淘汰了一批落后的公司和产品，亦有一批新的公司和产品崛起。当年DOS下的WPS如日中天，然而Windows一经推出，就难觅WPS的踪迹了。直到21世纪初，新版Windows下的金山WPS才姗姗而来。如果当年求伯君先生能及时推出新版的Windows下的WPS，那么中文办公软件的市场也许就不是今天这样了。

第二次换代发生在20世纪90年代末到21世纪初，Java和.NET这样的基于虚拟机的开发技术开始兴起。在DOS/Windows下，C/C++一直是主流的开发语言，而随着计算机软件的复杂度越来越高，C/C++语言开发效率低、查错难的弊端越来越突出，于是Java和.NET这样的基于虚拟机的开发技术日益受到开发者的青睐。这一波技术换代使Java的应用范围越来越广，也动摇了微软公司的技术在开发领域中无可替代的地位。

进入21世纪后，由于互联网技术，尤其是移动互联网技术的发展，诞生了一批大型互联网公司（谷歌、亚马逊、Facebook、百度、阿里巴巴、腾讯等）。在互联网发展的初期，这些公司在汹涌而来的大数据面前毫无准备，被折腾得筋疲力尽。痛定思痛后，以谷歌为代表的互联网新贵们，开发了许多优秀的技术（GFS、Bigtable、Hadoop等），完美而优雅地解决了大多数曾经令它们头疼的问题。继谷歌将其技术以论文的形式发表后，业界就开始模仿其设计和架构，相继开源了许多优秀的产品（如Apache Hadoop系列产品）。从此，大型分布式技术的大门被轰然推开，一个又一个大型的分布式应用（淘宝网上商城、京东网上商城、微信、谷歌搜索、Twitter社交平台、Facebook社交平台等）相继诞生。这就是以移动互联网和分布式计算的大规模应用为代表的第三次技术换代。无疑，这一次换代又

是一次各大公司实力的大洗牌。

第四次大的技术换代是什么？也许它正在进行（如果它是人工智能的话），也许还没有开始……

每一次大的技术换代，都会有公司和技术被无情地淘汰，从 Novell 到 Borland，到 Netscape，再到 Nokia，昔日霸主早已风光不再。但被淘汰的不仅有大公司，同时被淘汰的还有不少从事技术工作的工程师和架构师。因此，为了不被技术浪潮无情地拍倒，每一个 IT 从业者都应该对新技术心存敬畏，都应该关注技术的发展。

今天，对于在 IT 行业从事技术工作的人，无论是工程师、架构师还是管理者，也无论从事的工作是否与分布式相关，都应该了解分布式技术，因为总有一天，你会遇到它、接触它、使用它、理解它、完善它。

然而，分布式技术涉及的方面（存储、计算、框架、中间件等）是如此之多，且迄今为止尚未见到一本书对其进行概括和梳理，要想对分布式技术有全面的了解，特别是对初学者而言，何其难哉！

因此，本书试图对近年来涌现出的各种主流分布式技术做一个简要介绍，以使不太熟悉这个领域的读者能了解其概貌、原理和根源。

本书共分为以下 6 部分。

- 第一部分对典型的分布式系统的组成及其中每个组件的功能进行简要介绍，以使读者对分布式系统有一个总体了解。
- 第二部分介绍分布式系统的前端经常使用的 Web 框架、反向代理及负载均衡技术。
- 第三部分对分布式系统中经常使用的各种中间件技术逐一进行介绍，包括分布式同步服务中间件、关系型数据库访问中间件、分布式服务调用中间件、分布式消息服务中间件和分布式跟踪服务中间件。
- 第四部分介绍分布式文件系统、各种 NoSQL 数据库技术（基于键值对的 NoSQL 技术、基于列的 NoSQL 技术、基于文档的 NoSQL 技术、基于图的 NoSQL 技术）和 NewSQL 数据库系统。
- 第五部分对业界在构建大型分布式系统中的主要经验加以总结，使后来者避免重蹈覆辙。
- 第六部分介绍业界几个知名的大型分布式系统的主要设计思想和架构，包括谷歌搜索系统、淘宝网、阿里云和领英的社交应用。此外，还会探讨和思考分布式系统实现中的一些问题。

