

Introduction to Text Technology: Using XML in Natural Language Processing

Georg Rehm

Justus-Liebig-Universität Gießen
Institut für Germanistik
Angewandte Sprachwissenschaft und Computerlinguistik
Otto-Behaghel-Straße 10 D
35394 Gießen

Telefon: +49 641 99 29052
Fax: +49 641 99 29059

Georg.Rehm@uni-giessen.de
<http://www.uni-giessen.de/~g91063/>
<http://www.uni-giessen.de/germanistik/ascl/>

September 10, 2003

Abstract

This document covers the basic structure of the course “Introduction to Text Technology: Using XML in Natural Language Processing”, the topics of the individual sessions, requirements for a graded Schein, the most important websites and essential reading.

Introduction to Text Technology – Course Overview

- Week 1: 09/22/2003—09/26/2003
 1. **Introduction:**
 - Text Technology: Definitions and Related Fields.
 - A Visual Introduction to XML.
 - The History of Markup Languages.
 - Example of a Text Technology application.
 2. **Markup Languages I:**
 - HTML Crash Course.
 - Example: Rumelhart's Story Grammar as an XML Application.
 - Introduction to the Syntax of XML-DTDs – Part 1: Elements.
 3. **Markup Languages II:**
 - Introduction to the Syntax of XML-DTDs – Part 2: Attributes and Entities.
 - Linguistic Background of Document Grammars.
 - How to write Document Type Definitions.
 - Parsing XML Document Instances.
 - Basic XML Software: XML Parsers and XML Editors.
 4. **Practical Exercise I:**
 - Building and parsing DTDs and XML Document Instances.
 5. **Doing it with Style I:**
 - Cascading Style Sheets.
 - Introduction to XSL.
 - XPath as a Prerequisite for XSLT.
- Week 2: 09/29/2003—10/03/2003
 6. **Doing it with Style II:**
 - XSLT.
 7. **Practical Exercise II:**
 - Handling of a standalone XSLT processor.
 - Implementation of XSLT stylesheets
 8. **Related Standards:**
 - Beyond DTDs: XML Schema.
 - XML and Related Standards.
 9. **Glueing it All Together: Text Technology Applications**
 - Text Technology Applications.
 - Programming XML interfaces: SAX and DOM.
 10. **Coming Full Circle:**
 - The Future: The Semantic Web.
 - Topics we didn't (and couldn't) cover.
 - Disadvantages of XML.

Requirements for a Schein

In order to receive a graded Schein for this course, you can choose one of the following options. Regular attendance and especially practical work during the two exercises is mandatory.

1. If you already happen to

- be in the planning or even writing phase of the final thesis for the course you are reading and if your thesis could be related to Text Technology in one way or another, or if you
- work in an applied NLP- or Computational Linguistics research project:

Write a paper which introduces the topic of your thesis or the project you work on and describe its potential relations to Text Technology. Where could XML-related formalisms or languages be useful? Would it be possible to implement proprietary data structures in a Text Technology way? You should make the relations as concrete as possible, i. e., you should include, e. g., example data, XSLT stylesheets, architectures, Document Type Definitions or code. The paper should not exceed 20 pages.

Example: If your thesis deals with natural language parsing, you should discuss the representation of parse trees as XML document instances where the hierarchical structure of the parse tree is reflected by the XML structure of the instance and you could develop XML markup languages for the representation of the lexicon and the grammar.

2. If you have a background in programming:

Develop a piece of software that utilizes Text Technology methods and write a paper as its documentation. The function of this software is open to negotiation. The paper should not exceed 20 pages.

Example: Since the mid-90ies, corpora have become a crucial aspect of almost any NLP- or CL-related project, especially with regard to the evaluation of natural language software. Choose a language that some NL-software should process, then choose a website containing HTML documents of this language which could be processed by the abovementioned NL-software. Implement a robot and a wrapper which (a) automatically connect to this website, (b) identify webpages of interest, (c) wrap them up in simple XML code, (d) perform basic tokenization and (e) store these XML document instances in a corpus for further processing. You could even try to embed this corpus-building-process into a larger architecture, using NL-software already available at your university or department (taggers, parsers etc.). Another aspect would include the implementation of word frequency statistics, simple text summarization methods (e. g., with graphical representations based on XSLT transformations) and so on.

3. If you do not have a background in programming:

Write a paper which discusses Text Technology methods and approaches in concrete research projects and software prototypes, which process natural language input on *one* of the following levels of language description or fields of application respectively: (a) morphology, (b) syntax/grammar, (c) semantics, (d) pragmatics, (e) natural language generation, (f) summarization, (g) web document analysis, or (h) information retrieval. Limit yourself to at most three different approaches and describe their individual advantages and disadvantages. Suggest potential improvements of the discussed approaches. The paper should not exceed 20 pages.

