# Theory and Applications
# of Natural Language Processing

Series Editors:
Graeme Hirst (Textbooks)
Eduard Hovy (Edited volumes)
Mark Johnson (Monographs)

## Aims and Scope

The field of Natural Language Processing (NLP) has expanded explosively over the past decade: growing bodies of available data, novel fields of applications, emerging areas and new connections to neighboring fields have all led to increasing output and to diversification of research.

"Theory and Applications of Natural Language Processing" is a series of volumes dedicated to selected topics in NLP and Language Technology. It focuses on the most recent advances in all areas of the computational modeling and processing of speech and text across languages and domains. Due to the rapid pace of development, the diversity of approaches and application scenarios are scattered in an ever-growing mass of conference proceedings, making entry into the field difficult for both students and potential users. Volumes in the series facilitate this first step and can be used as a teaching aid, advanced-level information resource or a point of reference.

The series encourages the submission of research monographs, contributed volumes and surveys, lecture notes and textbooks covering research frontiers on all relevant topics, offering a platform for the rapid publication of cutting-edge research as well as for comprehensive monographs that cover the full range of research on specific problem areas.

The topics include applications of NLP techniques to gain insights into the use and functioning of language, as well as the use of language technology in applications that enable communication, knowledge management and discovery such as natural language generation, information retrieval, question-answering, machine translation, localization and related fields.

The books are available in printed and electronic (e-book) form:

* Downloadable on your PC, e-reader or iPad
* Enhanced by Electronic Supplementary Material, such as algorithms, demonstrations, software, images and videos
* Available online within an extensive network of academic and corporate R&D libraries worldwide
* Never out of print thanks to innovative print-on-demand services
* Competitively priced print editions for eBook customers thanks to MyCopy service http://www.springer.com/librarians/e-content/mycopy

For other titles published in this series, go to
www.springer.com/series/8899

Slav Petrov

# Coarse-to-Fine
# Natural Language Processing

Foreword by Eugene Charniak

🐴 Springer

Slav Petrov
Google
slav@petrovi.de

*Cover design*: deblik, Berlin

Printed on acid-free paper

*To my family*

# Foreword

Grammars for natural languages show how sentences (and their meaning) are built up out of smaller pieces. Syntactic parsing is the task of applying a grammar to a string of words (a sentence) in order to reconstruct this structure. For example, "The dog thought there was day-old food in his dish" has a sub-structure "there was day-old food in his dish" which in turn contains structures like "day-old food." Before we can build the meaning of the whole we must at least identify the parts from which it is built. This is what parsing gives us.

As with most all areas of natural-language processing (NLP) parsing research has greatly benefited from the statistical revolution — the process of absorbing statistical learning techniques into NLP that began about twenty five years ago. Prior to that time we had no parser that could, say, assign a plausible structure for every sentence in your local newspaper. Now you can download several good ones on the web.

From the outside the result has looked sort of like a Moore's law scenario. Every few years parsers got more accurate, or much more efficient, or both. From inside, however, things looked quite different. At more than one occasion we in the community had no idea where the next improvement would come from and some thought that we had, perhaps, reached the end of the road. The last time the improvement came from Slav Petrov and the ideas in this monograph. The embodiment of these ideas is the "Berkeley Parser."

The best parsers models are all "supervised," e.g., we have a corpus of sentences, in the case here the so-called "Penn tree-bank" where sentences have been analyzed by people so for each sentence has been broken down into a tree structure of components. A computer learns to parse new sentences by collecting statics from the training data that (we hope) reflect generalizations about a particular language, in this case English. We then recast the parsing problem as one of applied statistics and probability — find the most probable parse for the sentences according the the probabilities already obtained from the corpus.

To over simplify, until Slav's work, the best parsers could be thought of as word-based — rules should be based upon the words found in their examples. A paradigmatic case wold be, say, the use of the prepositional phrase "out of ..."

when talking about removing something by "washing", but not by, say, "sanding." Unfortunately the number of words in English is very large (really unbounded), so this data would be missing many crucial word-grammar combinations. In this cases the parser would "back off" and look from grammar rules ignoring the particular words in question.

The Berkeley parser, however bases rules not on words, but on sets of words. The "coarse to fine' of the title refers to the graularity of these sets. So the counter claim would be that "washing" is not unique here, but is rather one of a group of words that also include "scrubbing" and in some cases "flooding" (I flooded the cinder out of my eye). Unfortunately such groups can be quite idiosyncratic, so it might be that we are still better off at the word level. Indeed, the two methods can be thought of as two ends of a continuum, and perhaps future work can now combine the approaches. But until the Berkeley parser we did have a good concrete example of this second approach.

Furthermore, for anyone with a good machine learning background, once you see how this parser works, it makes immediate sense. Thus for people like me, at least, Slav's work is very easy to read. Perhaps I am not a "typical" person, but take it from me, there are a lot of papers in my research area that I do not find so easy.

Thus I strongly recommend Slav's work to you. It is major advance in the area of syntactic parsing, and a great advertisement for the superiority of the machine-learning approach to the field.

Brown University                                                      *Eugene Charniak*

# Preface

This book is based on my homonymous PhD thesis filed at the University of California, Berkeley in 2009. It has been updated to reference new work that has happened since then. It has also been reformatted to fit this paper size.

# Acknowledgements

# Contents

# List of Figures

# List of Tables