

NINE ALGORITHMS
THAT CHANGED
THE FUTURE



改变未来的九大算法

[美] 约翰·麦考密克◎著

John MacCormick

管策◎译



指尖上的精灵

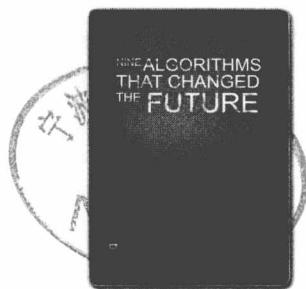
NLIC2970846616



中信出版社·CHINA CITIC PRESS

改变未来的九大算法

[美] 约翰·麦考密克 (John MacCormick) 著 管策 译



NLIC2970846616

图书在版编目 (CIP) 数据

改变未来的九大算法 / (美) 麦考密克著；管策译。—北京：中信出版社，2013.6

书名原文：Nine Algorithms That Changed The Future

ISBN 978-7-5086-3901-7

I. ①改… II. ①麦… ②管… III. ①数字技术－影响－经济发展－研究 IV. ①F201

中国版本图书馆 CIP 数据核字 (2013) 第 059002 号

Copyright © 2012 by Princeton University Press

Simplified Chinese edition © 2013 by China CITIC Press

All rights reserved.

No part of this book may be reproduced or transmitted in any
form or by any means, electronic or mechanical, including photocopying, recording
or by any information storage and retrieval system, without permission in writing
from the Publisher.

本书仅限中国大陆地区发行销售

改变未来的九大算法

著 者：[美] 约翰 · 麦考密克

译 者：管 策

策划推广：中信出版社（China CITIC Press）

出版发行：中信出版集团股份有限公司

（北京市朝阳区惠新东街甲 4 号富盛大厦 2 座 邮编 100029）

（CITIC Publishing Group）

承 印 者：三河市西华印务有限公司

开 本：787mm×1092mm 1/16

印 张：16 字 数：200 千字

版 次：2013 年 6 月第 1 版

印 次：2013 年 6 月第 1 次印刷

京权图字：01-2012-1063

广告经营许可证：京朝工商广字第 8087 号

书 号：ISBN 978-7-5086-3901-7 / F · 2873

定 价：39.00 元

版权所有·侵权必究

凡购本社图书，如有缺页、倒页、脱页，由发行公司负责退换。

服务热线：010-84849555 服务传真：010-84849000

投稿邮箱：author@citicpub.com

Nine Algorithms
That Changed the Future
目录

- v 序
- 001 第一章
前言
计算机日常运用的卓越思想有哪些
- 013 第二章
搜索引擎索引
在世界上最大的草垛中寻针
- 031 第三章
PageRank
让谷歌腾飞的技术
- 049 第四章
公钥加密
用明信片传输秘密
- 075 第五章
纠错码
自纠正的错误

第六章

图形识别

099 从经验中学习

第七章

数据压缩

127 有益无害

第八章

数据库

147 追求一致性的征程

第九章

数字签名

177 这个软件究竟由谁编写

第十章

205 **什么可以计算**

有些程序不可能存在

第十一章

结论

235 更多在你指尖的精灵

243 致 谢

245 注 释

第一章
前言
——计算机日常运用的卓越思想有哪些

ALGORITHMS
THAT CHANGED THE FUTURE

此乃小技……为诗之诀在于有气、有势、有情、
有韵、有起、有承、有转、有合。①

——威廉姆·莎士比亚，
《爱的徒劳》(*Love's Labour's Lost*)