在本书写作的过程中，参考了许多网上的资料和书籍，感谢这些资料的作者们。由于涉及的资料太多，无法将其作者一一列出，谨表歉意。

感谢人民邮电出版社的杨海玲编辑，没有你的辛勤付出，也不会有本书的顺利出版。

另外，还要感谢我的好朋友申天雷，感谢你在繁忙的工作之余抽出时间来全面审阅这本书。

最后，感谢我的家人。没有父母的养育，我不可能有机会掌握这些知识；没有我岳父母帮我看管孩子，没有我爱人的关心与支持，我不可能有时间来写作；没有我儿子那双充满期待的眼睛，我也不可能一直坚持把本书写完！

由于分布式系统涉及的内容实在太多，我不可能对其中每个细节都有很深刻和精确的理解，因此，书中错误及疏漏之处在所难免，欢迎批评指正。

李庆旭

2019年夏于北京

# 资源与支持

本书由异步社区出品，社区 (<https://www.epubit.com/>) 为您提供相关资源和后续服务。

## 配套资源

本书提供源代码下载，要获得相关配套资源，请在异步社区本书页面中单击 **配套资源**，跳转到下载界面，按提示进行操作即可。注意：为保证购书读者的权益，该操作会给出相关提示，要求输入提取码进行验证。

## 提交勘误

作者和编辑尽最大努力来确保书中内容的准确性，但难免会存在疏漏。欢迎您将发现的问题反馈给我们，帮助我们提升图书的质量。

当您发现错误时，请登录异步社区，按书名搜索，进入本书页面，单击“提交勘误”，输入勘误信息，单击“提交”按钮即可。本书的作者和编辑会对您提交的勘误进行审核，确认并接受后，您将获赠异步社区的 100 积分。积分可用于在异步社区兑换优惠券、样书或奖品。

详细设置 写书评 提交勘误

页码:  页内位置 (行号):  勘误内容:

勘误信息

字数统计

提交

## 扫码关注本书

扫描下方二维码，您将会在异步社区微信服务号中看到本书信息及相关的服务提示。



## 与我们联系

我们的联系邮箱是 [contact@epubit.com.cn](mailto:contact@epubit.com.cn)。

如果您对本书有任何疑问或建议，请您发邮件给我们，并在邮件标题中注明本书书名，以便我们更高效地做出反馈。

如果您有兴趣出版图书、录制教学视频，或者参与图书翻译、技术审校等工作，可以发邮件给我们；有意出版图书的作者也可以到异步社区在线提交投稿（直接访问 [www.epubit.com/selfpublish/submission](http://www.epubit.com/selfpublish/submission) 即可）。

如果您来自学校、培训机构或企业，想批量购买本书或异步社区出版的其他图书，也可以发邮件给我们。

如果您在网上发现有针对异步社区出品图书的各种形式的盗版行为，包括对图书全部或部分内容的非授权传播，请您将怀疑有侵权行为的链接发邮件给我们。您的这一举动是对作者权益的保护，也是我们持续为您提供有价值的内容的动力之源。

## 关于异步社区和异步图书

“异步社区”是人民邮电出版社旗下 IT 专业图书社区，致力于出版精品 IT 技术图书和相关学习产品，为译者提供优质出版服务。异步社区创办于 2015 年 8 月，提供大量精品 IT 技术图书和电子书，以及高品质技术文章和视频课程。更多详情请访问异步社区官网 <https://www.epubit.com>。

“异步图书”是由异步社区编辑团队策划出版的精品 IT 专业图书的品牌，依托于人民邮电出版社近 30 年的计算机图书出版积累和专业编辑团队，相关图书在封面上印有异步图书的 LOGO。异步图书的出版领域包括软件开发、大数据、AI、测试、前端、网络技术 etc。



异步社区



微信服务号

# 目录

## 第一部分 分布式系统概述

第 1 章 分布式系统概述	3	1.4.1 批处理分布式计算系统	9
1.1 分布式系统的组成	4	1.4.2 流处理分布式计算系统	10
1.2 分布式协调组件	6	1.4.3 混合系统	11
1.3 分布式存储系统	7	1.5 分布式系统中节点之间的关系	11
1.4 分布式计算系统	9		

