

ENGLISH CORPORA
AND
AUTOMATED GRAMMATICAL
ANALYSIS

英语语料库与自动语法分析

ALEX CHENGYU FANG



THE COMMERCIAL PRESS

H314
E1595

English Corpora and Automated Grammatical Analysis

Alex Chengyu Fang



上314

THE COMMERCIAL PRESS

2007 • Beijing

图书在版编目(CIP)数据

英语语料库与自动语法分析/方称宇著. —北京：
商务印书馆,2007

ISBN 978 - 7 - 100 - 05659 - 5

I . 英… II . 方… III . 英语—语料库—语法 IV . H314

中国版本图书馆 CIP 数据核字(2007)第 162671 号

所有权利保留。

未经许可，不得以任何方式使用。

YINGYÙ YÜLIAOKÙ YÙ ZIDÒNG YŪFĀ FÉNХÍ

英语语料库与自动语法分析

方称宇 著

商 务 印 书 馆 出 版

(北京王府井大街36号 邮政编码100710)

商 务 印 书 馆 发 行

商务印书馆上海印刷股份有限公司印刷

ISBN 978 - 7 - 100 - 05659 - 5

2007年11月第1版

开本 787×960 1/16

2007年11月上海第1次印刷

印张 15 1/4

印数 1 000 册

定价：38.00 元

For Lan

Preface

With this book, I would like to commemorate my father, an experimental scientist, who taught me so much about research, practice, and application, not by words but through deeds and demonstrated enthusiasm about novelties. This book could not have been completed without the unreserved support and care from Lan. I am grateful to my two children, Adrien and Elizabeth, for giving me so much joy and pleasure during the course of my research described in this book. I am grateful to all the others in the family, whose enduring love accompanied me throughout the years in London.

Much of the research described in this book was conducted during my association with the Survey of English Usage, University College London, when the International Corpus of English (ICE) was under way to record authentic uses of English worldwide. Pressure to process large quantities of such data for detailed linguistic description prompted me into this research, for which I am deeply indebted to the late Professor Sidney Greenbaum, director of the Survey and coordinator of the ICE project. This is an appropriate place to record my gratitude to him.

I would like to thank Dr Mark Huckvale of the Department of Phonetics and Linguistics, University College London, for his great patience, academic insight, and valuable comments on an earlier draft. His emphasis on practical applications of parsing technologies, especially in the area of linguistically informed speech processing, had an important influence on my research. Different parts of the book have been read by Mr Eric Atwell at Leeds University and Dr Jonathan Ginzburg at King's College London, to whom I would like to express sincere thanks. Of course, responsibility for all the shortcomings remains mine.

Credits must also go to Professor John Campbell of the Computer Science Department, University College London, and Professor Jan Svartvik of the Survey of Spoken English, University of Lund, for many an enlightening after-dinner discussion on various topics in Delancey

viii *Preface*

Street, Highfield, and Alyth Gardens. Mr Nigel Birch of the Engineering and Physical Sciences Research Council (UK) provided me with timely valuable support of the research.

Finally, I would like to acknowledge the support I have received from the Halliday Centre for Intelligent Applications of Language Studies as well as the Department of Chinese, Translation, and Linguistics of City University of Hong Kong. Not least, thanks go to Professor Housheng Qian of the Commercial Press for his great patience and kind assistance in the publication of this book.

October 2007

A.C.F.

前言

从 1993 年到 2005 年，我在伦敦大学学院 (University College London, 简称 UCL) 从事科研和教学工作。本书记载了我多年来在语料库语言学和计算语言学这两个领域的主要研究心得和成果。

上世纪 90 年代，是英国语料库语言学发展的黄金时期。伦敦的 Randolph Quirk 教授和 Sidney Greenbaum 教授、兰开斯特的 Geoffrey Leech 教授、伯明翰的 John Sinclair 教授都在进行语料库的开发工作。

