

ACL 2007



# ACL 2007

---

## The LAW Proceedings of The Linguistic Annotation Workshop

June 28-29, 2007  
Prague, Czech Republic

---



Production and Manufacturing by  
*Omnipress*  
2600 Anderson Street  
Madison, WI 53704  
USA

©2007 Association for Computational Linguistics

Order copies of this and other ACL proceedings from:

Association for Computational Linguistics (ACL)  
209 N. Eighth Street  
Stroudsburg, PA 18360  
USA  
Tel: +1-570-476-8006  
Fax: +1-570-476-0860  
[acl@aclweb.org](mailto:acl@aclweb.org)

## Preface

Welcome to The Linguistic Annotation Workshop (The LAW).

Linguistically annotated corpora play a major role in parsing, information extraction, question answering, machine translation and many other areas of computational linguistics, and provide an empirical testbed for theoretical linguistics research. This has led to a proliferation of annotation systems, frameworks, formats, and schemes. Recognition of the need to harmonize annotation practices and frameworks has become increasingly critical, as witnessed by numerous workshops dealing with different aspects of linguistic annotation over the past few years.

The LAW addresses all aspects of linguistic annotation in a single forum by merging two existing workshop series: NLPXML (Natural Language Processing and XML) and FLAC (Frontiers in Linguistically Annotated Corpora). The goals of the workshop include:

1. The exchange and propagation of research results with respect to the annotation, manipulation and exploitation of corpora, taking into account different applications and theoretical investigations in the field of language technology and research;
2. Working towards harmonization and interoperability from the perspective of the increasingly large number of tools and frameworks that support the creation, instantiation, manipulation, querying, and exploitation of annotated resources;
3. Working towards a consensus on all issues crucial to the advancement of the field of corpus annotation.

These proceedings include 11 long papers, 5 short papers, 4 demo descriptions and 8 posters selected by the program committee from 51 submissions for presentation at the workshop. In addition to these presentations, the workshop includes demonstrations of annotation tools, reports by working groups, and an open discussion session.

We would like to thank the members of the program committee for their timely reviews. We also thank the Workshops Chair and other organizers of ACL-2007 for their support. Finally, we congratulate Adriane Boyd, the winner of the Innovative Student Annotation Award for the paper *Discontinuity Revisited: An Improved Conversion to Context-Free Representations*.

Branimir Boguraev  
Nancy Ide  
Adam Meyers  
Shigeko Nariyama  
Manfred Stede  
Janyce Wiebe  
Graham Wilcock



# Organizers

## Workshop Chairs

Branimir Boguraev, IBM T. J. Watson Research Center  
Nancy Ide, Vassar College  
Adam Meyers, New York University  
Shigeko Nariyama, University of Melbourne  
Manfred Stede, University of Potsdam  
Janyce Wiebe, University of Pittsburgh  
Graham Wilcock, University of Helsinki

## Program Committee

David Ahn, University of Amsterdam  
Lars Ahrenberg, Linköping University  
Timothy Baldwin, University of Melbourne  
Francis Bond, NICT  
Kalina Bontcheva, University of Sheffield  
Paul Buitelaar, DFKI  
Jean Carletta, University of Edinburgh  
Key-Sun Choi, KAIST  
Christopher Cieri, Linguistic Data Consortium/University of Pennsylvania  
Hamish Cunningham, University of Sheffield  
David Day, MITRE Corporation  
Thierry Declerck, DFKI  
Ludovic Denoyer, LIP6 - University of Paris 6  
Tomaz Erjavec, Jozef Stefan Institute  
David Farwell, New Mexico State University  
Alex Chengyu Fang, City University of Hong Kong  
Chuck Fillmore, International Computer Science Institute, Berkeley  
Anette Frank, DFKI  
John Fry, San Jose State University  
Claire Grover, University of Edinburgh  
Jan Hajic, Charles University  
Ed Hovy, University of Southern California  
Baden Hughes, University of Melbourne  
Emi Izumi, NICT  
Tsai Jia-Lin, Tung Nan Institute of Technology  
Aravind Joshi, University of Pennsylvania  
Ewan Klein, University of Edinburgh  
Mounia Lalmas, University of London  
Mike Maxwell, University of Maryland  
Chieko Nakabasami, Toyo University  
Stephan Oepen, University of Oslo  
Kyonghee Paik, KLI

Martha Palmer, University of Colorado  
Antonio Pareja-Lora, Universidad Complutense de Madrid / OEG - UPM  
Manfred Pinkal, Saarland University  
James Pustejovsky, Brandeis University  
Owen Rambow, Columbia University  
Laurent Romary, MPG-INRIA  
Henry Thompson, University of Edinburgh  
Erik Tjong Kim Sang, University of Amsterdam  
Theresa Wilson, University of Edinburgh  
Nainwen Xue, University of Colorado