## 第二部分 分布式系统的前端构造技术

第 2 章 Web 框架的实现原理	15	2.4 Go 语言 Web 开发	30
2.1 Web 框架简介	16	2.4.1 Go 语言简介	31
2.2 PHP Web MVC 框架的工作原理	17	2.4.2 Go 语言 Web 开发	32
2.2.1 框架的入口	17	第 3 章 反向代理与负载均衡	33
2.2.2 URL 到控制器的映射	17	3.1 反向代理	33
2.2.3 如何将模型传给视图	18	3.1.1 Nginx	34
2.3 Java Web MVC 框架原理	19	3.1.2 Tengine	35
2.3.1 Java Servlet API 3.0	19	3.1.3 Varnish	35
2.3.2 框架的入口	20	3.2 负载均衡	36
2.3.3 Spring 4.0 框架	22	3.2.1 DNS 负载均衡	36
2.3.4 Java EE	29	3.2.2 硬件负载均衡	36
2.3.5 Spring 与 Java EE 的比较	30	3.2.3 软件负载均衡	37

## 第三部分 分布式中间件

第 4 章 分布式同步服务中间件	41	4.3.1 架构	44
4.1 分布式一致性协议	42	4.3.2 如何消除单点故障	45
4.2 分布式同步服务中间件简介	43	4.3.3 Chubby 对外提供的 API	45
4.3 分布式同步服务中间件的实现原理	44	4.3.4 数据库	46

4.3.5 Chubby 使用示例: 集群的主服务器选举	46	6.3 其他分布式服务调用中间件	67
4.4 其他分布式同步服务中间件	46	6.3.1 Protocol Buffers	67
4.4.1 Linux 心跳机制	47	6.3.2 gRPC	68
4.4.2 ZooKeeper	47	6.3.3 Thrift	69
4.4.3 iNexus	48	6.3.4 Motan	72
4.5 分布式同步服务的应用	48	6.3.5 sofa-pbrpc	73
<b>第 5 章 关系型数据库访问中间件</b>	<b>53</b>	6.4 分布式服务调用中间件的应用	73
5.1 数据库访问中间件的形式	54	<b>第 7 章 分布式消息服务中间件</b>	<b>75</b>
5.2 数据库访问中间件的工作原理	55	7.1 分布式消息服务中间件简介	75
5.3 著名的数据库访问中间件	56	7.2 分布式消息服务中间件的实现原理	77
5.3.1 MySQL 代理	57	7.2.1 消息模型	77
5.3.2 Cobar	58	7.2.2 架构	77
5.3.3 TDDL	59	7.3 其他分布式消息服务中间件	79
5.3.4 MyCAT	59	7.3.1 阿里巴巴 RocketMQ	79
5.3.5 Heisenberg	59	7.3.2 Apache Pulsar	80
5.4 数据库访问中间件的应用	60	7.4 分布式消息服务中间件的应用	82
5.4.1 使用 MySQL 代理实现读写数据库分离	60	7.4.1 秒杀系统中使用 Kafka 以削平峰值流量	83
5.4.2 研发自己的数据库中间件, 实现 MySQL 的分库分表	60	7.4.2 使用 Kafka 流实现消息推送	83
<b>第 6 章 分布式服务调用中间件</b>	<b>63</b>	<b>第 8 章 分布式跟踪服务中间件</b>	<b>85</b>
6.1 分布式服务调用中间件简介	63	8.1 分布式跟踪服务中间件的实现原理	85
6.2 分布式服务调用中间件的实现原理	64	8.2 其他分布式跟踪服务中间件	88
6.2.1 Dubbo 的架构	64	8.2.1 Twitter 的 Zipkin	88
6.2.2 Dubbo 中各组件的交互	65	8.2.2 Pinpoint	90
6.2.3 Dubbo 的实现及特点	66	8.2.3 阿里巴巴的 EagleEye	92
6.2.4 Dubbox	66	8.3 分布式跟踪服务中间件的应用	92

## 第四部分 分布式存储技术

<b>第 9 章 分布式文件系统</b>	<b>95</b>	<b>第 10 章 基于键值对的 NoSQL 数据库</b>	<b>107</b>
9.1 分布式文件系统的实现原理	96	10.1 NoSQL 数据库的 CAP 权衡	108
9.2 其他分布式文件系统	102	10.2 基于键值对的 NoSQL 数据库的实现	
9.3 分布式文件系统的应用	104		

