

Word Sense Disambiguation

Text, Speech and Language Technology

VOLUME 33

Series Editors

Nancy Ide, *Vassar College, New York*

Jean Véronis, *Université de Provence and CNRS, France*

Editorial Board

Harald Baayen, *Max Planck Institute for Psycholinguistics, The Netherlands*

Kenneth W. Church, *AT & T Bell Labs, New Jersey, USA*

Judith Klavans, *Columbia University, New York, USA*

David T. Barnard, *University of Regina, Canada*

Dan Tufis, *Romanian Academy of Sciences, Romania*

Joaquim Llisterri, *Universitat Autònoma de Barcelona, Spain*

Stig Johansson, *University of Oslo, Norway*

Joseph Mariani, *LIMSI-CNRS, France*

Eneko Agirre • Philip Edmonds
Editors

Word Sense Disambiguation Algorithms and Applications

 Springer

Eneko Agirre
University of the Basque Country
Basque Country
Spain

Philip Edmonds
Sharp Laboratories of Europe
Oxford
U.K.

ISBN 978-1-4020-4808-4

e-ISBN 978-1-4020-4809-2

Library of Congress Control Number: 2007938211

© 2007 Springer

No part of this work may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, microfilming, recording or otherwise, without written permission from the Publisher, with the exception of any material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work.

Printed on acid-free paper.

9 8 7 6 5 4 3 2 1

springer.com

Contents

Contributors	xiii
Foreword	xvii
Preface	xxi
1 Introduction	1
<i>Eneko Agirre and Philip Edmonds</i>	
1.1 Word Sense Disambiguation	1
1.2 A Brief History of WSD Research	4
1.3 What is a Word Sense?	8
1.4 Applications of WSD.....	10
1.5 Basic Approaches to WSD	12
1.6 State-of-the-Art Performance	14
1.7 Promising Directions	15
1.8 Overview of this Book.....	19
1.9 Further Reading	21
References	22
2 Word Senses	29
<i>Adam Kilgarriff</i>	
2.1 Introduction	29
2.2 Lexicographers	30
2.3 Philosophy	32
2.3.1 Meaning is Something You Do	32
2.3.2 The Fregean Tradition and Reification.....	33
2.3.3 Two Incompatible Semantics?.....	33
2.3.4 Implications for Word Senses.....	34
2.4 Lexicalization	35
2.5 Corpus Evidence.....	39
2.5.1 Lexicon Size	41
2.5.2 Quotations.....	42
2.6 Conclusion	43
2.7 Further Reading	44
Acknowledgments	45
References	45

3 Making Sense About Sense	47
<i>Nancy Ide and Yorick Wilks</i>	
3.1 Introduction	47
3.2 WSD and the Lexicographers	49
3.3 WSD and Sense Inventories	51
3.4 NLP Applications and WSD.....	55
3.5 What Level of Sense Distinctions Do We Need for NLP, If Any?	58
3.6 What Now for WSD?	64
3.7 Conclusion	68
References	68
 4 Evaluation of WSD Systems.....	 75
<i>Martha Palmer, Hwee Tou Ng, Hoa Trang Dang</i>	
4.1 Introduction	75
4.1.1 Terminology	76
4.1.2 Overview	80
4.2 Background.....	81
4.2.1 WordNet and Semcor	81
4.2.2 The <i>Line</i> and <i>Interest</i> Corpora.....	83
4.2.3 The DSO Corpus	84
4.2.4 Open Mind Word Expert	85
4.3 Evaluation Using Pseudo-Words.....	86
4.4 Senseval Evaluation Exercises	86
4.4.1 Senseval-1.....	87
Evaluation and Scoring.....	88
4.4.2 Senseval-2.....	88
English All-Words Task	89
English Lexical Sample Task	89
4.4.3 Comparison of Tagging Exercises.....	91
4.5 Sources of Inter-Annotator Disagreement	92
4.6 Granularity of Sense: Groupings for WordNet.....	95
4.6.1 Criteria for WordNet Sense Grouping.....	96
4.6.2 Analysis of Sense Grouping	97
4.7 Senseval-3.....	98
4.8 Discussion.....	99
References	102
 5 Knowledge-Based Methods for WSD.....	 107
<i>Rada Mihalcea</i>	
5.1 Introduction	107
5.2 Lesk Algorithm.....	108
5.2.1 Variations of the Lesk Algorithm.....	110
Simulated Annealing	110
Simplified Lesk Algorithm.....	111