Important Standards, Books and Papers

Standards: World Wide Web Consortium

- ADLER, SHARON; BERGLUND, ANDERS; CARUSO, JEFF; DEACH, STEPHEN; GRAHAM, TONY; GROSSO, PAUL; GUTENTAG, EDUARDO; MILOWSKI, ALEX; PARNELL, SCOTT; RICHMAN, JEREMY AND ZILLES, STEVE (2001): "Extensible Stylesheet Language (XSL) 1.0". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xsl1/>.
- ALTHEIM, MURRAY; BOUMPHREY, FRANK; DOOLEY, SAM; MCCARRON, SHANE; SCHNITZENBAUMER, SEBASTIAN AND WUGOFSKI, TED (2001): "Modularization of XHTML". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xhtml-modularization/>.
- ALTHEIM, MURRAY AND MCCARRON, SHANE (2001): "XHTML 1.1 – Module-based XHTML". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xhtml11/>.
- AYARS, JEFF; BULTERMAN, DICK; COHEN, AARON; DAY, KEN; HODGE, ERIK; HOSCHKA, PHILIPP; HYPHE, ERIC; JORDAN, MURIEL; KIM, MICHELLE; KUBOTA, KENICHI; LANPIER, ROB; LAYAÏDA, NABIL; MICHEL, THIERRY; NEWMAN, DEBBIE; VAN OSSENBRUGGEN, JACCO; RUTLEDGE, LLOYD; SACCOCIO, BRIDIE; SCHMITZ, PATRICK AND TEN KATE, WERNER (2001): "Synchronized Multimedia Integration Language (SMIL 2.0)". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/smil20/>.
- BAKER, MARK; ISHIKAWA, MASAYASU; MATSUI, SHINICHI; STARK, PETER; WUGOFSKI, TED AND YAMAKAMI, TOSHIHIKO (2000): "XHTML Basic". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xhtml-basic/>.
- BERGLUND, ANDERS; BOAG, SCOTT; CHAMBERLIN, DON; FERNANDEZ, MARY F.; KAY, MICHAEL; ROBIE, JONATHAN AND SIMÉON, JÉRÔME (2002): "XML Path Language (XPath) – Version 2.0". Technische Spezifikation (Working Draft), W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xpath20/>.
- BOS, BERT; LIE, HÅKON WIUM; LILLEY, CHRIS AND JACOBS, IAN (1998): "Cascading Style Sheets, Level 2 – CSS2 Specification". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/REC-CSS2/>.
- BOX, DON; EHNEBUSKE, DAVID; KAKIVAYA, GOPAL; LAYMAN, ANDREW; NIELSEN, NOAH MENDELSON HENRIK FRYSTYK; THATTE, SATISH AND WINER, DAVE (2000): "Simple Object Access Protocol (SOAP) 1.1". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/SOAP/>.
- BRAY, TIM; HOLLANDER, DAVE AND LAYMAN, ANDREW (1999): "Namespaces in XML". Technische Spezifikation, W3C (World Wide Web Consortium).
- BRAY, TIM; PAOLI, JEAN; SPERBERG-MCQUEEN, C. M. AND MALER, EVE (2000): "Extensible Markup Language (XML) 1.0 (Second Edition)". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/2000/REC-xml-20001006/>.
- BRICKLEY, DAN AND GUHA, R.V. (2003): "RDF Vocabulary Description Language 1.0: RDF Schema". Technische Spezifikation (Working Draft), W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/rdf-schema/>.
- CHRISTENSEN, ERIK; CURBERA, FRANCISCO; MEREDITH, GREG AND WEERAWARANA, SANJIVA (2001): "Web Services Description Language (WSDL) 1.1". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/wsdl>.
- CLARK, JAMES (1999): "XSL Transformations (Version 1.0)". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xslt/>.

- CLARK, JAMES AND DEROSE, STEVE (1999): "XML Path Language (XPath) – Version 1.0". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xslt/>.
- DEROSE, STEVE; MALER, EVE AND ORCHARD, DAVID (2001): "XML Linking Language (XLink) Version 1.0". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xlink/>.
- FALLSIDE, DAVID C.; THOMPSON, HENRY S.; BEECH, DAVID; MALONEY, MURRAY; MENDELSON, NOAH; BIRON, PAUL V. AND MALHOTRA, ASHOK (2001): "XML Schema". Technische Spezifikation, W3C (World Wide Web Consortium). W3C Recommendation. Besteht aus Part 0 (Primer), Part 1 (Structures), Part 2 (Datatypes). Online verfügbar: <http://www.w3.org/XML/Schema>.
- FERRAILOLO, JON; FUJISAWA, JUN AND JACKSON, DEAN (2003): "Scalable Vector Graphics (SVG) 1.1 Specification". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/SVG11/>.
- GROSSO, PAUL; MALER, EVE; MARSH, JONATHAN AND WALSH, NORMAN (2003): "XPointer Framework". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xptr-framework/>.
- LASSILA, ORA AND SWICK, RALPH R. (1999): "Resource Description Framework (RDF). Model and Syntax Specification". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/REC-rdf-syntax/>.
- MITRA, NILO; GUDGIN, MARTIN; HADLEY, MARC; MENDELSON, NOAH; MOREAU, JEAN-JACQUES AND NIELSEN, HENRIK FRYSTYK (2002): "SOAP Version 1.2". Technische Spezifikation, W3C (World Wide Web Consortium). W3C Recommendation. Besteht aus Part 0 (Primer), Part 1 (Messaging Framework), Part 2 (Adjuncts). Online verfügbar: <http://www.w3.org/2002/ws/>.
- PEMBERTON, STEVEN (2000): "XHTML 1.0: The Extensible Hypertext Markup Language". Technische Spezifikation, W3C (World Wide Web Consortium).
- PEMBERTON, STEVEN (2002): "XHTML 1.0: The Extensible Hypertext Markup Language (Second Edition)". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/xhtml1/>.
- RAGGETT, DAVE; HORS, ARNAUD LE AND JACOBOS, IAN (1999): "HTML 4.01 Specification". Technische Spezifikation, W3C (World Wide Web Consortium). Online verfügbar: <http://www.w3.org/TR/html401/>.

Standards: ISO

- ISO 10179 (1996): "Information Processing – Processing Languages – Document Style Semantics and Specification Language (DSSSL)". Internationaler Standard, International Organization for Standardization, Genf. Online verfügbar: <http://www.ornl.gov/sgml/wg8/>.
- ISO 10744 (1997): "Information Processing – Hypermedia/Time-Based Structuring Language (HyTime) – Second Edition". Internationaler Standard, International Organization for Standardization, Genf. Online verfügbar: <http://www.ornl.gov/sgml/wg8/>.
- ISO 15445 (2000): "ISO/IEC 15445: Information Technology – Document Description and Processing Languages – HyperText Markup Language (HTML)". Internationaler Standard, International Organization for Standardization, Genf. Online verfügbar: <http://pur1.org/NET/ISO+IEC.15445/15445.html>.
- ISO 639 (1998): "ISO/IEC 639: Codes for the representation of names of languages. Part 1 (1988), Part 2 (1998)". Internationaler Standard, International Organization for Standardization, Genf.

ISO 8879 (1986): “Information Processing – Text and Office Information Systems – Standard Generalized Markup Language”. Internationaler Standard, International Organization for Standardization, Genf.

ISO 8879 TC2 (1998): “ISO/IEC JTC1/SC34 N0029 – Document Description and Processing Languages – Technical Corrigendum 2 to ISO 8879:1986 (Annex K: WebSGML Adaptations; Annex L: Declaration for XML)”. Anhang/Revision eines internationalen Standards, International Organization for Standardization, Genf. Online verfügbar: <http://www.sgmlsource.com/8879/n0029.htm>.