原理	108	12.1.1 数据模型	131
10.2.1 谷歌的 LevelDB	108	12.1.2 自动分片	132
10.2.2 阿里巴巴的 Tair	111	12.1.3 副本	132
10.2.3 亚马逊的 Dynamo	111	12.1.4 索引	133
10.3 其他基于键值对的 NoSQL 数据库	115	12.1.5 查询路由	133
10.3.1 Memcached	115	12.2 其他基于文档的 NoSQL 数据库	133
10.3.2 Redis	116	12.2.1 CouchDB	133
10.3.3 Berkeley DB	117	12.2.2 RethinkDB	134
10.3.4 Facebook RocksDB	117	12.3 基于文档的 NoSQL 数据库的应用	135
10.3.5 Riak	118	<b>第 13 章 其他 NoSQL 数据库</b>	137
10.3.6 Voldemort	118	13.1 基于图的 NoSQL 数据库 Neo4j	138
10.4 基于键值对的 NoSQL 数据库的应用	118	13.1.1 数据模型	138
10.4.1 使用 Redis 缓存会话数据	118	13.1.2 图的存储	138
10.4.2 使用 Berkeley DB/LevelDB/RocksDB 构建自己的分布式存储系统	119	13.1.3 高可用性	139
10.4.3 使用 Berkeley DB/LevelDB/RocksDB 作为本地数据库	119	13.1.4 水平扩展	139
<b>第 11 章 基于列的 NoSQL 数据库</b>	121	13.2 多数据模型 NoSQL 数据库 OrientDB	139
11.1 基于列的 NoSQL 数据库的实现	121	13.2.1 基本概念	140
原理	121	13.2.2 图的表示	140
11.1.1 数据模型	121	13.2.3 节点与集群	140
11.1.2 架构	124	13.2.4 分片	141
11.2 其他基于列的 NoSQL 数据库	126	13.2.5 ACID 支持	141
11.2.1 Apache HBase	126	13.2.6 CAP 的权衡	142
11.2.2 Apache Cassandra	127	13.2.7 集群配置信息的管理	142
11.2.3 Baidu Tera	128	13.3 时间序列 NoSQL 数据库	142
11.3 基于列的 NoSQL 数据库的应用	128	<b>第 14 章 NewSQL 数据库</b>	143
11.3.1 HBase 用于数据分析系统	128	14.1 NewSQL 和 CAP 理论	144
11.3.2 HBase 用于存储呼叫记录	129	14.2 采用新架构的 NewSQL 系统	145
<b>第 12 章 基于文档的 NoSQL 数据库</b>	131	14.2.1 谷歌的 Megastore	145
12.1 基于文档的 NoSQL 数据库的实现	131	14.2.2 谷歌的 Spanner	146
原理	131	14.2.3 谷歌的 F1	147
		14.2.4 阿里巴巴的 OceanBase	148
		14.2.5 其他采用新架构的 NewSQL 数据库	152

## 第五部分 分布式系统的构建思想

第 15 章 云化 .....	157	16.4 实用主义的思想 .....	169
15.1 云化的技术基础 .....	157	16.4.1 搜索引擎作为查询工具 使用 .....	169
15.1.1 虚拟机技术 .....	157	16.4.2 阿里巴巴的 OceanBase 的 架构 .....	169
15.1.2 容器技术 .....	159	16.4.3 根据需要选择最适合的开发 工具和开发语言 .....	170
15.2 公有云能提供什么 .....	159	16.4.4 根据需要选择不同的存储 系统 .....	170
15.3 云化对软件架构的要求 .....	161	16.5 异步化以解耦并削平峰值 .....	170
第 16 章 分布式系统的构建思想 .....	163	16.6 最终一致性的思想 .....	171
16.1 一切都可能失败与冗余的思想 .....	163	16.7 微服务的思想 .....	171
16.1.1 如何避免单点故障 .....	164	16.8 MapReduce 的思想 .....	172
16.1.2 避免单点故障的具体做法 .....	165	16.9 服务跟踪的思想 .....	172
16.2 水平而不是垂直扩展的思想 .....	165	16.10 资源池化的思想 .....	173
16.2.1 数据的水平扩展 .....	166		
16.2.2 服务的水平扩展 .....	167		
16.2.3 数据中心的水平扩展 .....	167		
16.3 尽可能简单的思想 .....	168		