当时，Sidney Greenbaum 教授任 UCL 的英语用法调查中心 (Survey of English Usage) 主任，正在从事国际英语语料库 (The International Corpus of English) 的创建工作。100 万字的英国英语语料已经采集完毕，语法标码也已完成，但句法分析遇到不少困难。一是所用的句法分析系统不适用，每输入一个语句，常生成几十、上百、甚至上千棵句法树，然后再人工选取，十分耗时耗力。二是所用的形式语法不适用。当时的语法为英语书面语所写，而 100 万字的英国英语语料包含 60 万字的口语，所以几乎每天都要开会讨论一些语句的具体处理，语法的某些部分干脆需要重写，尤其是不同层次上的并列结构。尽管如此，最后还是有大约 30% 的语句，自动句法分析系统根本无法应付。

于是，Sidney Greenbaum 教授和我在 1994 年一同撰写了一份项目申请书，然后约见了英国工程及物理科学研究委员会 (Engineering and Physical Sciences Research Council) 的有关人员，其中包括 Nigel Birch 先生和 Mark Tatham 教授，提出了我们的研究设想。这份申请最后通过了委员会的评审，获得了一笔约 50 万英镑的资助，专门用于研制一个新的自动句法分析系统并重写一部新的、可用于英语口语分析的形式语法。

研究项目的主要思路就是将已经分析过的语料库变成一个句法知识库，从中提取短语结构语法规则，并通过基于实例的手段，在

x 前言

知识库中为待分析语句提取一棵最佳句法树。这样的句法分析机制涉及几个重要课题：首先需要一个高质量的自动词类标码系统，不仅能对大类进行判别，而且能对小类的细分进行快速、有效的精确分析，比如说动词的配价问题。然后，我们需要一个短语分析系统，将待分析语句处理成一个短语结构集，然后据此计算句法相似度，最终生成相应的句法树。这样一种句法分析途径，具有强劲、高效、精确和自动学习等特性，在对国际英语语料库及其他海量语料库的处理中得到广泛检测和验证。

本书对上述各个部分的研究进行了详细的描述，对系统的实际表现进行了深入的量化评测，并有专门章节来探讨句法分析的评测问题。除此之外，还探讨了介词短语的自动分析，特别是这类短语的句法功能的自动判定，因为这一研究和句法相似度分析有着密切的关系。同时，本书还就自动语法分析在语音合成及语音识别中的应用做了相应的介绍和说明，希望对读者能有所帮助。

我的不少朋友及同事都看过本书的初稿或部分章节，并提出过许多建议，在此表示感谢，特别是伦敦大学学院的 John Campbell 教授和 Mark Huckvale 博士、伦敦国王学院的 Jonathan Ginzburg 博士、利兹大学的 Eric Atwell 先生、瑞典隆德大学的 Jan Svartvik 教授及商务印书馆上海信息中心主任钱厚生教授。当然，我对书中的所有错误负全责，并恳请读者提出宝贵批评和建议。

最后，我以此书来缅怀先父对我的言传身教和恩师 Sidney Greenbaum 教授对我的栽培，并感谢家人对我的关心和支持。

方称宇
2007年10月

Abstract

This book describes novel research in the robust practical processing of English for detailed word class and syntactic annotations of natural texts. In particular, this book describes an automatic part-of-speech tagger and a robust parser, both of which demonstrate state-of-the-art performance in efficiency, accuracy, and robustness through their application for the annotation of a large corpus of transcribed speech and writing.

The first part of the book is a description of AUTASYS (An Automatic Text Analysis System of English), which combines both a statistical backbone and a set of heuristics. Its novelties include the mapping of the Lancaster-Oslo-Bergen (LOB) tag set, achieved through a bi-gram language model, to a syntactically rich set designed for the International Corpus of English (ICE), successfully achieved even without global syntactic information. This research was extensively tested and evaluated with corpus data. The accuracy was estimated at about 93%.