ISO 9070 (1991): “ISO/IEC 9070: Information Technology – SGML Support Facilities – Registration Procedures for Public Text Owner Identifiers”. Internationaler Standard, International Organization for Standardization, Genf.

ISO/IEC 13250 (2000): “Information Technology – Document Description and Processing Languages – Topic Maps”. Internationaler Standard, International Organization for Standardization, Genf. Online verfügbar: <http://www.ornl.gov/sgml/wg4/>.

Standards: Requests for Comments – RFCs

RFC 2376 (1998): “XML Media Types”. Network Working Group – Request for Comments (RFC). E. James Whitehead und Makoto Murata. Online verfügbar: <http://www.ietf.org/rfc/>.

RFC 2396 (1998): “Uniform Resource Identifiers (URI): Generic Syntax”. Network Working Group – Request for Comments (RFC). Tim Berners-Lee, Roy Fielding und Larry Masinter. Online verfügbar: <http://www.ietf.org/rfc/>.

RFC 2413 (1999): “Dublin Core Metadata for Resource Discovery”. Network Working Group – Request for Comments (RFC). Stuart L. Weibel, John A. Kunze, Carl Lagoze und Misha Wolf. Online verfügbar: <http://www.ietf.org/rfc/>.

RFC 2616 (1999): “Hypertext Transfer Protocol – HTTP/1.1”. Network Working Group – Request for Comments (RFC). Roy T. Fielding, James Gettys, Jeffrey C. Mogul, Henrik Frystyk Nielsen, Larry Masinter, Paul J. Leach und Tim Berners-Lee. Online verfügbar: <http://www.ietf.org/rfc/>.

RFC 2731 (1999): “Encoding Dublin Core Metadata in HTML”. Network Working Group – Request for Comments (RFC). John A. Kunze. Online verfügbar: <http://www.ietf.org/rfc/>.

RFC 2854 (2000): “The ‘text/html’ Media Type”. Network Working Group – Request for Comments (RFC). Dan Connolly und Larry Masinter. Online verfügbar: <http://www.ietf.org/rfc/>.

Text Technological Applications

CHEN, LI-QUN; XIE, XING; MA, WEI-YING; ZHANG, HONG-JIANG; ZHOU, HEQIN AND FENG, HUANQING (2003): “DRESS: A Slicing Tree Based Web Page Representation for Various Display Sizes”. In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.

COLING (1994): *COLING 94 – The 15th International Conference on Computational Linguistics*, Kyoto, Japan. Association for Computational Linguistics.

DCES (2003): “Dublin Core Metadata Element Set, Version 1.1: Reference Description”. Dublin Core Metadata Initiative, DCMI. Online verfügbar: <http://dublincore.org/documents/2003/02/04/dces/>.

DCMT (2003): “DCMI Metadata Terms”. Dublin Core Metadata Initiative, DCMI. Online verfügbar: <http://dublincore.org/documents/dcmi-terms/>.

DCTV (2003): “DCMI Type Vocabulary”. Dublin Core Metadata Initiative, DCMI. Online verfügbar: <http://dublincore.org/documents/dcmi-type-vocabulary/>.

- DUNLOP, DOMINIC (1995): "Practical Considerations in the Use of TEI Headers in a Large Corpus". *Computers and the Humanities* (29): pp. 85–98.
- FUJII, ATSUSHI AND ISHIKAWA, TETSUYA (2000): "Utilizing the World Wide Web as an Encyclopedia: Extracting Term Descriptions from Semi-Structured Texts". In: *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL-EACL 2000)*. pp. 488–495.
- FUJII, ATSUSHI AND ISHIKAWA, TETSUYA (2001): "Organizing Encyclopedic Knowledge based on the Web and its Application to Question Answering". In: *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics (ACL-EACL 2001)*. pp. 196–203.
- GUPTA, SUHIT; KAISER, GAIL; NEISTADT, DAVID AND GRIMM, PETER (2003): "DOM-based Content Extraction of HTML Documents". In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- IDE, NANCY (1994): "Encoding Standards for Large Text Resources: The Text Encoding Initiative". In: Coling (1994), pp. 574–578.
- KIM, SUN AND ZHANG, BYOUNG-TAK (2000): "Web-Document Retrieval by Genetic Learning of Importance Factors for HTML Tags". In: *Proceedings of the International Workshop on Text and Web Mining (PRICAI 2000)*, edited by Tan, Ah-Hwee and Yu, Philip S. Melbourne, pp. 13–23.
- KLABUNDE, RALF; CARSTENSEN, KAI-UWE; EBERT, CHRISTIAN; ENDRISS, CORNELIA; JEKAT, SUSANNE; LANGER, HAGEN AND SCHIEHLEN, MICHAEL (editors) (2003): *Computerlinguistik und Sprachtechnologie – Eine Einführung*. Heidelberg: Spektrum, 2nd edition. Erscheint.
- LABROU, YANNIS AND FININ, TIM (1999): "Yahoo! as an Ontology – Using Yahoo! Categories to Describe Documents". In: *Proceedings of the 8th International Conference on Information and Knowledge Management (CIKM '99)*. ACM Press, pp. 180–187.
- LEECH, GEOFFREY; GARSIDE, ROGER AND BRYANT, MICHAEL (1994): "CLAWS4: The Tagging of the British National Corpus". In: Coling (1994), pp. 622–628.
- LOBIN, HENNING (1998): "Intelligente Dokumente – Inhaltliche Strukturierung und flexible Hypertextualisierung von Multimedia-Dokumenten". In: *Forschung an der Universität Bielefeld*, Bielefeld: Universität Bielefeld, pp. 35–40.
- LOBIN, HENNING (1999a): "Intelligente Dokumente – Linguistische Repräsentation komplexer Inhalte für die hypermediale Wissensvermittlung". In: Lobin (1999b), pp. 155–178.
- LOBIN, HENNING (editor) (1999b): *Text im digitalen Medium – Linguistische Aspekte von Textdesign, Texttechnologie und Hypertext Engineering*. Wiesbaden: Westdeutscher Verlag.
- LOBIN, HENNING (2000): *Informationsmodellierung in XML und SGML*. Berlin, Heidelberg, New York etc.: Springer.
- LOBIN, HENNING (editor) (2001): *Sprach- und Texttechnologie in digitalen Medien – Proceedings der Frühjahrstagung der Gesellschaft für Linguistische Datenverarbeitung*, Justus-Liebig-Universität Gießen. Gesellschaft für linguistische Datenverarbeitung.
- LOBIN, HENNING AND LEMNITZER, LOTHAR (editors) (2003): *Texttechnologie – Anwendungen und Perspektiven*. Tübingen: Stauffenburg. Erscheint.
- LOBIN, HENNING AND STÜHRENBERG, MAIK (2003): "XML-strukturierte Learning Objects". In: *Sprache zwischen Theorie und Technologie – Festschrift für Wolf Papprotz zum 60. Geburtstag*, edited by Cyrus, Lea; Feddes, Hendrik; Schumacher, Frank and Steiner, Petra, Wiesbaden: Deutscher Universitäts-Verlag, Sprachwissenschaft, pp. 185–198.