① 选自朱生豪译本。——译者注

计算机科学中的伟大思想是如何诞生的？以下遴选部分思想进行介绍：

- 20世纪30年代，在第一台数字计算机发明以前，一位英国天才开创了计算机科学领域。之后，这位天才继续证明，不管未来建造的计算机运行多快、功能多强大、设计得多好，仍旧有一些问题是计算机不能解决的。
- 1948年，一位供职于电话公司的科学家发表了一篇论文，开创了信息理论领域。这位科学家的工作让计算机能以完美的精确度传输信息，即便大部分数据都因为被干扰而破坏。
- 1956年，一群学者在达特茅斯举行会议。这次会议的目标很清晰，也很大胆，那就是开创人工智能领域。在取得了许多重大成功以及经历了无数失望之后，我们仍期待出现一个真正的智能计算机程序。
- 1969年，IBM公司的一名研究人员发明了一种将信息组织进数据库中的先进方法。目前，绝大多数在线交易都使用该技术存储及检索信息。

• 1974 年，英国政府秘密通信实验室的研究人员发明了一种让计算机安全通信的方法，即另一台计算机可以查看在计算机之间交换的所有信息。这些研究人员为政府保密所限——不过幸运的是，三名美国专家独立开发并拓展了这项重大发明，为互联网上所有的安全通信打下了基础。

• 1996 年，两名斯坦福大学博士生决定联手搭建一个互联网搜索引擎。几年后，他们创办了谷歌公司——互联网时代的第一个数字巨头。

在我们享受 21 世纪技术惊人增长的同时，使用计算机设备——不管是现有最强大一组机器还是最新、最时尚的手持设备——都不可避免地要依赖计算机科学的基础思想，而这些思想都诞生于 20 世纪。想一想：你今天做过什么令人印象深刻的事情吗？好吧，这个问题的答案取决于你怎么看。也许你搜索了包含数十亿份文档的资料库，从中选出两到三份与你的需求最相关的文档？即便有能够影响所有电子设备的电磁干扰，存储或传输了数百万条信息，也没犯一点错误？你是否成功地完成了一次在线交易，即便同时有成千上万名消费者在访问同一个服务器？你是否在能够被其他数十台计算机嗅探到的线路中传输了一些机密信息（比如信用卡卡号）？你是否运用过压缩的魔力，将数兆的照片压缩成更易于管理的大小，以便在电子邮件中发送？你是否在手持设备上触发了人工智能，自动纠正你在手持设备的小巧键盘上输入的内容？

这些令人印象深刻的壮举都依赖于之前提到的伟大发现。然而，绝大多数计算机用户每天都会多次运用这些独创想法，却从没有意识到！本书旨在面向大众解释这些概念——我们每天使用的计算机科学的伟大思想。在解释每个概念时，我都假设读者不了解任何计算机科学的任何知识。

算法：指尖精灵的构件

到目前为止，我一直在谈计算机科学的伟大“思想”，但计算机科学家们将许多重要思想形容为“算法”。那么思想和算法之间有什么区别呢？究竟什么是算法？这一问题最简单的答案是，算法是一张精确的处方，按顺序详细列出了解决一个问题所需的具体步骤。我们小时候在学校学到的一种算法就是很好的例子：将两个大数字相加的算法。如下例所示。这个算法涉及一连串步骤，开始的步骤如下：“首先，将两个数的最末位数相加，写下结果的最末位数，将剩下的数放到左侧的下一栏；接着，将下一栏的数相加，再将除结果末位数之外的数字和前一栏余下的数相加……”。依此类推。

$$\begin{array}{r}
 4844978 \\
 +3745945 \\
 \hline
 \end{array}
 \quad
 \begin{array}{r}
 4844978 \\
 +3745945 \\
 \hline
 3
 \end{array}
 \quad
 \begin{array}{r}
 4844978 \\
 +3745945 \\
 \hline
 23
 \end{array}$$

将两个数字相加的算法的前两步。

请注意算法步骤近乎机械化的感觉。事实上，这是算法的关键特点之一：每一步都必须绝对精确，没有任何人类意图或推测掺杂其中。这样，每一个完全机械化的步骤才能被编入计算机。算法的另一个重要特点是，不管输入什么，算法总能运行。我们在学校学到的相加算法就拥有这一特性：不管你想把哪两个数相加，算法最终都会得出正确答案。比如，用这一算法将两个长达 1 000 位的数相加，你肯定能得到答案，尽管这需要相当长的时间。

