SHIGERU KATAGIRI

EDITOR

HANDBOOK OF
NEURAL
NETWORKS
FOR SPEECH PROCESSING

# Handbook of Neural Networks for Speech Processing

Shigeru Katagiri

*Editor*

AH

Artech House
Boston • London
www.artechhouse.com

**Cover design by Igor Valdman**

# Handbook of Neural Networks for Speech Processing

For a listing of recent titles in the *Artech House Signal Processing Library*, turn to the back of this book.

# Preface

Speech is one of the most important means of human communication. Human beings communicate with one another by speaking and hearing, and they have made great efforts to achieve freer and more convenient communication by developing speech telecommunications technology such as telephony. In recent years, such technology has greatly advanced. The dream of communication—speaking with anybody, anywhere, anytime—is now becoming a reality. The dream, however, continues to expand. A more economic telecommunications system is desired. Even language barriers are expected to be overcome. Morever, work is being carried out to develop human-machine interfaces that use speech as a communications medium. Indeed, satisfying these demands is of great significance. Global and borderless speech communications assist daily human life, promote mutual understanding among nations, and greatly contribute to the general progress of human society. In pursuing these challenges, a key goal is to develop a useful engineering model of the human speech processing mechanism. This handbook seeks to provide a comprehensive introduction to one of the most important approaches to this technological challenge, speech processing using artificial neural networks (ANNs).

ANNs and speech processing are both transdisciplinary research fields, each having a history of several decades and involving many disciplines such as physiology, physics, statistics, psychology, linguistics, and engineering. Accordingly, a single book cannot be expected to fully cover the entire subject of speech processing using ANNs.

Fortunately, many excellent textbooks have been published in both the ANNs and speech processing fields. Readers can find detailed expertise in such books. In this handbook, we shall concentrate on the results of emerging and challenging studies on using ANNs for speech processing. We shall strive to convey the excitement of developments on the research front. The handbook will guide its readers through the unfinished but promising technological strides made in this field.

The book's 13 chapters are divided into three parts. Part I, entitled "Fundamentals," consists of 4 chapters that introduce basic information about speech processing and ANN technologies. In this introductory section, however, our focus is on speech processing technologies; the fundamentals of ANNs are discussed in Parts II and III. Part II, "Current Issues in Speech Recognition," emphasizes the considerable vigor with which ANNs have been investigated in their application to speech recognition. This second part consists of Chapters 5 through 9. Finally, Part III, "Current Issues in Speech Signal Processing," discusses speech-related topics such as speaker recognition, voice conversion, speech coding, and speech enhancement.

The handbook is aimed at researchers, engineers, and graduate-level students who wish to study the fundamentals and practical applications of neural-network-based speech processing. For these readers, Part I provides the necessary basis for proceeding to the later chapters on applications; Parts II and III provide a comprehensive introduction to the research front of neural networks for speech processing, which can be useful for further study of topics elaborated in related technical books and journals. Despite the primary aim of providing an introductory text, the book can also be used by researchers who have experience in ANN and speech processing technologies. The two latter parts provide good archives of results of the ongoing research. Each chapter was written with the aim of being as self-contained as possible. The book can therefore be used for various types of study by various readers. However, we, the chapter contributors, have a common motivation: we hope that the handbook becomes a vehicle to stimulate research by scientists and engineers on the human mechanism of speech processing and its engineering embodiment, i.e., neural networks for speech processing.

The inspiration for this handbook came from the series editor, Professor Alexsander D. Poularikas, University of Alabama in Huntsville. I wish to express my sincere thanks to him for his careful planning and for asking me to serve as editor of this text. The value of the book derives from the considerable efforts of the contributors, each of whom

brought extensive experience to the task of writing the individual chapters. I would like to express my considerable gratitude to these authors. Finally, it would not have been possible to produce the book without the patience and professional assistance provided by the crew at Artech House, especially Mark Walsh, Alexia Rosoff, Traci Beane, and Barbara Lovenvirth. I wish to thank them for help in all phases of making this book a reality.

—Shigeru Katagiri

## Recent Titles in the Artech House
## Signal Processing Library

# Contents

## Part II   Current Issues in Speech Recognition

## 5   Discriminative Prototype-Based Methods for Speech Recognition   159