- MALER, EVE AND ANDALOUSSI, JEANNE EL (1996): *Developing SGML DTDs – From Text to Model to Markup*. Upper Saddle River: Prentice Hall.
- MARTELLI, SILVIA; CARROLL, JEREMY J. AND SIGNORE, ORESTE (2003): “Syntax for Semantic Enriching of Web Pages”. In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- MARTIN, PHILIPPE AND EKLUND, PETER (1999): “Embedding Knowledge in Web Documents”. *Computer Networks* (31): pp. 1403–1419.
- McKELVIE, DAVID; BREW, CHRIS AND THOMPSON, HENRY (1997): “Using SGML as a Basis for Data-Intensive NLP”. In: *Proceedings of Applied Natural Language Processing (ANLP) 97*. Association for Computational Linguistics, Washington D. C. Online verfügbar: <http://www.ltg.hcr.c.ed.ac.uk/~dmck/anlp97.ps>.
- MEHLER, ALEXANDER AND LOBIN, HENNING (editors) (2003): *Automatische Textanalyse – Systeme und Methoden zur Annotation und Analyse natürlichsprachlicher Texte*. Wiesbaden: Westdeutscher Verlag. Erscheint.
- MIKHEEV, ANDREI; GROVER, CLAIRE AND MOENS, MARC (1999): “XML Tools and Architecture for Named Entity Recognition”. *Markup Languages* 1 (3): pp. 89–113.
- MINNING, HOLGER; PUSKÁS, CSILLA AND SALISBURY, JUSTIN (2002): “Rhetorisches Parsing deutschsprachiger Texte”. In: *Proceedings of the 12th Student Conference on Computational Linguistics (TaCoS 2002)*, edited by Reitter, David. Potsdam. Online verfügbar: <http://www.ling.uni-potsdam.de/tacos/proceedings/>.
- MÖHR, WIEBKE AND SCHMIDT, INGRID (editors) (1999): *SGML und XML – Anwendungen und Perspektiven*. Berlin, Heidelberg, New York etc.: Springer.
- MYLLYMAKI, JUSSI (2001): “Effective Web Data Extraction with Standard XML Technologies”. In: *Proceedings of the 10th International World Wide Web Conference (WWW-10)*. Hong Kong, pp. 689–696.
- NAGAO, KATASHI AND HASIDA, KÔITI (1998): “Automatic Text Summarization Based on the Global Document Annotation”. In: *COLING 98 – The 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*. Association for Computational Linguistics, Montreal, Quebec, Kanada, volume 2, pp. 917–921. 2 Bände.
- PEPPER, STEVE AND MOORE, GRAHAM (2001): “XML Topic Maps (XTM) 1.0”. Technische Spezifikation, TopicMaps.Org. Online verfügbar: <http://www.topicmaps.org/xtm/1.0/>.
- POTOK, THOMAS E.; ELMORE, MARK T.; REED, JOEL W. AND SAMATOVA, NAGIZA F. (2002): “An Ontology-based HTML to XML Conversion Using Intelligent Agents”. In: *Proceedings of the 35th Hawaii International Conference on System Sciences (HICSS-35)*. Big Island, Hawaii: IEEE Computer Society.
- REHM, GEORG (1998): *Vorüberlegungen zur automatischen Zusammenfassung deutschsprachiger Texte mittels einer SGML- und DSSSL-basierten Repräsentation von RST-Relationen*. Magisterarbeit, Studiengang Computerlinguistik und Künstliche Intelligenz, Universität Osnabrück. Online verfügbar: <http://www.uni-giessen.de/~g91063/papers.shtml>.
- REHM, GEORG (2003a): “Das World Wide Web”. In: Klabunde et al. (2003). Erscheint.
- REHM, GEORG (2003b): “Texttechnologie und das World Wide Web – Anwendungen und Perspektiven”. In: Lobin and Lemnitzer (2003), pp. 433–464. Erscheint.
- REHM, GEORG (2003c): “Texttechnologische Grundlagen”. In: Klabunde et al. (2003). Erscheint.
- REHM, GEORG AND REINSCH, MARKUS (2001): “Die Chronik der “Chronik” – Über die Konvertierung und Weiterverarbeitung proprietär annotierter Daten”. In: *Sprach- und Texttechnologie in digitalen Medien – Proceedings der Frühjahrstagung der Gesellschaft für Linguistische Datenverarbeitung*, edited by Lobin, Henning. Gesellschaft für linguistische Datenverarbeitung, pp. 211–221.