Augmented Semantic Spaces	113
Summary	113
5.3 Semantic Similarity	114
5.3.1 Measures of Semantic Similarity.....	114
5.3.2 Using Semantic Similarity Within a Local Context	117
5.3.3 Using Semantic Similarity Within a Global Context.....	118
5.4 Selectional Preferences.....	119
5.4.1 Preliminaries: Learning Word-to-Word Relations	120
5.4.2 Learning Selectional Preferences	120
5.4.3 Using Selectional Preferences	122
5.5 Heuristics for Word Sense Disambiguation.....	123
5.5.1 Most Frequent Sense	123
5.5.2 One Sense Per Discourse.....	124
5.5.3 One Sense Per Collocation	124
5.6 Knowledge-Based Methods at Senseval-2	125
5.7 Conclusions	126
References	127

6 Unsupervised Corpus-Based Methods for WSD..... 133

Ted Pedersen

6.1 Introduction	133
6.1.1 Scope	134
6.1.2 Motivation	136
Distributional Methods.....	137
Translational Equivalence	139
6.1.3 Approaches	140
6.2 Type-Based Discrimination.....	141
6.2.1 Representation of Context	142
6.2.2 Algorithms	145
Latent Semantic Analysis (LSA).....	146
Hyperspace Analogue to Language (HAL).....	147
Clustering By Committee (CBC)	148
6.2.3 Discussion.....	150
6.3 Token-Based Discrimination.....	150
6.3.1 Representation of Context	151
6.3.2 Algorithms	151
Context Group Discrimination	152
McQuitty’s Similarity Analysis.....	154
6.3.3 Discussion.....	157
6.4 Translational Equivalence	158
6.4.1 Representation of Context	159
6.4.2 Algorithms	159
6.4.3 Discussion.....	160

6.5 Conclusions and the Way Forward.....	161
Acknowledgments	162
References	162

7 Supervised Corpus-Based Methods for WSD 167

Lluís Màrquez, Gerard Escudero, David Martínez, German Rigau

7.1 Introduction to Supervised WSD.....	167
7.1.1 Machine Learning for Classification	168
An Example on WSD.....	170
7.2 A Survey of Supervised WSD.....	171
7.2.1 Main Corpora Used	172
7.2.2 Main Sense Repositories	173
7.2.3 Representation of Examples by Means of Features.....	174
7.2.4 Main Approaches to Supervised WSD.....	175
Probabilistic Methods.....	175
Methods Based on the Similarity of the Examples	176
Methods Based on Discriminating Rules	177
Methods Based on Rule Combination.....	179
Linear Classifiers and Kernel-Based Approaches.....	179
Discourse Properties: The Yarowsky Bootstrapping Algorithm	181
7.2.5 Supervised Systems in the Senseval Evaluations	183
7.3 An Empirical Study of Supervised Algorithms for WSD.....	184
7.3.1 Five Learning Algorithms Under Study	185
Naïve Bayes (NB)	185
Exemplar-Based Learning (kNN)	186
Decision Lists (DL).....	187
AdaBoost (AB).....	187
Support Vector Machines (SVM).....	189
7.3.2 Empirical Evaluation on the DSO Corpus.....	190
Experiments.....	191
7.4 Current Challenges of the Supervised Approach.....	195
7.4.1 Right-Sized Training Sets.....	195
7.4.2 Porting Across Corpora	196
7.4.3 The Knowledge Acquisition Bottleneck.....	197
Automatic Acquisition of Training Examples.....	198
Active Learning.....	199
Combining Training Examples from Different Words	199
Parallel Corpora.....	200
7.4.4 Bootstrapping	201
7.4.5 Feature Selection and Parameter Optimization	202
7.4.6 Combination of Algorithms and Knowledge Sources	203
7.5 Conclusions and Future Trends	205

