

# Advances in Computational Linguistics

## Volume Editor:

Editor de Volumen

*Alexander Gelbukh*

Instituto Politécnico Nacional  
Centro de Investigación en Computación  
México 2009



## Preface

Computational linguistics is an interdisciplinary research area that combines the ideas and methods of linguistics, computer science, and artificial intelligence and has two-fold goal: on the one hand, to study human language by means of modern computational methods, and on the other hand, to develop computer programs capable of human-like activities related to understanding or producing texts or speech in human language, such as English or Chinese.

The most important technical applications of computational linguistics include information retrieval and information organization, machine translation, and natural language interfaces, among others. However, as in any science, the activities of the researchers are mostly concentrated on its internal art and craft, that is, on the solution of the problems arising in analysis or synthesis of natural language text or speech, such as syntactic and semantic analysis, disambiguation, or compilation of dictionaries and grammars necessary for such analysis.

This volume presents 25 original research papers written by 64 authors representing 23 different countries: Algeria, Brazil, Canada, China, Czech Republic, France, Germany, Greece, Hong Kong, India, Islamic Republic of Iran, Italy, Jordan, Lithuania, Macao, Mexico, Portugal, Russian Federation, Spain, Switzerland, Turkey, United Kingdom, and United States. The volume is structured in 8 thematic areas of both theory and applications of computational linguistics:

- Computational lexicography and lexical resources
- Morphology and syntax
- Semantics
- Anaphora and co-reference
- Text classification
- Text summarization
- Speech generation
- Applications

The papers included in this volume were selected on the base of rigorous international reviewing process out of 65 submissions considered for evaluation; thus the acceptance rate of this volume was 38%.

I would like to cordially thank all people involved in the preparation of this volume. In the first place I want to thank the authors of the published paper for their excellent research work that gives sense to the work of all other people involved, as well as the authors of rejected papers for their interest and effort. I also thank the members of the Editorial Board of the volume and additional reviewers for their hard work on reviewing and selecting the papers. I thank Sulema Torres, Ignacio Garcia Araoz, Oralia del Carmen Pérez Orozco, and Raquel López Alamilla for their valuable collaboration in preparation of this volume. The submission, reviewing, and selection process was supported for free by the EasyChair system, [www.EasyChair.org](http://www.EasyChair.org).

# Table of Contents

## Índice

Page/Pág.

### Computational Lexicography and Lexical Resources

Scaling to Billion-plus Word Corpora .....	3
<i>Jan Pomikálek, Pavel Rychlý and Adam Kilgarriff</i>	
Detecting and Grounding Terms in Biomedical Literature.....	15
<i>Kaarel Kaljurand, Fabio Rinaldi, Thomas Kappeler and Gerold Schneider</i>	
Automatic Word Clustering in Studying Semantic Structure of Texts.....	27
<i>Olga Mitrofanova</i>	
Semantically-Driven Extraction of Relations between Named Entities .....	35
<i>Caroline Brun and Caroline Hagège</i>	
Evaluation of Named Entity Extraction Systems.....	47
<i>Mónica Marrero, Sonia Sánchez-Cuadrado, Jorge Morato Lara and George Andreidakis</i>	
Building a Brazilian Portuguese Parallel Corpus of Original and Simplified Texts .....	59
<i>Helena M. Caseli, Tiago F. Pereira, Lucia Specia, Thiago A. S. Pardo, Caroline Gasperin and Sandra M. Aluisio</i>	

### Morphology and Syntax

Synchronized Morphological and Syntactic Disambiguation for Arabic .....	73
<i>Daoud Daoud</i>	
Parsing with Polymorphic Categorical Grammars .....	87
<i>Matteo Capelletti and Fabio Tamburini</i>	
A CCG-based System for Valence Shifting for Sentiment Analysis .....	99
<i>František Šimančík and Mark Lee</i>	
Classification of “Inheritance” Relations: a Semi-automatic Approach .....	109
<i>Ekaterina Lapshinova-Koltunski</i>	

### Semantics

Discovering Discourse Motifs in Instructional Dialog .....	123
<i>Juan M. Huerta</i>	
Ontology Oriented Computation of English Verbs Metaphorical Trait.....	135
<i>Zili Chen, Jonathan J. Webster, Ian I. Chow and Tianyong Hao</i>	
Automatic Formal Verification of Conceptual Model Documentation by Means of Self-organizing Map .....	143
<i>Algirdas Laukaitis and Olegas Vasilecas</i>	

## **Anaphora and Co-reference**

Coreference Resolution using Markov Logic Network .....	157
<i>Shujian Huang, Yabing Zhang, Junsheng Zhou and Jiajun Chen</i>	
A Ranking Approach to Persian Pronoun Resolution.....	169
<i>Nafiseh Sadat Moosavi and Gholamreza Ghassem-Sani</i>	

## **Text Classification**

Classification of Clinical Conditions: A Case Study on Prediction of Obesity and Its Co-morbidities .....	183
<i>Archana Bhattarai, Vasile Rus and Dipankar Dasgupta</i>	
Analysis of Stemming Alternatives and Dependency Pattern Support in Text Classification.....	195
<i>Levent Özgür and Tunga Güngör</i>	
Investigating Variations in Adjective Use across Different Text Categories .....	207
<i>Jing Cao and Alex Chengyu Fang</i>	
Multi-category Support Vector Machines for Identifying Arabic Topics.....	217
<i>Mourad Abbas, Kamel Smaili and Daoud Berkani</i>	

## **Text Summarization**

USUM: Update Summary Generation System .....	229
<i>C Ravindranath Chowdary and P Sreenivasa Kumar</i>	

## **Speech Generation**

Towards the Speech Synthesis of Raramuri: A Unit Selection Approach based on Unsupervised Extraction of Suffix Sequences.....	243
<i>Alfonso Medina Urrea, José Abel Herrera Camacho and Maribel Alvarado García</i>	
Pronunciation Rules in Portuguese Regional Speech (PORT REG) for Coarticulation Process.....	257
<i>Sara Candeias and Jorge Morais Barbosa</i>	

## **Applications**

A Maximum Entropy (ME) Based Translation Model for Chinese Characters Conversion .....	267
<i>Fai Wong, Sam Chao, Cheong Cheong Hao and Ka Seng Leong</i>	
Tracking Out-of-date Newspaper Articles.....	277
<i>Frederik Cailliau, Aude Giraudel and Béatrice Arnulphy</i>	
PtiClic: A Game for Vocabulary Assessment combining JeuxDeMots and LSA.....	289
<i>Mathieu Lafourcade and Virginie Zampa</i>	

# Scaling to Billion-plus Word Corpora

Jan Pomikálek<sup>1</sup>, Pavel Rychlý<sup>1</sup> and Adam Kilgarriff<sup>2</sup>

<sup>1</sup> Masaryk University, Brno, Czech Republic

<sup>2</sup> Lexical Computing Ltd, Brighton, UK

**Abstract.** Most phenomena in natural languages are distributed in accordance with Zipf's law, so many words, phrases and other items occur rarely and we need very large corpora to provide evidence about them. Previous work shows that it is possible to create very large (multi-billion word) corpora from the web. The usability of such corpora is often limited by duplicate contents and a lack of efficient query tools.

This paper describes BiWeC, a Big Web Corpus of English texts currently comprising 5.5b words fully processed, and with a target size of 20b. We present a method for detecting near-duplicate text documents in multi-billion-word text collections and describe how one corpus query tool, the Sketch Engine, has been re-engineered to efficiently encode, process and query such corpora on low-cost hardware.