- SAHUGUET, ARNAUD AND AZAVANT, FABIEN (1999): "Looking at the Web through XML Glasses". In: *Proceedings of the 4th International Conference on Cooperative Information Systems (CoopIS '99)*. Edinburgh: IEEE Computer Society Press, pp. 148–159.
- SEO, HEEKYOUNG AND CHOI, JAEYOUNG YANG JOONGMIN (2001): "Knowledge-based Wrapper Generation by Using XML". In: *Proceedings of the IJCAI-2001 Workshop on Adaptive Text Extraction and Mining*. Seattle.
- SIGLETOS, GEORGIOS; FARMAKIOTOU, DIMITRA; STAMATAKIS, KOSTAS; PALIOURAS, GEORGIOS AND KARKALETSIS, VANGELIS (2003): "Annotating Web pages for the needs of Web Information Extraction applications". In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- SPERBERG-MCQUEEN, C. M. AND BURNARD, LOU (editors) (2002): *TEI P4: Guidelines for Electronic Text Encoding and Interchange*. Text Encoding Initiative Consortium; Humanities Computing Unit, University of Oxford.
- SPERBERG-MCQUEEN, C. M. AND GOLDSTEIN, ROBERT F. (1994): "HTML to the Max – A Manifesto for Adding SGML Intelligence to the World-Wide Web". In: *Proceedings of the Second International WWW Conference – Mosaic and the Web*. Chicago. Online verfügbar: <http://www.uic.edu/~cmsmcq/htmlmax.html>.
- THIERRY, DECLERCK AND WITTENBURG, PETER (2001): "Xml and NLP: Their Interaction and their Role for HLT Applications". In: *Proceedings of the First NLP and XML Workshop*. Tokyo.
- ULE, TYLMAN AND HINRICH, ERHARD (2003): "Linguistische Annotation". In: Lobin and Lemnitzer (2003), pp. 217–243. Erscheint.
- ULE, TYLMAN AND MÜLLER, FRANK HENRIK (2003): "KaRoPars: Ein System zur linguistischen Annotation großer Text-Korpora des Deutschen". In: Mehler and Lobin (2003). Erscheint.
- VOLK, MARTIN (1998): "Markup of a Test Suite with SGML". In: *Linguistic Databases*, edited by Nerbonne, John, Cambridge: Cambridge University Press, number 77 in CSLI Lecture Notes, pp. 59–76. Online verfügbar: <http://www.ifi.unizh.ch/CL/volk/papers/SGMLMarkup.ps.gz>.
- VOLK, MARTIN (2000): "Scaling up: Using the WWW to resolve PP attachments and ambiguities". In: *KONVENS-2000 / Sprachkommunikation, Vorträge der gemeinsamen Veranstaltung 5. Konferenz zur Verarbeitung natürlicher Sprache*, edited by Zühlke, Werner and Schukat-Talamazzini, Ernst Günter. Ilmenau: VDE Verlag, pp. 151–155.
- VOLK, MARTIN (2001): "Exploiting the WWW as a corpus to resolve PP attachment ambiguities". In: *Proceedings of the Corpus Linguistics 2001 Conference*, edited by Rayson, Paul; Wilson, Andrew; McEnery, Tony; Hardie, Andrew and Khoja, Shereen. Lancaster, pp. 601–606.
- VOLK, MARTIN (2002): "Using the Web as Corpus for Linguistic Research". In: *Tähendusepüüdja. Catcher of the Meaning. A Festschrift for Professor Haldur Öim*, Publications of the Department of General Linguistics 3, University of Tartu. Online verfügbar: <http://www.ifi.unizh.ch/CL/volk/publications.html>.
- WITT, ANDREAS (1999): "SGML und Linguistik". In: Lobin (1999b), pp. 121–154.
- WITT, ANDREAS (2003): "Linguistische Informationsmodellierung mit XML". In: Mehler and Lobin (2003). Erscheint.
- WOLFF, CHRISTIAN (2003): "Systemarchitekturen – Aufbau texttechnologischer Anwendungen". In: Lobin and Lemnitzer (2003), pp. 165–192. Erscheint.

Corpora and the Web

- CAMPOS, JOÃO AND SILVA, MÁRIO J. (2001): "Versus: A Model for a Web Repository". In: *CRC'01 – 4ª Conferência de Redes de Computadores*. Covilhã.
- CHUNG, CHRISTINA YIP; GERTZ, MICHAEL AND SUNDARESAN, NEEL (2001): "Quixote: Building XML Repositories from Topic Specific Web Documents". In: *Fourth International Workshop on the Web and Databases (WebDB 2001)*, edited by Giansalvatore Mecca, Jérôme Siméon. Santa Barbara, pp. 103–108.
- COWIE, JIM; LUDOVIK, EVGENY AND ZACHARSKI, RON (1998): "An Autonomous, Web-based, Multilingual Corpus Collection Tool". In: *Proceedings of the International Conference on Natural Language Processing and Industrial Applications*. Moncton, pp. 142–148. Online verfügbar: <http://cr1.nmsu.edu/~raz/langrec/nlpia.htm>.
- DEWE, JOHAN; KARLGRÉN, JUSSI AND BRETAN, IVAN (1998): "Assembling a Balanced Corpus from the Internet". In: *Proceedings of the 11th Nordic Conference of Computational Linguistics*. Copenhagen, pp. 100–107.
- FAIRON, CÉDRICK (2000): "GlossaNet: Parsing a Web Site as a Corpus". *Linguisticae Investigationes* 22 (2): pp. 327–340.
- FLETCHER, WILLIAM H. (2001): "Concordancing the Web with KWICFinder". In: *Proceedings of the 3rd North American Symposium on Corpus Linguistics and Language Teaching*. Boston.
- GREFENSTETTE, GREGORY (1999): "The World Wide Web as a resource for example-based machine translation tasks". In: *Proceedings of the ASLIB Conference on Translating and the Computer*. London.
- HIRAI, JUN; RAGHAVAN, SRIRAM; GARCIA-MOLINA, HECTOR AND PAEPCKE, ANDREAS (2000): "WebBase: A Repository of Web Pages". In: *Proceedings of the 9th International World Wide Web Conference*. Amsterdam, pp. 277–293.
- IDE, NANCY AND VÉRONIS, JEAN (1994): "MULTEXT: Multilingual Text Tools and Corpora". In: *COLING 94 – The 15th International Conference on Computational Linguistics*. Association for Computational Linguistics, Kyoto, Japan, volume 1, pp. 588–592.
- JONES, ROSIE AND GHANI, RAYID (2000): "Automatically Building a Corpus for a Minority Language from the Web". In: *Proceedings of the Student Workshop at the 38th Annual Meeting of the Association for Computational Linguistics (ACL-2000)*. Hong Kong. Online verfügbar: <http://www.cs.cmu.edu/~webkb/>.
- KEHOE, ANDREW AND RENOUF, ANTOINETTE (2002): "WebCorp: Applying the Web to Linguistics and Linguistics to the Web". In: *Proceedings of the 11th International World Wide Web Conference (WWW 2002)*. Honolulu.
- KILGARRIFF, ADAM (2001): "Web as Corpus". In: *Proceedings of the Corpus Linguistics 2001 Conference*, edited by Rayson, Paul; Wilson, Andrew; McEnery, Tony; Hardie, Andrew and Khoja, Shereen. Lancaster, pp. 342–344.
- LI, FANG; SHENG, HUANYE AND WEISWEBER, WILHELM (2001): "World Wide Web – A Multilingual Language Resource". In: *Web Intelligence: Research and Development*, edited by Zhong, Ning; Yao, Yiyu; Liu, Jiming and Ohsuga, Setsuo. Berlin, Heidelberg, New York etc.: Springer, number 2198 in Lecture Notes in Artificial Intelligence, pp. 373–378.
- MORLEY, BARRY; RENOUF, ANTOINETTE AND KEHOE, ANDREW (2003): "Linguistic Research with XML/RDF-aware WebCorp tool". In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.