Acknowledgments	206
References	207
8 Knowledge Sources for WSD.....	217
<i>Eneko Agirre and Mark Stevenson</i>	
8.1 Introduction	217
8.2 Knowledge Sources Relevant to WSD	218
8.2.1 Syntactic	219
Part of Speech (KS 1)	219
Morphology (KS 2).....	219
Collocations (KS 3)	220
Subcategorization (KS 4).....	220
8.2.2 Semantic	220
Frequency of Senses (KS 5)	220
Semantic Word Associations (KS 6)	221
Selectional Preferences (KS 7)	221
Semantic Roles (KS 8).....	222
8.2.3 Pragmatic/Topical.....	222
Domain (KS 9).....	222
Topical Word Association (KS 10)	222
Pragmatics (KS 11).....	223
8.3 Features and Lexical Resources.....	223
8.3.1 Target-Word Specific Features.....	224
8.3.2 Local Features	225
8.3.3 Global Features.....	227
8.4 Identifying Knowledge Sources in Actual Systems	228
8.4.1 Senseval-2 Systems	229
8.4.2 Senseval-3 Systems	231
8.5 Comparison of Experimental Results	231
8.5.1 Senseval Results	232
8.5.2 Yarowsky and Florian (2002).....	233
8.5.3 Lee and Ng (2002).....	234
8.5.4 Martínez et al. (2002)	237
8.5.5 Agirre and Martínez (2001a)	238
8.5.6 Stevenson and Wilks (2001).....	240
8.6 Discussion.....	242
8.7 Conclusions	245
Acknowledgments	246
References	247
9 Automatic Acquisition of Lexical Information and Examples	253
<i>Julio Gonzalo and Felisa Verdejo</i>	
9.1 Introduction	253
9.2 Mining Topical Knowledge About Word Senses	254

9.2.1 Topic Signatures	255
9.2.2 Association of Web Directories to Word Senses.....	257
9.3 Automatic Acquisition of Sense-Tagged Corpora.....	258
9.3.1 Acquisition by Direct Web Searching	258
9.3.2 Bootstrapping from Seed Examples	261
9.3.3 Acquisition via Web Directories	263
9.3.4 Acquisition via Cross-Language Evidence	264
9.3.5 Web-Based Cooperative Annotation	268
9.4 Discussion.....	269
Acknowledgments	271
References	272
10 Domain-Specific WSD.....	275
<i>Paul Buitelaar, Bernardo Magnini, Carlo Strapparava, Piek Vossen</i>	
10.1 Introduction	275
10.2 Approaches to Domain-Specific WSD.....	277
10.2.1 Subject Codes	277
10.2.2 Topic Signatures and Topic Variation.....	282
Topic Signatures.....	282
Topic Variation	283
10.2.3 Domain Tuning.....	284
Top-down Domain Tuning	285
Bottom-up Domain Tuning.....	285
10.3 Domain-Specific Disambiguation in Applications	288
10.3.1 User-Modeling for Recommender Systems.....	288
10.3.2 Cross-Lingual Information Retrieval.....	289
10.3.3 The MEANING Project.....	292
10.4 Conclusions	295
References	296
11 WSD in NLP Applications	299
<i>Philip Resnik</i>	
11.1 Introduction	299
11.2 Why WSD?.....	300
Argument from Faith.....	300
Argument by Analogy.....	301
Argument from Specific Applications	302
11.3 Traditional WSD in Applications	303
11.3.1 WSD in Traditional Information Retrieval.....	304
11.3.2 WSD in Applications Related to Information Retrieval.....	307
Cross-Language IR.....	308
Question Answering.....	309
Document Classification	312
11.3.3 WSD in Traditional Machine Translation	313