## 1 Introduction

There's no data like more data, and one place to get more data almost without limit (for general English and some other languages and varieties) is the web. One way to use the web is to create a local corpus by downloading web pages: in [1] we argue that it is the optimal way to use the web for linguistic research. A number of corpora have been built in this way: Baroni and colleagues developed web corpora with nearly 2 billion words for German, Italian and English [2, 3] and have made them available for research as tar archives. Liu et al. [4] describe the creation of a 10 billion word corpus. In this paper we introduce BiWeC, a Big Web Corpus currently of 5.5b words, with a target size of 20b.

Very large corpora can be created on low cost hardware in a few person-months. Most of the steps have linear complexity and scale up well. The two outstanding issues we focus on in this paper are:

1. removing duplicate content
2. efficient querying.

The article is organized as follows. In section two our motivation for creating larger corpora is discussed and the advantages of using more data for various tasks is explained. Section three describes the process of creating BiWeC, with a focus on removing duplicate and near-duplicate documents. Section four deals with corpus processing and querying using the Sketch Engine corpus manager. Then we present some figures about BiWeC and outline future plans.

# Detecting and Grounding Terms in Biomedical Literature

Kaarel Kaljurand, Fabio Rinaldi, Thomas Kappeler and Gerold Schneider

Institute of Computational Linguistics, University of Zurich  
kalju@ifi.uzh.ch, rinaldi@ifi.uzh.ch, kappeler@bluewin.ch,  
gschneid@ifi.uzh.ch

**Abstract.** We present an approach towards the automatic detection of names of proteins, genes, species, etc. in biomedical literature and their grounding to widely accepted identifiers. The annotation is based on a large term list that contains the common expression of the terms, a normalization step that matches the terms with their actual representation in the texts, and a disambiguation step that resolves the ambiguity of matched terms. We describe various characteristics of the terms found in existing term resources and of the terms that are used in biomedical texts. We evaluate our results against a corpus of manually annotated protein mentions and achieve a precision of 57% and recall of 72%.

## 1 Introduction

The complexity of biological organisms and the success of biological research in describing them, have resulted in a large body of biological entities (genes, proteins, species, etc.) to be indexed, named and analyzed. Probably the most important entities are proteins. They are an essential part of an organism and participate in every process within cells. Most proteins function in collaboration with other proteins, and one of the research goals in molecular biology is to identify which proteins interact.

While the number of different proteins is large, the amount of their possible interactions and combinations is even larger. In order to record such interactions and represent them in a structured way, human curators who work for knowledge base projects, e.g. Molecular INTeraction database (MINT)<sup>1</sup>, Human Protein Reference Database (HPRD)<sup>2</sup>, IntAct<sup>3</sup> (see [4] for a detailed overview), carefully analyze published biomedical articles. As the body of articles is growing rapidly, there is a need for effective automatic tools to help curators in their work. Such tools must be able to detect mentions of biological entities in the text and tag them with identifiers that have been assigned by existing knowledge bases. As the names that are used to reference the proteins can be very ambiguous, there is a need for an effective ambiguity resolution.

<sup>1</sup> <http://mint.bio.uniroma2.it>

<sup>2</sup> <http://www.hprd.org/>

<sup>3</sup> <http://www.ebi.ac.uk/intact>

# Automatic Word Clustering in Studying Semantic Structure of Texts

Olga Mitrofanova

St. Petersburg State University, Faculty of Philology and Arts,  
Department of Mathematical Linguistics,  
Universitetskaya emb., 11  
199034 St. Petersburg, Russia  
{alkonost-om@yandex.ru}

**Abstract.** The purpose of the study is to prove that results of automatic word clustering (AWC) may contribute much in investigating semantic structure of texts and in evaluating plot complexity. Experiments were carried out for Russian texts, mainly stories and short novels. Data obtained in course of study allowed to formulate and verify several linguistic hypotheses.

**Keywords:** Automatic Word Clustering, Russian Corpora, Semantic Structure of Texts

## 1 Introduction

Formalization of text structure and quantitative evaluation of semantic relations between text units prove to be of considerable importance in various fields of natural language understanding: modelling plot structure, text summarization, evaluation of translation adequacy in parallel texts, automatic text indexing, classification of texts in corpora, etc. (for a detailed analysis cf. [1], [2]).

One of the procedures providing linguistic data on semantic structure of texts is automatic word clustering (AWC). It is assumed that AWC results help to reveal semantic structure of texts and to determine plot complexity. To prove this assumption, AWC procedure was carried out with the help of a specialized AWC toolkit based on word space model. Experimental procedure implied processing Russian texts, mainly stories and short novels. A set of key words describing major topics of the plot was assigned to each text, clusters of words with similar distributions were created for each key word. Data extracted from texts through AWC procedure admit thorough linguistic interpretation. Further comparison of cluster content and structure allowed to distinguish texts characterized by a plot including a dominating topic with a number of subtopics and texts characterized by a plot including a set of major (independent or correlating) topics.

# Semantically-Driven Extraction of Relations between Named Entities

Caroline Brun and Caroline Hagège

Xerox Research Centre Europe, 6 chemin de Maupertuis 38240 Meylan, France  
{Caroline.Brun, Caroline.Hagege}@xrce.xerox.com

**Abstract.** In this paper, we describe a method that automatically generates lexico-syntactic patterns which are then used to extract semantic relations between named entities. The method uses a small set of seeds, i.e. named entities that are a priori known to be in relation. This information can easily be extracted from encyclopedias or existing databases. From very large corpora we extract sentences that contain combinations of these attested entities. These sentences are then used in order to automatically generate, using a syntactic parser, lexico-syntactic patterns that links these entities. These patterns are then re-applied on texts in order to extract relations between new entities of the same type. Furthermore, the patterns that are extracted not only provide a way to spot new entities relations but also build a valuable paraphrase resource. An evaluation on the relation holding between an event, the place of the event occurrence and the date of the event occurrence has been carried out on French corpus and shows good results. We believe that this kind of methodology can be applied for other kinds of relation between named entities.

## 1 Introduction

In this paper we describe a system that extracts accurately semantic relations between named entities from raw text. Taking as input a small set of already known relations that can be extracted from encyclopedias or from databases, our system first learn from a large corpus a wide range of lexico-syntactic patterns conveying the desired semantic relation. These learned patterns are then further applied on texts, and as a result, new occurrences of the given semantic relations linking new entities are detected. As all the patterns extracted represent a comparable semantic situation, they can be considered as paraphrase patterns. These patterns can then be used both in generation and for information extraction tasks.

## 2 Related Work

Many research works on extraction of relation between entities have already been performed since this kind of information is useful for a wide range of applications of information extraction. For instance [8] describe an algorithm to extract relations between named entities and the resulting improvement of a question answering



# Evaluation of Named Entity Extraction Systems

Mónica Marrero, Sonia Sánchez-Cuadrado,  
Jorge Morato Lara and George Andreiadakis

Computer Engineering Department, University Carlos III of Madrid  
Av. de la Universidad 30, 28911 Leganés (Madrid), Spain  
mmarrero@inf.uc3m.es, ssanche@ie.inf.uc3m.es, jmorato@inf.uc3m.es, gand@ie.inf.uc3m.es

**Abstract.** The suitability of the algorithms for recognition and classification of entities (NERC) is evaluated through competitions such as MUC, CONLL or ACE. In general, these competitions are limited to the recognition of predefined entity types in certain languages. In addition, the evaluation of free applications and commercial systems that do not attend the competitions has been lightly studied. Shallowly studied have also been the causes of erroneous results. In this study a set of NERC tools are assessed. The assessment of the tools has consisted of: 1) the elaboration of a test corpus with typical and marginal types of entities; 2) the elaboration of a brief technical specification for the tools evaluated; 3) the assessment of the quality of the tools for the developed corpus by means of precision-recall ratios; 4) the analysis of the most frequent errors. The sufficiency of the technical characteristics of the tools and their evaluation ratios, presents an objective perspective of the quality and the effectiveness of the recognition and classification techniques of each tool. Thus, the study complements the information provided by other competitions and aids the choice or the design of more suitable NER tools for a specific project.