REHM, GEORG (2001): "korpust.html – Zur Sammlung, Datenbank-basierten Erfassung, Annotation und Auswertung von HTML-Dokumenten". In: *Proceedings of the GLDV Spring Meeting 2001*, edited by Lobin, Henning. Gesellschaft für linguistische Datenverarbeitung (Society for Computational Linguistics and Language Technology), Giessen, Germany, pp. 93–103. Online verfügbar: <http://www.uni-giessen.de/fb09/ascl/gldv2001/>.

RESNIK, PHILIP AND SMITH, NOAH A. (2002): "The Web as a Parallel Corpus". Technical Report UMIACS-TR-2002-61, University of Maryland. Online verfügbar: <http://www.umiacs.umd.edu/~resnik/pubs.html>.

Language Identification and the Web

GREFFENSTETTE, GREGORY AND NIOCHE, JULIEN (2000): "Estimation of English and non-English Language Use on the WWW". In: *Proceedings of RIAO'2000: Content-Based Multimedia Information Access*. Paris, pp. 237–246. Online verfügbar: <http://de.arXiv.org/abs/cs.CL/0006032>.

LANGER, STEFAN (2002): "Grenzen der Sprachenidentifizierung". In: *KONVENS 2002 – 6. Konferenz zur Verarbeitung natürlicher Sprache*, edited by Busemann, Stefan. Saarbrücken: Deutsches Forschungszentrum für Künstliche Intelligenz, pp. 99–106. DFKI Document D-02-01.

LAVOIE, BRIAN F. AND O'NEILL, EDWARD T. (1999): "How "World Wide" is the Web? Trends in the Internationalization of Web Sites". In: *Annual Review of OCLC Research 1999*, Dublin: OCLC Online Computer Library Center. Online verfügbar: <http://www.oclc.org>.

MUTHUSAMY, YESHWANT K. AND SPITZ, LAWRENCE (1998): "Automatic Language Identification". In: *Survey of the State of the Art in Human Language Technology*, edited by Cole, Ronald; Mariani, Joseph; Uszkoreit, Hans; Varile, Giovanni Battista; Zaenen, Annie and Zampolli, Antonio, Cambridge: Cambridge University Press, pp. 314–317. Online verfügbar: <http://cslu.cse.ogi.edu/HLsurvey/HLsurvey.html>.

REHM, GEORG (2003): "Ontologie-basierte Hypertextsorten-Klassifikation". In: *Automatische Textanalyse – Systeme und Methoden zur Annotation und Analyse natürlichsprachlicher Texte*, edited by Mehler, Alexander and Lobin, Henning, Wiesbaden: Westdeutscher Verlag. Erscheint.

Crawling the Web

CHAKRABARTI, SOUMEN; VAN DEN BERG, MARTIN AND DOM, BYRON (1999): "Focused crawling: A new approach to topic-specific Web resource discovery". *Computer Networks* 31 (11–16): pp. 1623–1640.

COWIE, JIM; LUDOVIK, EVGENY AND ZACHARSKI, RON (1998): "An Autonomous, Web-based, Multilingual Corpus Collection Tool". In: *Proceedings of the International Conference on Natural Language Processing and Industrial Applications*. Moncton, pp. 142–148. Online verfügbar: <http://cr1.nmsu.edu/~raz/langrec/nlpia.htm>.

HEYDON, ALLAN AND NAJORK, MARC (1999): "Mercator: A Scalable, Extensible Web Crawler". *World Wide Web* 2 (4): pp. 219–229.

O'NEILL, EDWARD T.; McCLAIN, PATRICK D. AND LAVOIE, BRIAN F. (1997): "A Methodology for Sampling the World Wide Web". In: *Annual Review of OCLC Research 1997*, Dublin: OCLC Online Computer Library Center. Online verfügbar: <http://www.oclc.org>.

RAGHAVAN, SRIRAM AND GARCIA-MOLINA, HECTOR (2001): "Crawling the Hidden Web". In: *Proceedings of the 27th International Conference on Very Large Databases (VLDB)*. pp. 129–138.

REHM, GEORG (2001): "korpus.html – Zur Sammlung, Datenbank-basierten Erfassung, Annotation und Auswertung von HTML-Dokumenten". In: *Proceedings of the GLDV Spring Meeting 2001*, edited by Lobin, Henning. Gesellschaft für linguistische Datenverarbeitung (Society for Computational Linguistics and Language Technology), Giessen, Germany, pp. 93–103. Online verfügbar: <http://www.uni-giessen.de/fb09/ascl/gldv2001/>.

THELWALL, MIKE (2001): "A Web Crawler Design for Data Mining". *Journal of Information Science* 27 (5): pp. 319–326.

TURAU, VOLKER (1998): "Web-Roboter". *Informatik Spektrum* 21 (3): pp. 159–160.

The Semantic Web

BERNERS-LEE, TIM (1999): *Weaving the Web – The Original Design and Ultimate Destiny of the World Wide Web by Its Inventor*. San Francisco: Harper San Francisco.

BERNERS-LEE, TIM; CAILLIAU, ROBERT; GROFF, JEAN-FRANÇOIS AND POLLERMANN, BERND (1992): "World-Wide Web: The Information Universe". *Electronic Networking: Research, Applications and Policy* 1 (2): pp. 52–58.

BERNERS-LEE, TIM; HENDLER, JAMES AND LASSILA, ORA (2001): "The Semantic Web". *Scientific American* 284 (5): pp. 34–43.

DILL, STEPHEN; EIRON, NADAV; GIBSON, DAVID; GRUHL, DANIEL; GUHA, RAMANATHAN; JHINGRAN, ANANT; KANUNGO, TAPAS; RAJAGOPALAN, SRIDHAR; TOMKINS, ANDREW; TOMLIN, JOHN A. AND ZIEN, JASON Y. (2003): "SemTag and Seeker: Bootstrapping the Semantic Web via Automated Semantic Annotation". In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.

