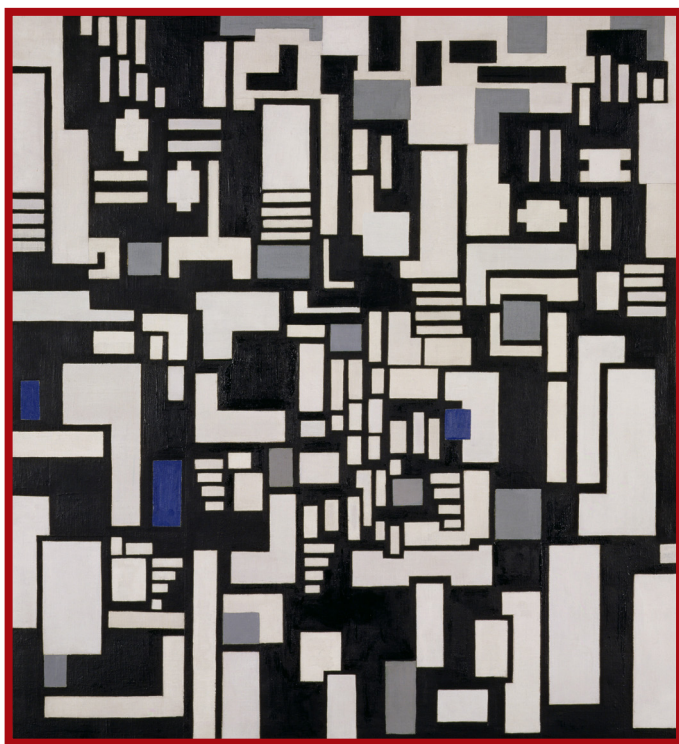*The Handbook of*
# Computational Linguistics and Natural Language Processing

*Edited by*

**Alexander Clark, Chris Fox, and Shalom Lappin**

Praise for *The Handbook of Computational Linguistics and Natural Language Processing*

"All in all, this is very well compiled book, which effectively balances the width and depth of theories and applications in two very diverse yet closely related fields of language research."

*Machine Translation*

"This *Handbook* is exceptionally broad and exceptionally deep in its coverage. The contributions, by noted experts, cover all aspects of the field, from fundamental theory to concrete applications. Clark, Fox and Lappin have performed a great service by compiling this volume."

*Richard Sproat, Oregon Health & Science University*

# Blackwell Handbooks in Linguistics

This outstanding multi-volume series covers all the major subdisciplines within linguistics today and, when complete, will offer a comprehensive survey of linguistics as a whole.

**Already published:**

*The Handbook of Child Language*
Edited by Paul Fletcher and Brian MacWhinney

*The Handbook of Phonological Theory, Second Edition*
Edited by John A. Goldsmith, Jason Riggle, and Alan C. L. Yu

*The Handbook of Contemporary Semantic Theory*
Edited by Shalom Lappin

*The Handbook of Sociolinguistics*
Edited by Florian Coulmas

*The Handbook of Phonetic Sciences, Second Edition*
Edited by William J. Hardcastle and John Laver

*The Handbook of Morphology*
Edited by Andrew Spencer and Arnold Zwicky

*The Handbook of Japanese Linguistics*
Edited by Natsuko Tsujimura

*The Handbook of Linguistics*
Edited by Mark Aronoff and Janie Rees-Miller

*The Handbook of Contemporary Syntactic Theory*
Edited by Mark Baltin and Chris Collins

*The Handbook of Discourse Analysis*
Edited by Deborah Schiffrin, Deborah Tannen, and Heidi E. Hamilton

*The Handbook of Language Variation and Change*
Edited by J. K. Chambers, Peter Trudgill, and Natalie Schilling-Estes

*The Handbook of Historical Linguistics*
Edited by Brian D. Joseph and Richard D. Janda

*The Handbook of Language and Gender*
Edited by Janet Holmes and Miriam Meyerhoff

*The Handbook of Second Language Acquisition*
Edited by Catherine J. Doughty and Michael H. Long

*The Handbook of Bilingualism and Multilingualism, Second Edition*
Edited by Tej K. Bhatia and William C. Ritchie

*The Handbook of Pragmatics*
Edited by Laurence R. Horn and Gregory Ward

*The Handbook of Applied Linguistics*
Edited by Alan Davies and Catherine Elder

*The Handbook of Speech Perception*
Edited by David B. Pisoni and Robert E. Remez

*The Handbook of the History of English*
Edited by Ans van Kemenade and Bettelou Los

*The Handbook of English Linguistics*
Edited by Bas Aarts and April McMahon

*The Handbook of World Englishes*
Edited by Braj B. Kachru; Yamuna Kachru, and Cecil L. Nelson

*The Handbook of Educational Linguistics*
Edited by Bernard Spolsky and Francis M. Hult

*The Handbook of Clinical Linguistics*
Edited by Martin J. Ball, Michael R. Perkins, Nicole Müller, and Sara Howard

*The Handbook of Pidgin and Creole Studies*
Edited by Silvia Kouwenberg and John Victor Singler

*The Handbook of Language Teaching*
Edited by Michael H. Long and Catherine J. Doughty

*The Handbook of Language Contact*
Edited by Raymond Hickey

*The Handbook of Language and Speech Disorders*
Edited by Jack S. Damico, Nicole Müller, Martin J. Ball

*The Handbook of Computational Linguistics and Natural Language Processing*
Edited by Alexander Clark, Chris Fox, and Shalom Lappin

*The Handbook of Language and Globalization*
Edited by Nikolas Coupland

*The Handbook of Hispanic Linguistics*
Edited by Manuel Díaz-Campos

*The Handbook of Language Socialization*
Edited by Alessandro Duranti, Elinor Ochs, and Bambi B. Schieffelin

*The Handbook of Intercultural Discourse and Communication*
Edited by Christina Bratt Paulston, Scott F. Kiesling, and Elizabeth S. Rangel

*The Handbook of Historical Sociolinguistics*
Edited by Juan Manuel Hernández-Campoy and Juan Camilo Conde-Silvestre

*The Handbook of Hispanic Linguistics*
Edited by José Ignacio Hualde, Antxon Olarrea, and Erin O'Rourke

*The Handbook of Conversation Analysis*
Edited by Jack Sidnell and Tanya Stivers

*The Handbook of English for Specific Purposes*
Edited by Brian Paltridge and Sue Starfield

# The Handbook of Computational Linguistics and Natural Language Processing

Edited by

*Alexander Clark, Chris Fox, and Shalom Lappin*

*For Camilla*
לאחיי דוד ודניאל, ולאחותי נעמי באהבה ובהומור

# Contents

# List of Figures