**Keywords:** Named entity extraction, named entity recognition and classification, information extraction, named entity extraction tools.

## 1 Introduction

There is currently a wide variety of named entities (NE) recognition systems. Competitive events are organized for the evaluation of NERC systems, in which the ability of identification and classification of the entities existing in a corpus is analyzed. Nevertheless, the competitions normally establish certain limitations such as:

- They focus on a limited group of NE types. This feature is quite variable due to the ambiguity in the use of the term *Named Entity* depending on the different forums or events. In the case of the MUC conferences, NEs were considered *personal names*, *organizations*, *locations* and at a later stage, *temporal entities* and *measurements* [1]. On the other hand, the CONLL-2002/2003 conferences defined the categories *person*, *organism*, *localization* and *miscellaneous* [2, 3]. The latter (*miscellaneous*), includes proper names of different nature with different categories: *gentilics*, *project names*, *team names*, etc. Finally, the ACE

# Building a Brazilian Portuguese Parallel Corpus of Original and Simplified Texts

Helena M. Caseli, Tiago F. Pereira, Lucia Specia, Thiago A. S. Pardo,  
Caroline Gasperin and Sandra M. Aluisio

Center of Computational Linguistics (NILC)/ Department of Computer Sciences, University of  
São Paulo, Av. Trabalhador São-Carlense, 400. 13560-970 - São Carlos/SP, Brazil  
helenacaseli@dc.ufscar.br, tiagofrepereira@yahoo.com.br, lspecia@icmc.usp.br,  
taspardo@icmc.usp.br, cgasperin@icmc.usp.br, sandra@icmc.usp.br

**Abstract.** In this paper we address the problem of building the necessary tools and resources for performing Brazilian Portuguese text simplification. We describe our efforts on the design and development of: (a) a XCES-based annotation schema, (b) an annotation edition tool, and (c) a portal to access parallel corpora of original-simplified texts. These contributions were intended to (i) allow the creation and public release of a corpus of original and simplified texts with two different versions of simplification (called here *natural* and *strong*), targeting two levels of functional illiteracy and (ii) register simplification decisions during the creation of such corpus. We also provide an analysis of the first corpus created using the resources presented here: 104 newspaper texts and their simplified versions, produced by an expert in text simplification.

**Keywords:** Text Simplification, Brazilian Portuguese, annotation standards, annotation edition tool.

## 1 Introduction

In Brazil, “letramento” (literacy) is the term used to designate people's ability to use written language to obtain and record information, express themselves, plan and learn continuously [1]. In Brazil, according to the index used to measure the literacy level of the population (*INAF - National Indicator of Functional Literacy*), a vast number of people belong to the so called *rudimentary* and *basic* literacy levels. These people are able to find explicit information in short texts (rudimentary level) and also process slightly longer texts and make simple inferences (basic level).

The PorSimples project (*Simplificação Textual do Português para Inclusão e Acessibilidade Digital*)<sup>1</sup> aims at producing text simplification tools for promoting digital inclusion and accessibility for people with such levels of literacy, and possibly other kinds of reading disabilities. More specifically, the goal is to help these readers to process documents available on the web. Additionally, it could help children learning to read texts of different genres or adults being alphabetized. Two tools are

---

<sup>1</sup> <http://caravelas.icmc.usp.br/wiki/index.php/Principal>

# Synchronized Morphological and Syntactic Disambiguation for Arabic

Daoud Daoud

Princess Sumaya University for Technology  
Daoud@batelco.jo

**Abstract.** In this paper, we present a unique approach to disambiguation Arabic using a synchronized rule-based model. This approach helps in highly accurate analysis of sentences. The analysis produces a semantic net like structure expressed by means of Universal Networking Language (UNL)- a recently proposed interlingua. Extremely varied and complex phenomena of Arabic language have been addressed.

**Keywords:** Arabic Language, Synchronized Model, Disambiguation, UNL

## 1 Introduction

Compared to French or English, Arabic as an agglutinative and highly inflected language shows its proper types of difficulties in morphological disambiguation, since a large number of its ambiguities come from both the stemming and the categorization of a morpheme while most of ambiguities in French or English are related to the categorization of a morpheme only.

Phrases and sentences in Arabic have a relatively free word. The same grammatical relations can have different syntactic structures. Thus, morphological information is crucial in providing signs for structural dependencies.

Arabic sentences are characterized by a strong tendency for agreement between its constituents, between verb and noun, noun and objective, in matters of numbers, gender, definitiveness, case, person etc. These properties are expressed by a comprehensive system of affixation.

Arabic uses a diverse system of prefixes, suffixes, and pronouns that are attached to the words, creating compound forms that further complicate text manipulation. Simultaneously, Arabic exhibits a large-scale ambiguity already at the word level, which means that there are multiple ways in which a word can be categorized or broken down to its constituent morphemes. This is further complicated by the fact that most vocalization marks (diacritics) are omitted in Arabic texts.

However, the morphological analysis of a word-form, and in particular its morphological segmentation, cannot be disambiguated without reference to context,

# Parsing with Polymorphic Categorical Grammars

Matteo Capelletti<sup>1</sup> and Fabio Tamburini<sup>2</sup>

<sup>1</sup> Lix, École Polytechnique - France  
`matteo.capelletti@elisanet.fi`

<sup>2</sup> DSLO, University of Bologna - Italy  
`fabio.tamburini@unibo.it`

**Abstract.** In this paper we investigate the use of polymorphic categorical grammars as a model for parsing natural language. We will show that, despite the undecidability of the general model, a subclass of polymorphic categorical grammars, which we call *linear*, is mildly context-sensitive and we propose a polynomial parsing algorithm for these grammars.

## 1 Introduction

The simplest model of a categorical grammar is the so called Ajdukiewicz–Bar-Hillel calculus of [2] and [4]. Syntactic categories are formed from a given set of atoms as *functions*  $a/b$  and  $b\backslash a$ , with  $b$  and  $a$  categories. The intuitive meaning of a syntactic category of the form  $a/b$  (resp.  $b\backslash a$ ) is that it looks for an *argument* of category  $b$  to its right (resp. left) to give a category of type  $a$ . The resulting grammar system is known to be context-free.

Contemporary categorical grammars in the style of Ajdukiewicz–Bar-Hillel grammars are called *combinatory categorical grammars*, see [25]. Such systems adopt other forms of composition rules which enable them to generate non-context-free languages, see [29; 28]. The other main tradition of categorical grammar, the type-logical grammars of [20; 18], stemming from the work of [15], adopt special kinds of structural rules, that enable the system to generate non-context-free languages.

Both approaches increase the generative power of the basic system by adding special kinds of rules. In this paper, we adopt a different strategy which consists in keeping the elementary rule component of Ajdukiewicz–Bar-Hillel grammar and in introducing *polymorphic* categories, that is syntactic categories that contain variables ranging over categories. The inference process will be driven by unification, rather than by simple identity of formulas. We will see two kinds of polymorphic categorical grammars, one that is Turing complete and another, resulting from a restriction on the first, which is mildly context-sensitive. This second system, which is obviously the most interesting one for linguistics, has some important advantages with respect to the aforementioned ones. In respect to TLG, the polymorphic system we define is *polynomial*, as we will prove by providing a parsing algorithm. In respect to CCG (and most known TLG), our system is not affected by the so called *spurious ambiguity* problem, that is the problem of generating multiple, semantically equivalent, derivations.