DOAN, ANHAI; MADHAVAN, JAYANT; DOMINGOS, PEDRO AND HALEVY, ALON (2002): "Learning to Map between Ontologies on the Semantic Web". In: *Proceedings of the 11th International World Wide Web Conference (WWW 2002)*. Honolulu.

KLEIN, MICHAEL; FENSEL, DIETER; VAN HARMELLEN, FRANK AND HORROCKS, IAN (2001): "The Relation Between Ontologies and XML Schemas". *Linköping Electronic Articles in Computer and Information Science* 6 (4). Online verfügbar: <http://www.ida.liu.se/ext/epa/cis/>.

Web Document Analysis

ASIRVATHAM, ARUL PRAKASH AND RAVI, KRANTHI KUMAR (2001): "Web Page Classification Based on Document Structure". Technical Report, International Institute of Information Technology, Hyderabad. Online verfügbar: http://www.iiit.net/stud_pub.htm.

BUYUKKOKTEN, ORKUT; GARCIA-MOLINA, HECTOR; PAEPCKE, ANDREAS AND WINOGRAD, TERRY (2000): "Power Browser: Efficient Web Browsing for PDAs". In: *Proceedings of the Conference on Human Factors in Computing Systems (CHI'00)*. The Hague.

CARCHIOLO, VINCENZA; LONGHEU, ALESSANDRO AND MALGERI, MICHELE (2002): "Extraction of Hidden Semantics from Web Pages". In: *Intelligent Data Engineering and Automated Learning – IDEAL 2002*, edited by Yin, Hujun; Allinson, Nigel; Freeman, Richard; Keane, John and Hubbard, Simon, Berlin, Heidelberg, New York etc.: Springer, number 2412 in Lecture Notes in Computer Science, pp. 117–122.

CHAN, MICHAEL AND YU, GIN (1999): "Extracting Web Design Knowledge: The Web De-Compiler". In: *IEEE International Conference on Multimedia Computing and Systems (ICMCS 1999)*. IEEE Computer Society, Florence, volume 2, pp. 547–552.

- CHEN, YU; MA, WEI-YING AND ZHANG, HONG-JIANG (2003): “Detecting Web Page Structure for Adaptive Viewing on Small Form Factor Devices”. In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- CUTLER, MICHAEL; SHIH, YUNGMING AND MENG, WEIYI (1997): “Using the Structure of HTML Documents to Improve Retrieval”. In: *USENIX Symposium on Internet Technologies and Systems (NSITS '97)*. Monterey, pp. 241–251.
- DiPASQUO, DAN (1998): “Using HTML Formatting to Aid in Natural Language Processing on the World Wide Web”. Senior Honors Thesis, School of Computer Science, Carnegie Mellon University. Online verfügbar: <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-11/www/wwkb/>.
- DORAI, GAUTHAM K. AND YACOOB, YASER (2002): “Embedded Grammar Tags: Advancing Natural Language Interaction on the Web”. *IEEE Intelligent Systems* 17 (1): pp. 48–53.
- JIANG, YUAN (1998): *Record-Boundary Detection in Web Documents*. Master’s thesis, Brigham Young University.
- REHM, GEORG (2003): “Hypertextsorten-Klassifikation als Grundlage generischer Informationsextraktion”. In: *Automatische Textanalyse – Systeme und Methoden zur Annotation und Analyse natürlichsprachlicher Texte*, edited by Mehler, Alexander and Lobin, Henning, Wiesbaden: Westdeutscher Verlag. Erscheint.
- SAHUGUET, ARNAUD AND AZAVANT, FABIEN (1999): “Looking at the Web through XML Glasses”. In: *Proceedings of the 4th International Conference on Cooperative Information Systems (CoopIS '99)*. Edinburgh: IEEE Computer Society Press, pp. 148–159.
- WALKER, DEREK (1999): “Taking Snapshots of the Web with a TEI Camera”. *Computers and the Humanities* 33 (1–2): pp. 185–192.

Wrapping and Information Extraction

- ABITEBOUL, SERGE; BUNEMAN, PETER AND SUCIU, DAN (2000): *Data on the Web – From Relations to Semistructured Data and XML*. San Francisco: Morgan Kaufmann.
- ADELBERG, BRAD (1998): “NoDoSE – A Tool for Semi-Automatically Extracting Structured and Semistructured Data from Text Documents”. In: *Proceedings of the 1998 ACM International Conference on Management of Data (SIGMOD'98)*. Seattle, pp. 283–294.
- ADELBERG, BRAD AND DENNY, MATT (1999): “Building Robust Wrappers for Text Sources”. Technical Report, Department of Computer Science, Northwestern University.
- BAUMGARTNER, ROBERT; FLESCA, SERGIO AND GOTTLÖB, GEORG (2001): “Visual Web Information Extraction with Lixto”. In: *Proceedings of the 27th VLDB Conference*. Rom, pp. 119–128.
- BOLLACKER, KURT; LAWRENCE, STEVE AND GILES, C. LEE (1998): “CiteSeer: An Autonomous Web Agent for Automatic Retrieval and Identification of Interesting Publications”. In: *Proceedings of the Second International Conference on Autonomous Agents*, edited by Sycara, Katia P. and Wooldridge, Michael. New York: ACM Press, pp. 116–123.
- BUTTLER, DAVID; LIU, LING AND PU, CALTON (2001): “A Fully Automated Object Extraction System for the World Wide Web”. In: *Proceedings of the 2001 International Conference on Distributed Computing Systems (ICDCS '01)*. Phoenix.
- CHAKRABARTI, SOUMEN (2003): *Mining the Web – Discovering Knowledge from Hypertext Data*. Amsterdam, Boston, London etc.: Morgan Kaufmann.
- EIKVIL, LINE (1999): “Information Extraction from World Wide Web – A Survey”. Technical Report 945, Norwegian Computing Center.