The book then describes the Survey Parser, which uses AUTASYS as a pre-processor for tokenised input tagged with the ICE tag set. It then produces either a complete parse that indicates both the syntactic category and the syntactic function for each unit of analysis or a partial analysis with syntactic functions or attachment undetermined. A novelty of this work is the use of a large context-free grammar automatically generated from the syntactically annotated ICE corpus. A second novelty is treating parsing as a search problem for a string that has been previously analysed. A major contribution by this work is the use of tag-feature combinations as the anchor of sub-trees for permissive derivations, a more generalised approach than other similarity-based parsing methodologies such as data-oriented parsing that typically make use of lexical anchors for sub-trees. Extensive tests and evaluations of this module were carefully conducted to determine the coverage, the labelling accuracy, and the attachment accuracy of the Survey Parser. Experiments demonstrated a combined accuracy of about 80%, com-

parable to the performance of some of the major systems reported in literature.

A special chapter was devoted to the study of the automatic determination of the syntactic functions of the prepositional phrase. A lexicalist approach was proposed that successfully resolved the attachment of 85% of the prepositional phrases without global parsing. As the final chapter concludes, the successful handling of prepositional phrases is expected to enhance significantly the performance of the Survey Parser. This parser will also benefit from a future study that aims at the automatic detection of clause boundaries.

Contents

Preface	vii
前言	ix
List of Figures	xvii
List of Tables.....	xix
Abstract	xxi
1. Introduction	1
1.1. What is Parsing?	1
1.2. The Introspective View.....	3
1.3. The Retrospective View	6
1.4. Data-Oriented Parsing	9
1.5. General Problems.....	12
1.6. The Proposed Research.....	13
1.6.1. Background to the Proposed Research.....	13
1.6.2. The Basic Approach of the Proposed Research	16
1.6.3. The Strengths and Novelties of the Proposed Approach.....	17
1.6.3.1. Automated Grammar Generation	18
1.6.3.2. De-Lexicalised Terminal Nodes.....	18
1.6.3.3. Global Parse with Subcategorisation Features	18
1.6.3.4. High-Quality Partial Parse.....	19
1.6.3.5. Intrinsic Ability to Learn	19
1.7. The Organisation of the Book.....	20
2. The Automatic Analysis of English Word Classes	22
2.1. An Overview of Word Class Tagging	22
2.2. Major Word Class Tagging Schemes	22
2.2.1. The Lancaster-Oslo/Bergen Tagging Scheme	23

2.2.1.1. The Lancaster-Oslo-Bergen Corpus.....	24
2.2.1.2. The Lancaster-Oslo-Bergen Tag Set	24
2.2.1.3. Summary	27
2.2.2. The International Corpus of English Tagging Scheme	28
2.2.2.1. The International Corpus of English	28
2.2.2.2. The International Corpus of English Tag Set	29
2.2.3. A Comparison of LOB and ICE.....	33
2.3. Word Class Tagging Methodologies	39
2.3.1. The Rule-Based Approach.....	39
2.3.2. The Probabilistic Approach	40
2.4. AUTASYS: A Hybrid Tagging System.....	42
2.4.1. A Probabilistic Approach Using the LOB Tag Set	42
2.4.1.1. The Tag Assignment Module.....	43
2.4.1.1.1. Tokenisation	44
2.4.1.1.2. The treatment of “.”	44
2.4.1.1.3. The treatment of “”	45
2.4.1.1.4. Sentence boundary markers.....	46
2.4.1.2. Orthographic Analysis	47
2.4.1.3. Lexicon Lookup	47
2.4.1.3.1. The lexicon	48
2.4.1.3.2. The coverage of the lexicon	49
2.4.1.4. Morphological Analysis	51
2.4.2. The Idiom Identification Module.....	52
2.4.3. The Probabilistic Tag Selection Module.....	53
2.4.3.1. The Bigram Probabilistic Matrix.....	53
2.4.3.2. Implementing Probabilistic Tag Selection	55
2.4.4. The Rule-Based Refinement Module.....	57
2.4.5. Empirical Evaluation.....	58
2.4.6. Permissive AUTASYS-LOB Disagreements.....	61
2.4.6.1. NNP-NPT	61
2.4.6.2. JJ-JJB.....	62
2.4.6.3. NNP-NPL	62
2.4.6.4. RB-NN	63