# A CCG-based System for Valence Shifting for Sentiment Analysis

František Simančík and Mark Lee

School of Computer Science, University of Birmingham, Birmingham, UK  
frantisek.simancik@worc.ox.ac.uk, m.g.lee@cs.bham.ac.uk

**Abstract.** The automatic classification of sentiment in text is becoming an important area of research. In this work, we present a linguistic system for sentence-level valence annotation. Our system uses the formalism of Combinatory Categorical Grammar to represent words as functions acting on their syntactic arguments, which provides a unified way of implementing various classes of valence shifters. We propose two simple semi-automatic methods for estimating the valence of individual terms based on the lexical relations of WordNet. We evaluate the system on the data generated for the Affective Text task of SemEval 2007 and show that it compares favourably with the systems participating in the task.

**Keywords:** sentiment analysis, valence annotation, valence shifters, headlines, combinatory categorical grammar.

## 1 Introduction

The number of opinion-rich resources such as discussions, blogs and review sites has been growing rapidly in recent years. As a result of this, there is a demand for tools capable of classifying texts not only by the topic but also the attitude and opinion they convey; giving rise to new areas in Natural Language Processing called Opinion Mining and Sentiment Analysis.

One of the most prominent tasks in the field is the classification of valence (positive/negative orientation). Researchers (Pang et al. [7], Kennedy and Inkpen [6] and others) have successfully applied supervised machine learning methods<sup>1</sup> to determine the valence of longer texts. These approaches rely on the availability of a large amount of human-tagged training data and, compared to linguistic methods, reveal very little about the nature of the connection between a text and the opinion it expresses.

---

<sup>1</sup>Naive Bayes Classifiers, Support Vector Machines, etc.

# Classification of “Inheritance” Relations: a Semi-automatic Approach

Ekaterina Lapshinova-Koltunski

IMS, Universität Stuttgart  
Azenbergstr.12  
70174 Stuttgart  
`katerina@ims.uni-stuttgart.de`

**Abstract.** This study describes a semi-automatic approach to the classification of “inheritance” relations between morphologically related predicates.

Predicates, such as verbs and nouns subcategorizing for a subclause, are automatically extracted from text corpora and are classified according to their subcategorisation properties. For this purpose, we elaborate a semi-automatic knowledge-rich extraction and classification architecture. Our aim is also to compare subcategorisation properties of morphologically related predicates, i.e. verbs and deverbal nouns.

In this work, we concentrate exclusively on the predicates with sentential complements, such as *dass*, *ob* and *w*-clauses (that, if and wh-clauses) in German, although our methods can be applied for other complement types as well.

## 1 Introduction

This paper describes a semi-automatic approach to the analysis of subcategorisation properties of morphologically related predicates, such as verbs and nouns. We classify predicates according to their subcategorisation properties by means of extracting them from German corpora along with their complements. In this work, we concentrate exclusively on sentential complements, such as *dass*, *ob* and *w*-clauses, although our methods can be also applied for other types of complements.

It is usually assumed that subcategorisation properties of nominalisations are taken over from their underlying verbs. However, our preliminary tests show that there exist different types of relations between them. Thus, our aim is to review the properties of morphologically related words and to analyse the phenomenon of “inheritance” of subcategorisation properties.

For this purpose, we elaborate a set of semi-automatic procedures, with the help of which we not only classify extracted units according to their subcategorisation properties, but also compare the properties of verbs and their nominalisations. Our aim is to serve NLP, especially such large symbolic grammar for deep processing as HPSG or LFG, which need detailed subcategorisation data for their lexicons and grammars.

# Discovering Discourse Motifs in Instructional Dialog

Juan M. Huerta

IBM T.J. Watson Research Center  
1101 Kitchawan Road, Yorktown Heights, NY 10598, USA  
huerta@us.ibm.com

**Abstract.** We propose a method to analyze conversational interaction using discourse motifs (sequence of labels). We focus specifically on instructional transactive discourse. We first describe the characteristics of transactive discourse, its relationship to other frameworks of instructional discourse, and introduce a refined taxonomy of transactive discourse. Based on this new taxonomy, we construct a set of classifiers to automatically label instructional dialog segments. After labeling, we search for salient patterns of discourse common to these chains of labels using Multiple EM for Motif Elicitation and Gapped Local Analysis of Motifs (which are two techniques available for DNA and protein motif discovery). From our analysis of a corpus of classroom data, a set of Transactive-Participatory-Coherent motifs emerge. This approach to interaction-motif discovery and analysis can find application in dialog and discourse analysis, pedagogical domains (e.g., assessment and professional development), automatic tutoring systems, meeting analysis, problem solving, etc.

## 1 Introduction

We focus on the analysis of classroom discourse particularly when the focus is on solving mathematical problems. While the analysis of classroom discourse and mathematical problem solving is useful in providing pedagogical insight into teaching practices (see for example Huerta (2008), Blanton (2008)), its analysis can also shed light into interaction mechanisms used in more general collaborative problem solving.

Research in human dialog has been approached from various viewpoints using frameworks and methodologies of analysis that have been tailored to address the specific requirements of these viewpoints (examples of relatively recent perspectives to dialog analysis include Stolcke (2000), Stent (2000) among others, and a good summary can be found in Moore (2003)).

More problem-solving specific frameworks have also been proposed to analyze planning-oriented and instructional dialog in the classroom (Linden (1995)). Additionally, there have also been other efforts in the manual analysis of classroom interaction from purely pedagogical and sociological perspectives (Blanton (2008), Mehan (1985), Stark (2002), Haussman (2003)). There has been also work focusing on specific theoretical frameworks of interaction and the correlation of their elements to individual learning (e.g., Haussman (2006) and elaborative discourse, Meyer (2002) and scaffolding and self-regulation) as well as development of discourse frameworks

# Ontology Oriented Computation of English Verbs Metaphorical Trait

Zili Chen, Jonathan J. Webster, Ian I. Chow and Tianyong Hao

Department of Chinese, Translation and Linguistics, City University of Hong Kong,  
81 Tat Chee Avenue, Kowloon Tong, KLN, Hong Kong

**Abstract.** Research on metaphor has generally focused on exploring its context-dependent behavior and function. This current study aims to testify the postulate of English verb's innate trait of Metaphor Making potential. This paper intends to carry out an in-depth case study of a group of English core verbs using WordNet and SUMO ontology. In order to operationalize the assessment of an English verb's metaphor making potential, a refined algorithm has been developed, and a program made to realize the computation. At last, it is observed that higher frequency verbs generally possess greater metaphor making potential; while a verb's metaphor making potential on the other hand is also strongly influenced by its functional categories. As a preliminary context-free experiment with metaphor, this research foresees the possibility of providing an annotation schema for critical discourse analysis and a new parameter for scaling the difficulty level of reading comprehension of English texts.

**Keywords:** ontological computation, English verbs, MMP

## 1 Introduction and Previous Work

Metaphorical computation continues to remain a significant challenge to NLP. Recent researches of it mainly fall into two categories: rule-based approaches and statistical-based approaches. Up to now, some achievements have been attained, among which knowledge representation based methods are predominant [1]. These methods mainly employ knowledge representation based ontologies, such as The Suggested Upper Merged Ontology (SUMO), as their working mechanism. However, those researches are all limited to the study of metaphor's behavior and function in different contexts.