11.3.4 Sense Ambiguity in Statistical Machine Translation.....	315
11.3.5 Other Emerging Applications.....	317
11.4 Alternative Conceptions of Word Sense.....	320
11.4.1 Richer Linguistic Representations.....	320
11.4.2 Patterns of Usage.....	321
11.4.3 Cross-Language Relationships.....	323
11.5 Conclusions.....	325
Acknowledgments.....	325
References.....	326
A Resources for WSD.....	339
A.1 Sense Inventories.....	339
A.1.1 Dictionaries.....	339
A.1.2 Thesauri.....	341
A.1.3 Lexical Knowledge Bases.....	341
A.2 Corpora.....	343
A.2.1 Raw Corpora.....	343
A.2.2 Sense-Tagged Corpora.....	345
A.2.3 Automatically Tagged Corpora.....	347
A.3 Other Resources.....	348
A.3.1 Software.....	348
A.3.2 Utilities, Demos, and Data.....	349
A.3.3 Language Data Providers.....	350
A.3.4 Organizations and Mailing Lists.....	350
Index of Terms.....	353
Index of Authors and Algorithms.....	361

Contributors

Eneko Agirre is an Associate Professor in the University of the Basque Country, where he is member of the IXA NLP group. He organized the Basque tasks for Senseval and coordinates the construction of the Basque WordNet and Semcor. Department of Computer Science, University of the Basque Country, Manuel de Lardizabal 1, E-20018 Donostia, Basque Country, Spain.

Paul Buitelaar is a Senior Researcher in the Language Technology Lab and co-chair of the Competence Center Semantic Web at DFKI (German Research Center for Artificial Intelligence) GmbH. He was organizer of several international workshops and has been an invited speaker at panels and workshops on topics in semantic annotation and ontology development. DFKI GmbH Language Technology Department, Stuhlsatzenhausweg 3, Saarbrücken, Germany.

Hoa Trang Dang is a Computer Scientist at the National Institute of Standards and Technology (NIST), where she coordinates evaluations of automatic question answering and summarization systems in TREC and DUC. National Institute of Standards and Technology, 100 Bureau Drive, Mailstop 8940, Gaithersburg, MD 20899-8940, U.S.A.

Philip Edmonds is a Research Scientist at Sharp Laboratories of Europe. He was chair of Senseval, 2001–2004, and is the author of the entry on lexical disambiguation in the *Elsevier Encyclopedia of Language and Linguistics, 2nd Ed.* Sharp Laboratories of Europe Limited, Oxford Science Park, Oxford OX4 4GB, United Kingdom.

Gerard Escudero is an Assistant Professor at the Universitat Politècnica de Catalunya. He was a participant in Senseval-2 and Senseval-3. He also participated in the MEANING project, funded by the EU. EUETIB, Urgell 187, E-08036 Barcelona, Catalonia, Spain.

Julio Gonzalo is an Assistant Professor at the UNED School of Computer Science. He is co-editor of the CLEF (Cross-Language Evaluation Forum) proceedings on multilingual information access published by Springer. In 2006, he is co-chair of the Programme Committee of the European

Conference on Advanced Research and Development for Digital Libraries. Dep. Lenguajes y Sistemas Informáticos, ETSI Informática – UNED, Ciudad Universitaria, c/ Juan del Rosal 16, 28040 Madrid, Spain.

Nancy Ide is a Professor of Computer Science at Vassar College and chair of the Computer Science Department. She is founder of the Text Encoding Initiative (TEI) and creator of the Corpus Encoding Standard. Currently she is directing the development of the American National Corpus. Department of Computer Science, Vassar College, 124 Raymond Avenue, Poughkeepsie, New York 12604-0520, U.S.A.

Adam Kilgarriff is Director of Lexical Computing Ltd. and Visiting Research Fellow at the University of Sussex, U.K. He works on both the theory and the practice, at the intersection of language corpora, language technologies and practical dictionary-making. Lexical Computing Ltd., 71 Freshfield Road, Brighton BN2 0BL, U.K.

Bernardo Magnini is a Senior Researcher at ITC-irst, where he coordinates the research group on Text Technologies. He is the local organizer co-chair of EACL-06, the 11th Conference of the European Chapter of the Association for Computational Linguistics. ITC-irst, Via Sommarive 18, I-38050, Povo-Trento, Italy.