2.4.7. Summary	65
2.5. A Rule-Based Approach towards LOB to ICE Translation.....	66
2.5.1. Solutions for Verbs	67
2.5.1.1. Auxiliary vs. Lexical	67
2.5.1.2. Monotransitive vs. Complex Transitive	67
2.5.1.3. Finite vs. Nonfinite.....	68
2.5.2. Closed Sets.....	70
2.5.3. Initial Results	71
2.5.4. Problems	71
2.5.5. Summary	74
3. The Automatic Induction of a Formal Grammar.....	76
3.1. Introduction.....	76
3.2. The ICE Parsing Scheme	77
3.3. Grammar Generation	79
3.3.1. Phrase Structure Rules	80
3.3.2. Phrase Structure Cluster Rules.....	82
3.4. Evaluating the Coverage.....	83
3.4.1. The Construction of the Training and Test Sets.....	84
3.4.2. The Number of Extracted Rules.....	84
3.4.3. The Coverage of Adjective Phrase Rules.....	85
3.4.4. The Coverage of Adverb Phrase Rules	85
3.4.5. The Coverage of Noun Phrase Rules	86
3.4.6. The Coverage of Verb Phrase Rules	87
3.4.7. The Coverage of Prepositional Phrase Rules	87
3.4.8. The Coverage of Phrase Structure Cluster Rules	88
3.4.9. Discussion	88
4. Robust Practical Analogy-Based Parsing.....	90
4.1. Introduction.....	90
4.1.1. Analogy-Based Parsing.....	90
4.2. An Overview of the Survey Parser	91
4.3. The Construction of the Syntactic Knowledge Base	93

4.3.1. Phrase Structure Rules	93
4.3.2. Phrase Structure Cluster Rules.....	96
4.3.2.1., Feature Constraints.....	97
4.3.2.2. Nonobligatory Elements.....	98
4.3.2.3. A Definition of Analogy.....	99
4.4. Parsing with Phrase Structure and Phrase Structure Cluster Rules	100
4.4.1. The Analysis of Word Classes	100
4.4.2. The Analysis of Phrases.....	102
4.4.3. The Analysis of Clauses.....	103
4.4.3.1. Partial Analysis	105
4.5. Some Initial Evaluation	106
4.5.1. Evaluating the Coverage of Phrase Structure Rules	106
4.5.2. Evaluating the Coverage of Phrase Structure Cluster Rules	108
4.5.3. Evaluating the Labelling Precision	109
4.5.4. Evaluating the Accuracy of Analysis.....	110
4.5.5. Evaluating the Processing Speed	111
4.6. Concluding Remarks	111
5. Extensive Evaluations of the Survey Parser.....	114
5.1. Introduction.....	114
5.2. Commonly Used Metrics	116
5.2.1. Labelled Match	116
5.2.2. Bracketed Match	118
5.2.3. Crossing-Brackets Rate.....	119
5.2.4. Summary	120
5.3. Evaluations with the NIST Scheme	123
5.3.1. Labelling Accuracy	123
5.3.1.1. Methodology	123
5.3.1.2. Evaluating the Scoring Program.....	127
5.3.1.3. Evaluating the Labelling Accuracy of Parser- Produced Trees	128