In line with Lakoff's view [2], "Metaphor allows us to understand one domain of experience in terms of another. This suggests that understanding takes place in terms of entire domains of experience and not in terms of isolated concepts", SUMO, an effort of the IEEE Standard Upper Ontology Working Group with the support of Teknowledge, contains terms chosen to cover all general domain concepts needed to represent world knowledge. Whereas Ahrens & Huang's research with SUMO and metaphor has focused on specific domain metaphors [3, 4], thus failing to make full use of SUMO's overall domain coverage.

Now that verb maintains the core for language processing, as believed by some



# Automatic Formal Verification of Conceptual Model Documentation by Means of Self-organizing Map

Algirdas Laukaitis and Olegas Vasilecas

Vilnius Gediminas Technical University , Sauletekio al. 11,  
LT-10223 Vilnius-40, Lithuania  
{algirdas.laukaitis,olegas}@fm.vgtu.lt

**Abstract.** By using background knowledge of the general and specific domains and by processing new natural language corpus experts are able to produce a conceptual model for some specific domain. In this paper we present a model that tries to capture some aspects of this conceptual modeling process. This model is functionally organized into two information processing streams: one reflects the process of formal concept lattice generation from domain conceptual model, and the another one reflects the process of formal concept lattice generation from the domain documentation. It is expected that similarity between those concept lattices reflects similarity between documentation and conceptual model. In addition to this process of documentation formal verification the set of natural language processing artifacts are created. Those artifacts then can be used for the development of information systems natural language interfaces. To demonstrate it, an experiment for the concepts identification from natural language queries is provided at the end of this paper.

**Key words:** Information systems engineering, formal concept analysis, IS documents self-organization, natural language processing.

## 1 Introduction

Software engineers spend hours in defining information systems (IS) requirements and finding common ground of understanding. The overwhelming majority of IS requirements are written in natural language supplemented with conceptual model and other semi-formal UML diagrams. The bridge in the form of semantic indexes between documents and conceptual model can be useful for more effective communication and model management. Then, an integration of the natural language processing (NLP) into information system documentation process is an important factor in meeting challenges for methods of modern software engineering.

Reusing natural language IS requirement specifications and compiling them into formal statements has been an old challenge [1], [14]. Kevin Ryan claimed that NLP is not mature enough to be used in requirements engineering [13] and

# Coreference Resolution using Markov Logic Network

Shujian HUANG<sup>1</sup>, Yabing ZHANG<sup>1</sup>, Junsheng ZHOU<sup>1,2</sup> and Jiajun CHEN<sup>1</sup>

<sup>1</sup> State Key Laboratory for Novel Software Technology,  
Nanjing University, Nanjing 210093, China

<sup>2</sup> Department of Computer Science, Nanjing Normal University,  
Nanjing, Jiangsu, 210097, China  
{huangsj, zhangyb, zhoujs, chenjj}@nlp.nju.edu.cn

**Abstract.** Most previous work treats the solution for pronouns and noun phrases either in two separate processes or in a single process. We argue that resolving them in two processes may result in the loss of potential useful information for each process. However, resolving them in a single process is also problematic. These two types of mentions have very different characteristics in some commonly used features. Current models cannot catch those differences and thus the two types may interfere with each other. In this paper, we propose a modeling strategy using Markov logic networks (MLNs) which can explicitly discriminate the two types in one single process. Experiments on ACE2005 Chinese dataset show that our modeling using MLNs, together with the correlation clustering technique, brings significant improvements to the task.

**Key words:** Coreference Resolution, Markov Logic Networks

## 1 Introduction

Coreference resolution (CR) has drawn a lot of attentions over the past decade, especially since McCarthy[1], Cardie and Wagstaff[2] introduced machine learning techniques into this field. It plays an important role in understanding complex texts and is widely used in a lot of applications such as question answering[3], summarization[4], etc. The strong relation with other popular topics such as entity resolution in database and citation analysis[5] makes it more attractive.

Pronoun and noun phrase are two major types of mentions in CR. There are two strategies for the resolution of them. One strategy tends to split the resolution of pronouns and noun phrases into two separate processes (*Separate Strategy*). Some works focus on just pronoun resolution, aiming to find the right antecedent for each pronoun[6–9]. Denis and Baldridge subdivide mentions into five categories such as third person pronouns, speech pronouns, etc. Then, specialized models are proposed for each individual type[10]. However, we argue that just considering pairwise relation between pronoun and each of its antecedent candidates does not make full use of the information among those candidate

© A. Gelbukh (Ed.)

*Advances in Computational Linguistics.*

*Research in Computing Science 41, 2009, pp. 157-168*

*Received 10/11/08*

*Accepted 15/12/08*

*Final version 04/02/09*

# A Ranking Approach to Persian Pronoun Resolution

Nafiseh Sadat Moosavi and Gholamreza Ghassem-Sani

Department of Computer Engineering, Sharif University of Technology, Iran  
n\_moosavi@ce.sharif.edu, sani@sharif.edu

**Abstract.** Coreference resolution is an essential step toward understanding discourses, and it is needed by many NLP tasks such as machine translation, question answering, and summarization. Pronoun resolution is a major and challenging subpart of coreference resolution, in which only the resolution of pronouns is considered. Classification approaches have been widely used for coreference/pronoun resolution, but it has been shown that ranking approaches outperform classification approaches in a variety of fields such as English pronoun resolution (Denis and Baldridge, 2007), question answering (Ravichandran, 2003), and tagging/parsing (Collins and Duffy, 2002; Charniak and Johnson, 2005). The strength of ranking is in its ability to consider all candidates at once and selecting the best one based on the model, while existing classification methods consider at most two candidate responses at a time. Persian and its varieties are spoken by more than 71 million people, and it has some characteristic that make parsing and other related processing of Persian more difficult than those of English. In this paper, we have evaluated maximum entropy ranker on Persian pronoun resolution and compared the results with that of four base classifiers.

**Keywords:** Natural Language Processing, Machine Learning, Ranking, Classification, Pronoun Resolution, Persian.

## 1 Introduction

The final goal of natural language processing (NLP) is that computers understand human languages. Different NLP research areas such as part of speech (POS) tagging, word sense disambiguation (WSD), and grammatical parsing concentrate only on a partial solution of this ultimate goal. All of these are required for a computer to understand a natural language.

NLP tasks can be divided into micro-tasks and macro-tasks. Micro-tasks focus on a word level processing or a sentence level processing such as WSD and parsing. On the other hand, macro-tasks include tasks which do a document level processing such as information retrieval and document classification. Before the introduction of machine learning approaches in NLP, higher level tasks such as semantic processing needed a variety of lower level tasks such as POS tagging and parsing. However, the use of machine learning methods may make it possible to obtain enough statistical in-

# Classification of Clinical Conditions: A Case Study on Prediction of Obesity and Its Co-morbidities

Archana Bhattarai, Vasile Rus and Dipankar Dasgupta

Department of Computer Science, The University of Memphis,  
209 Dunn Hall  
Memphis, TN 38152-3240, USA  
{abhattachar, vrus, dasgupta}@memphis.edu

**Abstract.** We investigate a multiclass, multilabel classification problem in medical domain in the context of prediction of obesity and its co-morbidities. Challenges of the problem not only lie in the issues of statistical learning such as high dimensionality, interdependence between multiple classes but also in the characteristics of the data itself. In particular, narrative medical reports are predominantly written in free text natural language which confronts the problem of predominant synonymy, hyponymy, negation and temporality. Our work explores the comparative evaluation of both traditional statistical learning based approach and information extraction based approach for the development of predictive computational models. In addition, we propose a scalable framework which combines both the statistical and extraction based methods with appropriate feature representation/selection strategy. The framework leads to reliable results in making correct classification. The framework was designed to participate in the second i2b2 Obesity Challenge.