对于把算法定义为一张精确、机械化的处方的说法，你也许会略感好奇。这张处方究竟要有多精确？要进行哪些基本操作？比如，在上面的相加算法中，简单地说一句“把两个数相加”是不是就行了？还是说我

们要在加法表上列出所有个位数字？这些细节看起来也许有点乏味，甚至会显得有点学究气，但其实离真相不远了：这些问题的真正答案正处于计算机科学的核心，并且也和哲学、物理学、神经科学以及遗传学有联系。有关算法究竟是什么的深层问题都归结于一个前提——也就是众所周知的邱奇—图灵论题（Church–Turing Thesis）。我们将在第十章重温这些问题，届时我们还将讨论计算的理论极限，以及邱奇—图灵论题的一些方面。同时，将算法比作一张非常精确的处方这一非正式概念效果会非常好。

现在我们知道算法是什么，但算法和计算机有什么联系呢？关键在于，计算机需要用非常精确的指令编程。因此，在能让计算机为我们解决某个特定问题之前，我们需要为那个问题开发一个算法。在数学和物理学等其他学科中，重要的结果通常是由一个方程式获得的。（著名的例子包括勾股定理 $a^2+b^2=c^2$ ，或爱因斯坦的质量守恒定理 $E=mc^2$ 。）相反，计算机科学的伟大思想通常是形容如何解决一个问题——当然，是使用一种算法。因此，本书的主要目的是，解释让计算机成为你的个人精灵的东西——计算机每天使用的伟大算法。

一个伟大的算法由什么构成？

这会引出一个刁钻的问题：什么才是真正伟大的“算法”？潜在的候选算法清单相当大，但我用几条基本标准缩减了用于本书的候选算法清单。第一条，也是最重要的一条标准是，伟大的算法要被普通计算机用户每天用到。第二条重要的标准是，伟大的算法应该能处理具体的现实问题，如压缩一个特定文件或通过一个噪声链接精确地传输文件。对于已经了解部分计算机科学的读者而言，下面的文字框解释前面两大标准的部分后果。

第一条标准——要被普通计算机用户每天用到——排除了主要由计算机专业人士使用的算法，如编译器和程序验证技术。第二条标准——针对某个特定问题的具体程序——排除了许多作为计算机科学本科课程核心内容的伟大算法，如排序算法（快速排序）、图形算法（迪杰斯特拉最短路径算法）、数据结构（哈希表）。这些算法的伟大性毋庸置疑，而且很轻易地就满足了第一条标准，因为普通用户使用的绝大多数应用程序都会反复应用这些算法。但这些算法太通用了：它们能用于解决众多问题。在本书中，我决定要专注于解决特定问题的算法，因为对于普通计算机用户而言，这些算法能让他们拥有更清晰的动机。

一些和本书选取算法有关的额外细节。本书读者无须具备计算机科学的任何知识。但如果读者具备计算机科学背景知识，这个文字框会解释为何这类读者之前偏好的许多内容没有出现在本书中。

第三个标准是，算法主要和计算机科学理论相关。这排除了主要和计算机硬件——如CPU、监视器以及网络——有关的技术。这条标准也减轻了对基础设施——如互联网——设计的重视。为什么我要着重于计算机科学理论？部分原因是由于公众对计算机科学认知的不平衡：有一种广泛的观点认为，计算机科学基本上就是编程（如“软件”）和设备设计（如“硬件”）。事实上，最优美的计算机科学思想中有许多是十分抽象的，并不属于以上任意一类。我希望通过着重于这些理论思想，让更多人将计算机科学的本质作为一门知识学科来理解。

你也许已经注意到了，我列出的标准可能会遗漏一些伟大的算法，但却从一开始就避免了定义伟大这个极其麻烦的问题。针对这一问题，我依赖于自己的直觉。在本书中说明的每一个算法中，其核心都是一个让整件事情奏效的精巧把戏。对我而言，当这个把戏显露出来时，那个“惊叹”