## Table of Contents

<i>GrAF: A Graph-based Format for Linguistic Annotations</i> Nancy Ide and Keith Suderman .....	1
<i>Efficient Annotation with the Jena ANnotation Environment (JANE)</i> Katrin Tomanek, Joachim Wermter and Udo Hahn .....	9
<i>Mining Syntactically Annotated Corpora with XQuery</i> Gosse Bouma and Geert Kloosterman .....	17
<i>Associating Facial Displays with Syntactic Constituents for Generation</i> Mary Ellen Foster .....	25
<i>An Annotation Type System for a Data-Driven NLP Pipeline</i> Udo Hahn, Ekaterina Buyko, Katrin Tomanek, Scott Piao, John McNaught, Yoshimasa Tsuruoka and Sophia Ananiadou .....	33
<i>Discontinuity Revisited: An Improved Conversion to Context-Free Representations</i> Adriane Boyd .....	41
<i>Usage of XSL Stylesheets for the Annotation of the Sámi Language Corpora.</i> Saara Huhmarniemi, Sjur N. Moshagen and Trond Trosterud .....	45
<i>Criteria for the Manual Grouping of Verb Senses</i> Cecily Jill Duffield, Jena D. Hwang, Susan Windisch Brown, Dmitriy Dligach, Sarah E. Vieweg, Jenny Davis and Martha Palmer .....	49
<i>Semi-Automated Named Entity Annotation</i> Kuzman Ganchev, Fernando Pereira, Mark Mandel, Steven Carroll and Peter White .....	53
<i>Querying Multimodal Annotation: A Concordancer for GeM</i> Martin Thomas .....	57
<i>Annotating Chinese Collocations with Multi Information</i> Ruifeng Xu, Qin Lu, Kam-Fai Wong and Wenjie Li .....	61
<i>Computing Translation Units and Quantifying Parallelism in Parallel Dependency Treebanks</i> Matthias Buch-Kromann .....	69
<i>Adding Semantic Role Annotation to a Corpus of Written Dutch</i> Paola Monachesi, Gerwert Stevens and Jantine Trapman .....	77
<i>A Search Tool for Parallel Treebanks</i> Martin Volk, Joakim Lundborg and Maël Mettler .....	85
<i>Annotating Expressions of Appraisal in English</i> Jonathon Read, David Hope and John Carroll .....	93

<i>Active Learning for Part-of-Speech Tagging: Accelerating Corpus Annotation</i>	
Eric Ringger, Peter McClanahan, Robbie Haertel, George Busby, Marc Carmen, James Carroll, Kevin Seppi and Deryle Lonsdale .....	101
<i>Combining Independent Syntactic and Semantic Annotation Schemes</i>	
Marc Verhagen, Amber Stubbs and James Pustejovsky .....	109
<i>XARA: An XML- and Rule-based Semantic Role Labeler</i>	
Gerwert Stevens .....	113
<i>ITU Treebank Annotation Tool</i>	
Gülşen Eryiğit .....	117
<i>Two Tools for Creating and Visualizing Sub-sentential Alignments of Parallel Text</i>	
Ulrich Germann .....	121
<i>Building Chinese Sense Annotated Corpus with the Help of Software Tools</i>	
Yunfang Wu, Peng Jin, Tao Guo and Shiwen Yu .....	125
<i>Annotating a Japanese Text Corpus with Predicate-Argument and Coreference Relations</i>	
Ryu Iida, Mamoru Komachi, Kentaro Inui and Yuji Matsumoto .....	132
<i>Web-based Annotation of Anaphoric Relations and Lexical Chains</i>	
Maik Stührenberg, Daniela Goecke, Nils Diewald, Alexander Mehler and Irene Cramer .....	140
<i>Standoff Coordination for Multi-Tool Annotation in a Dialogue Corpus</i>	
Kepa Joseba Rodríguez, Stefanie Dipper, Michael Götze, Massimo Poesio, Giuseppe Riccardi, Christian Raymond and Joanna Rąbiega-Wiśniewska .....	148
<i>PoCoS - Potsdam Coreference Scheme</i>	
Olga Krasavina and Christian Chiarcos .....	156
<i>Multiple-step Treebank Conversion: From Dependency to Penn Format</i>	
Cristina Bosco .....	164
<i>Experiments with an Annotation Scheme for a Knowledge-rich Noun Phrase Interpretation System</i>	
Roxana Girju .....	168
<i>IGT-XML: An XML Format for Interlinearized Glossed Text</i>	
Alexis Palmer and Katrin Erk .....	176
<i>Shared Corpora Working Group Report</i>	
Adam Meyers, Nancy Ide, Ludovic Denoyer and Yusuke Shinyama .....	184
<i>Discourse Annotation Working Group Report</i>	
Manfred Stede, Janyce Wiebe, Eva Hajicova, Brian Reese, Simone Teufel, Bonnie Webber and Theresa Wilson .....	191



# Workshop Program

**Thursday, June 28, 2007**

## **Session 1: Introduction and Long Papers**

- 14:30–14:55 Introduction to the Workshop
- 14:55–15:20 *GrAF: A Graph-based Format for Linguistic Annotations*  
Nancy Ide and Keith Suderman
- 15:20–15:45 *Efficient Annotation with the Jena ANnotation Environment (JANE)*  
Katrin Tomanek, Joachim Wermter and Udo Hahn
- 15:45–16:15 Break