**Keywords:** Text classification, Information Extraction, natural language processing

## 1 Introduction

Medical informatics is chiefly inductive and information intensive science where observation and analysis of comprehensive clinical data can lead to complex and powerful evidence based decision support systems. One of the primary goals of these automated systems is to make information more accessible, representative and meticulous in a quick span[4]. Furthermore, they have gained increased importance in the recent years as it can even outperform a human expert in some cases in diagnosing diseases as the process is highly subjective and fundamentally depends on the experiences of the assessor and his/her interpretation on the information[4]. Conversely, most medical institutions are still keeping a large amount of medical data in narrative form resulting in huge volume of potential information with limited or no utility and accessibility. An effort to exploit this data poses multiple challenges as it involves processing free text data with the presence of acronyms, synonyms, negation

# Analysis of Stemming Alternatives and Dependency Pattern Support in Text Classification

Levent Özgür and Tunga Güngör

Department of Computer Engineering, Boğaziçi University,  
Bebek, 34342 Istanbul, Turkey  
{ozgurlev, gungort}@boun.edu.tr

**Abstract.** In this paper, we study text classification algorithms by utilizing two concepts from Information Extraction discipline; dependency patterns and stemmer analysis. To the best of our knowledge, this is the first study to fully explore all possible dependency patterns during the formation of the solution vector in the Text Categorization problem. The benchmark of the classical approach in text classification is improved by the proposed method of pattern utilization. The test results show that support of four patterns achieves the highest ranks, namely, *participle modifier*, *adverbial clause modifier*, *conjunctive* and *possession modifier*. For the stemming process, we benefit from both morphological and syntactic stemming tools, Porter stemmer and Stanford Stemmer, respectively. One of the main contributions of this paper is its approach in stemmer utilization. Stemming is performed not only for the words but also for all the extracted pattern couples in the texts. Porter stemming is observed to be the optimal stemmer for all words while the raw form without stemming slightly outperforms the other approaches in pattern stemming. For the implementation of our algorithm, two formal datasets, Reuters - 21578 and National Science Foundation Abstracts, are used.

**Key words:** Text Classification, Dependency Patterns, Stemmer Analysis, Information Extraction

## 1 Introduction

Text Classification (TC) is a learning task, where pre-defined category labels are assigned to documents based on the likelihood suggested by a training set of labelled documents.

Most of the approaches used in this problem study it in bag-of-words (bow) form, where only the words in the text are analyzed by some machine learning algorithms for TC [1]. In this approach, documents are represented by the widely used vector-space model, introduced by Salton et al. [2]. In this model, each document is represented as a vector  $d$ . Each dimension in the vector  $d$  stands for a distinct term (word) in the term space of the document collection.

© A. Gelbukh (Ed.)  
Advances in Computational Linguistics.  
Research in Computing Science 41, 2009, pp. 195-206

Received 17/11/08  
Accepted 15/12/08  
Final version 03/02/09

# Investigating Variations in Adjective Use across Different Text Categories

Jing Cao<sup>1</sup> and Alex Chengyu Fang<sup>2</sup>

Department of Chinese, Translation and Linguistics  
City University of Hong Kong  
Hong Kong SAR, China

<sup>1</sup>cjing3@student.cityu.edu.hk, <sup>2</sup>acfang@cityu.edu.hk

**Abstract.** Adjectives are an informative but understudied linguistic entity with good potentials in sentiment analysis, text classification and automatic genre detection. In this article, we report an investigation of the variations in adjective use across different text categories represented in a sizable corpus. In particular, we report the distribution of adjectives across a range of categories grouped together as academic prose in the British National Corpus. We shall measure inter-category similarity in the use of adjectives and demonstrate with empirical data that adjectives are an effective differentia of text categories or domains, at least in terms of arts and sciences as the two major sub-categories within academic prose.

**Key Words:** corpus, text category, adjective, similarity, BNC

## 1 Introduction

Adjectives are an informative but understudied linguistic entity [1, 2], drawing more and more attention within the research community. Focus has been mostly on the semantic aspect of adjectives for practical research in sentiment analysis applicable to automatic evaluations of email communication [3], blogs [4] and customer reviews in [5]. Studies in this respect typically focus on evaluative adjectives [6] and size adjectives [7]. In addition to the semantic approach, adjectives are also used for purposes of text categorization and genre detection in [8]. In this respect, [2] and [9] have generally shown with corpus evidence that adjectives occur more often in written texts than in spoken ones, and more frequently in informative writing than in imaginative writing. According to [8], ‘the literature suggests that adjectives and adverbs will vary by genre because of their unique patterns of usage in text’ (p. 4).

This paper describes one of the recent attempts to study adjectives from the perspectives of text categorization and genre detection. In particular, we investigate the variations of adjective use across various types of academic writing selected from a large-sized corpus. We attempt to ascertain whether adjective-based indices will be able to classify texts in such a way that conforms to manual classification. As we shall show in this article with empirical data, adjectives do differ by text categories and therefore appear to be an important differentia of text categories. More importantly,

# Multi-Category Support Vector Machines for Identifying Arabic Topics

Mourad Abbas, Kamel Smaili and Daoud Berkani

CRSTDLA, Speech Processing Lab.  
1 rue D. E. Alafghani, Algeria  
INRIA-LORIA, Parole team  
B.P. 101, Villers les Nancy, France  
Polytechnic School, Signal and Communication Lab.  
10 rue H. Badi, Algeria  
m\_abbas04@yahoo.fr, kamel.smaili@loria.fr, dberkani@enp.edu.dz

**Abstract.** It is known that Support Vector Machines were designed for binary classification. Nevertheless, it would be fruitful to extend this operation to what is called Multi-category classification. That is why Multi-category Support Vector Machines (MSVM) become nowadays the current subject of several serious researches, aiming to achieve high levels of multi-category classification tasks. This technique has been assessed recently in some fields as text categorization, Cancer classification, etc. We should notify that experiments which have been realized until now using MSVM are limited to small data sets, since its computation is more expensive. In this paper we are interested in the use of this method, for the first time in topic identification. The experiments conducted concern topic identification of Arabic language. The corpora are extracted from Alwatan newspaper. Achieved results lead to an improvement of MSVM performance in comparison to the baseline SVM method. Nevertheless, SVM still outperforms MSVM when using larger sizes of the vocabulary.

## 1 Introduction

The main objective of topic identification is to assign one or several topic labels to a flow of textual data. Labels are chosen from a set of topics fixed a priori. Talking about topics conduct us to clarify the definition of a topic. In [1], each keyword is considered as a topic. Whereas in other works, topics are more sophisticated corresponding to specific subject, for example politics and sports [2]. In our case, we are dealing with six topics: Culture, Religion, Economy, Local news, International news and sports.

Topic identification is used in several areas: to adapt language models for speech recognition and for machine translation, to focus on a specific use for search engines,...etc. In spontaneous speech recognition process the vocabulary has to be as large as possible. Enlarging the vocabulary increases the search space and consequently could reduce system's performance.