- GAO, XIAOYING AND STERLING, LEON (1999): “AutoWrapper: Automatic Wrapper Generation for Multiple Online Services”. In: *Proceedings of Asia Pacific Web Conference 1999 (APWeb99)*. pp. 61–70.
- GUPTA, SUHIT; KAISER, GAIL; NEISTADT, DAVID AND GRIMM, PETER (2003): “DOM-based Content Extraction of HTML Documents”. In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- LIU, LING; PU, CALTON AND HAN, WEI (2000): “XWRAP: An XML-Enabled Wrapper Construction System for Web Information Sources”. In: *Proceedings of the International Conference on Data Engineering (ICDE)*. pp. 611–621.
- MYLLYMAKI, JUSSI (2001): “Effective Web Data Extraction with Standard XML Technologies”. In: *Proceedings of the 10th International World Wide Web Conference (WWW-10)*. Hong Kong, pp. 689–696.
- PARADIS, FRANÇOIS (2000): “Information Extraction and Gathering for Search Engines: The Taylor Approach”. RIAO (Recherche d’Informations Assistée par Ordinateur), Paris, France.
- SEO, HEEKYOUNG AND CHOI, JAEYOUNG YANG JOONGMIN (2001): “Knowledge-based Wrapper Generation by Using XML”. In: *Proceedings of the IJCAI-2001 Workshop on Adaptive Text Extraction and Mining*. Seattle.
- SODERLAND, STEPHEN (1997): “Learning to Extract Text-Based Information from the World Wide Web”. In: *Proceedings of the 3rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-97)*, edited by Heckerman, David; Mannila, Heikki and Pregibon, Daryl. Newport Beach: AAAI Press, pp. 251–254.

Searching the Web

- BRIN, SERGEY AND PAGE, LAWRENCE (1998): “The Anatomy of a Large-Scale Hypertextual Web Search Engine”. *Computer Networks and ISDN Systems* 30 (1–7): pp. 107–117.
- CHAKRABARTI, SOUMEN; DOM, BYRON; KUMAR, S. RAVI; RAGHAVAN, PRABHAKAR; RAJAGOPALAN, SRIDHAR; TOMKINS, ANDREW; KLEINBERG, JON M. AND GIBSON, DAVID (1999): “Hypersearching the Web”. *Scientific American* 280 (6): pp. 54–60.
- HU, WEN-CHEN; CHEN, YINING; SCHMALZ, MARK S. AND RITTER, GERHARD X. (2001): “An Overview of World Wide Web Search Technologies”. In: *Proceedings of the 5th World Multiconference on Systemics, Cybernetics and Informatics (SCI 2001)*. Orlando.
- HUANG, LAN (2000): “A Survey on Web Information Retrieval Technologies”. Technical Report, Experimental Computer Systems Laboratory (ECSL), State University of New York (SUNY). Online verfügbar: <http://www.cs.sunysb.edu/~lanhuang/>.

Automatic Summarization in the Web

- BUYUKKOKTEN, ORKUT; GARCIA-MOLINA, HECTOR AND PAEPCKE, ANDREAS (2001a): “Accordion summarization for end-game browsing on PDAs and cellular phones”. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM Press, pp. 213–220.
- BUYUKKOKTEN, ORKUT; GARCIA-MOLINA, HECTOR AND PAEPCKE, ANDREAS (2001b): “Seeing the Whole in Parts: Text Summarization for Web Browsing on Handheld Devices”. In: *Proceedings of the 10th International World Wide Web Conference (WWW 2001)*. Hong Kong.
- CORSTON-OLIVER, SIMON (1998): *Computing Representations of the Structure of Written Discourse*. Ph.D. thesis, University of California, Santa Barbara. Zugleich: Microsoft Research Technical Report, MSR-TR-98-15.

- LEONG, H.; KAPUR, S. AND DE VEL, O. (1997): "Text Summarisation for Knowledge Filtering Agents in Distributed Heterogeneous Environments". In: Mahesh (1997b), pp. 87–94. Online verfügbar: <http://crl.nmsu.edu/users/mahesh/aaai-web-nlp-symposium/proceedings.html>.
- LOK, SIMON AND KAN, MIN-YEN (2003): "Employing Natural Language Summarization and Automated Layout for Effective Presentation and Navigation of Information Retrieval Results". In: *Proceedings of the 12th International World Wide Web Conference (WWW 2003)*. Budapest.
- MAHESH, KAVI (1997a): "Hypertext Summary Extraction for Fast Document Browsing". In: Mahesh (1997b), pp. 95–103. Online verfügbar: <http://crl.nmsu.edu/users/mahesh/aaai-web-nlp-symposium/proceedings.html>.
- MAHESH, KAVI (editor) (1997b): *Working Notes of the AAAI Spring Symposium: Natural Language Processing for the World Wide Web*, Stanford. American Association for Artificial Intelligence. Online verfügbar: <http://crl.nmsu.edu/users/mahesh/aaai-web-nlp-symposium/proceedings.html>.
- REHM, GEORG (1998): *Vorüberlegungen zur automatischen Zusammenfassung deutschsprachiger Texte mittels einer SGML- und DSSSL-basierten Repräsentation von RST-Relationen*. Magisterarbeit, Studiengang Computerlinguistik und Künstliche Intelligenz, Universität Osnabrück. Online verfügbar: <http://www.uni-giessen.de/~g91063/papers.shtml>.

Natural language processing helps computer to understand human language as it is spoken. Real world use of natural languages such as English, Hindi, German, French etc doesn't have a formulated grammar. Over the years there have been many advancements in Natural language processing. NLP Terminologies : Lets understand the Basic Terminologies used in NLP : Tokenization, Corpus or Corpora, Stemming, Bag of Words, Stop Words, Tf-idf, Disambiguation, Topic Models , Word Boundaries. Tokenization : Tokenization is a process to split longer strings into smaller pieces. Large documents can be tokenized into paragraphs, Paragraphs can be tokenized into sentences and sentences can be tokenized into phrases, words or letters. Corpus or Corpora Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit by Steven Bird Paperback \$38.99. In Stock. Ships from and sold by Amazon.com. "Instructors looking for a good introductory text to use in a course devoted to computational linguistics should consider this book as a strong candidate if they wish to emphasize a linguistics approach. The book is interestingly written with many insightful discussions, and it is the only (introductory) computational linguistics textbook that looks at the field from a linguist's point of view." NLTK (Natural Language Toolkit) is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to many corpora and lexical resources. Also, it contains a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning. Best of all, NLTK is a free, open source, community-driven project. We use NLTK to show some basics of the natural language processing field. For the examples below, we assume that we have imported the NLTK toolkit. We can do this like this: `import nltk`.