Lluís Màrquez is an Associate Professor at the Polytechnical University of Catalunya (UPC). He organized two shared tasks on semantic role labeling at the Conference on Natural Language Learning (CoNLL) in 2004 and 2005, and led the team that organized the Catalan and Spanish lexical sample tasks at Senseval-3. In 2006, he will be the co-chair of the CoNLL conference. Despatx S120 - Edifici Omega, Campus Nord UPC, C/ Jordi Girona Salgado 1-3, E-08034 Barcelona, Catalonia, Spain.

David Martínez is a post-doc researcher in the NLP group of the University of Sheffield. Natural Language Processing Group, Department of Computer Science, University of Sheffield, Sheffield, S1 4DP, United Kingdom.

Rada Mihalcea is an Assistant Professor of Computer Science at University of North Texas. She is the president of ACL SIGLEX and was a co-chair of Senseval-3. Department of Computer Science, University of North Texas, PO Box 311366, Denton, TX 76203, U.S.A.

Hwee Tou Ng is an Associate Professor of Computer Science at the National University of Singapore. He is on the editorial board of *Computational Linguistics*, was program co-chair of the ACL-2005 conference, and has served on the program committees of many past conferences including ACL, SIGIR, AAI, and IJCAI. Department of Computer Science, School of Computing, National University of Singapore, 3 Science Drive 2, Singapore 117543.

Martha Palmer is an Associate Professor in the Departments of Linguistics and Computer Science and a Faculty Fellow of the Institute of Cognitive Science at the University of Colorado at Boulder. She has been a member of the Advisory Committee for the DARPA TIDES program, chair of ACL SIGLEX and ACL SIGHAN, and is currently Past-President of the Association for Computational Linguistics. Department of Linguistics, 295 UCB - Hellems 295, Boulder, CO 80309, U.S.A.

Ted Pedersen is an Associate Professor in the Department of Computer Science at the University of Minnesota, Duluth. He is the recipient of a National Science Foundation (NSF, USA) Faculty Early Career Development (CAREER) Award. Department of Computer Science, 1114 Kirby Drive, University of Minnesota, Duluth, MN 55812, U.S.A.

Philip Resnik is an associate professor at the University of Maryland, College Park, in the Department of Linguistics and the Institute for Advanced Computer Studies. He is on the editorial board of *Cognitive Linguistics*. 1401 Marie Mount Hall, University of Maryland, College Park, MD 20742, U.S.A.

German Rigau is an Associate Professor in the Department of Computer Science of the Basque Country University. He coordinated the EU's 5th framework MEANING project. He has also participated in Senseval-2 and Senseval-3. Department of Computer Science, University of the Basque Country, Manuel de Lardizabal 1, E-20018 Donostia, Basque Country, Spain.

Mark Stevenson is a Lecturer in Computer Science at the University of Sheffield. He is author of the monograph *Word Sense Disambiguation: Combining Knowledge Sources for Sense Resolution* (2003) based on his Ph.D. thesis. Natural Language Processing Group, Department of Computer Science, University of Sheffield, Sheffield, S1 4DP, United Kingdom.

Carlo Strapparava is a Senior Researcher at ITC-irst in the Communication and Cognitive Technologies Division. He is author of over ninety published papers on topics including artificial intelligence, natural language processing, and word sense disambiguation. ITC-irst, Via Sommarive, 18 I-38050, Povo-Trento, Italy.

M. Felisa Verdejo is Full Professor and head of the department Lenguajes y Sistemas Informáticos (LSI) at National Distance Learning University (UNED). She has been involved in several large-scale EU-funded projects such as EuroWordNet and CLEF. Dep. Lenguajes y Sistemas Informáticos, ETSI Informática – UNED, Ciudad Universitaria, c/ Juan del Rosal 16, 28040 Madrid, Spain.

Piek Vossen is CTO of Irion Technologies. He worked on several EU projects: Acquilex, Sift, EAGLES, EuroWordNet, EuroTerm, BalkaNet and MEANING and most recently on an American project to develop an Arabic wordnet. He is also founder and president of the Global Wordnet Association (GWA). Irion Technologies BV., Delftechpark 26, 2628 XH Delft, PO Box 2849, 2601 CV Delft, The Netherlands.