## 第六部分 大型分布式系统案例研究及分析

第 17 章 大型分布式系统案例研究 .....	177	18.5.2 Java 语言的插件模式 实现 .....	203
17.1 案例研究之谷歌搜索系统 .....	177	18.5.3 采用专用语言的插件模式 实现 .....	205
17.2 案例研究之淘宝网 .....	182	18.6 关于分布式服务调用中间件的 实现 .....	205
17.3 案例研究之阿里云 .....	185	18.7 动态链接还是静态链接 .....	206
17.4 案例研究之领英 .....	189	18.8 无所不用其极的压榨性能手段 .....	206
第 18 章 关于分布式系统设计的思考 .....	197	18.8.1 编译后代码的原生态化 .....	206
18.1 大型互联网公司架构的共性 .....	197	18.8.2 定制的 Linux 内核 .....	207
18.2 为何大型互联网公司的架构如此 相似 .....	198	18.8.3 定制的 Java 虚拟机 .....	208
18.3 关于分布式监控系统 .....	199	18.8.4 定制的 MySQL .....	208
18.4 Linux 系统调用 epoll() .....	200	参考文献 .....	209
18.5 关于插件设计模式的实现 .....	201	后记 .....	211
18.5.1 C/C++ 语言的动态库形式的 实现 .....	202		

# 第一部分

第 1 章

## 分布式系统概述

之所以需要有分布式系统，最根本的原因还是单机的计算和存储能力不能满足系统的需要，但要把成百上千台计算机组织成一个有机的系统绝非易事。

这一部分会对典型的分布式系统的组成及其每个组件的功能做简要介绍，以便读者对分布式系统有一个总体的了解。

## 1. 2 的 特 殊 形 式

### 1. 2. 1 的 特 殊 形 式

1. 2. 1. 1 的 特 殊 形 式

1. 2. 1. 1. 1

1. 2. 1. 1. 2 的 特 殊 形 式

1. 2. 1. 1. 2. 1

1. 2. 1. 1. 2. 2 的 特 殊 形 式

1. 2. 1. 1. 2. 2. 1

1. 2. 1. 1. 2. 2. 2 的 特 殊 形 式

1. 2. 1. 1. 2. 2. 2. 1

1. 2. 1. 1. 2. 2. 2. 2

1. 2. 1. 1. 2. 2. 2. 2. 1

1. 2. 1. 1. 2. 2. 2. 2. 2

1. 2. 1. 1. 2. 2. 2. 2. 2. 1

1. 2. 1. 1. 2. 2. 2. 2. 2. 2

1. 2. 1. 1. 2. 2. 2. 2. 2. 2. 1

1. 2. 1. 1. 2. 2. 2. 2. 2. 2. 2

## 1. 2. 2 的 特 殊 形 式

### 1. 2. 2. 1 的 特 殊 形 式

1. 2. 2. 1. 1 的 特 殊 形 式

1. 2. 2. 1. 2 的 特 殊 形 式

1. 2. 2. 1. 3 的 特 殊 形 式

1. 2. 2. 1. 4 的 特 殊 形 式

### 1. 2. 2. 2 的 特 殊 形 式

1. 2. 2. 2. 1 的 特 殊 形 式

1. 2. 2. 2. 2 的 特 殊 形 式

1. 2. 2. 2. 3 的 特 殊 形 式

1. 2. 2. 2. 4 的 特 殊 形 式

1. 2. 2. 2. 5 的 特 殊 形 式

1. 2. 2. 2. 6 的 特 殊 形 式

1. 2. 2. 2. 7 的 特 殊 形 式

1. 2. 2. 2. 8 的 特 殊 形 式

1. 2. 2. 2. 9 的 特 殊 形 式

1. 2. 2. 2. 10 的 特 殊 形 式

1. 2. 2. 2. 11 的 特 殊 形 式

1. 2. 2. 2. 12 的 特 殊 形 式

# 第 1 章

## 分布式系统概述

1999 年 8 月 6 日，CNN 报道了一起 eBay 网站的事故：从 7:30 开始，整个网站崩溃，一直持续了 9 个多小时。下午 5:30 后，技术人员开始进行系统恢复，但搜索功能依然不能使用。

2011 年 4 月 21 日至 22 日，亚马逊 EC2 (Elastic Computer Cloud) 服务出现大面积事故，导致数以千计的初创公司受到影响，而且造成大约 11 小时的历史数据永久性丢失。

