# 分布式算法导论

## （第二版）

# Introduction to Distributed Algorithms

## Second Edition

## 英文版

[荷] Gerard Tel 著

# 分布式算法导论
## （第二版）
### （英文版）

# Introduction to Distributed Algorithms
## Second Edition

［荷］　Gerard　Tel　著

# 内 容 简 介

分布式算法 20 多年来一直是倍受关注的主流方向。本书第二版不仅给出了算法的最新进展，还深入探讨了与之相关的理论知识。这本教材适合本科高年级和研究生使用，同时，本书所覆盖的广度和深度也十分适合从事实际工作的工程师和研究人员参考。书中重点讨论了点对点消息传递模型上的算法，也包括计算机通信网络的实现算法。其他重点讨论的内容包括分布式应用的控制算法（如波算法、广播算法、选举算法、终止检测算法、匿名网络的随机算法、快照算法、死锁检测算法、同步系统算法等），还涉及了利用分布式算法实现容错计算。第二版新增的关于方向感和故障检测器的内容都代表了当今最新技术发展水平，为在这些方向上从事研究的人员提供了很好的帮助。

# 出 版 说 明

21世纪初的5至10年是我国国民经济和社会发展的重要时期，也是信息产业快速发展的关键时期。在我国加入WTO后的今天，培养一支适应国际化竞争的一流IT人才队伍是我国高等教育的重要任务之一。信息科学和技术方面人才的优劣与多寡，是我国面对国际竞争时成败的关键因素。

当前，正值我国高等教育特别是信息科学领域的教育调整、变革的重大时期，为使我国教育体制与国际化接轨，有条件的高等院校正在为某些信息学科和技术课程使用国外优秀教材和优秀原版教材，以使我国在计算机教学上尽快赶上国际先进水平。

电子工业出版社秉承多年来引进国外优秀图书的经验，翻译出版了"国外计算机科学教材系列"丛书，这套教材覆盖学科范围广、领域宽、层次多，既有本科专业课程教材，也有研究生课程教材，以适应不同院系、不同专业、不同层次的师生对教材的需求，广大师生可自由选择和自由组合使用。这些教材涉及的学科方向包括网络与通信、操作系统、计算机组织与结构、算法与数据结构、数据库与信息处理、编程语言、图形图像与多媒体、软件工程等。同时，我们也适当引进了一些优秀英文原版教材，本着翻译版本和英文原版并重的原则，对重点图书既提供英文原版又提供相应的翻译版本。

在图书选题上，我们大都选择国外著名出版公司出版的高校教材，如Pearson Education培生教育出版集团、麦格劳-希尔教育出版集团、麻省理工学院出版社、剑桥大学出版社等。撰写教材的许多作者都是蜚声世界的教授、学者，如道格拉斯·科默（Douglas E. Comer）、威廉·斯托林斯（William Stallings）、哈维·戴特尔（Harvey M. Deitel）、尤利斯·布莱克（Uyless Black）等。

为确保教材的选题质量和翻译质量，我们约请了清华大学、北京大学、北京航空航天大学、复旦大学、上海交通大学、南京大学、浙江大学、哈尔滨工业大学、华中科技大学、西安交通大学、国防科学技术大学、解放军理工大学等著名高校的教授和骨干教师参与了本系列教材的选题、翻译和审校工作。他们中既有讲授同类教材的骨干教师、博士，也有积累了几十年教学经验的老教授和博士生导师。

在该系列教材的选题、翻译和编辑加工过程中，为提高教材质量，我们做了大量细致的工作，包括对所选教材进行全面论证；选择编辑时力求达到专业对口；对排版、印制质量进行严格把关。对于英文教材中出现的错误，我们通过与作者联络和网上下载勘误表等方式，逐一进行了修订。

此外，我们还将与国外著名出版公司合作，提供一些教材的教学支持资料，希望能为授课老师提供帮助。今后，我们将继续加强与各高校教师的密切联系，为广大师生引进更多的国外优秀教材和参考书，为我国计算机科学教学体系与国际教学体系的接轨做出努力。

电子工业出版社

# 教材出版委员会

# Preface

Distributed systems and distributed information processing have received considerable attention in the past few years, and almost every university offers at least one course on the design of distributed algorithms. There exist a large number of books about principles of distributed systems; see for example Tanenbaum [Tan96] or Sloman and Kramer [SK87], but these concentrate on architectural aspects rather than on algorithms. Since the first edition of this book, other texts on distributed algorithms have been published by Barbosa [Bar96], Lynch [Lyn96], and Attiya and Welch [AW98].