Yorick Wilks is Professor of Computer Science at the University of Sheffield. He is author of numerous articles and six books including *Electric Words: Dictionaries, Computers and Meanings* (1996 with Brian Slator and Louise Guthrie). He is a Fellow of the American and European Associations for Artificial Intelligence, and on the boards of some fifteen AI-related journals. Department of Computer Science, University of Sheffield, Sheffield S1 4DP, United Kingdom.

Foreword

Graeme Hirst

University of Toronto

Of the many kinds of ambiguity in language, the two that have received the most attention in computational linguistics are those of word senses and those of syntactic structure, and the reasons for this are clear: these ambiguities are overt, their resolution is seemingly essential for any practical application, and they seem to require a wide variety of methods and knowledge-sources with no pattern apparent in what any particular instance requires.

Right at the birth of artificial intelligence, in his 1950 paper “Computing machinery and intelligence”, Alan Turing saw the ability to understand language as an essential test of intelligence, and an essential test of language understanding was an ability to disambiguate; his example involved deciding between the generic and specific readings of the phrase *a winter’s day*. The first generations of AI researchers found it easy to construct examples of ambiguities whose resolution seemed to require vast knowledge and deep understanding of the world and complex inference on this knowledge; for example, *Pharmacists dispense with accuracy*. The disambiguation problem was, in a way, nothing less than the artificial intelligence problem itself. No use was seen for a disambiguation method that was less than 100% perfect; either it worked or it didn’t. Lexical resources, such as they were, were considered secondary to non-linguistic common-sense knowledge of the world.

And because the methods that were developed required a resource whose eventual existence was merely hypothesized – a knowledge base containing everything a typical adult knows – and because there were no test data available, it was not possible to empirically test them or quantitatively evaluate them or their underlying ideas in any serious way. Rather, systems and methods were presented like theorems whose truth or correctness could be demonstrated by a rational argument bolstered by hand-waving and a ‘toy’

demonstration: a knowledge source would be built for a few words and facts, and the system would be run on a few “interesting” constructed examples to show that it did “the right thing”. This approach to evaluation was quite normal in the milieu in which this research was carried out and didn’t seem to worry anyone at the time: computational linguistics had not yet achieved its empirical orientation.

Contemporary approaches have turned all that upside-down. Statistical and machine-learning methods and methodologies that have been adopted in the last decade have revolutionized our view of ambiguity resolution. It is now understood that imperfect methods that rely on rich lexical resources but limited additional knowledge have great use in the world; and that systems must undergo rigorous evaluation. The present volume demonstrates this in particular for word sense disambiguation – both the strengths and the inherent limitations of these approaches.¹ In particular, contemporary methods are less ambitious and have lower expectations. Unlike the earlier research, they don’t worry about case roles, about helping a parser with attachment decisions, or about working with a semantic interpretation process aimed at a deep level of “understanding”. Rather than aiming for a complete solution and hypothesizing a resource that this necessitates, they rely on an existing resource and try to see how much can be done with it. And yet they still have enormous application in NLP (see Chap. 11).

One issue that has remained constant is what kinds of information in the text may be drawn upon as cues for disambiguation, and how near in the text to the target word those cues should be. In my own early work (Hirst 1987), restrictions on communication between disambiguating processes arose from two competing principles: any particular word or structural cue for disambiguation has quite a limited sphere of influence, and yet almost anything in a text or discourse is potentially a cue for disambiguation (cf. McRoy 1992). In contemporary systems, the analogous dilemma is in the choice of features and the window size (see Chap. 8).

The other thing that hasn’t changed is how hard the lexical disambiguation problem is. Many sophisticated systems struggle merely to reach the modest accuracy of simple baseline algorithms such as that of Lesk (1986) (see Chap. 5) or just choosing the most frequent sense. But what is a poor computer to do when humans themselves frequently disagree on what the correct answer is supposed to be (see Chaps. 2–4)?