2013 年 4 月 27 日，《大掌门》游戏的开发商玩蟹科技 CEO 叶凯在微博上吐槽，“我们在阿里云上用了 20 多台机器。半年时间，出现过 1 次所有机器全部断电，2 次多个硬盘突然只读，3 次硬盘 I/O 突然变满……”。

2013 年 12 月 28 日，春运第一天，铁道部首次推出了网上订票系统，但很快就出现许多用户无法访问、响应缓慢甚至串号等事故。

2017 年 9 月 17 日，谷歌的网盘服务 Drive 出现故障，成千上万用户受到影响。

上面的这几起事故，当时都闹得沸沸扬扬，不仅给受影响的用户带来了很大的损失，也极大地影响了厂商的形象。事实上，几乎每一家互联网公司的后台系统都曾经不止一次地经历过这样或那样的尴尬时刻。可以这样说，几乎每一家互联网公司的后台架构都是在发现问题、解决问题的循环中发展起来的。

即便是执分布式系统技术牛耳的谷歌，在 2017 年 9 月，也出现过分布式系统的故障。可见，开发并维护一个成功的分布式系统是多么不易！

最早得到广泛应用的分布式系统是诞生于 20 世纪 70 年代的以太网。尽管分布式系统存在的历史已经有近半个世纪，然而其大规模的发展和应用则是 2000 年以后的事情。

21 世纪以来，随着雅虎、谷歌、亚马逊、eBay、Facebook、Twitter 等众多互联网公司的崛起，其用户量以及要处理的数据量迅速增长，远远超过了传统的计算机系统能够处理的范围，因此，以谷歌为代表的互联网公司提出了许多新技术（如 HDFS、Bigtable、

MapReduce 等)。以 BAT 为代表的中国互联网公司，也在 21 世纪整体崛起，在初期借鉴美国公司技术的基础上，他们也自行开发了许多新的技术（如淘宝的管理海量小文件的分布式存储系统 TFS、阿里巴巴开源的分布式调用框架 Dubbo、阿里巴巴开源的数据库中间件 Cobar 等）。

为了解决分布式系统中的各种各样的问题，各大互联网公司开发了各种各样的技术，当然，这也促进了当今分布式系统技术领域的飞速发展。为了存储大量的网站索引，谷歌设计了 GFS 分布式文件存储系统和基于列存储的 Bigtable NoSQL 数据库系统；为了计算 PageRank 算法中的页面 rank 值，谷歌又设计了 MapReduce 分布式计算系统；为了方便其分布式系统中不同主机间的协调，谷歌还设计了 Chubby 分布式锁系统；为了解决不同语言实现的组件间的通信问题，Facebook 设计了 Thrift；为了解决大量消息的快速传递问题，领英设计了 Kafka……这个列表可以很长很长。

为了“压榨”分布式系统中每个组件的性能，人们已经不再仅仅满足于在程序库（如网络编程库 Netty、内存管理库 TCMalloc 等）、程序框架（如 Spring）等“略显浅薄”的地方提高，而是已经渗透到了硬件（如谷歌为其计算中心专门设计了计算机）、网络（如 SDN）、操作系统（如各大互联网公司定制的 Linux 内核）、语言（如谷歌设计的 Go 语言）、数据库系统（如各种 NoSQL 系统）、算法（如人工智能领域的突飞猛进）等各种计算机基础领域。

毫无疑问，我们处于计算机技术发展最为迅猛的时代。在这个如火如荼的时代里，许多尘封多年的计算机技术（如人工智能、分布式系统、移动计算、虚拟计算等），一改往日不温不火的模样，在互联网这片广袤的土地上如日中天，发展迅速。

今天的计算机领域，已经与 20 年前大为不同。20 年前，只需要对操作系统、数据库、网络、编译等领域有深刻的理解，再熟练掌握几门计算机语言，了解一些常见的软件架构（客户服务器架构、管道架构、分层架构等）和软件工程（主要是瀑布模型）的知识，基本上就能胜任大多数软件开发工作了。而今天，仅了解这些基础知识已经远远不够，因为在近 20 年内，人类创造了太多的新技术，而这些新技术又大都起源并服务于分布式计算领域。