时刻，会让解释这些算法成为令人愉悦的经历，我希望你也能有此感受。因为我会用到“把戏”(trick)这个词很多次，需要指出的是，我并非指那些卑劣或骗人的把戏——那种孩子可能会用在弟弟或妹妹身上的把戏。相反，本书中的把戏类似于交易诀窍或魔术：为达成目标而采用的聪明技巧，否则目标很难或不可能达成。

因此，根据直觉，我选出了自认为是计算机科学世界中最精巧、最神奇的把戏。在英国数学家高德菲·哈罗德·哈代(G. H. Hardy)的《一个数学家的辩白》(*A Mathematician's Apology*)中，作者试图向公众解释数学家从事数学的原因：“美是第一道测试：丑陋的数学在这个世界中无永存之地。”这道美的测试也适用于计算机科学中蕴含的理论思想。因此，选取在本书中出现的算法的最后一条标准，就是哈代的——也许可以这么称呼——美的测试：希望我至少能成功地向读者展示部分美——我在每个算法中感觉到的美。

接下来谈谈我选择展示的这些算法。搜索引擎的巨大影响，也许是算法技术影响所有计算机用户最明显的例子，我自然也将部分互联网搜索的核心算法收入了本书中。第二章描述了搜索引擎如何使用索引寻找与请求的文件，而第三章则解释了网页排名(PageRank)算法——谷歌公司为保证匹配度最高的文件出现在搜索结果列表顶部的原始算法。

即便我们不经常想这件事情，绝大多数人也能意识得到，为提供出人意料的强大搜索结果，搜索引擎使用着一些深邃的计算机科学思想。相反，其他一些伟大的算法也经常被用到，但计算机用户对此甚至都没有意识到。第四章描述的公钥加密(public key cryptography)就是这样一种算法。用户每次访问一个安全网站(地址以https而非http开头)，用户都会用到公钥加密的一个方面——也就是众所周知的密钥交换(key exchange)——来展开一段安全对话。第四章解释了密钥交换过程的实现原理。

第五章的主题是纠错码(error correcting codes)，这是我们经常使用

但却没有意识到的另一类算法。事实上，纠错码极有可能是有史以来唯一一个使用次数最频繁的伟大算法。纠错码可以让计算机识别并纠正正在储存或传输数据中出现的错误，而不必依靠备份或再次传输。纠错码无处不在：它们被用于所有硬盘驱动器、众多网络传输、CD和DVD，甚至还存在于一些计算机的内存。不过，纠错码的能力太强了，以至于我们意识不到它们存在。

第六章稍微有点特殊，介绍了图形识别算法（pattern recognition algorithm）。图形识别算法也能进入伟大的计算机科学思想榜单，但却违背了第一条标准：要被普通计算机用户每天用到。图形识别属于计算机识别高度可变信息——如笔迹、讲话和人脸——的技术。事实上，在21世纪的第一个十年，绝大多数日常计算并没有用到这些技术。但在2011年，图形识别的重要性急剧增大：配备小型屏幕键盘的移动设备需要自动纠错，平板设备必须识别手写输入，而且所有这些设备（特别是智能手机）越来越趋向于语音操作。一些网站甚至使用图形识别来决定向用户展示哪种广告。另外，我对图形识别也有偏好，因为它是我的研究领域。因此，第六章描述了3种最有趣、最成功的图形识别技术：最近邻分类器（nearest-neighbor classifier）、决策树（decision tree）以及神经网络（neural network）。

第七章讨论了压缩算法。压缩算法组成了另一组使计算机变成我们指尖精灵的伟大思想。计算机用户的确会时不时地直接进行压缩，也许是為了节省磁盘空间，也许是为了缩减照片容量，以便用电子邮件寄出。不过在私底下，压缩使用的频率要更高：我们根本没有意识到，我们的下载或上传也可以通过压缩以节省带宽，而数据中心通常会压缩消费者的数据以降低成本。电子邮件提供商提供给你的5GB空间，经压缩后很有可能只占据电子邮件提供商5GB空间的很小一部分。