5.3.2. Bracketing Accuracy	128
5.3.2.1. A Linear Representation of the Hierarchical Structure	131
5.3.2.2. Evaluating the Scoring Program.....	132
5.3.2.3. Evaluating the Bracketing Accuracy of Parser- Produced Trees	133
5.3.3. Labelling and Bracketing Accuracy Scores Combined.....	134
5.3.3.1. A Description	135
5.3.3.2. Empirical Scores by the Survey Parser	136
5.3.3.3. A Comparison with other Reports.....	137
5.3.3.3.1. A Comparison with Keller (2003).....	137
5.3.3.3.2. A Comparison with Plaehn (2004).....	137
5.3.3.3.3. A Comparison with Henderson (2004).....	138
5.3.3.3.4. Summary	138
5.4. Summary.....	138
6. The Resolution of Prepositional Phrases.....	141
6.1. Introduction.....	141
6.2. PPs in Contemporary British English	145
6.2.1. Prepositions	145
6.2.2. Prepositional Complement.....	146
6.2.3. The Syntactic Functions of Prepositional Phrases	148
6.3. Data and Resources Used in the Experiment.....	149
6.4. Scope of Experiment.....	150
6.5. Test Data.....	151
6.6. Lexical Database.....	152
6.7. Rules and their Coverage.....	154
6.7.1. Prepositional Phrases as Adjective Phrase Postmodifiers....	155
6.7.1.1. Test Data	155
6.7.1.2. Attachment Rules	156
6.7.1.3. Morphological Analysis	157
6.7.1.4. Coverage.....	158
6.7.2. Prepositional Phrases as Noun Phrase Postmodifiers	159

6.7.2.1. Test Data	160
6.7.2.2. Attachment Rules	160
6.7.2.3. Morphological Analysis	161
6.7.2.4. Coverage of the Rules	161
6.7.3. Prepositional Phrases as Adverbials.....	162
6.8. Overall Success Rate and Discussion	162
7. Conclusions and Further Work	165
7.1. Conclusions	165
7.2. Applications of AUTASYS and the Survey Parser	167
7.2.1. AUTASYS	167
7.2.2. The Survey Parser	168
7.2.2.1. SpeechMaker.....	169
7.2.2.2. Correlation between Tone Units and Syntax.....	169
7.3. Future Work.....	171
7.3.1. Implementation of Automated Prepositional Phrase Attachment	171
7.3.2. Automated Clause Boundary Detection and Attachment	172
References	175
Appendix A: A List of LOB Tags	187
Appendix B: A List of ICE Tags.....	192
Appendix C: A List of AUTASYS Idioms.....	200
Appendix D: A List of ICE Parsing Symbols	204
Appendix E: A List of ICE Prepositions in Descending Frequency Order	207
Appendix F: A Distributional Profile of ICE-GB Prepositions ...	211
Index	214

List of Figures

Figure 1.	A syntactic tree for [1].....	2
Figure 2.	A syntactic tree for [2].....	4
Figure 3.	A syntactic tree with rule probabilities for T_1	7
Figure 4.	A syntactic tree with rule probabilities for T_2	8
Figure 5.	Fragment t obtained from a syntactic tree T	11
Figure 6.	Step 1: The substitution of Mary into T to produce T_1	11
Figure 7.	Step 2: The substitution of sports into T_1 to produce T_2	11
Figure 8.	The transitivity types in ICE.....	30
Figure 9.	The major components in the probabilistic model	42
Figure 10.	The major components in the tag assignment module.....	43
Figure 11.	A flowchart for the processing of “.”	44
Figure 12.	A flowchart for the processing of “”	46
Figure 13.	The subclasses of common nouns in LOB	61
Figure 14.	The syntactic tree of [76].....	78
Figure 15.	A NP tree without functional label.....	81
Figure 16.	A global tree for [j].....	83
Figure 17.	The derivational tree of [82].....	83
Figure 18.	The number of extracted rules	84
Figure 19.	The coverage of AJP rules.....	85
Figure 20.	The coverage of AVP rules	86
Figure 21.	The coverage of NP rules	86
Figure 22.	The coverage of VP rules	87
Figure 23.	The coverage of PP rules	87
Figure 24.	The coverage of CL rules	88
Figure 25.	The major components in the Survey Parser	92
Figure 26.	The tree structure for Example [83]	94
Figure 27.	A flowchart of analogy identification by the Survey Parser	99
Figure 28.	Subcategorisation of verbs	101
Figure 29.	Input string as a PS cluster with associated sub-trees	103
Figure 30.	The final analysis of [84].....	104