Although it is an edited volume, this book is not an anthology of “recent advances” papers by individual authors on their own research, requiring

¹ A similar revolution has occurred in parsing and structural disambiguation; see Manning and Schütze (2000, Chaps. 11–12) for an overview.

each reader to synthesize a view of the overall situation in a research topic. Rather, editors Agirre and Edmonds have enlisted the leading researchers of the field to do the hard work. Each chapter of this book presents an overview and synthesis of one facet of current research. The result is a clear and well-organized presentation of the state of the art in word sense disambiguation that can be read, like a textbook, from start to finish. I commend it to you.

Graeme Hirst is the author of Semantic Interpretation and the Resolution of Ambiguity (Cambridge University Press, 1987), which presents an integrated theory of lexical disambiguation, structural disambiguation, and semantic interpretation.

References

- Hirst, Graeme. 1987. *Semantic Interpretation and the Resolution of Ambiguity*. Cambridge University Press.
- Lesk, Michael. 1986. Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from an ice cream cone. *Proceedings of SIGDOC-86: 5th International Conference on Systems Documentation*, Toronto, Canada, 24–26.
- Manning, Christopher D. & Hinrich Schütze. 2000. *Foundations of Statistical Natural Language Processing*. Cambridge, MA: The MIT Press.
- McRoy, Susan. 1992. Using multiple knowledge sources for word sense discrimination. *Computational Linguistics*, 19(1):1–30.
- Turing, Alan M. 1950. Computing machinery and intelligence. *Mind*, 59:433–460. Reprinted in: Stuart Shieber, ed. 2004. *The Turing Test: Verbal Behavior as the Hallmark of Intelligence*. Cambridge, MA: The MIT Press.

Preface

Word sense disambiguation is a core research problem in computational linguistics, which was recognized at the beginning of the scientific interest in machine translation and artificial intelligence. And yet no book has been fully devoted to review the wide variety of approaches to solving the problem. The time is right for such a book.

This book had its genesis over five years ago when Nancy Ide, series co-editor of then Kluwer's, now Springer's, *Text, Speech, and Language Technology* series, approached us with the project. Word sense disambiguation is an active and quickly progressing research field, so we thought it far more beneficial to the research community if we were to enlist the main experts to each give their own view of the field.

Being the first major book on the topic, and with the hope of it becoming the definitive reference, we endeavoured to fashion a coherent, consistent, critical, and readable survey of the current state of the art. We started by sketching an overview of the main topics that should be covered, and then approached experts in the field with desiderata for each chapter. We requested that authors give a general overview of their topic and proceed with a thorough exposition of the theory, methodology, algorithms, critical analysis, experimentation, results, and open issues. We are indebted to all of the authors, who worked with us most patiently.

The manuscript has taken time to produce, having been through numerous reviews and revisions along the way. Many difficult decisions were made in the attempt to best embrace all of the important research in the field, and to keep up with new developments. We apologize if we have missed something.

Please visit the book website, www.wsdbook.org, for the latest information updates, and a book search interface.

Word sense disambiguation is a fascinating topic; we hope you enjoy reading this book as much as we did creating it!

Acknowledgments

First all, it is the chapter authors who created the book; we thank them all for exceeding their remits and for their patience during the lengthy reviewing cycles. We owe the existence of this book to Nancy Ide, the series co-editor.

A few people gave us encouragement and feedback at various stages about the content and organization of the book. We are grateful to Robert Dale, David Farwell, Graeme Hirst, Eduard Hovy, Inderjeet Mani, Pete Whitelock, and David Yarowsky.

Three anonymous reviewers helped us get around the weak points. We also thank the team at Kluwer and Springer for their support: Tamara Welschot, Jacqueline Bergsma, Helen van der Stelt, and Jolanda Voogd (Associate Publishing Editor).

Phil was supported by Sharp Laboratories of Europe, and Eneko by the IXA research group.

Finally, this book is dedicated to our families, who had to sacrifice their time with us for the sake of this book.

Phil Edmonds and Eneko Agirre
27 January 2006