## **Session 2: Long Papers**

- 16:15–16:40 *Mining Syntactically Annotated Corpora with XQuery*  
Gosse Bouma and Geert Kloosterman
- 16:40–17:05 *Associating Facial Displays with Syntactic Constituents for Generation*  
Mary Ellen Foster
- 17:05–17:30 *An Annotation Type System for a Data-Driven NLP Pipeline*  
Udo Hahn, Ekaterina Buyko, Katrin Tomanek, Scott Piao, John McNaught,  
Yoshimasa Tsuruoka and Sophia Ananiadou

## **Demonstration and Poster Session**

- 17:30–18:30 Demonstrations and Posters  
(Listed at end of program)

**Friday, June 29, 2007**

**Session 3: Working Groups**

- 09:30–10:00 Shared Corpora Working Group Report  
Adam Meyers, Nancy Ide, Ludovic Denoyer and Yusuke Shinyama
- 10:00–10:45 Panel Session on Discourse Annotation  
Manfred Stede, Eva Hajicova, Brian Reese, Simone Teufel, Bonnie Webber and  
Theresa Wilson
- 10:45–11:15 Break

**Session 4: Short Papers**

- 11:15–11:30 *Discontinuity Revisited: An Improved Conversion to Context-Free Representations*  
Adriane Boyd
- 11:30–11:45 *Usage of XSL Stylesheets for the Annotation of the Sámi Language Corpora.*  
Saara Huhmarniemi, Sjur N. Moshagen and Trond Trosterud
- 11:45–12:00 *Criteria for the Manual Grouping of Verb Senses*  
Cecily Jill Duffield, Jena D. Hwang, Susan Windisch Brown, Dmitriy Dligach,  
Sarah E. Vieweg, Jenny Davis and Martha Palmer
- 12:00–12:15 *Semi-Automated Named Entity Annotation*  
Kuzman Ganchev, Fernando Pereira, Mark Mandel, Steven Carroll and Peter White
- 12:15–12:30 *Querying Multimodal Annotation: A Concordancer for GeM*  
Martin Thomas
- 12:30–14:30 Lunch

**Friday, June 29, 2007 (continued)**

**Session 5: Long Papers**

- 14:30–14:55 *Annotating Chinese Collocations with Multi Information*  
Ruifeng Xu, Qin Lu, Kam-Fai Wong and Wenjie Li
- 14:55–15:20 *Computing Translation Units and Quantifying Parallelism in Parallel Dependency Treebanks*  
Matthias Buch-Kromann
- 15:20–15:45 *Adding Semantic Role Annotation to a Corpus of Written Dutch*  
Paola Monachesi, Gerwert Stevens and Jantine Trapman
- 15:45–16:15 Break

**Session 6: Long Papers**

- 16:15–16:40 *A Search Tool for Parallel Treebanks*  
Martin Volk, Joakim Lundborg and Maël Mettler
- 16:40–17:05 *Annotating Expressions of Appraisal in English*  
Jonathon Read, David Hope and John Carroll
- 17:05–17:30 *Active Learning for Part-of-Speech Tagging: Accelerating Corpus Annotation*  
Eric Ringger, Peter McClanahan, Robbie Haertel, George Busby, Marc Carmen, James Carroll, Kevin Seppi and Deryle Lonsdale

**Discussion Session**

- 17:30–18:30 Open Discussion

## Demonstrations and Posters

### Demonstrations

*Combining Independent Syntactic and Semantic Annotation Schemes*

Marc Verhagen, Amber Stubbs and James Pustejovsky

*XARA: An XML- and Rule-based Semantic Role Labeler*

Gerwert Stevens

*ITU Treebank Annotation Tool*

Gülşen Eryiğit

*Two Tools for Creating and Visualizing Sub-sentential Alignments of Parallel Text*

Ulrich Germann

### Posters

*Building Chinese Sense Annotated Corpus with the Help of Software Tools*

Yunfang Wu, Peng Jin, Tao Guo and Shiwen Yu

*Annotating a Japanese Text Corpus with Predicate-Argument and Coreference Relations*

Ryu Iida, Mamoru Komachi, Kentaro Inui and Yuji Matsumoto

*Web-based Annotation of Anaphoric Relations and Lexical Chains*

Maik Stührenberg, Daniela Goecke, Nils Diewald, Alexander Mehler and Irene Cramer

*Standoff Coordination for Multi-Tool Annotation in a Dialogue Corpus*

Kepa Joseba Rodríguez, Stefanie Dipper, Michael Götze, Massimo Poesio, Giuseppe Riccardi, Christian Raymond and Joanna Rabiega-Wiśniewska

*PoCoS - Potsdam Coreference Scheme*

Olga Krasavina and Christian Chiarcos

*Multiple-step Treebank Conversion: From Dependency to Penn Format*

Cristina Bosco

**Demonstrations and Posters (continued)**

*Experiments with an Annotation Scheme for a Knowledge-rich Noun Phrase Interpretation System*

Roxana Girju

*IGT-XML: An XML Format for Interlinearized Glossed Text*

Alexis Palmer and Katrin Erk