第八章讲到了数据库中运用的一些基础算法。这一章侧重为实现一致性——指一个数据库中的关系不互相冲突——而采用的聪明技巧。没

有这些精巧的技术，我们的绝大部分在线生活（包括网络购物以及通过Facebook之类的社交网站进行互动）就会消亡于众多计算机错误中。这一章解释了一致性真正的问题是什么，以及计算机科学家是如何解决这一问题的。前提是不牺牲我们所期望的在线系统拥有的高效性。

在第九章，我们会了解理论计算机科学无可争议的瑰宝之一：数字签名。乍看之下，用数字形式“签署”一份电子文档似乎不可能。你也许会想，这种签名必须由数字信息组成，而任何想要伪造签名的人都可以毫不费力地拷贝这些信息。这一悖论的解决方案，就是计算机科学取得的最令人瞩目的成就之一。

第十章采取了截然不同的视角：与其描述一个已经存在的伟大算法，我们不如去了解一个假如存在则必然会伟大的算法。不过我们会震惊地发现，这个特别伟大的算法不可能存在。这表明计算机解决问题的能力存在一些绝对极限，而我们将简单地从哲学和生物学角度探讨这一结果的应用。

第十一章我们会总结伟大算法的一些共性，花些时间畅想未来会怎样。会有更多伟大算法出现吗？或者说，我们已经发现了所有的伟大算法？

在此，不得不提前说一下本书的风格。任何科普作品都必须清楚地告知来源，但引用会破坏文本的流畅性，并让读者产生学术的感觉。由于可读性和易读性是本书的首要目标，所以本书正文不会出现引用。不过，我清楚地记录了所有来源，并在本书末尾的“来源和延伸阅读”板块中列出，并时不时附上拓展评论。这个板块还列出了一些额外材料，以便感兴趣的读者能去寻找更多和计算机科学中伟大算法有关的东西。

既然提前说了本书的风格，我还要谈谈本书书名中采取的少量诗化。本书无疑是革命性的，但真的有九种算法吗？这一说法值得探讨，因为要取决于有多少算法被算作单独算法。让我们来算下“九”是怎么来的。除了前言和结论两章外，本书还有九章，每一章都介绍了对一种计算任务产生革命性影响的算法，例如加密、压缩、图形识别。因此，书名中的“九

大算法”实际上指的是处理这九种任务的九类算法。

为什么我们要关注这些伟大的算法？

希望对这些迷人思想的快速总结能让你渴望深入了解它们的运行方式。不过，也许你仍然在思考：本书的终极目标是什么？让我简短地说下本书的真正目的。这本书绝不是一本问答式操作手册。在读完本书后，你不会成为计算机安全方面的专家，也不会成为人工智能或其他领域的专家。你也许能学到一些有用的技能，这倒是真的。比如：你会对如何检查“安全”网站凭证以及“已签名”软件包了解更多；你能针对不同任务在有损和无损压缩之间做出明智选择；而且通过理解搜索引擎索引和排名技术的某些方面，你能更高效地使用搜索引擎。

在读完本书后，你不会成为一名更加熟练的计算机用户。但你会更加珍视每天在所有计算设备上不停使用的美的思想。

为什么这是件好事？我用类比的方式来说明。我肯定不是一位天文学专家——事实上，我在这个项目上相当无知，我想知道更多。但每当我注视夜空，我知道的少量天文学知识增强了我对这一经验的享受。有时，我对自己看到事物的理解，让我产生了一种满足和惊奇的感觉。希望在读完本书后，你在使用计算机时也能经常获得同样的满足和惊奇之感，这也是我殷切的希望。你将真正珍视我们时代最常见、最神秘的黑盒子：你的个人电脑，你指尖的精灵。