It has been remarked that algorithms are the backbone of every computer application; therefore a text devoted solely to distributed algorithms seems to be justified. The aim of this book is to present a large body of theory about distributed algorithms, which has been developed over the past twenty years or so. This book can be used as a textbook for a one- or two-semester course on distributed algorithms; the teacher of a one-semester course may select topics to his own liking.

The book will also provide useful background and reference information for professional engineers and researchers working with distributed systems.

**Exercises.** Each chapter (with the exception of Chapters 1 and 13) ends with a list of exercises and small projects. The projects usually require the reader to develop a small but non-trivial extension or application of the material treated in the chapter, and in most cases I do not have a "solution". If the reader succeeds in working out one of these small projects, I would be pleased to have a copy of the result.

A list of answers (sometimes partial) to most of the exercises is available for teachers; it can be obtained from the author or by anonymous ftp.

**Corrections and suggestions.** If the reader finds errors or omissions in this book, please inform the author (preferably by electronic mail). All constructive criticism, including suggestions for more exercises, is most welcome.

# 目 录 概 览

# Contents

# Chapter 1

# Introduction: Distributed Systems

This chapter gives reasons for the study of distributed algorithms by briefly introducing the types of hardware and software systems for which distributed algorithms have been developed. By a distributed system we mean all computer applications where several computers or processors cooperate in some way. This definition includes wide-area computer communication networks, but also local-area networks, multiprocessor computers in which each processor has its own control unit, and systems of cooperating processes.

The different types of distributed system and the reasons why distributed systems are used are discussed in Section 1.1. Some examples of existing systems will be given. The main topic of this book, however, is not what these systems look like, or how they are used, but how they can be made to work. And even that topic will be further specialized towards the treatment of the *algorithms* used in the systems.

Of course, the entire structure and operation of a distributed system is not fully understood by a study of its algorithms alone. To understand such a system fully one must also study the complete architecture of its hardware and software, that is, the partition of the entire functionality into modules. Also, there are many important questions related to properties of the programming languages used to build the software of distributed systems. These subjects will be discussed in Section 1.2.

It is the case, however, that there are already in existence excellent books about distributed systems, which concentrate on the architectural and language aspects; see, e.g., Tanenbaum [Tan96], Sloman and Kramer [SK87], Bal [Bal90], Coulouris and Dollimore [CD88] or Goscinski [Gos91]. As already mentioned, the present text concentrates on algorithms for distributed systems. Section 1.3 explains why the design of distributed algorithms dif-

fers from the design of centralized algorithms, sketches the research field of distributed algorithms, and outlines the remainder of the book.

## 1.1 What is a Distributed System?

In this chapter we shall use the term "distributed system" to mean an interconnected collection of autonomous computers, processes, or processors. The computers, processes, or processors are referred to as the *nodes* of the distributed system. (In the subsequent chapters we shall use a more technical notion, see Definition 2.6.) To be qualified as "autonomous", the nodes must at least be equipped with their own private control; thus, a parallel computer of the single-instruction, multiple-data (SIMD) model does not qualify as a distributed system. To be qualified as "interconnected", the nodes must be able to exchange information.

As (software) processes can play the role of nodes of a system, the definition includes software systems built as a collection of communicating processes, even when running on a single hardware installation. In most cases, however, a distributed system will at least contain several processors, interconnected by communication hardware.

More restrictive definitions of distributed systems are also found in the literature. Tanenbaum [Tan96], for example, considers a system to be distributed only if the existence of autonomous nodes is *transparent* to users of the system. A system distributed in this sense behaves like a virtual, stand-alone computer system, but the implementation of this transparency requires the development of intricate distributed control algorithms.