## 1.1 分布式系统的组成

一个大型的分布式系统虽然非常复杂，但其设计目标却往往是非常简单的，例如，京东和淘宝这样的电商，其设计目标是卖东西；谷歌和百度这样的搜索引擎，其设计目标是帮助大家在网上找相关的内容；Facebook 和微信这样的社交应用，其设计目标是方便大家相互联系并分享自己生活中的点点滴滴。

如前文所述，之所以需要有分布式系统，最根本的原因还是单机的计算和存储能力不能满足系统的需要。但要把成百上千台计算机组织成一个有机的系统，绝非易事。在人类

社会中，其实也一样，找到 1000 个人容易，但要把这 1000 个人组织成一只能战斗的军队可就没那么简单了。

一个典型的分布式系统如图 1-1 所示。

- 分布式系统大都有一个 Web 前端，用户可以通过浏览器随时随地访问，当然，前端也可以是运行在 Windows/Linux 上的桌面程序或者运行在手机上的应用。
- 分布式系统还要有后端支撑。分布式系统的后端大都是基于 Linux 的集群<sup>①</sup>。之所以采用 Linux，一是因为开源操作系统成本低，二是因为开源软件可以定制。
- 就像人类社会需要有一定的组织和管理一样，为了组成一个集群，在单机的操作系统之上，还需要集群管理系统。在集群管理系统中，一个非常重要的组件是分布式协调组件，用来协调不同机器之间的工作。这些协调系统大都基于一些著名的分布式一致性协议（如 Paxos、Raft 等）。有些超大型的后端还拥有专门的集群操作系统，这些系统不仅有分布式协调功能，还有资源的分配与管理功能。
- 为了满足大规模数据的存储需要<sup>②</sup>，需要有能够存储海量数据的后端存储系统。
- 为了满足大规模数据的计算需要<sup>③</sup>，还需要有能够分析海量数据的后端计算系统。

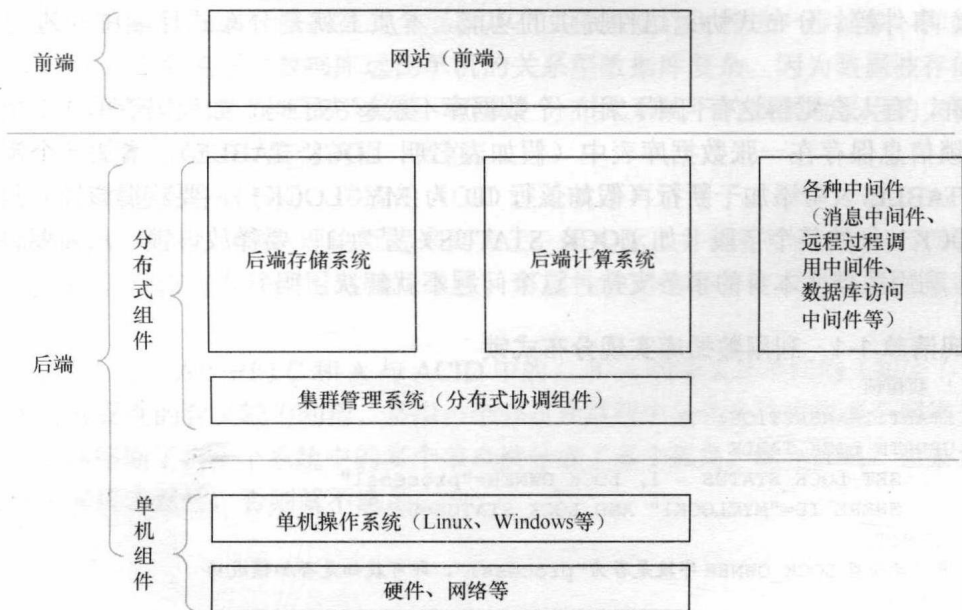


图 1-1 一个典型的分布式系统

① 所谓集群，是指采用同样或类似配置的许多台机器，为了达到一个共同的目的而组成的系统。

② 像谷歌和百度这样的公司，因为索引的页面量非常庞大，需要很大的存储空间。微信和 Facebook 这样的社交应用亦然。

③ 例如谷歌，其 PageRank 算法就需要很大的计算量；再如京东和淘宝这样的电商，其商品推荐系统也需要很大的计算量。