A language model is one of the knowledge sources which is used by a speech

# USUM: Update Summary Generation System

C Ravindranath Chowdary and P Sreenivasa Kumar

Department of Computer Science and Engineering  
Indian Institute of Technology Madras  
Chennai 600 036, India  
`{chowdary,psk}@cse.iitm.ac.in`

**Abstract.** *Huge amount of information is present in the World Wide Web and a large amount is being added to it frequently. A query-specific summary of multiple documents is very helpful to the user in this context. Currently, few systems have been proposed for query-specific, extractive multi-document summarization. If a summary is available for a set of documents on a given query and if a new document is added to the corpus, generating an updated summary from the scratch is time consuming and many a times it is not practical/possible. In this paper we propose a solution to this problem. This is especially useful in a scenario where the source documents are not accessible. We cleverly embed the sentences of the current summary into the new document and then perform query-specific summary generation on that document. Our experimental results show that the performance of the proposed approach is good in terms of both quality and efficiency.*

## 1 Introduction

Currently, the World Wide Web is the largest source of information. Huge amount of data is present on the Web and large amount of data is added to the web constantly. Often the information pertaining to a topic is present across several web pages. It is a tedious task for the user to go through all these documents as the number of documents available on a topic will range from tens to thousands. It will be of great help for the user if a query specific multi-document summary is generated. Summary generation can be broadly divided as abstractive and extractive. In abstractive summary generation, the abstract of the document is generated. The summary so formed need not have exact sentences as present in the document. In extractive summary generation, important sentences are extracted from the document. The generated summary contains all such extracted sentences arranged in a meaningful order. In this paper, generated summaries are extractive. Summary can be generated either on a single document or on several documents. In multi-document summary generation, other issues like time, ordering of extracted sentences, scalability etc. will arise.

Summary can be either generic or query specific. In generic summary generation, the important sentences from the document are extracted and the sentences

© A. Gelbukh (Ed.)

*Advances in Computational Linguistics.*

*Research in Computing Science 41, 2009, pp. 229-240*

*Received 10/11/08*

*Accepted 15/12/08*

*Final version 05/02/09*



# Towards the Speech Synthesis of Raramuri: A Unit Selection Approach based on Unsupervised Extraction of Suffix Sequences

Alfonso Medina Urrea,<sup>1</sup> José Abel Herrera Camacho<sup>2</sup>  
and Maribel Alvarado García<sup>3</sup>

<sup>1</sup> GIL-II, Universidad Nacional Autónoma de México  
04510 Coyoacán, DF, MEXICO  
`amedinau@ii.unam.mx`

<sup>2</sup> FI, Universidad Nacional Autónoma de México  
04510 Coyoacán, DF, MEXICO  
`abelh@verona.fi-p.unam.mx`

<sup>3</sup> Escuela Nacional de Antropología e Historia  
14030 Tlalpan, DF, MEXICO  
`marvarado1978@yahoo.com.mx`

**Abstract.** This work deals with the design of a synthesis system to provide an audio database for Raramuri or Tarahumara, a Yuto-Nahua language spoken in Northern Mexico. In order to achieve the most natural speech possible, the synthesis system is proposed which uses a unit selection approach based on function words, suffix sequences (derivational and inflectional morphemes) and diphones of the language. In essence, the unknown suffix units were extracted from a corpus and recorded, along diphones and function words, in order to build the audio database that provides data for Text-to-Speech synthesis.

## 1 Introduction

The ultimate objective of Text-to-Speech (TTS) synthesis systems is to create applications which listeners, and users in general, cannot easily determine whether the speech he or she is hearing comes from a human or a synthesizer.

Synthesized speech can be produced by concatenating recorded units (waveforms) selected from a large, single-speaker speech database. The primary motivation for using a database with a large number of units that covers wider prosodic and spectral characteristics, gives us the great benefit to produce a synthesized speech that sounds more natural than those produced by systems that use a small set of controlled units (*e.g.* diphones) [1]. There is a paradigm for achieving high-quality synthesis that uses a large corpus of recorded speech units; it is called *unit-selection synthesis*. Unit selection is a method in which we can concatenate waveforms from different linguistic structures such as sentences, words, syllables, triphones, diphones and phones. Due to the increasing computer's storage capacity, we are able to create a corpus of prerecorded

# Pronunciation Rules in Portuguese Regional Speech (PORT REG) for Coarticulation Process

Sara Candeias<sup>1</sup> and Jorge Morais Barbosa<sup>2</sup>

<sup>1</sup> Instituto de Telecomunicações, Department of Computers and Electrical Engineering,  
University of Coimbra, PORTUGAL

<sup>2</sup> Departement of Portuguese Language, Faculty of Letters, University of Coimbra,  
PORTUGAL

saracandeias@co.it.pt, jbarbosa@ci.uc.pt

**Abstract.** This paper describes one aspect of an ongoing work to incorporate pronunciation variability in the Portuguese (PORT) speech system. This work focuses on the linguistic rules to improve the grapheme-(multi)phone transcription algorithm that will be implemented. Portuguese ‘Beira Interior’ regional speech (PORT-BI REG) is considered to be in the realm of coarticulation (post-lexical) phenomena. A set of linguistic rules for most of the common vowel transformation in an utterance (vocalic segments at both the left and right edges of the word) is presented. The analysis focuses on the distinctive features that originate vowel sound challenges in connected speech. The results are interesting from the point of view of setting up models to reconstruct a grapheme-phone transcription algorithm for Portuguese multi-pronunciation speech systems. We propose that the linguistic documentation of Portuguese minority speech can be an optimal start for Portuguese speech system development process, too.

**Keywords:** Text-to-Speech; coarticulation (phonology); structural analysis (linguistic features); pronunciation instruction (phonetic).

## 1 Introduction

Several frameworks have been proposed for the grapheme-to-phone transcription module for Portuguese language, such as [2, 3, 12]. However, the problem with the Portuguese regional speech under development is the shortage of speech and text corpora. This is one of the reasons why their linguistic structure has been very poorly investigated, especially at linguistic levels such as phonetics. The applications of the Portuguese speech system are mainly based on standard Portuguese language and on isolated word recognition. It is well known that the sequence of phones spoken by a human speaker is not the same sequence as that which derives from the phonetic transcription of a word in isolation. Coarticulation (post-lexical) rules must be included in the course of phonetic transcription. In order to obtain a more natural speech, these rules must be applied to varying sequences of phones. Several methods can be used to elicit grapheme-to-phoneme rules from pre-existing lexicons. However, these automatic techniques do not cope very well with the concurrent multi-

# A Maximum Entropy (ME) Based Translation Model for Chinese Characters Conversion

Fai Wong, Sam Chao, Cheong Cheong Hao and Ka Seng Leong

Faculty of Science and Technology of University of Macau,  
Av. Padre Tomás Pereira S.J., Taipa, Macao  
{derekfw, lidiase}@umac.mo

**Abstract.** As the growth of exchange activities between four regions of cross strait, the problem to correctly convert between Traditional Chinese (TC) and Simplified Chinese (SC) is getting important and attention from many people, especially in business organizations and translation companies. Different from the approaches of many conventional code conversion systems, which rely on various levels of human constructed knowledge (from character set to semantic level) to facilitate the translation purpose, this paper proposes a Chinese conversion model based on Maximum Entropy (ME), a Machine Learning (ML) technique. This approach uses tagged corpus as the only information source for creating the conversion model. The constructed model is evaluated with selected ambiguous characters to investigate the recall rate as well as the conversion accuracy. The experiment results show that the proposed model is comparable to the state of the art conversion system.

**Keywords:** Maximum Entropy, Machine Learning, Natural Language Processing, Chinese translation, Traditional Chinese, Simplified Chinese.