### *1.1.1 Motivation*

Distributed computer systems may be preferred over sequential systems, or their use may simply be unavoidable, for various reasons, some of which are discussed below. This list is not meant to be exhaustive. The choice of a distributed system may be motivated by more than one of the arguments listed below, and some of the advantages may come as a spin-off after the choice has been made for another reason. The characteristics of a distributed system may vary also, depending on the reason for its existence, but this will be discussed in more detail in Subsections 1.1.2 through 1.1.6.

(1) *Information exchange.* The need to exchange data between different computers arose in the sixties, when most major universities and companies started to have their own mainframe computer. Cooperation between people of different organizations was facilitated by

the exchange of data between the computers of these organizations, and this gave rise to the development of so-called *wide-area networks* (WANs). ARPANET, the predecessor of the current Internet, went on-air in December, 1969. A computer installation connected in a wide-area network (sometimes called a *long-haul network*) is typically equipped with everything a user needs, such as backup storage, disks, many application programs, and printers.

Later computers became smaller and cheaper, and soon each single organization had a multitude of computers, nowadays often a computer for each person (a personal computer or workstation). In this case also the (electronic) exchange of information between personnel of one organization already required that the autonomous computers were connected. It is even not uncommon for a single person or family to have multiple computers in the home, and connect these in a small personal home-network.

(2) *Resource sharing.* Although with cheaper computers it became feasible to equip each employee of an organization with a private computer, the same is not the case for the peripherals (such as printers, backup storage, and disk units). On this smaller scale each computer may rely on dedicated servers to supply it with compilers and other application programs. Also, it is not effective to replicate all application programs and related file resources in all computers; apart from the waste of disk space, this would create unnecessary maintenance problems. So the computers may rely on dedicated nodes for printer and disk service. A network connecting computers on an organization-wide scale is referred to as a *local-area network* (LAN).

The reasons for an organization to install a network of small computers rather than a mainframe are cost reduction and extensibility. First, smaller computers have a better price–performance ratio than large computers; a typical mainframe computer may perform 50 times faster than a typical personal computer, but cost 500 times more. Second, if the capacity of a system is no longer sufficient, a network can be made to fit the organization's needs by adding more machines (file servers, printers, and workstations). If the capacity of a stand-alone system is no longer sufficient, replacement is the only option.

(3) *Increased reliability through replication.* Distributed systems have the potential to be more reliable than stand-alone systems because they have a *partial-failure* property. By this it is meant that some nodes of the system may fail, while others are still operating correctly and can

take over the tasks of the failed components. The failure of a stand-alone computer affects the entire system and there is no possibility of continuing the operation in this case. For this reason distributed architectures are a traditional concern in the design of highly reliable computer systems.

A highly reliable system typically consists of a two, three, or four times replicated uniprocessor that runs an application program and is supplemented with a voting mechanism to filter the outputs of the machines. The correct operation of a distributed system in the presence of failures of components requires rather complicated algorithmical support.

(4) *Increased performance through parallelization.* The presence of multiple processors in a distributed system opens up the possibility of decreasing the turn-around time for a computation-intensive job by splitting the job over several processors.

Parallel computers are designed specifically with this objective in mind, but users of a local-area network may also profit from parallelism by shunting tasks to other workstations.

(5) *Simplification of design through specialization.* The design of a computer system can be very complicated, especially if considerable functionality is required. The design can often be simplified by splitting the system into modules, each of which implements part of the functionality and communicates with the other modules.

On the level of a single program modularity is obtained by defining abstract data types and procedures for different tasks. A larger system may be defined as a collection of cooperating processes. In both cases, the modules may all be executed on a single computer. But it is also possible to have a local-area network with different types of computers, one equipped with dedicated hardware for number crunching, another with graphical hardware, a third with disks, etc.

### 1.1.2 Computer Networks

By a computer network we mean a collection of computers, connected by communication mechanisms by means of which the computers can exchange information. This exchange takes place by sending and receiving messages. Computer networks fit our definition of distributed systems. Depending on the distance between the computers and their ownership, computer networks are called either *wide-area networks* or *local-area networks*.