## 1 Introduction

Modern Chinese typically involves two main dialects of writing, Traditional Chinese (TC) and Simplified Chinese (SC). In Chinese computing, these two systems adapt different coding schema for the computer to process the corresponding Chinese information. Traditional Chinese uses Big5 encoding while Simplified Chinese uses GB. For a Simplified Chinese document to be opened and read in a computer with Traditional Chinese operating system, conversion from Simplified Chinese encoding system into Traditional Chinese encoding is necessary for the purpose that the document can be further processed under the Traditional Chinese computer environment, and vice versa. As addressed by Wang [1] in the meeting of the 4<sup>th</sup> Chinese Digitization Forum, although there are many conversion systems implemented and available in the market, neither one of them can produce the conversion result with satisfaction. Reviewing the nature of this problem, Simplified Chinese is actually a simpler version of Traditional Chinese. It differs in two ways from the Traditional writing system: 1) a reduction of the number of strokes per character and 2) the reduction of the number of characters in use that is two different

© A. Gelbukh (Ed.)  
*Advances in Computational Linguistics.*  
*Research in Computing Science 41, 2009, pp. 267-276*

*Received 06/11/08*  
*Accepted 11/12/08*  
*Final version 02/02/09*

# Tracking Out-of-date Newspaper Articles

Frederik Cailliau<sup>1,2</sup>, Aude Giraudel<sup>3</sup> and Béatrice Arnulphy<sup>4</sup>

<sup>1</sup> Sinequa Labs, 12 Rue d'Athènes, F-75009 Paris

<sup>2</sup> LIPN-Univ. de Paris-Nord, 99 av. Jean-Baptiste Clément, F-93430 Villetaneuse

<sup>3</sup> Syllabs, 15 rue Jean-Baptiste Berlier, F-75013 Paris

<sup>4</sup> LIMS-CNRS, BP 133, F-91403 Orsay\*

cailliau@sinequa.com giraudel@syllabs.com beatrice.arnulphy@limsi.fr

**Abstract.** Local newspapers rely on their local correspondents to bring you the freshest news of your home town. Some of the articles written by these correspondents are not published immediately, but are put aside to be placed in a later edition. One of the challenges the editor in chief is confronted with, is to publish only up-to-date information about a local event. In this paper we present a system that tracks out-of-date newspaper articles to prevent their publishing. It firstly dates the event the article is talking about. The date detection grammar is written in a single but complex finite state automaton based on linguistic pattern matching. Secondly, it computes an absolute date for each relative date. A freshness score can then be deduced from the time difference between the extracted and the publication date. The system has been tested on several French local newspaper corpora. Baseline for the temporal extraction is our standard date extraction that was developed for general purposes.

**Keywords:** extraction of temporal information, named entities, events, information retrieval, newspaper articles

## 1 Introduction

Local daily newspapers rely on their network of correspondents to cover local events. Once revised by a journalist, these texts are able to be published as newspaper articles. Everyday the editor in chief validates or assembles the pages that will be published in the next edition. One of the challenges he is confronted with, is to validate only those articles covering hot news or coming events and to replace out-of-date articles by more recent ones. Unfortunately there is no easy or quick way to perform this task. To judge whether the article talks about a recent or a coming event, the editor has to read most of the article, since the day of writing is not a reliable indicator, when given. Even if most of the articles may be short, this activity is too time-expensive to be executed in the short delay before going to press. The editor's

---

\* Some of the work in this article was done in collaboration with Béatrice Arnulphy during her internship at Sinequa in the summer of 2006 as a graduate student of the French university *Université de Provence*.

# PtiClic: A Game for Vocabulary Assessment combining JeuxDeMots and LSA

Mathieu Lafourcade<sup>1</sup> and Virginie Zampa<sup>2</sup>

<sup>1</sup> LIRMM-TAL, Univ. Montpellier 2 France  
mathieu.lafourcade@lirmm.fr

<sup>2</sup> LIDILEM-DIP Univ. Grenoble 3 France  
virginie.zampa@u-grenoble3.fr

**Abstract.** One interesting path for designing software for vocabulary acquisition and assessment could be games. In this article we present PtiClic, a lexical game based on the principles behind JeuxDeMots and combined with LSA. Such a game can foster the interest of young people in developing lexical skills in either a tutored or an open environment.

**Keywords:** Lexical acquisition, lexical network, serious games, LSA, JeuxDeMots

## 1 Introduction

Developing software for vocabulary acquisition and/or assessment in general, and for young people in particular, is a risky business. There are several difficulties. First, we have to be able to design an activity that can foster an interest in learning which is not easy in the case of children. Then, the underlying dictionaries or lexical databases can prove tremendously hard to develop, especially if we want to go beyond just a list of words with parts-of-speech and venture into the realm of lexical functions and the relations intertwining terms.

Automated acquisition of lexical or functional relations between terms is necessary in a large number of tasks in Natural Language Processing (NLP) outscoping largely Technology-Enhanced Learning of language. These relations that we find generally in thesauruses or ontologies can of course be revealed in a manual way; for example, one of the oldest thesauruses is Roget's, its current version being (Kipfer, 2001), or the most famous lexical network is Wordnet (Miller, 1990). Such relations can be also determined computationally from corpora of texts, for example (Robertson and Spark Jones, 1976), (Lapata and Keller, 2005), or (Landauer et al 1998) in which statistical studies on the distributions of words are made. Moreover, many applications of NLP require information of various natures, like synonymy or antonymy, but also relations of hyperonymy / hyponymy, holonymy / meronymy etc. The building of such relations, when done manually by experts, requires resources

## Author Index

Abbas, Mourad	217	Lafourcade, Mathieu	289
Aluisio, Sandra M.	59	Lapshinova-Koltunski, E.	109
Alvarado García, Maribel	243	Laukaitis, Algirdas	143
Andreadakis, George	47	Lee, Mark	99
Arnulphy, Béatrice	277	Marrero, Mónica	47
Berkani, Daoud	217	Medina Urrea, Alfonso	243
Bhattarai, Archana	183	Mitrofanova, Olga	27
Brun, Caroline	35	Morais Barbosa, Jorge	257
Cailliau, Frederik	277	Morato Lara, Jorge	47
Candeias, Sara	257	Özgür, Levent	195
Cao, Jing	207	Pardo, Thiago A. S.	59
Capelletti, Matteo	87	Pereira, Tiago F.	59
Caseli, Helena M.	59	Pomikálek, Jan	3
Chao, Sam	267	Ravindranath Chowdary, C	229
Chen, Jiajun	157	Rinaldi, Fabio	15
Chen, Zili	135	Rus, Vasile	183
Chengyu Fang, Alex	207	Rychlý, Pavel	3
Cheong Hao, Cheong	267	Sadat Moosavi, Nafiseh	169
Chow, Ian I.	135	Sánchez-Cuadrado, Sonia	47
Daoud, Daoud	73	Schneider, Gerold	15
Dasgupta, Dipankar	183	Seng Leong, Ka	267
Gasperin, Caroline	59	Simančík, František	99
Ghassem-Sani, Gholamreza	169	Smaili, Kamel	217
Giraudel, Aude	277	Specia, Lucia	59
Güngör, Tunga	195	Sreenivasa Kumar, P	229
Hagège, Caroline	35	Tamburini, Fabio	87
Hao, Tianyong	135	Vasilecas, Olegas	143
Herrera Camacho, José A.	243	Webster, Jonathan J.	135
Huang, Shujian	157	Wong, Fai	267
Huerta, Juan M.	123	Zampa, Virginie	289
Kaljurand, Kaarel	15	Zhang, Yabing	157
Kappeler, Thomas	15	Zhou, Junsheng	157
Kilgarrieff, Adam	3		